

Thèse de Doctorat

de l'Université Sorbonne Paris Cité

Préparée à l'Université Paris Diderot

Equipe de Biologie Computationnelle et Biomathématiques à l'IJM
Laboratoire de Probabilités, Statistique et Modélisation

**Cloning Algorithms:
from Large Deviations to Population Dynamics**

Esteban GUEVARA HIDALGO

Dirigée par Khashayar PAKDAMAN
Co-dirigée par Vivien LECOMTE

Présentée et soutenue publiquement à Paris le 1 Juin 2018

Président du jury:	Mme.	Leticia Cugliandolo	Université Pierre et Marie Curie
Rapporteur:	M.	Juan P. Garrahan	University of Nottingham
Rapporteur:	M.	Hugo Touchette	Stellenbosch University
Examineur:	M.	Julien Tailleur	Université Paris Diderot
Examineur:	M.	Vivien Lecomte	Université Grenoble-Alpes
Directeur de thèse:	M.	Khashayar Pakdaman	Université Paris Diderot



Cloning Algorithms: from Large Deviations to Population Dynamics

Abstract: Population dynamics provides a numerical tool allowing for the study of rare events by means of simulating a large number of copies of the system, supplemented with a selection rule that favours the rare trajectories of interest. The cloning algorithm allows the estimation of a large deviation function (LDF) of additive observables in Markov processes. However, such algorithms are plagued by finite simulation time t and finite population size N_c effects that can render their use delicate. First, using a non-constant population approach, we analyze the small- N_c effects in the initial transient regime. These effects play an important role in the numerical determination of LDF. We show how to overcome these effects by introducing a time delay in the evolution of populations, additional to the discarding of the initial regime of the population growth where these discreteness effects are strong. Then, the study of the finite- t and finite- N_c scalings in the LDF evaluation is done using two different versions of the algorithm, in discrete and continuous-time. We show that these scalings behave as $1/N_c$ and $1/t$ in the large- N_c and large- t asymptotics respectively. Moreover, we show that one can make use of this convergence speed in order to extract the asymptotic behavior in the infinite- t and infinite- N_c limits resulting in a better LDF estimation. These scalings are later generalized and evidence of a breakdown for large-size systems is presented.

Keywords : Rare Events, Large Deviations, Population Dynamics Algorithms

Algorithmes de Clonage: des Grandes Déviations à la Dynamique des Populations

Résumé: La dynamique des populations fournit un outil numérique qui permet l'étude des événements rares grâce à la simulation d'un grand nombre de copies du système. Le processus est muni d'une règle qui favorise les trajectoires rares d'intérêt. La méthode de l'algorithme de clonage permet l'estimation de la fonction de grandes déviations (en anglais, LDF) pour les observables additives pour les processus de Markov. Cependant, cette méthode doit être soigneusement utilisée car il existe des effets de temps de simulation t finie et de taille de population N_c finie. Premièrement, nous analysons les effets de petit N_c dans un régime transitoire initial en utilisant une approche de population non constante. Ces effets jouent un rôle important dans la détermination numérique de la LDF. Pour surmonter ces effets, nous avons introduit un délai dans l'évolution des populations, en plus de l'exclusion du régime initial de la croissance de la population où ces effets sont forts. Ensuite, l'étude des lois d'échelle de t et N_c finie dans l'évaluation de LDF est faite en utilisant deux versions différentes de l'algorithme, en temps discret et en temps continu. Nous montrons que ces échelles se comportent comme $1/N_c$ et $1/t$ dans les régimes asymptotiques de grand N_c et de grand- t respectivement. En outre, nous montrons qu'il est possible d'utiliser cette vitesse de convergence pour extraire le comportement asymptotique des limites de t et N_c infinis, fournissant ainsi une meilleure estimation de la LDF. Enfin, ces lois d'échelles sont généralisées et les indications de leurs limites dans les systèmes de grandes dimensions sont présentées.

Mots clés : Evénements Rares, Grandes Déviations, Dynamique des Populations

Acknowledgements

Many thanks to Khashayar Pakdaman and Vivien Lecomte for their support and discussions. Thanks to the team of Biologie Computationnelle et Biomathématiques at Institut Jacques Monod where this thesis was developed. Special thanks to the Ecuadorian State and the Secretaría Nacional de Educación Superior, Ciencia, Tecnología e Innovación, SENESCYT, which supported financially this PhD program and the French Republic for welcoming me. Also to my friends and family in Ecuador who were always present. Finally, I thank my friends in France, especially, my girlfriend who was my companion during these years making the stay in Paris not only easier but also wonderful.

*A mi mamá,
mis hermanos
y mis perros*

Résumé

L'occurrence d'**événements rares** peut grandement contribuer à l'évolution des systèmes physiques en raison de leur effets dramatiques. La théorie des grandes déviations fournit un ensemble d'outils qui permettent leur traitement [1–3]. Ces probabilités et fluctuations ont la propriété de décroître exponentiellement en fonction d'un paramètre (comme le temps ou la température). Cela signifie que lorsque le paramètre se devient plus grand, l'événement devient moins probable [4]. D'un point de vue pratique, la théorie des grandes déviations peut être vue comme une collection de méthodes qui permettent de déterminer si un principe de grandes déviations existe pour une variable aléatoire donnée et pour déterminer sa fonction de taux ou “fonction des grandes déviations” (en anglais, LDF).

Seulement dans quelques cas simples, il est possible d'obtenir des expressions exactes et expressions explicites pour la fonction de taux [5, 6]. Pour la plupart des processus stochastiques, l'évaluation de ces fonctions est faite en utilisant des approches analytiques et des méthodes numériques [1–3, 7, 8]. Ils vont de “importance sampling method” [9], “adaptive multilevel splitting” [10] à “transition path sampling” [11–14] et des algorithmes “go with the winner” [15, 16] aux méthodes de dynamique des populations [7, 17] en temps discret [18] ou temps continu [19]. Ces méthodes ont été généralisées à de nombreux contextes [20–24]. En physique, ceux-ci sont de plus en plus utilisés dans l'étude des systèmes complexes, par exemple dans l'étude des fluctuations réelles des modèles de transport [25–27], glasses [12], du repliement de protéines [28] et des réseaux de signalisation [29, 30]. Mathématiquement, la procédure revient à déterminer la fonction des grandes déviations associée à la distribution d'une observable dépendant de la trajectoire, qui peut être reformulée à son tour en la détermination de l'état fondamental d'un opérateur linéaire [31], une question commune à la physique statistique et à la physique quantique [32].

Dans cette thèse, nous accordons une attention particulière à algorithmes basés sur **dynamique des populations** [7, 17–19, 33] afin d'étudier les trajectoires rares en biaisant exponentiellement leur probabilité. Dans ce contexte, la procédure numérique introduite par Giardinà, Kurchan et Peliti [18] surmonte la difficulté d'observer les fluctuations d'une observable (dont la probabilité diminue exponentiellement dans le temps) pour les chaînes de Markov à temps discret. La fonction des grandes déviations peut être obtenue comme la plus grande valeur propre d'une matrice d'évolution d'une dynamique modifiée [17, 18] qui peut être calculé numériquement [5, 6, 34] seulement pour les petits systèmes car la matrice d'évolution est exponentiellement grande dans la taille du système. Ensuite, une modification de cette procédure a été proposée [19, 35] pour lequel les problèmes de discrétisation lié à l'approche originale [18] sont contournés avec une approche directe en temps continu.

L'évolution du système a été représentée par une dynamique de population du type de diffusion Monte Carlo [32]. Cet algorithme a été appliqué pour calculer avec succès les grandes déviations du courant total dans le processus d'exclusion symétrique et asymétrique [36, 37], et de l'activité dans le processus de contact [38]. Aussi, pour analyser la dynamique [39, 40] des modèles cinétiquement contraints (KCM) [41–55], des “glassy systems” [56–58] à travers les statistiques de trajectoires de la dynamique montrant que ces modèles présentent une transition dynamique de premier ordre entre les phases dynamiques active et inactive. Il a également été utilisé pour étudier les symétries dans les fluctuations loin de l'équilibre [59] et dans les modèles de transport [21, 22, 60]. Ces études permettent non seulement de tester les prédictions de l'hydrodynamique fluctuante [21, 61], mais aussi les limites de la méthode elle-même [22]. Il a été également suggéré [7] que la méthode pourrait être appliquée pour étudier en détail de possible future et l'évolution passée des systèmes planétaires, et aussi l'auto-organisation de la stabilité de notre système solaire.

L'idée de la dynamique de populations est de traduire l'étude d'une classe de trajectoires rares (par rapport à une contrainte globale déterminée) dans l'évolution de plusieurs copies de la dynamique originale, avec un processus de sélection local-dans-temps rendant l'occurrence des trajectoires rares typiques dans la population évoluée. La distribution de la classe des trajectoires rares dans la dynamique originale est liée à la croissance (ou décroissance) exponentielle de la population des clones du système et la LDF peut être estimée à partir de leur taux de croissance. Les procédures numériques visant à simuler efficacement des événements rares, en utilisant un schéma de dynamique de population sont communément appelées **algorithmes de clonage**. Dans de tels algorithmes, les copies du système sont évoluées en parallèle et celles qui montrant le comportement rare d'intérêt sont multipliées itérativement [7, 16–19, 23, 32, 33, 62–70].

Les différentes versions de l'algorithme de clonage utilisées dans cette thèse sont détaillées dans le [Introduction](#). Là, nous partons de la construction de l'équation maîtresse, de sa solution et de son interprétation. Nous définissons le principe des grandes déviations pour certains observables \mathcal{O} ce qui peut être interprété comme la probabilité d'observer une valeur atypique de celle observables après une longue échelle de temps. La fonction de taux de ce principe des grandes déviations correspond à la LDF et c'est un équivalent dynamique de l'entropie intensive dans l'ensemble microcanonique [1]. Il code non seulement les fluctuations gaussiennes mais aussi les fluctuations non-gaussiennes (ou les grandes déviations) de l'observable \mathcal{O}/t qui peut être obtenu par son expansion au-delà de l'ordre quadratique. Dans la limite du temps infini, le LDF peut ne pas être analytique, ce qui peut être interprété comme une signature d'hétérogénéités dynamiques (transition de phase dynamique) [71, 72].

Le problème de la détermination de la fonction de taux est en général une tâche difficile dans l'ensemble microcanonique, on préfère donc aller à l'ensemble dynamique canonique (ou espace de Laplace). Au lieu de fixer la valeur de l'observable \mathcal{O} afin de déterminer la LDF on introduit un paramètre s (intensif dans le temps) qui biaise le poids statistique des histoires et fixe la valeur moyenne de \mathcal{O} , de sorte que $s \neq 0$ favorise ses valeurs non-typiques. Le paramètre s implique une modification (exponentielle) du poids statistique des histoires du système. Valeurs pour $s = 0$ correspondent aux moyennes de l'état stable de \mathcal{O} . Pendant ce temps, les valeurs des $s \neq 0$ favorisent les histoires avec des valeurs non-typiques de l'observable \mathcal{O} .

Pour des raisons pratiques, il est convenient de calculer la fonction génératrice des cumulants (par ses sigles en anglais CGF) au lieu de LDF (qui sont reliées par une transformée de Legendre), ce que nous calculons en pratique tout au long de la thèse. Nous montrons comment estimer le CGF à partir de l’interprétation de la dynamique des populations modifiée ou à partir de la plus grande valeur propre de l’opérateur d’évolution modifiée. Dans le premier cas, l’équation d’évolution temporelle qui décrit la dynamique modifiée peut être interprétée pas comme l’évolution d’un système unique, mais comme une dynamique des populations sur un grand nombre N_c des copies du système qui évolue de manière couplée [18, 33]. C’est-à-dire, comme un processus stochastique avec des taux de transition fourni par un mécanisme de sélection où un clone du système est copié s’il est rare ou tué sinon.

Nous détaillons également les approches avec population totale **non-constante** et **constante** de l’algorithme de clonage et les estimateurs CGF obtenus à partir de ceux-ci. Nous expliquons comment pour l’approche de la population totale constante, un uniforme élagage/clonage est appliqué au-dessus de la dynamique de clonage afin d’éviter l’explosion ou la disparition exponentielle de la population. Alors que la dernière version est évidemment plus “computer-friendly”, l’ancienne version présente des caractéristiques intéressantes: Premièrement, il est directement lié à l’évolution des systèmes biologiques (sauts stochastiques représentant des mutations, les règles de sélection étant interprétées comme une pression darwinienne); Seconde, le uniforme élagage/clonage de la population, bien que non biaisé, induit des corrélations dans la dynamique que l’on pourrait vouloir éviter; Enfin, dans certaines situations où les taux de sélection sont très fluctuants, l’algorithme de population constante ne peut pas être utilisé dans la pratique en raison des effets de population finie (la population étant éliminée par un seul clone), et on doit recourir à la non-constante. Un exemple de la mise en oeuvre de cette version peut être trouvé dans Ref. [73].

À la fin de l’introduction, nous présentons les exemples de modèles utilisés pour notre analyse: une simple dynamique d’annihilation-crédation à deux états, et un processus de contact sur un treillis périodique unidimensionnel. Le premier système (chapitres: II, III, IV et V) a été choisi pour sa simplicité et la possibilité de comparer les prédictions numériques avec les valeurs exactes de CGF. D’autre part, le processus de contact (chapitres: IV, V et VI) est utilisé pour étendre l’analyse et vérifier les résultats sur a (plus complexe) système de “many body” où la dépendance avec la taille du système peut également être analysée. Dans les deux cas, nous considérons l’activité dynamique K [12, 39, 40, 74–82] comme l’additif observable \mathcal{O} (I.19) d’intérêt. L’expression analytique du CGF est obtenue (lorsque cela est possible) en résolvant la plus grande eigenvalue de l’opérateur modifié comme discuté dans la Sec. I.6.1.

Dans le chapitre II: **Discreteness Effects in Population Dynamics** [P1], nous appliquons l’algorithme de population non-constante afin d’analyser numériquement les effets dus à la petite taille de population dans le régime transitoire initial sur un modèle simple d’annihilation-crédation (Sec. I.8.1) où sa mise en oeuvre et ses propriétés peuvent être examinées dans les moindres détails. Au cours du régime transitoire initial de l’évolution des populations, il y a une grande distribution des temps où la première série de sauts se produisent. Cela signifie que les fluctuations au moment initial produisent que certaines populations restent dans leur état initial beaucoup plus longtemps que d’autres, produisant un écart dans leur évolution individuelle. Cela induit un décalage relatif qui dure pour toujours.

Ces effets jouent un rôle important en particulier pour la détermination de la fonction de grandes déviations. L'estimation de CGF provient de la détermination du taux de croissance d'une log-population moyenne (Sec. I.7.1.1) sur de nombreuses réalisations de la dynamique.

Afin de réaliser cette moyenne de manière systématique, nous définissons une procédure que nous avons appelée **merging**. Toutefois, cette moyenne est fortement dépendante non seulement du nombre de réalisations, et sur la taille de la population initiale mais aussi sur le temps (ou population) de coupure considéré pour arrêter leur évolution. C'est-à-dire, en limitant l'évolution de nos populations jusqu'à un maximum de temps T_{\max} (ou population N_{\max}) ce qui n'est pas "assez grand", la population moyenne (et la détermination du CGF) peut être influencée par **effets de discrétion aux temps initiaux**, causés par une petite taille de la population. Nous avons proposé comme alternative afin de surmonter l'influence des effets de discrétisation se débarrasser des régions des populations où ces effets sont présents. Autrement dit, couper le régime transitoire initial des populations. Dans ce cas, nous voyons que la moyenne des populations est limitée à un intervalle qui peut être très petit et cela peut induire une mauvaise estimation de CGF.

En complément, nous avons trouvé un moyen de souligner les effets du régime de croissance exponentielle dans la détermination du CGF en utilisant le fait que les log-populations après un temps assez long deviennent parallèles (Fig. II.2(a)) et qu'une fois que les populations ont surmonté le régime des effets de discrétion, la distance entre elles devient constante (Fig. II.2(b)) et ces effets ne sont plus forts (Sec. II.3). D'autre part, nous soutenons en Sec. II.4.1 que ces effets de discrétion initiale ou le décalage initiale entre les populations pourrait être compensée en effectuant sur les populations un déplacement du temps (Eq. (II.2)). Cette procédure est choisie de façon à chevaucher les évolutions de la population dans leur régime de grande durée large-time régime (Fig. II.4(b)). Ceci avec un rejet des régimes initiaux dans l'évolution de la population surpasse l'influence des effets de discrétion améliorant l'estimation de CGF.

Nous montrons que c'est vrai, indépendamment de la méthode utilisée pour calculer le taux de croissance de la population moyenne, comme le Fig. II.7. En outre, il est montré que si en plus, nous effectuons la transformation temporaire, l'estimation de CGF est encore améliorée et plus proche de la valeur théorique (Sec. II.5.2). Ce résultat est confirmé plus tard en Sec. II.5.3 en calculant la distance relative des estimateurs numériques à la valeur théorique et leurs erreurs. Comme peut être observé dans Fig. II.8, l'écart de la valeur théorique est plus élevé pour les valeurs de s proches de 0, mais est plus petit après la "correction du temps" pour presque chaque valeur de s . De même pour l'estimateur d'erreur (Fig. II.9). De plus, nous étudions les propriétés de ces retards. Nos résultats numériques supportent également un comportement "loi de puissance" dans le temps de la variance de retard. Par ailleurs, la distribution des délais prend une forme universelle, après avoir rééchelonné la variance à 1.

Le chapitre II [P1] est structuré comme suit: en Sec. II.2 nous décrivons les problèmes liés à la moyenne des réalisations distincts, que nous quantifions en Sec. II.3. En Sec. II.4 nous proposons la méthode pour augmenter l'efficacité de l'algorithme de dynamique de la population en appliquant un temps de retard dépendant de la réalisation, et nous présentons les résultats de son application en Sec. II.5. Nous caractérisons numériquement la distribution de ces temps de retard en Sec. II.6. Nos conclusions et perspectives sont réunis en Sec. II.7.

En dehors des approches de population constante, les mécanismes de sélection dans l'algorithme de clonage peuvent être mis en œuvre de différentes manières. L'un d'eux, avec chaque évolution des copies des copies (**Continuous-Time**) [7, 17, 19] ou bien, pour chaque intervalle de temps pré-fixé (**Discrete-Time**) [18]. Les différences importantes entre les deux techniques sont discutées dans Secs. III.2.4.2 et IV.5.

L'algorithme proposé par Giardinà et al. [7, 17–19, 33, 70] (un technique à temps discrete) est utilisé pour évaluer numériquement le CGF d'additif (ou extensif dans le temps) observables dans les processus de Markov [1, 83]. Le CGF est obtenu comme le taux de croissance exponentiel que la population présenterait si elle n'était pas maintenue constante. Il a été appliqué à de nombreux systèmes physiques, y compris les systèmes chaotiques, la dynamique vitreuse et les modèles de gaz en treillis sans équilibre, a permis l'étude de nouvelles propriétés, telles que le comportement des respirateurs dans la chaîne de Fermi-Pasta-Ulam-Tsingou [33], transitions de phase dynamiques dans des modèles cinétiquement contraints [39], et un principe d'additivité pour les processus d'exclusion simples [60, 84]. Sous cette approche, le correspondant estimateur CGF n'est valide que dans les limites du temps de simulation infinie t et taille de la population infinie N_c . La stratégie habituelle suivie pour obtenir ces limites est d'augmenter le temps de simulation et la taille de la population jusqu'à ce que la moyenne de l'estimateur sur plusieurs réalisations ne dépende pas de ces deux paramètres, jusqu'à des incertitudes numériques.

Bien que la méthode a été largement utilisée, il y a eu moins d'études axées sur la justification analytique de l'algorithme. De plus, il introduit deux paramètres supplémentaires en considération: la taille de la population N_c et le temps de simulation t , tous les deux affectent considérablement la précision de l'estimation de CGF. Même si l'on croit heuristiquement que l'estimateur LDF converge vers le résultat correct à mesure que le nombre de copies N_c augmente, il n'y a pas de preuve de sa convergence. Relatif à ce manque de preuve, bien que nous utilisons l'algorithme en supposant sa validité, nous n'avons aucune idée de la vitesse à laquelle l'estimateur converge en $N_c \rightarrow \infty$. Nous discutons de cette convergence en effectuant une étude analytique en temps discret dans le chapitre III et en utilisant une approche numérique en temps continu dans le chapitre IV. Il est important de remarquer que les deux versions de l'algorithme (temps discret et continu) diffèrent sur un point crucial qui fait qu'une extension de l'analyse développée au chapitre III ne peut pas être faite directement pour comprendre le cas de temps continu dans chapitre IV.

Dans le chapitre III: **Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: I. Analytical Study using a Birth-Death Process** [P2], nous discutons de cette convergence en définissant deux types d'erreurs numériques. Premièrement, pour un nombre fini et fixe de clones N_c , en faisant la moyenne sur un grand nombre de réalisations, l'estimateur CGF converge vers une valeur incorrecte, qui est différente du résultat de grande déviation souhaité. Nous appelons cette déviation de la valeur correcte, **erreurs systématiques**. Par rapport à ces erreurs, nous considérons également les fluctuations de la valeur estimée. Plus précisément, pour une valeur fixe de N_c , les résultats obtenus dans différentes réalisations sont répartis autour de cette valeur incorrecte. Nous appelons les erreurs associées à ces fluctuations **erreurs stochastiques**. Bien que les deux erreurs soient importantes dans les simulations numériques, la dernière peut conduire cet algorithme à produire de mauvais résultats. Par exemple, l'erreur systématique croît exponentiellement à mesure que la température diminue [85].

Pour étudier ces erreurs, nous avons utilisé une description de processus de naissance-mort [86, 87] par l’algorithme de dynamique de population comme il est expliqué ci-dessous: Nous nous concentrons sur les systèmes physiques décrits par une dynamique de Markov [7, 18, 19] avec un nombre fini d’états M . Nous dénotons par i ($i = 0, 1, \dots, M - 1$) les états du système. Ce processus de Markov a sa propre dynamique stochastique, décrite par les taux de transition $w(i \rightarrow j)$. Dans les algorithmes de dynamique de population, afin d’étudier ses trajectoires rares, on prépare N_c copies du système, et simule ces copies en fonction de (i) la dynamique de $w(i \rightarrow j)$ (suivi indépendamment pour toutes les copies) et (ii) la étape de “clonage” dans laquelle l’ensemble des copies est directement manipulé, i.e., certaines copies sont éliminées pendant que d’autres sont multipliées (See Table III.1). Formellement, la dynamique des populations représente pour une *unique* copie du système, un processus qui ne préserve pas la probabilité, comme mentionné dans Sec. I.6.2. Ce fait a motivé l’étude des processus auxiliaires [88], des processus efficaces [89] et driven processes [90] pour construire une dynamique modifiée (et leurs approximations [91]) qui préserve la probabilité.

Au chapitre III, nous formulons explicitement la **méta-dynamique** des copies elles-mêmes en utilisant un processus stochastique de naissance-mort qui préserve la probabilité, et il nous permet d’étudier les erreurs numériques de l’algorithme lors de l’évaluation CGF. Nous considérons la dynamique des copies comme un processus stochastique de naissance-mort dont l’état est noté par $n = (n_0, n_1, n_2, \dots, n_{M-1})$, où $0 \leq n_i \leq N_c$ représente le nombre de copies qui sont dans l’état i dans l’ensemble des copies. Nous introduisons explicitement les taux de transition décrivant les dynamiques de n , que nous désignons par $\sigma(n \rightarrow \tilde{n})$. Nous montrons que la dynamique décrite par ces taux de transition conduit en général à l’estimation CGF correcte du système original $w(i \rightarrow j)$ dans la limit $N_c \rightarrow \infty$. Nous montrons aussi que les erreurs systématiques sont de l’ordre $\mathcal{O}(1/N_c)$, alors que les erreurs numériques sont de l’ordre $\mathcal{O}(1/(\tau N_c))$ (où τ is an averaging duration). Ce résultat contraste nettement avec les méthodes Monte-Carlo standard, où les erreurs systématiques sont toujours 0. La formulation développée au chapitre III [P2] nous donne la possibilité de calculer exactement les expressions des coefficients de convergence, comme nous le faisons en Sec. III.4 sur un exemple simple.

Chapitre III [P2] est structuré comme suit. Nous définissons d’abord le problème CGF au début de Sec. III.2, ensuite, nous formulons le processus de naissance-mort utilisé pour décrire l’algorithme en Sec. III.2.1. En utilisant ce processus de naissance-mort, nous démontrons que l’estimateur de l’algorithme converge vers la correct fonction des grandes desviations en Sec. III.2.2. À la fin de cette section, en Sec. III.2.3, nous discutons de la vitesse de convergence de cet estimateur (les erreurs systématiques) et nous dérivons son échelle $\sim 1/N_c$. In Sec. III.3, nous passons aux erreurs stochastiques. Pour discuter de cela, nous introduisons la LDF de l’estimateur, à partir de laquelle nous dérivons que la vitesse de convergence des erreurs stochastiques est proportionnelle à $1/(\tau N_c)$. Dans la section suivante, Sec. III.4, nous introduisons un modèle simple à deux états, auquel nous appliquons les formulations développées dans les sections précédentes. Nous dérivons les expressions exactes des erreurs systématiques en Sec. III.4.1 et des erreurs stochastiques en Sec. III.4.2. Ensuite, en Sec. III.4.3, nous proposons un autre grand estimateur de déviation et enfin, en Sec. III.5, nous résumons les résultats obtenus.

À partir de cette formulation, nous avons déduit les scalings finies en N_c et t des erreurs systématiques de l'estimateur CGF, montrant que ceux-ci se comportent comme $1/N_c$ et $1/t$ dans le grand- N_c et grand- t asymptotiques respectivement. En principe, connaissant l'échelle *a priori* signifie que la limite asymptotique de l'estimateur dans le $t \rightarrow \infty$ et $N_c \rightarrow \infty$ limites peuvent être interpolées à partir des données à fini t et N_c . Toutefois, si cette idée est réellement utile ou non est une question non triviale, comme il y a toujours une possibilité que les valeurs de début des \mathbf{N}_c^{-1} - et \mathbf{t}^{-1} -scalings sont trop volumineux pour les utiliser.

Dans le chapitre **IV: Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: II. Numerical Approach in Continuous Time [P3]**, nous considérons une version continue dans le temps des algorithmes de dynamique des populations [17, 19]. Nous montrons numériquement que on peut en effet faire usage de ces propriétés afin de concevoir une méthode originale et simple qui prenne en compte les échelles exactes du corrections de t et N_c finies afin de fournir des estimateurs CGF significativement meilleurs (**scaling method**) dans l'application à un système avec des interactions à plusieurs corps (un processus de contact). Nous soulignons que les deux versions de l'algorithme diffèrent sur un point crucial qui fait qu'une extension de l'analyse développée au chapitre III [P2] ne peut pas être fait directement afin de comprendre le cas en temps continu (Sec. IV.5). Nous soulignons donc que l'observation de ces échelles eux-mêmes est également non triviale.

Nous notons que les scalings qui régissent la convergence aux limites des temps temps et de taille infinies (avec corrections en $1/N_c$ et en $1/t$) doivent être pris en compte correctement: en effet, en tant que lois de puissance, elles ne présentent pas de taille et de temps caractéristiques au-delà desquelles les corrections seraient négligeables. La situation est très similaire à l'étude de la force de dépinçage critique dans les driven random manifolds: la force critique présente une correction de 1 sur la taille du système [92] qui doit être considéré correctement afin d'extraire sa valeur réelle. Génériquement, de telles échelles fournissent également un critère de convergence aux régimes asymptotiques de l'algorithme: il faut confirmer que l'estimateur CGF présente des corrections (premièrement) $1/t$ et (en second lieu) dans $1/N_c$ par rapport à une valeur asymptotique afin de s'assurer que cette valeur représente une évaluation correcte de la CGF.

Le chapitre IV [P3] est organisé comme suit. En Sec. IV.3.1 nous étudions le comportement de l'estimateur CGF en fonction du temps d'observation (pour une population fixe N_c) et nous voyons comment sa limite de temps infinie peut être extraite à partir des données numériques. En Sec. IV.3.2 nous analysons le comportement de l'estimateur en augmentant le nombre de clones (pour un temps de simulation final donné) et la limite de taille infinie de l'estimateur CGF. Sur la base de ces résultats, nous présentons en Sec. IV.4 a une méthode qui nous permet d'extraire les limit infinies de temps et taille de l'estimator de la fonction des grandes déviations à partir d'une analyse d'échelle à taille finie et à temps fini. En Sec. IV.5, nous discutons la difficulté d'une approche analytique de l'algorithme en temps continu. Enfin, nos conclusions sont faites en Sec. IV.6.

Afin de compléter la discussion principale effectuée au chapitre IV [P3], en chapitre V [P3] nous étudions les fluctuations de l'estimateur CGF (défini dans la version continue dans le temps). Ceci est fait en étudiant sa distribution et sa dépendance avec le temps de simulation et le nombre de clones. Compatible avec le théorème de la limite centrale, nous montrons

comment un redimensionnement approprié de l'estimateur CGF produit un effondrement des distributions dans une distribution standard normale pour différentes valeurs de N_c et des temps de simulation. De plus, nous discutons dans Sec. V.3 une autre façon de le définir, qui a déjà été introduit dans Sec. III.4.3 pour la version à temps discret.

L'analyse d'échelle de t et N_c finies dans l'évaluation de CGF a été réalisée suivant deux approches différentes: un analytique, au chapitre III [P2], en utilisant une version à temps discret de l'algorithme de dynamique des populations [18], et numérique, dans le chapitre IV [P3], en utilisant une version à temps continu [17, 19]. Dans les deux cas, les erreurs systématiques de ces échelles ont été trouvés à se comporter comme $1/t$ et $1/N_c$ dans les asymptotiques de t et N_c grandes respectivement. De plus, il a été montré comment ces propriétés d'échelle peuvent être utilisées pour améliorer l'estimation CGF par la mise en ouvre d'une scaling method (Sec. IV.4.1). Ceci a été fait en considérant que le comportement asymptotique de l'estimateur dans $t \rightarrow \infty$ et $N_c \rightarrow \infty$ limits peut être interpolé à partir des données obtenues à partir de simulations à temps de simulation et nombre de clones **finies et relativement petites**. Cependant, la validité de ces échelles et l'efficacité de la méthode n'ont été prouvées que dans les cas où le nombre de sites L (où la dynamique se produit) était petit: une simple dynamique d'annihilation-création à deux états (Sec. I.8.1) (dans un site) et un processus de contact unidimensionnel (Sec. I.8.2) (avec $L = 6$ sites).

En chapitre VI: **Breakdown of the Finite-Time and Finite- N_c Scalings in the Large- L Limit** [P4], nous complétons les résultats présentés dans les chapitres III [P2] et IV [P3] en étendant l'analyse à un processus de contact de grand- L Afin de le faire, nous redéfinissons ces échelles de manière plus générale. Nous supposons le comportement de l'estimateur CGF décrit par un $t^{-\gamma_t}$ -scaling (Eq. (VI.1)) et un $N_c^{-\gamma_{N_c}}$ -scaling (Eq. (VI.2)). Cette redéfinition nous permet de vérifier dans les systèmes de grand- L si effectivement $\gamma_t \approx 1$, $\gamma_{N_c} \approx 1$ et si les termes $\chi_\infty^{(N_c)}$ et χ_∞^∞ représentent les limites en $t \rightarrow \infty$ et $N_c \rightarrow \infty$ de l'estimateur CGF.

Ceci est fait d'abord en Sec. VI.3.1 où nous avons considéré un processus de contact avec $L = 100$ sites et deux valeurs représentatives du paramètre s ($s = -0.1$ et $s = 0.2$). Bien que le t^{-1} -scaling et le N_c^{-1} -scaling ont été prouvés à tenir pour $s = 0.1$, ce n'était pas le cas pour $s = 0.2$. Comment ce changement d'échelle est-il produit en fonction du paramètre s est présenté en détail dans Sec. VI.3.2 où les exposants $\gamma_t(s)$ et $\gamma_{N_c}(s)$ sont caractérisés. En particulier, pour $\gamma_{N_c}(s)$, nous avons été en mesure de distinguer trois étapes dans son comportement, où, le N_c^{-1} -scaling était valide jusqu'à $s = s^*$, puis γ_{N_c} diminue à 0 at $s = s^{**}$ et enfin, il devient négatif pour $s > s^{**}$. En Sec. VI.3.3 nous montrons comment ces échelles affectent la détermination de la limite infini en t et N_c de l'estimateur CGF. Cela se produit parce que le scaling method reposait sur la validité du t^{-1} - et N_c^{-1} -scalings. Comme pour $L = 100$ ce n'est pas le cas, il est possible de voir comment les différents estimateurs correspondaient les uns aux autres jusqu'à $s = s^*$ à partir de laquelle ils divergent jusqu'à $s = s^{**}$ où il y a une discontinuité. Cette analyse est étendue au plan $s - L$ en Sec. VI.4 où les exposants γ_t et γ_{N_c} ont été calculé pour une grille de valeurs des paramètres (s, L) . Leur caractérisation est faite en introduisant une dépendance du s' , s^* et s^{**} avec le nombre de sites précédemment défini en Sec. VI.3 ainsi que l'utilisation du nombre de zéros de l'exposant $\gamma_{N_c}^{(L)}(s)$ afin de caractériser les différents groupes de L .

Le chapitre VI [P4] est organisé comme suit: La généralisation du scalings du t et N_c fini du CGF pour systèmes avec grand L est fait en Sec. VI.2.2. Nous utilisons ces résultats dans Sec. VI.3 où nous vérifions la validité du t^{-1} - et N_c^{-1} -scalings (Sec. VI.3.1), leur comportement dans la dynamique modifié par s (Sec. VI.3.2) ainsi que l’applicabilité du scaling method (Sec. VI.3.3) pour un processus de contact avec $L = 100$ sites. Cette analyse est généralisée en Sec. VI.4 où nous caractérisons les échelles de t et N_c fini de la CGF dans le plan $s - L$. Avant de présenter nos conclusions en Sec. VI.6, nous discutons des effets de la transition de phase dynamique sur les scalings en Sec. VI.5.

Alternativement aux méthodes mentionnées au début de cette thèse, on peut faire usage d’une approche complètement différente afin d’étudier des événements rares. Ceci est l’étude empirique des modèles qui se cachent derrière les **données** correspondant à certains phénomènes naturels ou sociaux (e.g., tremblements de terre, marchés boursiers, météo, épidémies, etc). Dans un contexte de séries temporelles financières, ces modèles sont connus comme **stylized facts** ou **seasonalities** [93–99] et les rares événements d’intérêt pourraient correspondre, par exemple, à des chutes de marché ou à des bulles financières [100, 101]. Ces stylized facts ont la caractéristique d’être communs et persistants sur différents marchés, périodes de temps et actifs, éventuellement [99] parce que les marchés fonctionnent en synchronisation avec les activités humaines qui laissent une trace dans les séries temporelles financières.

Cependant, l’utilisation de la “bonne horloge” pourrait être d’une importance primordiale lorsqu’il s’agit de propriétés statistiques et les modèles pourraient varier en fonction si nous utilisons des données quotidiennes ou “intra-day data” et “event time”, temps commercial ou des intervalles de temps arbitraires (e.g., $T = 1, 5, 15$ minutes, etc). Par exemple, il est un fait bien connu que les distributions empiriques des rendements financiers et log-returns sont “fat tailed” [102, 103]. Cependant, comme on augmente l’échelle de temps du fat-tail la propriété devient moins prononcée et la distribution approche la forme gaussienne [104]. Comme l’a indiqué dans Ref. [96], le fait que la forme de la distribution change avec le temps indique clairement que le processus aléatoire sous-jacent aux prix doit avoir une structure temporelle non triviale.

Dans un travail précédent, Allez et al. [99] a établi plusieurs nouveaux stylized facts concernant les intra-day seasonalities de la dynamique des stocks individuels et transversaux. Cette dynamique est caractérisée par l’évolution des moments de ses retours au cours d’une journée type. Basé sur les travaux de Allez et al. [99] et Kaisoji [100], au chapitre VII, nous effectuons une analyse statistique sur les rendements et les prix relatifs des CAC 40 et S&P 500. Nous analysons les **intra-day seasonalities** de la dynamique du individuels et transversal stocks en le caractérisant par l’évolution des moments des retours (et des prix relatifs) au cours d’une journée typique. Pour “single stock intra-day seasonalities” nous nous référons au comportement moyen des moments des retours (et prix relatifs) d’un stock moyen dans une journée moyenne. De même, la cross-sectional intra-day seasonality n’est pas plus que le comportement moyen d’un moment d’index. Nous présentons ces saisonnalités pour les retours (Figs. VII.2 et VII.3) et prix relatifs (Figs. VII.7 et VII.8) et comparé la moyenne des stocks de la volatilité des stocks individuels [$\sigma_\alpha(k)$], la moyenne temporelle de la cross-sectional volatility $\langle \sigma_d(k, t) \rangle$ et la valeur absolue moyenne du equi-weighted index $\langle |\mu_d| \rangle$ (Figs. VII.4 et VII.9).

Notamment, dans le cas des retours, ces modèles dépendent réellement de la taille de la boîte. This fact is well illustrated with 5 différentes valeurs de la taille de la boîte à travers de Fig. VII.11 pour les volatilités et Fig. VII.12 pour le kurtosis dans lequel son inversé U-pattern est évident au moment nous considérons petites tailles de boîte. Dans le cas des prix relatifs, les volatilités présentent également le même type de intra-day pattern (Fig. VII.9), mais contrairement aux retours, il est indépendant de la taille de la boîte, et l'indice que nous considérons, mais caractéristique pour chaque indice. Nous suggérons dans Sec. VII.6 comment cette indépendance de taille de boîte le intra-day patterns en prix relatif pourrait être utilisé pour caractériser **atypical days** pour les index et **anomalous behaviors** en stocks. Ceci est présenté dans Figs. VII.13 et VII.14 où nous avons présenté nos intra-day seasonalites pour le (a) moyenne et (b) la volatilité en bleu et les respectifs les cross-sectional moments pour 3 jours (et les moments de stock unique pour 3 stocks) pris au hasard en bleu clair et nous avons vu comment le comportement moyen de leurs moments se déplacent avec nos intra-day patterns ce qui n'était pas le cas pour la journée 11 et le stock 228. Comme cette thèse est axée sur l'algorithme de clonage, nous avons préféré laisser cette étude dans le dernier chapitre VII: **Intra-day Seasonalities in High Frequency Financial Time Series** [P0].

Comme déjà suggéré par le placement de citations à côté des chapitres, séparé de **Introduction**, où nous établissons nos définitions, le reste de cette thèse est basée sur les résultats qui sont apparus dans **Publications** produit pendant ce programme de doctorat. La recherche actuelle et prospective, ainsi que quelques questions ouvertes, sont présentées dans le **Perspectives** après le **Conclusion**.

Contents

	Résumé	v
	Preface	xxi
I	Introduction	1
I.1	LDT: From Boltzmann to Cloning Algorithms	1
I.2	Discrete and Continuous Master Equation	5
I.2.1	Conservation of Probability and Equilibrium States	6
I.3	Master Equation Matrix Form	7
I.3.1	Conservation of Probability and Equilibrium States Revisited	7
I.4	Solution of the Master Equation	8
I.4.1	Interpretation	8
I.5	Large Deviations of Time-Extensive Observables	9
I.6	The s -modified Dynamics	11
I.6.1	ψ as the Largest Eigenvalue of W_s	11
I.6.2	A Mutation-Selection Mechanism	12
I.7	Continuous-Time Population Dynamics	13
I.7.1	The Cloning Algorithm	13
I.7.1.1	Non-Constant Population Approach	14
I.7.1.2	Constant-Population Approach	14
I.8	Example Models	14
I.8.1	Annihilation-Creation Dynamics	15
I.8.2	Contact Process	15
II	Discreteness Effects in Population Dynamics	17
II.1	Introduction	17
II.2	Average Population and the LDF	17
II.2.1	Populations Merging	18
II.2.2	Discreteness Effects at Initial Times	18
II.3	Parallel Behavior in Log-Populations	20
II.3.1	Distance between Populations	20
II.3.2	Properties of $D(\hat{N}_i, \hat{N}_j)$	20
II.4	Time Correction in the Evolution of Populations	21
II.4.1	Time Delay Correction	21
II.5	$\Psi(s)$ Before and After the Time Delay	23

II.5.1	Numerical Estimators for $\psi(s)$	23
II.5.2	Comparison between “Bulk” and “Fit” Estimators of $\psi(s)$	24
II.5.3	Relative Distance and Estimator Error	25
II.6	Time Delay Properties	26
II.7	Discussion	27
III	Finite-Time and Finite-Size Scalings. I. Discrete Time	31
III.1	Introduction	31
III.2	Birth-Death Process and the Population Dynamics Algorithm	33
III.2.1	Transition Matrices and the Population Dynamics Algorithm	35
III.2.1.1	Original Dynamics: \mathcal{T}	35
III.2.1.2	Cloning Part: \mathcal{C}	36
III.2.1.3	Maintaining Part: \mathcal{K}	37
III.2.1.4	Total Transition: \mathcal{KCT}	38
III.2.2	Large Deviation Results in the $N_c \rightarrow \infty$ Asymptotics	38
III.2.2.1	Estimator of the Large Deviation Function	38
III.2.2.2	Connection between the Distribution Functions of the Population and of the Original System	39
III.2.2.3	Justification of the Convergence of the Large Deviation Estimator as Population Size becomes Large	40
III.2.3	Systematic Errors due to Finite N_c : Convergence Speed of the Large Deviation Estimator as $N_c \rightarrow \infty$	41
III.2.4	Remarks	42
III.2.4.1	Relaxing the Condition $dt = \Delta t$	42
III.2.4.2	A Continuous-Time Cloning Algorithm	42
III.3	Stochastic Errors: Large Deviations of the Population Dynamics	43
III.4	Example: A Simple Two-State Model	45
III.4.1	Systematic Errors	45
III.4.2	Stochastic Errors	46
III.4.3	A Different Large Deviation Estimator	47
III.5	Discussion	48
IV	Finite-Time and Finite-Size Scalings. II. Continuous Time	49
IV.1	Introduction	49
IV.2	CGF Estimator: Constant-Population Approach	50
IV.3	Finite-Time and Finite- N_c Behavior	50
IV.3.1	Finite-Time Scaling	50
IV.3.2	Finite- N_c Scaling	52
IV.4	Finite-Time and Finite- N_c Scaling Method	53
IV.4.1	The Scaling Method	54
IV.4.2	Application to the Contact Process	55
IV.5	Issues on an Analytical Approach	57
IV.6	Conclusions	58

V	Fluctuations of CGF Estimator	59
V.1	Central Limit Theorem	59
V.2	Logarithmic Distribution of CGF Estimator	61
V.3	A Different CGF Estimator	62
VI	Finite-Time and Finite-N_c Scalings in the Large-L Limit	65
VI.1	Introduction	65
VI.2	Finite Scalings of the Large Deviation Function Estimator	66
VI.2.1	Large-Time and Large- N_c Limit	66
VI.2.2	Scalings in the Large- L Limit	66
VI.2.2.1	Determination of the Exponents γ_t & γ_{N_c}	67
VI.3	Finite Scalings for a Large- L Contact Process	68
VI.3.1	Finite-Time and Finite- N_c Scalings	68
VI.3.2	Exponents Characterization & s -Dependence	69
VI.3.3	Implementation of the Scaling Method	69
VI.4	L -Dependence of the Finite Scalings	71
VI.4.1	Characterization of the Exponent $\gamma_t(s, L)$	71
VI.4.2	Characterization of the Exponent $\gamma_{N_c}(s, L)$	72
VI.5	Dynamical Phase Transition and Finite-Scalings	73
VI.6	Conclusion	75
VII	Intra-day Seasonalities in Financial Time Series	77
VII.1	Introduction	77
VII.2	Definitions	77
VII.2.1	Single Stock Properties	79
VII.2.2	Cross-Sectional Stock Properties	79
VII.3	Intra-day Seasonalities for Returns	80
VII.3.1	Single Stock Intra-day Seasonalities	80
VII.3.2	Cross-Sectional Intra-day Seasonalities	81
VII.3.3	U-Pattern Volatilities	82
VII.3.4	Intra-day Seasonalities in the Stock Correlation	82
VII.4	Intra-day Seasonalities for Relative Prices	84
VII.4.1	Single Stock Intra-day Seasonalities	84
VII.4.2	Cross-Sectional Intra-day Seasonalities	84
VII.4.3	C-Pattern Volatilities	85
VII.5	Intra-day Patterns and Bin Size	86
VII.6	Intra-day Abnormal Patterns	88
VII.7	Discussion	90
VIII	Conclusion	91
IX	Perspectives	93
	Publications	95
	Bibliography	97

List of Figures

II.1	Log-Populations	19
II.2	Log-Populations Distance	20
II.3	Average Distance between Log-Populations	21
II.4	Time-Delayed Log-Populations	22
II.5	Delayed Log-Populations Variance	23
II.6	Numerical CGF Estimations	24
II.7	“Bulk” and “Fit” Estimations	25
II.8	Estimator Relative Distance	26
II.9	Estimator Error	26
II.10	Time Delay Variance	27
II.11	Distribution of Time Delays	28
III.1	Schematic Picture of the Population Dynamics Algorithm	32
IV.1	Time Evolution of the CGF Estimator	51
IV.2	Population-Size Evolution of the CGF Estimator	53
IV.3	CGF Estimator Distance	54
IV.4	Continuous-Time CGF Estimator	55
IV.5	Projection of the CGF Estimator	56
IV.6	Relative Systematic Error	57
V.1	Distribution of the CGF estimator	59
V.2	Collapse of the CGF Estimator Distribution	60
V.3	Logarithmic Distribution	62
V.4	Two Different CGF Estimators	63
VI.1	Finite Scalings in the Large-L Limit	67
VI.2	CGF Estimator in the Large-L Limit	68
VI.3	Finite Population-Size Scaling in the Large-L Limit	70
VI.4	CGF Estimators in the Large-L Limit	71
VI.5	Finite-Time Scaling on the Plane s-L	72
VI.6	Finite Population-Size Scaling on the Plane s-L	74
VII.1	Intra-Day Asynchronous Financial Time Series	78
VII.2	Single Stock Intra-Day Seasonalities: Returns	80
VII.3	Cross-Sectional Intra-day Seasonalities: Returns	81

VII.4	U-Pattern Volatilities	82
VII.5	Largest Eigenvalues Structure	83
VII.6	Top Eigenvalues	84
VII.7	Single Stock Intra-day Seasonalities: Relative Prices	85
VII.8	Cross-Sectional Intra-day Seasonalities: Relative Prices	86
VII.9	C-Pattern Volatilities	87
VII.10	Time Average of the Cross-Sectional Skewness	87
VII.11	Bin Size Dependence in the U-Pattern Volatilities	88
VII.12	Bin Size Dependence in the Inverted U-Pattern Kurtosis	88
VII.13	Atypical Day	89
VII.14	Anomalous Stock Behavior	90

Preface

The occurrence of **rare events** can vastly contribute to the evolution of physical systems because of their potential dramatic effects. Their understanding has gathered a strong interest and, focusing on stochastic dynamics, a large variety of numerical methods have been developed to study their properties [1, 7, 8]. They range from importance sampling [9], adaptive multilevel splitting [10] to transition path sampling [11–14] and from “go with the winner” algorithms [15, 16] to discrete-time [18] or continuous-time [19] population dynamics [7, 17]. These methods have been generalized to many contexts [20–24]. In Physics, those are being increasingly used in the study of complex systems, for instance in the study of current fluctuations in models of transport [25–27], glasses [12], protein folding [28] and signalling networks [29, 30]. Mathematically, the procedure amounts to determining a large deviation function (LDF) associated to the distribution of a given trajectory-dependent observable, which in turns can be reformulated in finding the ground state of a linear operator [31], a question common to both statistical and quantum physics [32].

In this thesis, we will give particular attention to **population dynamics algorithms** [7, 17–19, 33] which aim at studying rare trajectories by exponentially biasing their probability. The idea of population dynamics is to translate the study of a class of rare trajectories (with respect to a determined global constraint) into the evolution of several copies of the original dynamics, with a local-in-time selection process rendering the occurrence of the rare trajectories typical in the evolved population. The distribution of the class of rare trajectories in the original dynamics is related with the exponential growth (or decay) of the population of clones of the system and LDF can be estimated from its growth rate.

The numerical procedures aimed at simulating rare events efficiently, using a population dynamics scheme are commonly referred as **cloning algorithms**. In such algorithms, copies of the system are evolved in parallel and the ones showing the rare behavior of interest are multiplied iteratively [7, 16–19, 23, 32, 33, 62–70]. Some of the limitations and associated improvements of the population dynamics algorithms have been studied in [22, 85, 105, 106].

Two versions of such algorithms exist: the **non-constant** and the **constant total population** approaches. For the last one, a uniform pruning/cloning is applied on top of the cloning dynamics so as to avoid the exponential explosion or disappearance of the population. While the later version is obviously more computer-friendly, the former version presents interesting features: First, it is directly related to the evolution of biological systems (stochastic jumps representing mutations, selection rules being interpreted as Darwinian pressure); second, the uniform pruning/cloning of the population, although unbiased, induces correlations in the dynamics that one might want to avoid; last, in some situations where the selection

rates are very fluctuating, the constant-population algorithm cannot be used in practice because of finite-population effects (population being wiped out by a single clone), and one has to resort to the non-constant one. An example of the implementation of this version can be found in Ref. [73].

In chapter II: **Discreteness Effects in Population Dynamics** [P1], we apply the non-constant population algorithm in order to analyze numerically the small population-size effects in the initial transient regime. These effects play an important role for the numerical determination of large deviation functions of additive observables for stochastic processes. The LDF estimation, in this case, reduces to the determination of the growth rate of a population, averaged over many realizations of the dynamics. However, this averaging is highly dependent not only on the number of realizations, and on the initial population size but also on the cut-off time (or population) considered to stop their numerical evolution. This may result in an over-influence of **discreteness effects at initial times**, caused by small population size. We show how to overcome these effects by introducing a (realization-dependent) time delay in the evolution of populations, additional to the discarding of the initial transient regime of the population growth where these discreteness effects are strong. We show that the improvement in the estimation of the large deviation function comes precisely from these two main contributions.

Apart from the population-constraint approaches we just mentioned, the selection mechanisms within the cloning algorithm can be implemented in different ways. One of them, along with each evolution of the copies (**Continuous-Time**) [7, 17, 19] or alternatively, for each pre-fixed time-interval (**Discrete-Time**) [18]. The important differences between both techniques are discussed in Secs. III.2.4.2 and IV.5.

The algorithm proposed by Giardinà et al. [7, 17–19, 33, 70] (a discrete-time approach) is used to evaluate numerically the LDF of additive (or “time-extensive”) observables in Markov processes [1, 83]. The LDF is obtained as the exponential growth rate that the population would present if it was not kept constant. It has been applied to many physical systems, including chaotic systems, glassy dynamics and non-equilibrium lattice gas models, and it has allowed the study of novel properties, such as the behavior of breathers in the Fermi-Pasta-Ulam-Tsingou chain [33], dynamical phase transitions in kinetically constrained models [39], and an additivity principle for simple exclusion processes [60, 84]. Under this approach, the corresponding LDF estimator is in fact valid only in the limits of infinite simulation time t and infinite population size N_c . The usual strategy that is followed in order to obtain those limits is to increase the simulation time and the population size until the average of the estimator over several realizations does not depend on those two parameters, up to numerical uncertainties.

While the method has been used widely, there have been less studies focusing on the analytical justification of the algorithm. Moreover, it introduces two additional parameters into consideration: the population size N_c and the simulation time t , both of which affect considerably the accuracy of the LDF estimation. Even though it is heuristically believed that the LDF estimator converges to the correct result as the number of copies N_c increases, there is no proof of this convergence. Related to this lack of proof, although we use the algorithm by assuming its validity, we do not have any clue how fast the estimator converges as $N_c \rightarrow \infty$. We discuss this convergence performing an analytical study in discrete time in

chapter III and using a numerical approach in continuous time in chapter IV. It is important to remark that the two versions of the algorithm (discrete- and continuous-time) differ on a crucial point which implies that an extension of the analysis developed in chapter III cannot be done straightforwardly in order to comprehend the continuous-time case in chapter IV.

In chapter III: **Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: I. Analytical Study using a Birth-Death Process [P2]**, in order to study the numerical errors of this algorithm, we explicitly devise a stochastic birth-death process that describes the time evolution of the population probability. From this formulation, we derived the finite- N_c and finite- t scalings of the systematic errors of the LDF estimator, showing that these behave as $1/N_c$ and $1/t$ in the large- N_c and large- t asymptotics respectively. In principle, knowing the scaling *a priori* means that the asymptotic limit of the estimator in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits may be interpolated from the data at finite t and N_c . However, whether this idea is actually useful or not is a non-trivial question, as there is always a possibility that onset values of the N_c^{-1} - and t^{-1} -scalings are too large to use these scalings.

In chapter IV: **Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: II. Numerical Approach in Continuous Time [P3]**, we consider a continuous-time version of the population dynamics algorithms [17, 19]. We show numerically that one can indeed make use of these properties in order to devise an original and simple method that takes into account the exact scalings of the finite- t and finite- N_c corrections in order to provide significantly better LDF estimators (**scaling method**). We study the fluctuations of the standard estimator in chapter V [P3] and additionally, we discuss an alternative way of defining the LDF estimator. However, the validity of these scalings and the method efficiency is proved in chapter IV only in cases for which the number of sites L (where the dynamics occurs) was small: a simple two-states annihilation-creation dynamics (in one site) and a one-dimensional contact process [38, 107, 108] (with $L = 6$ sites).

In chapter VI: **Breakdown of the Finite-Time and Finite- N_c Scalings in the Large- L Limit [P4]**, we complement the results presented in chapter IV by extending the analysis of the finite-scalings of the LDF to a large- L contact process. The dependence of these scalings with the number of sites is analyzed by introducing the exponents γ_t and γ_{N_c} . The generalized $t^{-\gamma_t}$ - and $N_c^{-\gamma_{N_c}}$ -scalings allow to characterize the behavior in the large- L limit where we verify that t^{-1} and N_c^{-1} -scalings are no longer valid.

Alternatively to the methods mentioned at the beginning of this introduction, one can make use of a completely different approach in order to study rare events. This is the empirical study of the patterns that hide behind the **data** corresponding to some natural or social phenomena (e.g., earthquakes, stock markets, weather, epidemics, etc). In a financial time series context, these patterns are known as **stylized facts** or **seasonalities** [93–99] and the rare events of interest could correspond, for example, to market crashes or financial bubbles [100, 101]. These properties have the characteristic of being common and persistent across different markets, time periods and assets possibly [99] because markets operate in synchronization with human activities which leave a trace in the financial time series.

Following specially the works by Allez et al. [99] and Kaisoji [100], in chapter VII we perform a statistical analysis over the returns and relative prices of the CAC 40 and the S&P 500. We analyze the **intra-day seasonalities** of single and cross-sectional stock dynamics by characterizing it by the evolution of the moments of the stock returns (and relative prices) during a typical day. We show the bin-size and index independence for the case of the relative prices but not for the returns. However, we suggest how this fact could be used in order to characterize **atypical days** for indexes and **anomalous behaviours** of stocks. As this thesis is focused on the cloning algorithm, we have preferred to leave this study in the last chapter VII: **Intra-day Seasonalities in High Frequency Financial Time Series** [P0].

As already suggested by the placement of citations next to the chapters, apart from the **Introduction**, where we establish our definitions, the rest of this thesis is based on results that have appeared in **Publications** produced during this PhD program. The current and prospective research, as well as some open questions, are presented in the **Perspectives** after the **Conclusion**.

«Le secret de la liberté est d'éclairer les hommes,
comme celui de la tyrannie est de les retenir dans l'ignorance»

Maximilien Robespierre

I – Introduction

I.1 Large Deviation Theory: From Boltzmann to Cloning Algorithms

The theory of large deviations deals with probabilities of rare events [1–3]. These probabilities or fluctuations have the characteristic of decaying exponentially as a function of some parameter (like the time or the temperature) meaning that, as the parameter becomes larger, the event becomes less probable [4]. They are of important interest in many fields like statistics, queuing theory, finance, engineering and in equilibrium and non-equilibrium statistical physics. From a practical point of view, large deviation theory can be seen as a collection of methods which allow to determine if a large deviation principle exists for a given random variable and to determine its respective rate (or large deviation) function (LDF).

The first large deviation result is due to Boltzmann in 1877 [109, 110]. He showed how the relative entropy expresses the asymptotic behavior of multinomial probabilities presenting the entropy as a bridge between the microscopic level, of physical interactions, and a macroscopic one, where the physics laws are formulated. This constituted a probabilistic interpretation of the Second Law of Thermodynamics [109] and the basis which led to the development of the classical equilibrium statistical mechanics [111]. Ellis [110] describes this interpretation as “a revolutionary moment in human culture during which both statistical mechanics and the theory of large deviations were born”.

Some large deviation results like Cramér’s theorem [112] (who initiated a mathematical theory of large deviations in the 30’s), Chebyshev’s inequality [83] and the Sanov’s theorem [113], were also anticipated by Boltzmann [109, 110]. However, there was not a unified or general framework that dealt with them until the 60’s and 70’s when this theory was developed by Donsker and Varadhan [114–118] and by Freidlin and Wentzell [119].

In some cases, the large deviation principle can be determined directly from the probability distribution of a random variable. This is done by deriving a large deviation approximation using Stirling’s or other asymptotic formulae. However, a more general result was provided by the Gärtner-Ellis theorem [120, 121] which is the product of a result proved by Gärtner [120] and later generalized by Ellis [110, 121–124] which explicitly refers to the construction of the currently adopted large deviation principle. This was inspired from the work of Varadhan [118]. However, meanwhile the Gärtner-Ellis theorem is used to prove the existence of a large deviation principle and the determination of the corresponding rate function from the knowledge of the scaled cumulant generating function (CGF), the Varadhan’s theorem is used to calculate the CGF knowing the rate function. Moreover, the contraction principle [117] introduced by Donsker and Varadhan allows to compute a rate function from the knowledge of another rate function. A direct application of the Gärtner-Ellis theorem

or of the contraction principle allows to formulate a large deviation principle for many problems like sums of (binary, symmetric Levy, totally skewed Levy, etc) i.i.d. random variables (Cramér theorem [112]), random vectors (Sanov's theorem [113]), Markov processes [114–118], among others.

Donsker and Varadhan defined three levels of large deviation results [124]. Level-1 is the level of sample means, Level-2 is the level of empirical distributions and Level-3 is the level of the empirical processes. The latter distributions can be derived from the former using the contraction principle [117]. For example, the Level-2 rate function of Markov chains can be derived by contracting the large deviations of the pair empirical matrix [1].

Large deviation theory has been suggested to be a generalization of the Central Limit Theorem [125, 126] because it provides information not only about the small but also its large fluctuations of a random variable far away from its typical values. It is also considered that it extends the Law of Large Numbers [127] providing information of how fast a random variable converges in probability to its mean. In fact, the existence of a Law of Large Numbers for a random variable is a good sign that there also holds a large deviation principle and also it can be used as a departure point [128, 129].

Some physicists consider large deviation theory as a natural generalization of the entropy-probability relation fully exploited by Einstein in his theory of thermodynamic fluctuations [130, 131]. According to this theory, the probabilities can be expressed in terms of entropy functions. This fact allows to use large deviation theory to understand the foundations of statistical mechanics. In this way, large deviation theory explains for example why the entropy and free energy are related through a Legendre transform and why equilibrium states can be calculated via extremum principles (maximum entropy for the microcanonical ensemble and minimum free energy for the canonical ensemble) generalizing them to arbitrary macrostates and arbitrary many-particle systems [1]. On the other hand, the well-known maximum entropy principle of Jaynes [132–134] can be obtained by considering the Level-2 large deviations of systems of independent particles. Einstein fluctuation formula was used by Varadhan and Donsker [114–118] as the basis of the standard theory for static equilibrium fluctuations. Additionally, Onsager and Machlup [135–138] used it in order to propose a reformulation of linear fluctuation theory about equilibrium.

The implementation of large deviation techniques for studying the equilibrium properties of many-particle systems described at a probabilistic level by statistical mechanical ensembles has its roots in the work of Ruelle [139], Lanford [127], and especially Ellis [110, 122, 124]. Ellis is considered to provide (in Ref. [124]) the most complete framework in which the large deviation theory is introduced to physics stressing in the connections between probability, large deviations and equilibrium statistical mechanics. The first work on large deviations and statistical mechanics is attributed to Lanford [127] which uses concepts from large deviation theory to explain the fact that while matter is extremely complicated at microscopic level, it can be described at the macroscopic level by a small number of parameters. Moreover, using a large deviation approach on the ensembles in statistical mechanics, the study of equilibrium states and their fluctuations can be reduced to the study of properly defined rate functions (entropy functions) [1, 124, 140–145]. Many other links between statistical mechanics and large deviations also has been discussed by Lewis, Pfister, and Sullivan [140–145], as well as Oono [146], Amann [147] and in several reviews [1–3, 110, 122].

Behind the application of large deviation theory to equilibrium statistical physics lies the idea that outcomes of a macrostate involving n particles should concentrate in probability around certain stable or equilibrium values even though the state of the particles is described by a random variable. In many cases the outcomes satisfy a large deviation principle due to the probability of observing a departure from these equilibrium values is exponentially small with n . Thus, in order to describe the macrostate of a large many-particle system it is only necessary to know its equilibrium values [1, 2, 124, 140–147].

Additionally to the equivalences already mentioned, we have that the thermodynamic limit is a large deviation limit, and the free energy is the equivalent of a scaled cumulant generating function [1–3]. Moreover, the Legendre transform which connects the entropy and free energy in thermodynamics is nothing but the Legendre-Fenchel transform connecting the rate function and the scaled cumulant generating function in the Gärtner-Ellis theorem [120, 121] and in Varadhan’s theorem [124]. The equilibrium properties of mean field models can be studied as Level-2 or directly at the Level-1 of large deviations. Maximum entropy principles have been applied successfully to these models which consider all-to-all coupling between particles like the Curie-Weiss model [122, 124, 148] and its parent model, the Potts model [122, 149–151], the Blume-Emery-Griffiths model [152–154], the mean-field Hamiltonian model [155], as well as mean-field versions of the spherical model [156, 157], and the ϕ^4 model [158–160].

The microcanonical and canonical ensembles differ from each other in the way their respective microstates are weighted. In the microcanonical ensemble, the control parameter is the energy (or the mean energy), and the microstates are distributed with the same probabilistic weight if they have the same value of control parameter. On the other hand, in the canonical ensemble, the control parameter is the inverse temperature, the probability measure is the Gibbs measure and the rate functions are the macrostate free energies (which are the basis of the Ginzburg–Landau theory of phase transitions [161]). Some examples of results derived in the canonical ensemble can be found in Refs. [110, 122, 124, 152]. The thermodynamic equivalence (or non-equivalence) between the microcanonical and canonical ensembles is related to the concavity of the entropy. This comes from the fact that the free energy can always be obtained as the Legendre–Fenchel transform of the entropy, but the entropy can be obtained as the Legendre–Fenchel transform of the free energy only when the entropy is concave. Moreover, the Gärtner-Ellis theorem can be reformulated (in a physical way) as : “If there is no first-order phase transition in the canonical ensemble, then the microcanonical entropy is the Legendre transform of the canonical free energy” [1]. Examples of models with non-concave entropies are the mean-field Blume-Emery-Griffiths model [152–154], the mean-field Potts model [151, 162], some models of plasmas [163] and 2D turbulence [164–166], as well as models of gravitational systems [167, 168]. This thermodynamic equivalence is translated in terms of Gibbs’s entropy and Boltzmann’s entropy in the thermodynamic limit, where the Gibbs entropy is equal (up to a constant) to the Boltzmann entropy evaluated at the equilibrium mean energy value [1, 169].

Large deviation theory is becoming the standard formalism to study non-equilibrium systems [1, 170], modelled in general by stochastic differential equations or Markov processes [171]. They have the characteristic of evolving dynamically in time or to be maintained in out-of-equilibrium steady states under the application of an external force. It has been suggested that large deviation theory provides the proper basis for building a theory of non-equilibrium systems [146, 172]. This requires, of course the inclusion of the time

in the large deviation analysis and the consideration that we do not know the underlying probability distribution states as the concept of ensemble is not defined for non-equilibrium systems. In spite of this, many large deviation principles have been derived for example for Markovian models of interacting particles [36, 108, 170, 173] such as the exclusion process, the zero-range process and their different variants [5, 25, 84, 174–178] in some cases at the level of density field [176, 179–182] or at the level of current [5, 84, 177].

The so called fluctuation theorem [183] and other non-equilibrium work relations [184] concern the large deviations of work [185]. It states that the probability of observing an entropy production opposite to that dictated by the second law of thermodynamics decreases exponentially. It was first proposed and tested numerically in 1993 [186], the first mathematical proof was in 1994 [187] and it was verified experimentally in 2002 [188]. Gallavotti and Cohen [189, 190] used these results as a basis to proof a fluctuation theorem for the entropy rate of chaotic deterministic systems. This was extended later to general Markov processes by Kurchan [191], Lebowitz and Spohn [192], and Maes [193]. These results have inspired several experimental studies of fluctuation relations that appear for example, particles immersed in fluids [188, 194], electrical circuit [195, 196], granular media [197–201], turbulent fluids [202, 203], and the effusion of ideal gases [204], among other systems.

Other applications of large deviation theory are related to multifractals [205–208], chaotic systems [209–212], disordered systems, and quantum systems. Multifractal analysis can be seen as a large deviation theory of self-similar measures [213–215]. Dynamical systems often give rise to large deviation principles without a perturbing noise. Their study in the context of chaotic systems and ergodic theory is the subject of the so-called thermodynamic formalism [205, 216] developed by Ruelle [217, 218] and Sinai [219, 220]. This formalism introduces the topological pressure and the structure function which play the role of the CGF implying a direct connection between dynamical systems and large deviation theory [216, 221–224]. Additionally, large deviation principles can be obtained when studying disordered and quantum systems, for example, in random walks in random environments [225–228], spin glasses [229–231], boson gases [232–234], quantum gases [235, 236], and quantum spin systems [237–239].

In this point, it is important to remark that only in few simple cases is it possible to obtain exact and explicit expressions for the rate functions [5, 6]. For most stochastic processes, the evaluation of these functions is done by using analytical approximations and numerical methods [1–3]. They range from importance sampling [9], adaptive multilevel splitting [10] to transition path sampling [11–14] and “go with the winner” algorithms [15, 16]. Kurchan and his collaborators generalized a procedure used previously to study rare events in chemical reactions [240–242] in order to compute large deviation functions in dynamical systems [33], discrete-time [17, 18] and continuous-time [17, 19] population dynamics [7], being generalized then to many contexts [20–24].

The numerical procedure introduced by Giardinà, Kurchan and Peliti [18] overcomes the difficulty of observing the fluctuations of an observable (whose probability decreases exponentially in time) for discrete-time Markov chains. It was known that the large deviation function can be obtained as the largest eigenvalue of a evolution matrix of a modified dynamics [17, 18] which can be computed numerically [5, 6, 34] specially for small systems as the evolution matrix is exponentially large in the system size. Later, a modification of the procedure was proposed [19, 35] for which the time discretization issues of the original approach [18] are bypassed with a direct continuous-time approach. The evolution of the sys-

tem was represented by a population dynamics of the type of the diffusion Monte Carlo [32]. This **cloning algorithm** was applied to successfully compute the large deviations of the total current in the symmetric and asymmetric exclusion process [36, 37], and of the activity in the contact process [38].

Among its applications, Garrahan et al. [39, 40] analyzed the dynamics of kinetically constrained models [41–55] of glassy systems [56–58] by analyzing the statistics of trajectories of the dynamics. They showed that these models exhibit a first-order dynamical transition between active and inactive dynamical phases. It also has been used to study symmetries in fluctuations far from equilibrium [59] and in transport models [21, 22, 60]. These studies allow not only to test the predictions of fluctuating hydrodynamics [21, 61], but also the limits of the method itself [22]. It also has been suggested [7] that the method could be applied to study in detail the possible future and past evolution of planetary systems, and also the self-organization of the stability of our solar system.

In this chapter, we introduce the cloning algorithm. This method will be used through the thesis in order to analyze the issues previously mentioned in the [Preface](#). We start from the construction of the master equation, its solution and interpretation. Then, we introduce the large deviations of additive observables and the s -modified dynamics. We show how to estimate these large deviations from the population dynamics interpretation of the modified dynamics or from the largest eigenvalue of the modified evolution equation. Finally, we present the example models used for this analysis: a simple two-state annihilation-creation dynamics, and a contact process on a one-dimensional periodic lattice.

I.2 Discrete and Continuous Master Equation

Consider a system whose dynamics occurs in jumps between configurations. We denote the set of available configurations $\{C\}$ and the transition rates between them $W(C \rightarrow C')$, setting $W(C \rightarrow C) = 0$. We are interested in describing the probability for the system to be in configuration C at time t , that we denote $P(C, t)$. In order to do that, we start from the following considerations: During the time interval dt the system either stays in the same configuration C or changes configuration to C' . Thus, the transition probability p between configurations can be expressed in terms of dt and W as

$$p(C \rightarrow C') = dt W(C \rightarrow C') \quad \forall C' \neq C, \quad (\text{I.1})$$

$$p(C \rightarrow C) = 1 - dt \sum_{C'} W(C \rightarrow C'). \quad (\text{I.2})$$

The tendency of the dynamics to leave from a configuration C to any other is captured in the escape rate $r(C)$, defined as

$$r(C) = \sum_{C'} W(C \rightarrow C') \quad (\text{I.3})$$

which appears in the second term of Eq. (I.2). The probability of being in configuration C at time $t + dt$ is none other than the probability of being at configuration C given that the system was at configuration C' at time t , plus the probability of having remained at configuration C between time t and $t + dt$, which can be expressed as

$$P(C, t + dt) = P(C, t + dt | C', t) + P(C, t + dt | C, t), \quad (\text{I.4})$$

where each term in Eq. (I.4) is given by

$$P(C, t + dt | C', t) = \left(\sum_{C'} dt W(C' \rightarrow C) \right) P(C', t), \quad (\text{I.5})$$

$$P(C, t + dt | C, t) = \left(1 - dt \sum_{C'} W(C \rightarrow C') \right) P(C, t). \quad (\text{I.6})$$

The transition rules (I.1) and (I.2) are equivalent to the discrete evolution equation (I.4) which after replacing Eqs. (I.5) and (I.6) reads

$$P(C, t + dt) = \sum_{C'} \left[dt W(C' \rightarrow C) P(C', t) + \left(1 - dt W(C \rightarrow C') \right) P(C, t) \right]. \quad (\text{I.7})$$

Equation (I.7) is also known as the discrete master equation. The second right term of Eq. (I.7) ensures probability conservation as $\sum_C P(C, t + dt) = \sum_C P(C, t) = 1$. Taking the limit $dt \rightarrow 0$ in Eq. (I.7) and replacing Eq. (I.3), we obtain its analogous version in continuous time

$$\partial_t P(C, t) = \sum_{C'} \left[W(C' \rightarrow C) P(C', t) - r(C) P(C, t) \right]. \quad (\text{I.8})$$

I.2.1 Conservation of Probability and Equilibrium States

The probability $P(C, t)$ is conserved at all times, i.e.,

$$\partial_t \sum_C P(C, t) = 0. \quad (\text{I.9})$$

The steady state solution P_{st} of Eq. (I.8), which is obtained from $\partial_t P(C, t) = 0$, verifies the global balance condition

$$\sum_{C'} W(C \rightarrow C') P_{\text{st}}(C) = \sum_{C'} W(C' \rightarrow C) P_{\text{st}}(C'),$$

for all C . If the steady state also satisfies the detailed balance condition

$$W(C \rightarrow C') P_{\text{eq}}(C) = W(C' \rightarrow C) P_{\text{eq}}(C'), \quad (\text{I.10})$$

for all C and C' , then the steady state is an equilibrium state of the system, i.e., $P_{\text{st}} = P_{\text{eq}}$. The last condition implies that there is no current of probability in the steady state, and that the dynamics starting from P_{eq} is reversible. This can be seen (using Eq. (I.10)) from

$$W(C_0 \rightarrow C_1) \dots W(C_{K-1} \rightarrow C_K) P_{\text{eq}}(C_0) = W(C_K \rightarrow C_{K-1}) \dots W(C_1 \rightarrow C_0) P_{\text{eq}}(C_K),$$

where the probability density of the history $C_0 \rightarrow \dots \rightarrow C_K$ is the same as its time-reversed history $C_K \rightarrow \dots \rightarrow C_0$.

I.3 Master Equation Matrix Form

In order to study the properties of the master equation it is convenient to introduce the following vector and operator notations [243]: Consider an orthonormal vector space of basis $|C\rangle$ with scalar product $\langle C'|C\rangle = \delta_{CC'}$ where $\langle C|$ is the transpose of $|C\rangle$. A vector is denoted as $|v\rangle = \sum_C v_C |C\rangle$ with $v_C = \langle C|v\rangle$. An operator $\mathbb{A} = \sum_{CC'} \mathbb{A}_{CC'} |C\rangle\langle C'|$ can be represented in matrix form of elements $\mathbb{A}_{CC'} = \langle C|\mathbb{A}|C'\rangle$. Using this notation, the master equation (I.8) takes the linear form

$$\partial_t |P(t)\rangle = \mathbb{W}|P(t)\rangle \quad (\text{I.11})$$

for the probability vector

$$|P(t)\rangle = \sum_C P(C, t) |C\rangle.$$

The master operator \mathbb{W} is a matrix of elements

$$(\mathbb{W})_{CC'} = W(C' \rightarrow C) - r(C) \delta_{CC'}, \quad (\text{I.12})$$

i.e.,

$$(\mathbb{W})_{CC'} = \begin{cases} W(C \rightarrow C') & \text{if } C \neq C' \\ -r(C) & \text{if } C = C'. \end{cases}$$

The diagonal elements of matrix (I.12) correspond to waiting times between jumps when the system stays in the same configuration.

I.3.1 Conservation of Probability and Equilibrium States Revisited

Using the vector notation introduced above, the conservation probability (I.9) reads as $\sum_C (\mathbb{W})_{CC'} = 0$ which with $\langle -| = \sum_C \langle C|$ becomes

$$\langle -|\mathbb{W} = 0, \quad (\text{I.13})$$

meaning that the vector $\langle -|$ is a left eigenvector of \mathbb{W} (of eigenvalue 0). On the other hand, the global balance condition reads

$$\mathbb{W}|P_{\text{st}}\rangle = 0,$$

where the vector $|P_{\text{st}}\rangle$ is the right eigenvector of \mathbb{W} (also of eigenvalue 0). As \mathbb{W} and \mathbb{W}^T have the same spectrum, the conservation of probability ensures the existence of a steady state. Moreover, all the eigenvalues of \mathbb{W} are of real part negative. The detailed balance condition (I.10) can be written in terms of a diagonal operator \hat{P}_{eq} of elements $P_{\text{eq}}(C)$ as

$$\mathbb{W}\hat{P}_{\text{eq}} = \hat{P}_{\text{eq}}\mathbb{W}^T. \quad (\text{I.14})$$

By multiplying Eq. (I.14) by the left by $\langle -|$ and using Eq. (I.13) we verify that $|P_{\text{eq}}\rangle$ is indeed a steady state through $\langle -|\hat{P}_{\text{eq}} = \langle P_{\text{eq}}|$. Moreover from Eq. (I.14) we also have

$$\hat{P}_{\text{eq}}^{-1/2} \mathbb{W} \hat{P}_{\text{eq}}^{1/2} = \hat{P}_{\text{eq}}^{1/2} \mathbb{W}^T \hat{P}_{\text{eq}}^{-1/2},$$

implying that $\mathbb{W}^{\text{sym}} = \hat{P}_{\text{eq}}^{-1/2} \mathbb{W} \hat{P}_{\text{eq}}^{1/2}$ is a symmetric operator (self-adjoint in fact) and it can be diagonalized in an orthonormal basis. Given that \mathbb{W} and \mathbb{W}^{sym} have the same spectrum, the last equation also implies that \mathbb{W} has real eigenvalues.

I.4 Solution of the Master Equation

Given the initial condition $|P_0\rangle = |P(t=0)\rangle = \sum_C P_0(C)|C\rangle$, Eq. (I.11) has as solution

$$|P(t)\rangle = e^{t\mathbb{W}}|P(0)\rangle = \sum_{n \geq 0} \frac{t^n}{n!} \mathbb{W}^n |P(0)\rangle. \quad (\text{I.15})$$

However, as the matrix \mathbb{W} (I.12) has diagonal elements different from zero, Eq. (I.15) does not allow a description as a sum over trajectories of successively different visited configurations. Thus, in order to get rid of the diagonal terms of matrix \mathbb{W} (I.12), it is convenient to transform the master equation (I.8) by defining

$$Q(C, t) = e^{-tr(C)} P(C, t). \quad (\text{I.16})$$

Equation (I.16) verifies

$$\partial_t |Q(t)\rangle = \mathbb{W}_Q(t) |Q(t)\rangle, \quad (\text{I.17})$$

where

$$(\mathbb{W}_Q(t))_{CC'} = W(C' \rightarrow C) e^{t(r(C')-r(C))}.$$

The solution of Eq. (I.17) is given by

$$|Q(t)\rangle = \tau_{exp} \left\{ \int_0^t dt' \mathbb{W}_Q(t') \right\} |Q(0)\rangle,$$

where τ_{exp} is the time-ordered exponential:

$$\tau_{exp} \left\{ \int_0^t dt' \mathbb{W}_Q(t') \right\} = \sum_{K \geq 0} \int_{t_0}^t dt_1 \int_{t_1}^t dt_2 \dots \int_{t_{K-1}}^t dt_K \mathbb{W}_Q(t_K) \dots \mathbb{W}_Q(t_1)$$

with $K \in \mathcal{N}$. Finally, $P(C, t)$ is obtained coming back to Eq. (I.16) as

$$\begin{aligned} P(C, t) &= \sum_{K \geq 0} \sum_{C_0 \dots C_K} \int_{t_0}^t dt_1 \int_{t_1}^t dt_2 \dots \int_{t_{K-1}}^t dt_K \\ &\times r(C_0) e^{-(t_1-t_0)r(C_0)} \dots r(C_{K-1}) e^{-(t_K-t_{K-1})r(C_{K-1})} \times e^{-(t-t_K)r(C_K)} \quad (\text{I.18}) \\ &\times \frac{W(C_0 \rightarrow C_1)}{r(C_0)} \dots \frac{W(C_{K-1} \rightarrow C_K)}{r(C_{K-1})} \times P(C_0, t=0) \end{aligned}$$

with $C_K = C$.

I.4.1 Interpretation

Equation (I.18) allows us to have a better visualization of the process described by the master equation (I.8) [86]. The dynamics occurs in jumps between configurations with transition rates $W(C \rightarrow C')$ on a time window $[0, t]$. The histories of the systems are described by

$$(\vec{C}, \vec{t}) = (C_0, t_0; C_1, t_1; \dots; C_k, t_k; \dots; C_K = C, t_K),$$

where K is the total number of actual jumps between successively distinct configurations denoted by $\vec{C} = (C_0, C_1, \dots, C_k, \dots, C_K = C)$ and $\vec{t} = (t_0, t_1, \dots, t_k, \dots, t_K)$ is the increasing sequence of times at which the jumps occurs (e.g., at time t_k , the system jumps from configuration C_{k-1} to C_k). The factors

$$\frac{W(C_{k-1} \rightarrow C_k)}{r(C_{k-1})}$$

in Eq. (I.18) represent the probability of jumping from configuration C_{k-1} to C_k . Hence, the probability of a history of configurations \vec{C} is given by

$$P(\vec{C}) = \prod_{k=0}^{K-1} \frac{W(C_k \rightarrow C_{k+1})}{r(C_k)}.$$

The system rests in a configuration C_k during an interval $\Delta t_k = t_{k+1} - t_k$ which is distributed with a Poisson law of parameter $r(C_k)$

$$\rho(r(C_k)) = r(C_k) e^{-\Delta t_k r(C_k)}.$$

Thus, the probability density of instants $\{t_k\}_{1 \leq k \leq K}$ of the change of configuration is

$$\prod_{k=0}^{K-1} r(C_k) e^{-\Delta t_k r(C_k)}.$$

Meanwhile the probability of not jumping between times t_K and t is $e^{-(t-t_K) r(C_K)}$. Thus, if we fix the initial configuration C_0 , the probability of the path is given by

$$\mathbb{P}(C_1, t_1; \dots; C_K, t_K \mid C_0, t_0; t) = \prod_{k=0}^{K-1} \frac{W(C_k \rightarrow C_{k+1})}{r(C_{k-1})} \times \prod_{k=0}^{K-1} r(C_k) e^{-\Delta t_k r(C_k)} \times e^{-(t-t_K)}.$$

I.5 Large Deviations of Time-Extensive Observables

Once we have defined our system, we are now interested in the distribution of history-dependent observables and its fluctuations. These dynamical observables are defined as a sum along the history of small contributions for transitions between successive configurations during a time interval $[0, t]$. In general, they are of the form

$$\mathcal{O} = \sum_{k=0}^{K-1} a(C_k, C_{k+1}) + \int_0^t dt' b(C(t')), \quad (\text{I.19})$$

where $C(t')$ is the state of the system at time t' : when $t_k \leq t' < t_{k+1}$, $C(t') = C_k$ ($k = 0, 1, 2, \dots, K-1$) with $t_0 = 0$. The functions a and b describe elementary increments: a accounts for quantities associated with transitions (of state), whereas b does for static quantities. We commonly refer to observables of the form (I.19) for $b = 0$ to ‘type-A’, meanwhile to those for which $a = 0$ to ‘type-B’ observables [40]. A simple example of observables of this form is the dynamical activity K [12, 39, 40, 74–82], which is the number of configuration changes on the time interval $[0, t]$ (in this case one has $a(C, C') = 1$ and

$b \equiv 0$ in Eq. (I.19)). Another, is the current of particles Q [61, 244–250] in a one-dimensional lattice gas, where the value of the observable Q is incremented or decremented at each time a particle jumps to the left or right. This kind of observables contrasts from the static ones which depend only on the configuration of the system at a given time.

The probability density of being in configuration C at time t having observed a value \mathcal{O} of observable is denoted by $P(C, \mathcal{O}, t)$ and is related through the probability distribution of \mathcal{O} at time t , $P(\mathcal{O}, t)$, by

$$P(\mathcal{O}, t) = \sum_C P(C, \mathcal{O}, t).$$

This probability distribution scales as

$$P(\mathcal{O}, t) \sim e^{t\pi(\mathcal{O}/t)} \quad (\text{I.20})$$

in the infinite-time limit. Equation (I.20) is known as the large deviation principle for observable \mathcal{O} [1]. It can be interpreted as the probability of observing an atypical value of observable \mathcal{O} after a large-time scale. The rate function $\pi(\mathcal{O}/t)$ is a dynamical equivalent of the intensive entropy in the microcanonical ensemble and it is known as the large deviation function [1]. It encodes not only the Gaussian but also the non-Gaussian fluctuations (or large deviations) of the observable \mathcal{O}/t which can be obtained by an expansion beyond the quadratic order of the function $\pi(\mathcal{O}/t)$. In the infinite-time limit the function $\pi(\mathcal{O}/t)$ may not be analytic which can be interpreted as a signature of dynamical heterogeneities (dynamical phase transition) [71, 72].

The problem of the determination of the rate function $\pi(\mathcal{O}/t)$ is in general a difficult task, one thus prefers to go to the dynamical canonical ensemble or Laplace space. Instead of fixing the value of the observable \mathcal{O} in order to determine $\pi(\mathcal{O}/t)$ one introduces a parameter s (intensive in time) which biases the statistical weight of histories and fixes the average value of \mathcal{O} , so that $s \neq 0$ favors its non-typical values. In order to do that, we introduce the dynamical partition function (or moment generating function)

$$Z(s, t) = \langle e^{-s\mathcal{O}} \rangle, \quad (\text{I.21})$$

where $\langle \cdot \rangle$ is the expected value with respect to trajectories of duration t . Since the observable \mathcal{O} is additive and the system is described by a Markov process, $Z(s, t)$ satisfies at large times the scaling

$$Z(s, t) \sim e^{t\psi(s)} \quad (\text{I.22})$$

for $t \rightarrow \infty$. The growth rate of $Z(s, t)$ with respect to time, $\psi(s)$ is known as the scaled cumulant generating function (CGF) which fulfils the role of a dynamical free energy. It allows to recover the large-time limit of the cumulants of \mathcal{O} as derivatives of $\psi(s)$ in $s = 0$ from

$$\lim_{t \rightarrow \infty} \frac{1}{t} \langle \mathcal{O}^k \rangle_c = (-1)^k \partial_s^k \psi(s) |_{s=0},$$

where $\langle \mathcal{O}^k \rangle_c$ is the k^{th} cumulant of \mathcal{O} . The cumulative generating function $\psi(s)$ and the large deviation function $\pi(\mathcal{O}/t)$ are related through the Legendre transform

$$\psi(s) = \max_{\hat{o}} [\pi(\hat{o}) - s\hat{o}],$$

where $\hat{o} = \mathcal{O}/t$. If π is convex, i.e., $\pi''(\hat{o}) \leq 0$ [124], then

$$\pi(\hat{o}) = \min_s [\psi(s) + s\hat{o}].$$

I.6 The s -modified Dynamics

As mentioned before, the parameter s involves a (exponential) modification on the statistical weight of the histories of the system. Within this s parametrized ensemble, averages of the observable \mathcal{O} (I.19) defined as

$$\langle \mathcal{O} \rangle_s = \frac{\langle \mathcal{O} e^{-s\mathcal{O}} \rangle}{\langle e^{-s\mathcal{O}} \rangle}$$

for $s = 0$ correspond to the steady state averages of \mathcal{O} . Meanwhile, values of $s \neq 0$ favours histories with non-typical values of observable \mathcal{O} . The s -modified dynamics can be obtained taking the Laplace transform of the probability distribution $P(C, \mathcal{O}, t)$

$$\hat{P}(C, s, t) = \int d\mathcal{O} e^{-s\mathcal{O}} P(C, \mathcal{O}, t).$$

This Laplace transform allows to recover the moment generating function (I.21) as

$$Z(s, t) = \sum_C \hat{P}(C, s, t).$$

The probability $\hat{P}(C, s, t)$ satisfies an s -modified master equation for its time-evolution [40]

$$\partial_t |\hat{P}(t)\rangle = \mathbb{W}_s |\hat{P}(t)\rangle, \quad (\text{I.23})$$

where the s -modified master operator \mathbb{W}_s is given by

$$(\mathbb{W}_s)_{CC'} = W_s(C' \rightarrow C) - r_s(C) \delta_{CC'} + \delta r_s(C) \delta_{CC'}, \quad (\text{I.24})$$

where $\delta r_s(C)$ is defined as

$$\delta r_s(C) = r_s(C) - r(C) - sb(C), \quad (\text{I.25})$$

and $r(C)$ is the escape rate (I.3). On the other hand, $W_s(C \rightarrow C')$ and $r_s(C)$ can be seen as s -modified transition and escape rates, respectively,

$$\begin{aligned} W_s(C \rightarrow C') &= e^{-sa(C, C')} W(C \rightarrow C'), \\ r_s(C) &= \sum_{C'} W_s(C \rightarrow C'). \end{aligned} \quad (\text{I.26})$$

The cumulative generating function ψ in Eq. (I.22) can be determined from this s -modified dynamics as the maximum eigenvalue of the matrix \mathbb{W}_s (I.24) or also, by simulating Eq. (I.23) using a population dynamics algorithm (or cloning algorithm). Both of them are discussed below.

I.6.1 ψ as the Largest Eigenvalue of \mathbb{W}_s

Similarly as we saw for Eq. (I.11), equation (I.23) has as solution

$$|\hat{P}(t)\rangle = e^{t\mathbb{W}_s} |\hat{P}(0)\rangle. \quad (\text{I.27})$$

Matrix \mathbb{W}_s can be written in terms of its left $\langle L_n|$ and right $|R_n\rangle$ eigenvectors and their respective eigenvalues λ_n , with $\lambda_0 > \lambda_1 > \dots$, as

$$\mathbb{W}_s = \sum_n \lambda_n(s) |R_n\rangle \langle L_n|.$$

In the large-time limit, the exponential in Eq. (I.27) is dominated by the largest eigenvalue $\lambda_0(s)$, so that

$$e^{t\mathbb{W}_s} = |R_0\rangle \langle L_0| e^{t\lambda_0(s)} + \dots$$

Thus, $|\hat{P}(t)\rangle$ in the large-time limit also scales as

$$|\hat{P}(t)\rangle = e^{t\mathbb{W}_s} |\hat{P}(0)\rangle \sim |R_0\rangle e^{t\lambda_0(s)} \langle L_0| P_0\rangle + \dots$$

which is equivalent to

$$\hat{P}(C, s, t) \sim R_0(C, s) e^{t\lambda_0(s)}.$$

Thus, in the large-time limit

$$Z(s, t) = \sum_C \hat{P}(C, s, t) \sim e^{t\lambda_0(s)}$$

from which we can see that the maximum eigenvalue of matrix \mathbb{W}_s (I.24) corresponds to the cumulative generating function $\psi(s)$.

I.6.2 A Mutation-Selection Mechanism

Following the procedure used in Sec. I.4 for the master equation, the solution of its s -modified version (I.23) is given by

$$\begin{aligned} P(C, s, t) &= \sum_{K \geq 0} \sum_{C_0 \dots C_K} \int_0^t dt_1 \int_{t_1}^t dt_2 \dots \int_{t_{K-1}}^t dt_K \\ &\times r_s(C_0) e^{-(t_1-t_0) r_s(C_0)} \dots r_s(C_{K-1}) e^{-(t_K-t_{K-1}) r_s(C_{K-1})} \times e^{-(t-t_K) r_s(C_K)} \\ &\times \frac{W_s(C_0 \rightarrow C_1)}{r_s(C_0)} \dots \frac{W_s(C_{K-1} \rightarrow C_K)}{r_s(C_{K-1})} \\ &\times e^{(r_s(C_0)-r(C_0)) (t_1-t_0)} \dots e^{(r_s(C_{K-1})-r(C_{K-1})) (t_K-t_{K-1})} \times e^{(r_s(C_K)-r(C_K)) (t-t_K)} \\ &\times P(C_0, s, 0) \end{aligned}$$

which have been written conveniently in order to introduce the terms

$$Y(C_k) = e^{(r_s(C_k)-r(C_k)) \Delta t(C_k)}, \quad (\text{I.28})$$

where $\Delta t(C_k)$ is the time spent in the configuration C_k . Contrarily to the original operator \mathbb{W} (I.12), the s -modified operator \mathbb{W}_s (I.24) does not conserve probability (since $\delta r_s(C) \neq 0$), implying that there is no obvious way to simulate Eq. (I.23). However, this time-evolution equation can be interpreted not as the evolution of a single system, but as a population dynamics on a large number N_c of copies of the system which evolve in a coupled way [18, 33].

More precisely, reading the operator of the modified master equation (I.23) as in Eq. (I.24), we find that this equation can be seen as a stochastic process of transition rates $W_s(C \rightarrow C')$ (I.26) supplemented with a selection mechanism of rates $\delta r_s(C)$ (I.25), where a copy of the system in configuration C is copied at rate $\delta r_s(C)$ if $\delta r_s(C) > 0$ or killed at rate $|\delta r_s(C)|$ if $\delta r_s(C) < 0$. As detailed below, an estimator for the CGF $\psi(s)$ may be recovered from the exponential growth (or decay) rate of a population evolving with these rules.

I.7 Continuous-Time Population Dynamics

The basic idea of the population dynamics algorithm consists in preparing N_c copies of the system (or clones) and in evolving them according to the transition rates $W_s(C \rightarrow C')$ given by Eq. (I.26). During this evolution some copies are repeatedly multiplied or eliminated according to a selection mechanism of rates $\delta r_s(C)$ (I.25). The mutation-selection mechanism described above can be performed in a number of ways. One of them consists in keeping the total number of clones constant and another, in leaving this population of clones to grow (or decrease) in time (as implemented in chapter II). Additionally, the selection mechanism can be implemented along with each evolution of the copies (**Continuous-Time**) [7, 17, 19] or for each pre-fixed time-interval (**Discrete-Time**) [18]. This last one is implemented (in a constant population fashion) in chapter III, while the continuous-time version is used in the majority of the manuscript. An explanation about important differences between continuous and discrete-time techniques can be found in Secs. III.2.4.2 and IV.5. A detailed description of the continuous-time approaches is presented below.

I.7.1 The Cloning Algorithm

We consider N_c copies (or clones) of the system. The dynamics is continuous in time: for each copy, the actual changes of configuration occur at times (which we call ‘evolution times’) which are separated by intervals whose duration is distributed exponentially. At a given step of the algorithm, we denote by $\mathbf{t} = \{t^{(i)}\}_{i=1, \dots, N_c}$ the set of the future evolution times of all copies and by $C = \{C_i\}_{i=1, \dots, N_c}$ the configurations of the copies. Their initial configurations do not affect the resulting scaled cumulant generating function in the large-time limit. However, for the concreteness of the discussion, without loss of generality, we assume that these copies have the same configuration C_0 at $\mathbf{t} = 0$. The cloning algorithm is constituted of the repetition of the following procedures:

1. Find the clone whose next evolution time is the smallest among all the clones, i.e., $j = \operatorname{argmin}_i t^{(i)}$.
2. Compute $y_j = \lfloor Y(C_j) + \epsilon \rfloor$, where the cloning factor $Y(C_j)$ (I.28) is defined as $e^{\Delta t(C_j) \delta r_s(C_j)}$, $\Delta t(C_j)$ is the time spent by the clone j in the configuration C_j since its last configuration change, and ϵ is a random number uniformly distributed on $[0, 1]$.
3. If $y_j = 0$, remove this copy from the ensemble, and if $y_j > 0$, make $y_j - 1$ new copies of this clone.
4. For each of these copies (if any), the state C_j is changed independently to another state C'_j , with probability $W_s(C_j \rightarrow C'_j)/r_s(C_j)$.

5. Choose a waiting time $\Delta t^{(j)}$ from an exponential law of parameter $r_s(C'_j)$ for each of these copies. Their next change of configuration will occur at the evolution time $t^{(j)} + \Delta t^{(j)}$.

I.7.1.1 Non-Constant Population Approach

The repetition of this procedure will result (after an enough time) in an exponential growth (or decay) of the number of clones. We can keep track of the different changes in the number of clones and the times where these changes occur and denote by $N(s, t)$ the time-dependent population. The CGF estimator can be computed from the slope in time of the log-population $\hat{N}(s, t) = \log N(s, t)$, which constitutes an evaluation of the population growth rate. This can be done in different ways, for example by fitting $\hat{N}(s, t)$ by $p_t = \Psi(s)t + p_0$ where $\Psi(s)$ is the CGF estimator or also using

$$\Psi(s) = \frac{1}{T_{\max} - T_{\min}} \log \left(\frac{N_{\max}}{N_{\min}} \right), \quad (\text{I.29})$$

where N_{\max} and N_{\min} are the maximum and minimum values for $N(s, t)$ and T_{\max} and T_{\min} their respective times. This approach is implemented in chapter II, where we discuss the discreteness effects in populations dynamics and their influence in the CGF estimation.

I.7.1.2 Constant-Population Approach

In order to keep the total number of copies constant, we add to the procedure described above an additional step

6. We choose randomly and uniformly: (i) a clone k , $k \neq j$ and we copy it (if $y_j = 0$), or (ii) $y_j - 1$ clones and we erase them (if $y_j > 1$).

Thus, the CGF estimator $\Psi_s^{(N_c)}$ can be obtained from the exponential growth rate that the population would present if it was not kept constant [7]. More precisely, this estimator is defined as

$$\Psi_s^{(N_c)} = \frac{1}{t} \log \prod_{i=1}^{\mathcal{K}} X_i, \quad (\text{I.30})$$

where $X_i = (N_c + y_i - 1)/N_c$ are the ‘‘growth’’ factors at each step j of the procedure described above, and \mathcal{K} is the total number of configuration changes in the full population up to time t (which has not to be confused with K). This growth rate can also be computed from a linear fit over the reconstructed log-population and the initial transient regime, where the discreteness effects are present, can be discarded in order to obtain a better estimation. This approach is used in chapters IV, V and VI.

I.8 Example Models

In the next chapters, the cloning algorithm is implemented in order to obtain an estimation of the CGF $\psi(s)$. This is done with two specific models: a simple two-state annihilation-creation dynamics, and a contact process on a one-dimensional periodic lattice [19, 38, 107, 108]. The first system (chapters: II, III, IV and V) was chosen for its simplicity and the possibility of

comparing the numerical predictions with the exact values of $\psi(s)$. On the other hand, the contact process (chapters: IV, V and VI) is used to extend the analysis and to verify the results on a (more complex) many body system where the dependence with the size of the system can be also analyzed. In both cases, we consider the dynamical activity K [12, 39, 40, 74–82] as the additive observable \mathcal{O} (I.19). The analytical expression of the CGF $\psi(s)$ is obtained (when possible) by solving the largest eigenvalue of the operator \mathbb{W}_s (I.24) as discussed in Sec. I.6.1.

I.8.1 Annihilation-Creation Dynamics

The dynamics occurs in one site where the only two possible configurations C are either 0 or 1. The transition rates between configurations are given by

$$\begin{aligned} W(0 \rightarrow 1) &= c, \\ W(1 \rightarrow 0) &= 1 - c, \end{aligned}$$

where $c \in [0, 1]$. Eq. (I.8) for this process becomes

$$\partial_t \begin{pmatrix} P(0, t) \\ P(1, t) \end{pmatrix} \begin{pmatrix} -c & 1 - c \\ c & -1 + c \end{pmatrix} \begin{pmatrix} P(0, t) \\ P(1, t) \end{pmatrix}. \quad (\text{I.31})$$

As we mentioned before, one advantage of considering this process for our analysis is that the large deviation function for the activity can be determined analytically. The large-time cumulant generating function $\psi_K(s) = \lim_{t \rightarrow \infty} \frac{1}{t} \log \langle e^{-sK} \rangle$ corresponds to the maximum eigenvalue of the matrix \mathbb{W}_s (I.24)

$$\mathbb{W}_s = \begin{pmatrix} -c & (1 - c) e^{-s} \\ c e^{-s} & -1 + c \end{pmatrix} \quad (\text{I.32})$$

which results in
$$\psi_K(s) = -\frac{1}{2} + \frac{1}{2} \left(1 - 4c(1 - c)(1 - e^{-2s}) \right)^{1/2}. \quad (\text{I.33})$$

Equation (I.33) will allow us to assess the quality of our numerical results. The inverse of the difference between the eigenvalues of \mathbb{W}_s

$$t_{\text{gap}} = \frac{1}{\sqrt{1 - 4c(1 - c)(1 - e^{-2s})}} \quad (\text{I.34})$$

allows us to define the typical convergence time t_{gap} for the large-time behavior for Eq. (I.31) which is analyzed in Sec. II.4.1.

I.8.2 Contact Process

This system consists in a one-dimensional lattice with L sites and periodic boundary conditions [38, 107, 108]. Each position i is occupied by a spin which is either in the state $n_i = 0$ or $n_i = 1$. The configuration C is then constituted by the states of these spins, i.e., $C = (n_i)_{i=1}^L$. The transition rates for this process are given by

$$\begin{aligned} W(n_i = 1 \rightarrow n_i = 0) &= 1, \\ W(n_i = 0 \rightarrow n_i = 1) &= \lambda(n_{i-1} + n_{i+1}) + h, \end{aligned}$$

where λ and h (spontaneous rate of creation) are positive constants. This model is an example of contact processes [38], which have been studied in many contexts especially to model the spread of infection diseases [251]. Within this context, the state $n_i = 1$ is used to represent a sick individual and λ can be seen as a infection rate. The corresponding CGF develops a singularity as $L \rightarrow \infty$, showing a dynamical phase transition [19, 35, 74, 252]. The contact process is a model of the directed percolation universality class and its scaling properties have been discussed extensively [252–254].

II – Discreteness Effects in Population Dynamics

II.1 Introduction

In the present chapter [P1], we focus on the **non-constant population approach** of the cloning algorithm (as described in Sec. I.7.1.1), that we study numerically for the simple annihilation-creation (Sec. I.8.1) model where its implementation and its properties can be examined in great details. As we mentioned in the introduction, the cloning algorithm results (as time goes to infinity) in an exponential growth (for $s < 0$) or decay (for $s > 0$) of the number of clones. As we will see later, the “discreteness effects” in the evolution of our populations are strong at initial times. That is why the determination of the large deviation function using this algorithm is constrained not only to the parameters (c, s) , the initial number of clones N_c and the number of realizations R but also to the final time (or the maximum population) until which the process evolves in the numerical procedure. In Sec. II.2 we describe issues related to the averaging of distinct runs, that we quantify in Sec. II.3. In Sec. II.4 we propose a new method to increase the efficiency of the population dynamics algorithm by applying a realization-dependent time delay, and we present the results of its application in Sec. II.5. We characterize numerically the distribution of these time delays in Sec. II.6. Our conclusions and perspectives are gathered in Sec. II.7.

II.2 Average Population and the LDF

In order to obtain an accurate estimation of $\psi(s)$, we should average several realizations of the procedure described in Sec. I.7.1.1. To perform this average, we will define below a procedure that we have called **merging** which will allow us to determine in a systematic way the average population from which we can obtain this estimation that we denote $\Psi(s)$. Noteworthy, this erroneously could be seen as obtaining $\Psi(s)$ from the growth rate of the average (or equivalently the sum) of several runs of the population dynamics. This procedure would be incorrect since it amounts to performing a single run of the total population of the different runs, with a dynamics that would partition the total population into *non-interacting* sub-populations, while, as described in Sec. I.7.1.1, the population dynamics induces effective interactions among the whole set of copies inside the population. In fact, the right way of performing this numerical estimate comes from computing $\Psi(s)$ from the average growth rate of several runs of the population, i.e., from taking the average $\langle \log N(s, t) \rangle$ of the slopes of several $\log N(s, t)$ instead of the slope of $\log \langle N(s, t) \rangle$. The two results differ in general since

$\langle \log N(s, t) \rangle \neq \log \langle N(s, t) \rangle$. One can expect that the two results become equivalent in the large N_c limit as the distribution of growth rate should become sharply concentrated around its average value; however, they are different in the finite N_c regime that we are interested in. This alternative way of defining the CGF estimator is discussed deeply in Secs. III.4.3 and V.3.

II.2.1 Populations Merging

Let us consider J populations: $\mathcal{N} = \{N_1(s, t), N_2(s, t), \dots, N_J(s, t)\}$. In order to compute the average population $\langle \mathcal{N} \rangle$ defined as $\langle \mathcal{N} \rangle = \langle N_j(s, t) \rangle_{j=1}^J$, we introduce a procedure that we have called **merging** (of populations) which is described below.

Given $N_i(s, t)$ and $N_j(s, t)$ the result of merging these two populations $\mathcal{M}(N_i, N_j)$ is another population $N_{ij} = N_i + N_j$ which represents the total number of clones for each time where a change in population for N_i and N_j has occurred. If $\langle N_{ij} \rangle$ is the average population for N_i and N_j , this is related to the merged population through $\langle N_{ij} \rangle = \frac{N_{ij}}{2}$. If we add, for example, to our previous result another population N_k , the result $\mathcal{M}(N_{ij}, N_k)$ is related to the average by $\mathcal{M}(N_{ij}, N_k) = N_{ij} + N_k = N_i + N_j + N_k = N_{ijk} = 3\langle N_{ijk} \rangle$. These merging procedure can be repeated for each of the populations in \mathcal{N} so that

$$\mathcal{M}[\mathcal{N}] = \mathcal{M}(\mathcal{M}(\mathcal{M}(\dots(\mathcal{M}(\mathcal{M}(N_1, N_2), N_3), N_4) \dots), N_{J-1}), N_J)$$

is the result of systematically merging all the populations in \mathcal{N} . The average population $\langle \mathcal{N} \rangle$ can be recovered from $\mathcal{M}[\mathcal{N}]$ as

$$\langle \mathcal{N} \rangle = (1/J)\mathcal{M}[\mathcal{N}].$$

Similarly, in the case of log-populations ($\hat{N}_j(s, t) = \log N_j(s, t)$), the average $\langle \hat{N} \rangle = \langle \hat{N}_j(s, t) \rangle_{j=1}^J$ is obtained from merging all the log-populations in $\hat{N} = \{\hat{N}_1(s, t), \hat{N}_2(s, t), \dots, \hat{N}_J(s, t)\}$. The estimator $\Psi(s)$ is then computed from the slope of $\langle \hat{N} \rangle$ with $\langle \hat{N} \rangle = (1/J)\mathcal{M}[\hat{N}]$.

II.2.2 Discreteness Effects at Initial Times

Issues can emerge in the determination of $\Psi(s)$ (I.29) which are not only related to the dependence of the method in N_c (the initial number of clones) and J (the number of populations). At initial times there is a wide distribution of times at which the first series of jumps occurs. This means that fluctuations at initial times induce that some populations remain in their initial states longer than others, producing an effective delay compared to other populations that evolve faster in their initial regime. From a practical point of view, this can induce that the numerical determination of $\Psi(s)$ becomes a slow and inefficient task. One way of dealing with this issue comes from restricting the evolution of \mathcal{N} up to a maximum time T_{\max} or a maximum population N_{\max} . However, this implies that if T_{\max} or N_{\max} are not long enough, the determination of $\Psi(s)$ will be strongly affected by the behavior of \mathcal{N} at initial times. We now discuss two issues that are encountered in the numerical evaluation of the CGF estimator: (i) the influence of how the dynamics is halted; and (ii) the role of initial population in the initial regime.

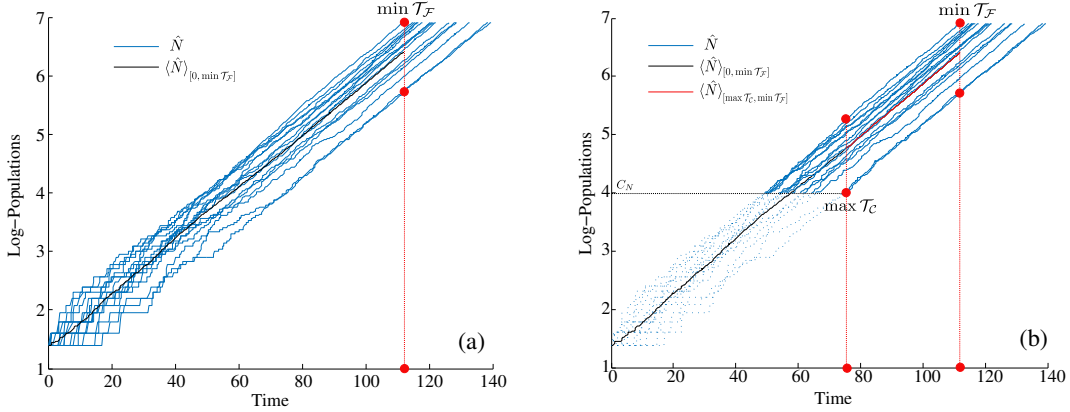


Figure II.1: Log-populations as function of time (blue). Their evolution has been restricted up to a maximum (log) population value. **(a)** The average log-population $\langle \hat{N} \rangle$ (black) is made in the interval $[0, \min \mathcal{T}_{\mathcal{F}}]$, where all the populations are defined. **(b)** After a cut in populations C_N (in order to eliminate the initial discreteness effects), the average log-population (red) that represents the new \hat{N} is defined only in the interval $[\max \mathcal{T}_{\mathcal{C}}, \min \mathcal{T}_{\mathcal{F}}]$.

Let us call $\mathcal{T}_{\mathcal{F}} = \{t_1^{\mathcal{F}}, \dots, t_j^{\mathcal{F}}\}$ the set of final times of \mathcal{N} , with $t_j^{\mathcal{F}} \leq T_{\max}, \forall j \in \{1, \dots, J\}$. Note that $t_j^{\mathcal{F}}$ depends on j whenever the simulation is stopped at N_{\max} (as in Fig. II.1) or at T_{\max} . This is due to the fact that the algorithm is continuous in time and the last $\Delta t(C)$ does not exactly lead to T_{\max} . We say that the average population $\langle \mathcal{N} \rangle$ **represents** \mathcal{N} only if the average is made in the interval $[0, \min \mathcal{T}_{\mathcal{F}}]$ where all the populations are defined. In other words, the average population in this interval takes into consideration all the populations while for times $t \geq \min \mathcal{T}_{\mathcal{F}}$ some populations have stopped evolving. This phenomenon is especially evident when considering a maximum population limit N_{\max} for the evolution of the populations (Fig. II.1(a)). As a consequence, $\langle \mathcal{N} \rangle$ depends on the distribution of final times of \mathcal{N} which are not necessarily equal to T_{\max} .

An alternative that can be considered in order to overcome the influence of initial discreteness effects in the determination of $\Psi(s)$ is to get rid of the initial transient regime where these effects are present. In other words, to cut the initial time regime of our populations. Let us call $C_N \geq \log N_c$ the initial cut in log-populations and equivalently $C_t \geq 0$ the initial cut in times. $\mathcal{T}_{\mathcal{C}} = \{t_1^{\mathcal{C}}, \dots, t_j^{\mathcal{C}}\}$ is the distribution of times at $C_{t,N}$. In that case, similarly as we analyzed before, the average population $\langle \mathcal{N} \rangle$ represents \mathcal{N} only if the average is made in the interval $[\max \mathcal{T}_{\mathcal{C}}, \min \mathcal{T}_{\mathcal{F}}]$ which can be in fact very small and could result in a bad approximation of $\Psi(s)$ (Fig. II.1(b)).

As we will see in the next section, the log-populations after a long enough time become parallel, i.e., once the populations have surpassed the discreteness effects regime, the distance between them is constant. We will use this fact in order to propose a method which allows us to overcome the problems we have described in this section. Throughout this chapter, we consider for our simulations $c = 0.3$, $N_c = 2^2$, $N_{\max} = 10^3$, $J = 2^8$ and $s \in [-0.3, 0]$.

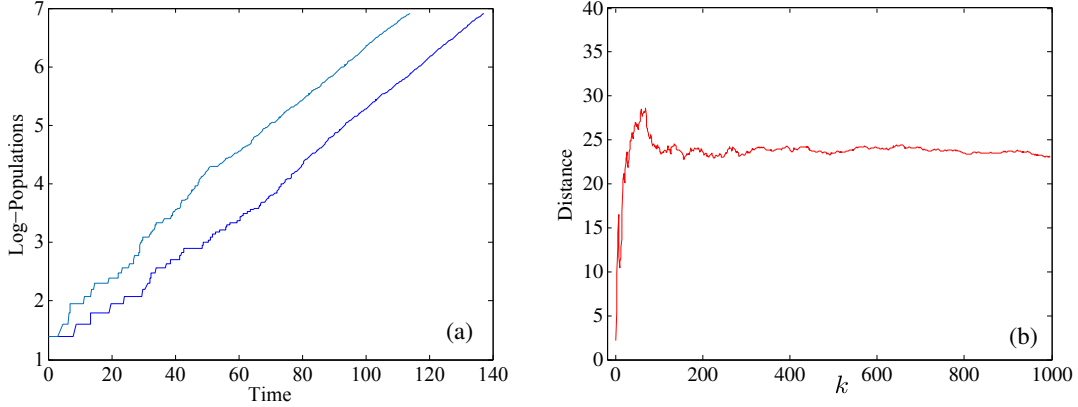


Figure II.2: Evolution of two log-populations \hat{N}_i, \hat{N}_j as function of time and the distance $D(\hat{N}_i, \hat{N}_j)$ between them as defined in Eq. (II.1). **(a)** Log-populations after a long enough time become parallel. **(b)** Once the populations have overcome the initial discrete population regime, the distance between them becomes constant. ($s = -0.1$).

II.3 Parallel Behavior in Log-Populations

II.3.1 Distance between Populations

Given $N_i(s, t)$ and $N_j(s, t)$, we define the distance between these populations at N^* (with $N^* \in N_i$ and $N^* \in N_j$), as

$$D(N_i, N_j)(N^*) = \left| \left(t_j(N^*) + \frac{\Delta t_j(N^*)}{2} \right) - \left(t_i(N^*) + \frac{\Delta t_i(N^*)}{2} \right) \right|, \quad (\text{II.1})$$

where $\Delta t_k(N^*)$ is the time interval $N_k(s, t)$ spent at N^* and $t_k(N^*)$ is the time where $N_k(s, t)$ changes to N^* . Evidently, there are cases where $N^* \notin N_i$ but $N^* \in N_j$, $N^* \in N_i$ but $N^* \notin N_j$ and $N^* \notin N_i$ and $N^* \notin N_j$. However, $D(N_i, N_j)(N^*)$ for these cases can also be computed. The last analysis (and definitions) is also valid for log-populations. These distances, $D(N_i, N_j)(N^*)$ and $D(\hat{N}_i, \hat{N}_j)(N^*)$ enjoy interesting properties that we discuss below.

II.3.2 Properties of $D(\hat{N}_i, \hat{N}_j)$

In Fig. II.2, we show two log-populations and the distance between them. These log-populations after a long enough time become parallel (Fig. II.2(a)), i.e., once the populations have overcome the discreteness effects regime, the distance between them becomes constant (Fig. II.2(b)). The region where the distance between populations is constant characterizes the exponential regime of the populations growth, i.e., the region where the discreteness effects are not strong anymore.

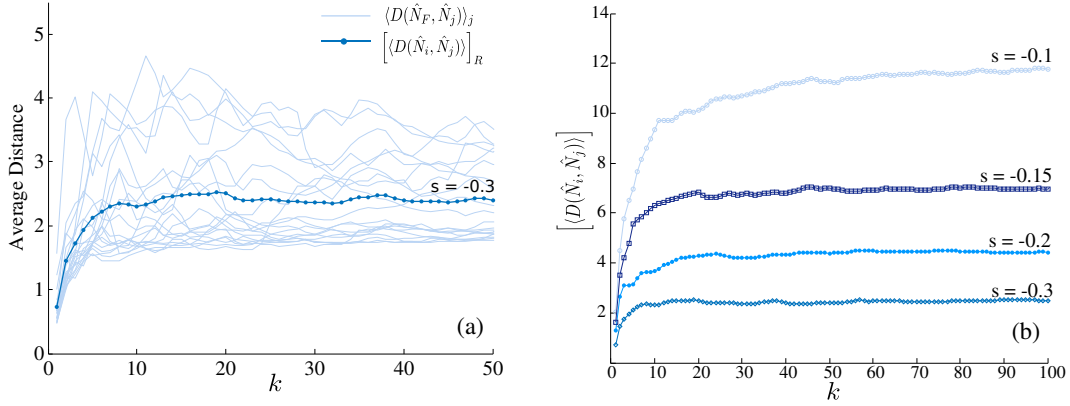


Figure II.3: **(a)** Average distance of the log-populations in \hat{N} with respect to a reference one for $R = 20$ realizations (light blue) and their average (dark blue). **(b)** Average distance between populations for several values of the parameter s . How fast a population “exits” from the discreteness effects regime depends on s : s closer to zero corresponds to a slower population growth and hence a longer discreteness regime.

If we consider some population $\hat{N}_F \in \hat{N}$ as reference, using the definitions above, it is possible to determine the distance $D(\hat{N}_F, \hat{N}_j)$ between \hat{N}_F and the rest of populations in $\hat{N} = \{\hat{N}_1, \hat{N}_2, \dots, \hat{N}_J\}$. In Fig. II.3(a) we show their average $\langle D(\hat{N}_F, \hat{N}_j) \rangle_j$ in light blue and its average over $R = 20$ realizations $\left[\langle D(\hat{N}_i, \hat{N}_j) \rangle \right]_R$ in dark blue. The parameter s characterizes atypical behaviors of the unbiased dynamics (as we mentioned in Sec. I.6), and this induces a dependence in s of the population growth. A population with a large value of s corresponds to a large deviation of K . Also, as it is clearly illustrated in Fig. II.3(b), the time of entrance of the system into a regime free of discreteness effects depends on s .

II.4 Time Correction in the Evolution of Populations

Based on the results we just illustrated, we propose a method to improve the estimation of $\psi(s)$ and reduce the influence of the initial transient regime we described in Sec. II.2.2. We aim at giving more weight to the exponential regime in the determination of $\psi(s)$. As detailed below, this can be done through a time delay in the evolution of populations.

II.4.1 Time Delay Correction

Consider J populations \mathcal{N} , their respective log-populations \hat{N} and their distribution of final times $\mathcal{T}_{\mathcal{F}} = \{t_1^{\mathcal{F}}, \dots, t_J^{\mathcal{F}}\}$. We define as **delay** $\Delta\tau_j$ of \hat{N}_j (with respect to a fixed reference population $\hat{N}_F \in \hat{N}$) the time interval

$$\Delta\tau_j = t_F^{\mathcal{F}} - t_j^{\mathcal{F}}$$

such that, if $\Delta\tau_j < 0$, \hat{N}_j is ahead with respect to \hat{N}_F , and if $\Delta\tau_j > 0$, \hat{N}_j is delayed with respect to \hat{N}_F . This lag can be compensated by performing on \hat{N}_j the time translation

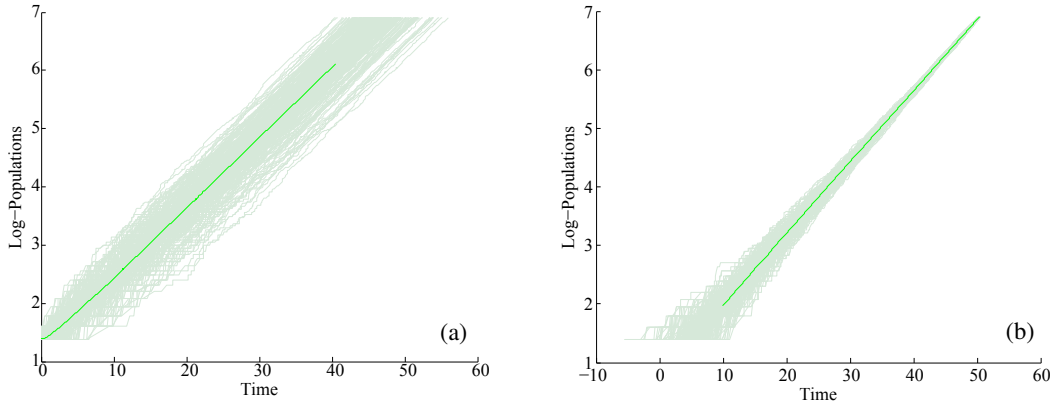


Figure II.4: **(a)** Log-populations, **(b)** time-delayed log-populations, and their average (dark green). The fluctuations at initial times produce a gap in the evolution of individual populations inducing a relative shift that lasts forever. This is compensated by delaying the populations in time, as explained in Sec. II.4.1. ($s = -0.25$).

$$\hat{N}_j^{\text{new}} = \hat{N}_j(s, t + \Delta\tau_j) \quad (\text{II.2})$$

which produces that \hat{N}_j^{new} and \hat{N}_F share not only the final population N_{max} , but also the same final time $t_F^{\mathcal{F}}$. Moreover, considering also the fact that log-populations are parallel at large times, this procedure produces that \hat{N}_j^{new} and \hat{N}_F overlap in a **free of discreteness effects region**. The result of performing this transformation to all the populations in \hat{N} is shown in Fig. II.4.

Fig. II.4 also illustrates many of the points we have discussed up to now. One of them is related to the “wide” distribution of final times, i.e., $\min \mathcal{T}_{\mathcal{F}}$ and $\max \mathcal{T}_{\mathcal{F}}$ can be very distant one from each other. This along with the fact that the average population depends on $\min \mathcal{T}_{\mathcal{F}}$ makes that the determination of $\Psi(s)$ omits a considerable region where the populations have already entered the exponential regime. This implies precisely that more weight is given to the initial discreteness effects than to the exponential regime. These effects are in fact present up to relatively long times which means that if we would like to get rid of the region where discreteness effects are strong by cutting the populations, the determination of $\Psi(s)$ would be restricted to the interval $[\max \mathcal{T}_{\mathcal{C}}, \min \mathcal{T}_{\mathcal{F}}]$. By applying precisely this time delay correction to \hat{N} we solve these two problems. First, we give more importance precisely to the region where the population growth is exponential. Second, we omit naturally the very first initial times of the evolution of our populations.

The inverse of the difference between the two largest eigenvalues of \mathbb{W}_s (I.32), t_{gap} (Eq. (I.34)) allows us to define the typical convergence time to the large time behavior for Eq. (I.31) (as we mentioned in Sec. I.8.1). A crucial remark is that, as observed numerically, the duration before the population enters into the exponential regime is in fact larger than the time scale given by the gap: for instance, for the parameters used to obtain Fig. II.4, from Eq. (I.34) one has $t_{\text{gap}} \approx 0.804$. The understanding of the duration of this discreteness effects regime would require a full analysis of the finite-population dynamics which are not fully understood. We propose in this section a numerical procedure to reduce its influence.

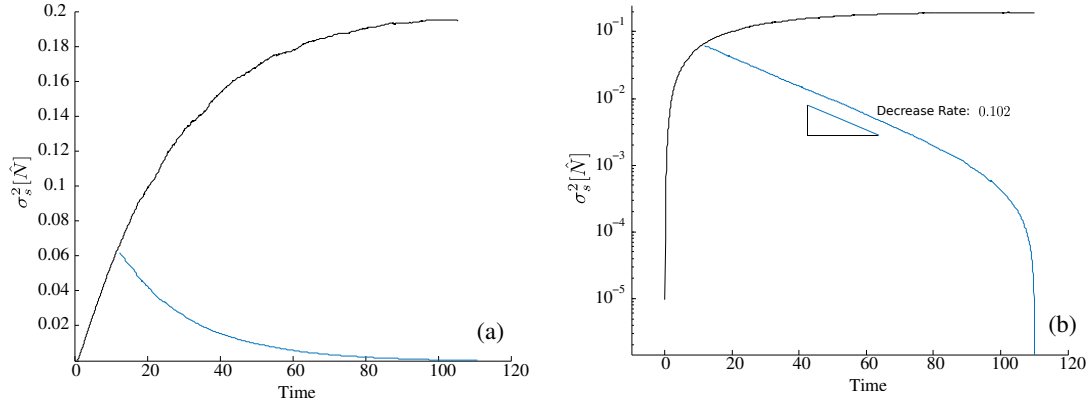


Figure II.5: **(a)** Variance of the log-populations (black) and the delayed log-populations (blue) as a function of time. The variance of log-populations increases (or decreases, after the time transformation) as function of time. **(b)** Log-population variance in semi-log scale. ($s = -0.1$).

Log-Population Variance

As can be seen from Fig. II.4, and as it is verified in Fig. II.5, the variance of log-populations (black) increases as a function of the time, faster during the transient regime, and slower during the exponential growth regime until the variance becomes constant. After the time-delay correction, the variance of the delayed log-population (blue) decreases to zero as a function of time. The s -dependent decrease rate is shown in Fig. II.5(b).

II.5 $\Psi(s)$ Before and After the Time Delay

The CGF estimator $\Psi(s)$ can be recovered from the slope in time of the logarithm of the average population (see Sec. II.2.1). We also mentioned in Sec. II.2.2, that an alternative we can consider to overcome the discreteness effects would be to eliminate the initial transient regime where these effects are strong. The improvement in the estimation of the analytical CGF $\psi(s)$ (I.33) comes precisely from these two main contributions, the time delaying of populations and the discarding of the initial transient regime of the populations. We denote $\Psi_{\text{num}}(s)$ the numerical estimator which is obtained from the slope of the logarithm of the average population (computed from merging several populations that have been generated using the cloning algorithm). On the other hand, $\Psi_{\tau}(s)$ is obtained through a time delay procedure over \hat{N} , as described above. These two numerical estimations are in fact averages over R realizations and over their last γ values. The approach followed in order to compute $\Psi_{\text{num}}(s)$ and $\Psi_{\tau}(s)$ are computed is explained below.

II.5.1 Numerical Estimators for $\psi(s)$

Let us call $\Psi_*(C_N)$ an estimation of ψ (by some method $(*) \in \{\text{num}, \tau\}$) as a function of the cut in log-population C_N . If we consider C_N as $C_N = \{C_N^1, \dots, C_N^\Gamma\}$ a set of Γ cuts, $\Psi_*(C_N)$

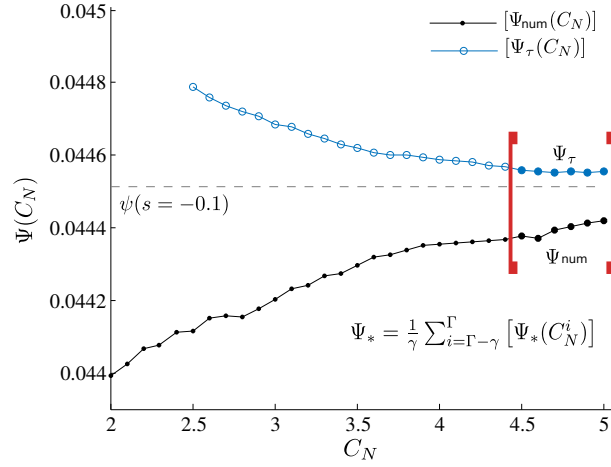


Figure II.6: Numerical estimations of $\psi(s = -0.1)$ as a function of the cut C_N in (log) population. $[\Psi_\tau(C_N)]$ is shown in blue and $[\Psi_{\text{num}}(C_N)]$ in black for $R = 40$. The numerical estimations Ψ_{num} and Ψ_τ are computed from an average of $[\Psi_*(C_N^i)]$ over its last $\gamma = 6$ values. The subscript “*” denotes “num” or “ τ ”.

is in fact $\Psi_*(C_N) = \{\Psi_*(C_N^1), \dots, \Psi_*(C_N^\Gamma)\}$. If $[\Psi_*(C_N^i)]$ is an average over R realizations,

$$[\Psi_*(C_N^i)] = \frac{1}{R} \sum_{r=1}^R \Psi_*^r(C_N^i)$$

our numerical estimation (for a given s) is then computed from an average of $[\Psi_*(C_N^i)]$ over its last γ values, i.e.,

$$\Psi_*(s) = \frac{1}{\gamma} \sum_{i=\Gamma-\gamma}^{\Gamma} [\Psi_*(C_N^i)] = \frac{1}{\gamma R} \sum_{i=\Gamma-\gamma}^{\Gamma} \sum_{r=1}^R \Psi_*^r(C_N^i)$$

as is shown in Fig. II.6. More details of the determination of these estimators are given in the subsection below.

II.5.2 Comparison between “Bulk” and “Fit” Estimators of $\psi(s)$

The estimators defined in the last subsection can be obtained from the “bulk” slope (Fig. II.7(a)) given by Eq. (I.29) and from the affine fit of the average log-population by $p_t = \Psi(s)t + p_0$ (Fig. II.7(b)), as explained in Sec. I.7.1.1. Fig. II.7 shows the average over $R = 40$ realizations of the numerical estimators $\Psi_{\text{num}}(C_N)$ and $\Psi_\tau(C_N)$ as a function of the cut in log-population for $s = -0.1$. As before, $[\Psi_\tau(C_N)]$ is shown in blue and $[\Psi_{\text{num}}(C_N)]$ (without the “time delay”) is shown in black. As we already mentioned, the estimation for ψ becomes better if we discard the initial transient regime where the discreteness effects are strong.

The black curves in Fig. II.7 represent the standard way of estimating ψ which comes from the slope of the average log-population, shown in dark green in Fig. II.4(a) for one realization. We can observe the effect of discarding the initial transient regime of these

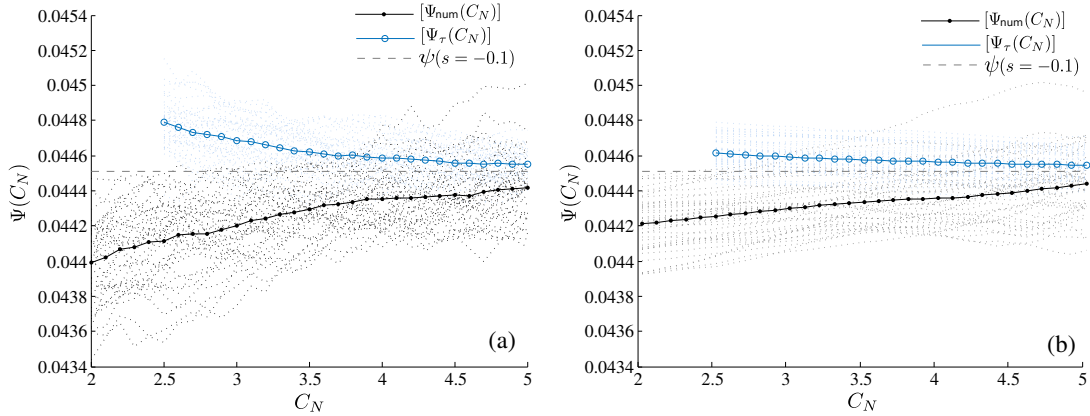


Figure II.7: Average over $R = 40$ realizations of the numerical estimators $\Psi_{\text{num}}(C_N)$ and $\Psi_{\tau}(C_N)$ as a function of the cut in log-population for $s = -0.1$. The CGF estimations were obtained from (a) a “Bulk” slope and from (b) a “Fit” slope. The estimation for ψ becomes better if we discard the initial transient regime where discreteness effects are strong.

populations by cutting systematically this curve and computing $\Psi_{\text{num}}(C_N)$ from the growth rate $\Psi(s)$ computed on the interval $[C_N, N_{\text{max}}]$. Independently if $\Psi_{\text{num}}(C_N)$ is computed from the “bulk” slope or by the “fit” slope, for appropriate values of C_N , $\Psi_{\text{num}}(C_N)$ becomes closer to the theoretical value. Additionally to this result, we can add the time correction or delay proposed in Sec. II.4.1 and the estimation $\Psi_{\tau}(C_N)$, shown by blue curves in Fig. II.7, is closer to the theoretical value than $\Psi_{\text{num}}(C_N)$ for all C_N .

Once we have proved that the estimation of ψ becomes better when we discard the initial times where the discreteness effects are strong and when we perform a “time delay” over our populations in order to give more weight to the final regime of our populations, the question that remains is related to what we should consider as $\Psi_{\text{num}}(s)$ and $\Psi_{\tau}(s)$. As we showed in Fig. II.6, $\Psi_{\text{num}}(s = -0.1)$ and $\Psi_{\tau}(s = -0.1)$ are computed from an average over the last γ values of $[\Psi_{\text{num}}(C_N)]$ and $[\Psi_{\tau}(C_N)]$. Below, we repeat this procedure and compute these estimators for several values of s , $s \in [-0.3, -0.05]$. The improvement in the determination of the CGF is measured through the relative distance of the numerical estimations with respect to the theoretical values and their errors.

II.5.3 Relative Distance and Estimator Error

The relative distance

$$D(\psi(s), \Psi_*(s)) = \left| \frac{\psi(s) - \Psi_*(s)}{\psi(s)} \right|$$

between the estimator $\Psi_*(s)$ and its theoretical value $\psi(s)$ is shown in Fig. II.8. These distances were also computed from the “bulk” and the “fit” slope and with (blue) and without time delay (black). As we can observe, the deviation from the theoretical value is larger for values of s close to 0, but is smaller after the “time correction” for almost every value of s . Fig. II.9 presents the estimator error for $\psi(s)$ defined as

$$\epsilon = \frac{\sigma_{\Psi_*}}{\sqrt{R}},$$

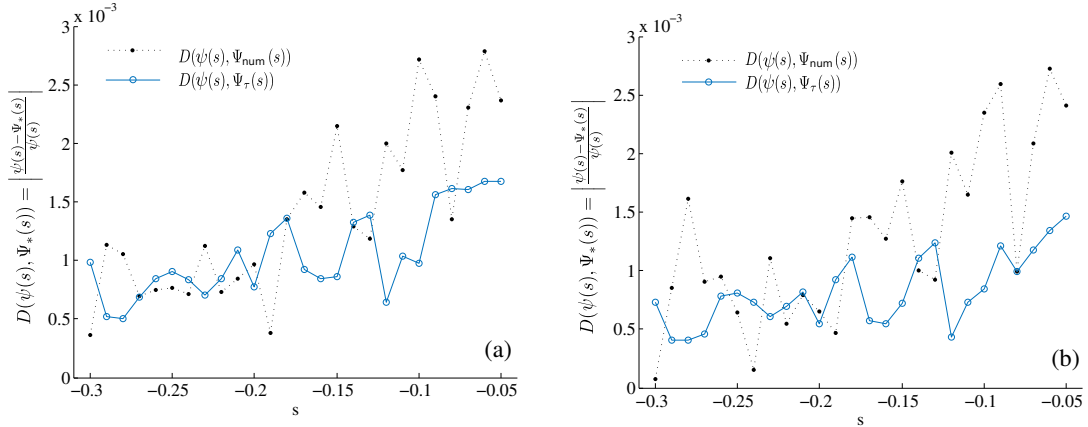


Figure II.8: Relative distance $D(\psi(s), \Psi_*(s))$ between the estimator $\Psi_*(s)$ and its theoretical value $\psi(s)$. The deviation from the theoretical value is larger for values of s close to 0, but is smaller after the “time delay correction” for almost every value of s . The CGF estimators were obtained from (a) a “Bulk” and from (b) a “Fit” slope.

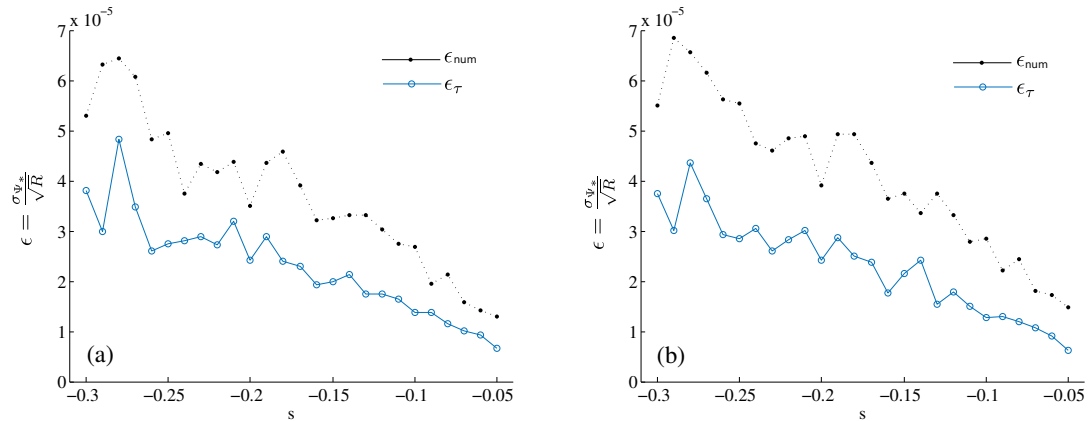


Figure II.9: Estimator error for $\psi(s)$, ϵ_{num} (black) and ϵ_τ (blue). The estimator error decreases as s approaches to 0 (for both, (a) “Bulk”, (b) and “Fit” slopes) and it is always smaller for $\Psi_\tau(s)$ for any value of s .

where R is the number of realizations and σ_{Ψ_*} is the standard deviation of $\Psi_*(s)$. Similarly as in previous results, the estimator error decreases as s approaches to 0 (for both slopes) and it is always smaller for $\Psi_\tau(s)$ for any value of s .

II.6 Time Delay Properties

Here, we analyze the properties of the distribution of time delays $\Delta\tau(s) = \{\Delta\tau_1(s), \dots, \Delta\tau_J(s)\}$ which has been centered with respect to its mean. In Fig. II.10(a), we show its variance $\sigma_s^2[\Delta\tau]$. The dispersion of time delays is large for values of s close to 0 and decreases quickly

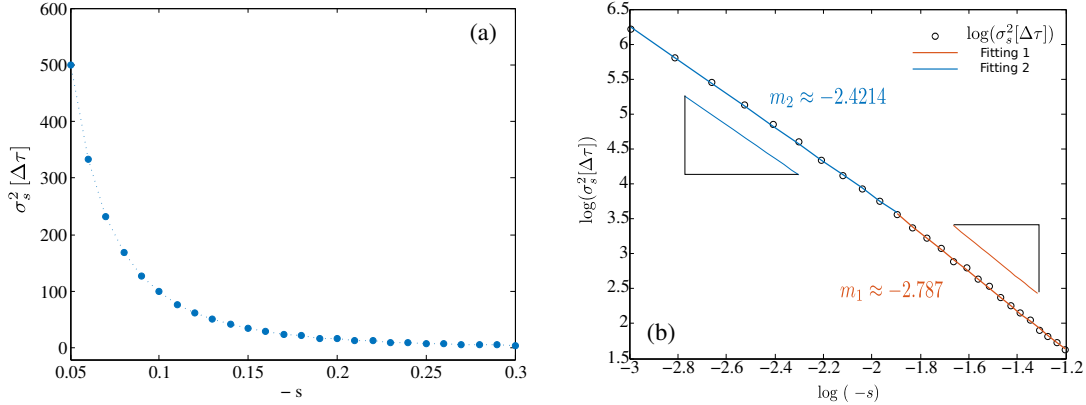


Figure II.10: **(a)** Time delay variance $\sigma_s^2[\Delta\tau]$. The dispersion of the time delays is large for values of s close to 0 and decreases rapidly as $-s$ increases. **(b)** Time delay variance regimes, one characterized with $m_1 \approx -2.877$ ($s \in [-0.15, -0.3]$) and the other with $m_2 \approx -2.4214$ ($s \in [-0.05, -0.15]$).

as $-s$ increases. This is understood by observing that the typical growth rate $[r_s(C) - r(C)]$ of the cloning algorithm goes to zero as $s \rightarrow 0$ inducing a longer transient regime between the small and large population regimes. When we plot the variance in log-log scale, as in Fig. II.10(b), we can observe two linear regimes, one characterized by an exponent $m_1 \approx -2.877$ ($s \in [-0.15, -0.3]$) and the other by $m_2 \approx -2.4214$ ($s \in [-0.05, -0.15]$). They correspond to power-law behaviors in time of the variance of the delays, which remain to be understood.

This dependence of the dispersion of time delays with s can be better seen in the distribution of time delays $P_s(\Delta\tau)$ shown in Fig. II.11 for various values of s . This distribution is wider for values of s closer to zero (Fig. II.11(a)). However if we rescale the distributions of time delays by their respective σ_s , as shown in Fig. II.11(b), the distributions become independent of s as $P_s(\Delta\tau) = \sigma_s[\Delta\tau] \hat{P}\left(\frac{\Delta\tau}{\sigma_s[\Delta\tau]}\right)$. This provides a strong numerical evidence supporting the existence of a universal distribution \hat{P} .

II.7 Discussion

In this chapter, we analyzed the discreteness effects at initial times in population dynamics. During the initial transient regime of the evolution of populations, there is a wide distribution of times at which the first series of jumps occurs. This means that fluctuations at initial times produce that some populations remain in their initial states for much longer than others, producing a gap in their individual evolution. This induces a relative shift that lasts forever. These effects play an important role specially for the determination of the large deviation function which may be obtained from the growth rate of the average log-population (Sec. I.7.1.1).

However, in Sec. II.2.2 we saw how by restricting the evolution of our populations up to a maximum time T_{\max} (or population N_{\max}) which is not “large enough”, the average population (and $\Psi(s)$) is strongly affected by the behavior of \mathcal{N} at initial times. We proposed

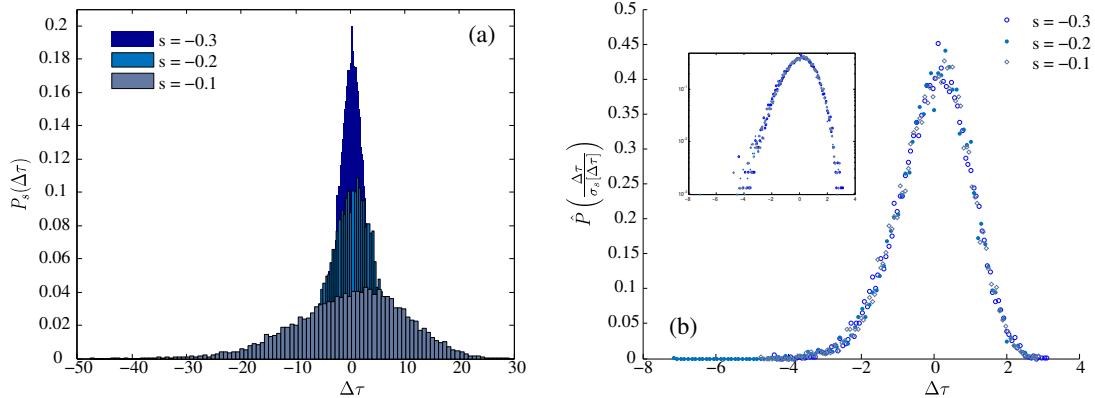


Figure II.11: **(a)** Distribution of time delays for different values of s . The dispersion of time delays is wider for values of s closer to zero. **(b)** Rescaled distribution of time delays $\hat{P}\left(\frac{\Delta\tau}{\sigma_s|\Delta\tau|}\right)$. The distribution of time delays depends only on their σ_s as $P_s(\Delta\tau) = \sigma_s[\Delta\tau] \hat{P}\left(\frac{\Delta\tau}{\sigma_s|\Delta\tau|}\right)$.

as an alternative to overcome the influence of initial discreteness effects to get rid of the regions of the populations where these effects are present. In other words, to cut the initial transient regime of the populations. In that case, we saw that the average of populations is restricted to the interval $[\max \mathcal{T}_C, \min \mathcal{T}_F]$ which can be in fact very small and this can induce a poor estimation of $\psi(s)$ (Fig. II.1(b)).

Complementary to this, we found a way of emphasizing the effects of the exponential growth regime in the determination of $\psi(s)$ by using the fact that log-populations after a long enough time become parallel (Fig. II.2(a)) and that once the populations have overcome the discreteness effects, the distance between them becomes constant (Fig. II.2(b)) and the discreteness effects are not strong anymore (Sec. II.3). We argued in Sec. II.4.1 that this initial discreteness effects or initial “lag” between populations could be compensated by performing over the populations a time translation (Eq. (II.2)). This time delay procedure is chosen so as to overlap the population evolutions in their large-time regime (Fig. II.4(b)). The improvement in the estimation of ψ comes precisely from these two main contributions, the time delaying of populations and the discarding of the initial transient regime of the populations.

We showed how the numerical estimations for the CGF are improved as the initial transient regime of the populations are discarded (independently of the method used to compute the growth rate of the average population, see Fig. II.7). Also, it was shown that if additionally, we perform the time delay procedure, the estimation of ψ is improved even more and closer to the theoretical value (Sec. II.5.2). This result was confirmed later in Sec. II.5.3 by computing the relative distance of the numerical estimators with respect to the theoretical value and their errors. As we observed in Fig. II.8, the deviation from the theoretical value is higher for values of s close to 0, but is smaller after the “time correction” for almost every value of s . Similarly for the error estimator (Fig. II.9).

Our numerical study was performed on a simple system, and we hope it can be extended to more complex phenomena. However, there remain open questions even for the system we have studied. The duration of the initial discrete-population regime could be understood from an analytical study of the population dynamics itself. Our numerical results also support a power-law behavior in time of the variance of the delays. Furthermore, it appeared that the distribution of the delays takes a universal form, after rescaling the variance to one. Those observations open questions for future studies.

III – Finite-Time and Finite-Size Scalings in the Evaluation of Large-Deviation Functions:

I. Analytical Study using a Birth-Death Process

III.1 Introduction

Cloning algorithms are numerical procedures aimed at simulating rare events efficiently, using a population dynamics scheme. In such algorithms, copies of the system are evolved in parallel and the ones showing the rare behavior of interest are multiplied iteratively [7, 16–19, 23, 32, 33, 62–70] (See Fig. III.1). One of these algorithms proposed by Giardinà et al. [7, 17–19, 33, 70] is used to evaluate numerically the cumulant generating function CGF (a large deviation function, LDF) of additive (or “time-extensive”) observables in Markov processes [1, 83]. While the method has been used widely, there have been less studies focusing on the analytical justification of the algorithm. Even though it is heuristically believed that the LDF estimator converges to the correct result as the number of copies N_c increases, there is no proof of this convergence. Related to this lack of the proof, although we use the algorithm by assuming its validity, we do not have any clue how fast the estimator converges as $N_c \rightarrow \infty$.

In this chapter [P2], we discuss this convergence defining two types of numerical errors. First, for a fixed finite N_c , averaging over a large number of realizations, the CGF estimator converges to an incorrect value, which is different from the desired large deviation result. We call this deviation from the correct value, **systematic errors**. Compared with these errors, we also consider the fluctuations of the estimated value. More precisely, for a fixed value of N_c , the results obtained in different realizations are distributed around this incorrect value. We call the errors associated to these fluctuations **stochastic errors**. Although both errors are important in numerical simulations, the former one can lead this algorithm to produce wrong results. For example as seen in Ref. [85], the systematic error grows exponentially as a temperature decreases (or generically in the weak noise limit of diffusive dynamics).

To study these errors, we employ a birth-death process [86, 87] description of the population dynamics algorithm as explained below: We focus on physical systems described by a Markov dynamics [7, 18, 19] with a finite number of states M . We denote by i

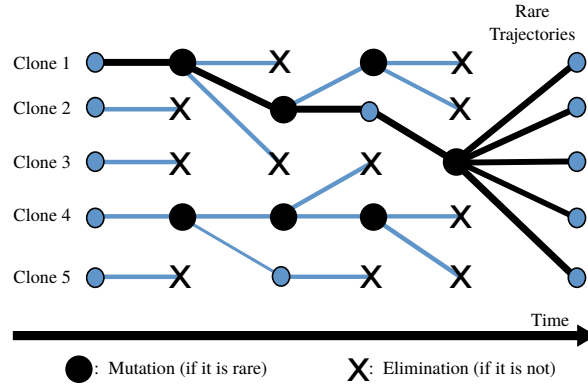


Figure III.1: Schematic picture illustrating the principle of the population dynamics algorithm. ‘Clones’ (or copies) of the system are prepared and they evolve following a mutation-and-selection process, maintaining the total population constant.

($i = 0, 1, \dots, M - 1$) the states of the system. This Markov process has its own stochastic dynamics, described by the transition rates $w(i \rightarrow j)$. In population dynamics algorithms, in order to study its rare trajectories, one prepares N_c copies of the system, and simulate these copies according to (i) the dynamics of $w(i \rightarrow j)$ (followed independently by all copies) and (ii) ‘cloning’ step in which the ensemble of copies is directly manipulated, i.e., some copies are eliminated while some are multiplied (See Table III.1). Formally, the population dynamics represents, for a *single* copy of the system, a process that does not preserve probability, as mentioned in Sec. I.6.2. This fact has motivated the studies of auxiliary processes [88], effective processes [89] and driven processes [90] to construct modified dynamics (and their approximations [91]) that preserve probability. Different from these methods, in this chapter, we formulate explicitly the **meta-dynamics** of the copies themselves by using a stochastic birth-death process which preserves probability, and it allows us to study the numerical errors of the algorithm when evaluating LDF. We consider the dynamics of the copies as a stochastic birth-death process whose state is denoted by $n = (n_0, n_1, n_2, \dots, n_{M-1})$, where $0 \leq n_i \leq N_c$ represents the number of copies which are in state i in the ensemble of copies. We explicitly introduce the transition rates describing the dynamics of n , which we denote by $\sigma(n \rightarrow \tilde{n})$. We show that the dynamics described by these transition rates lead in general to the correct LDF estimation of the original system $w(i \rightarrow j)$ in the $N_c \rightarrow \infty$ limit. We also show that the systematic errors are of the order $\mathcal{O}(1/N_c)$, whereas the numerical errors are of the order $\mathcal{O}(1/(\tau N_c))$ (where τ is an averaging duration). This result is in clear contrast with standard Monte-Carlo methods, where the systematic errors are always 0. The formulation developed in this chapter provides us the possibility to compute exactly the expressions of the convergence coefficients, as we do in Sec. III.4 on a simple example. The analytical analysis presented here [P2] is supplemented with a thorough numerical study in the next chapter IV [P3]. There, we employ an intrinsically different cloning algorithm, which is the **continuous-time** population dynamics algorithm, that cannot be studied by the methods presented here (see Secs. III.2.4.2 and IV.5). We show in chapter IV that the validity of the scaling that we derive analytically here is very general, we make use of the convergence speed to propose a simple interpolation technique demonstrating in practice its efficiency in the

	Population dynamics algorithm	Birth-death process describing the population dynamics
State of the system	i $(i = 0, 1, \dots, M - 1)$	$n = (n_0, n_1, \dots, n_{M-1})$ $(0 \leq n_i \leq N_c \text{ with } \sum_i n_i = N_c)$
Transition rates	$w(i \rightarrow j)$ Markov process on states i	$\sigma(n \rightarrow \tilde{n})$ Markov process on states n
Numerical procedure for rare-event sampling	Prepare N_c clones and evolve those with a mutation-selection procedure	Described by the dynamics of rates $\sigma(n \rightarrow \tilde{n})$

Table III.1: Correspondence between the population dynamics and the birth-death process to describe it.

evaluation of the LDF, irrespectively of the details of the population dynamics algorithm.

The present chapter is structured as follows. We first define the LDF problem in the beginning of Sec. III.2, and then formulate the birth-death process used to describe the algorithm in Sec. III.2.1. By using this birth-death process, we demonstrate that the estimator of the algorithm converges to the correct large deviation function in Sec. III.2.2. At the end of this section, in Sec. III.2.3, we discuss the convergence speed of this estimator (the systematic errors) and derive its scaling $\sim 1/N_c$. In Sec. III.3, we turn to stochastic errors. For discussing this, we introduce the large deviation function of the estimator, from which we derive that the convergence speed of the stochastic errors is proportional to $1/(\tau N_c)$. In the next section, Sec. III.4, we introduce a simple two-state model, to which we apply the formulations developed in the previous sections. We derive the exact expressions of the systematic errors in Sec. III.4.1 and of the stochastic errors in Sec. III.4.2. Then, in Sec. III.4.3, we propose another large deviation estimator and finally, in Sec. III.5, we summarize the results obtained.

III.2 Birth-Death Process and the Population Dynamics Algorithm

As explained in the introduction of this chapter (also see Table III.1), the state of the population is $n = (n_0, n_1, \dots, n_{M-1})$, where n_i represents the number of clones in the state i . The total population is preserved: $\sum_i n_i = N_c$. Below, we introduce the transition rates of the dynamics between the occupations n , $\sigma(n \rightarrow \tilde{n})$ that describe corresponding large deviations of the original system whose dynamics is given by the rates $w(i \rightarrow j)$.

As the **original system**, we consider the continuous-time Markov process in a discrete-time representation. By denoting by dt the time step, the transition matrix $R_{j,i}$ for time evolution of the state i is described as

$$R_{j,i} = \delta_{i,j} + dt \left[w(i \rightarrow j) - \delta_{i,j} \sum_k w(i \rightarrow k) \right], \quad (\text{III.1})$$

where we set $w(i \rightarrow i) = 0$. The probability distribution of the state i , $p_i(t)$, evolves in time as $p_i(t + dt) = \sum_j R_{i,j} p_j(t)$. In the $dt \rightarrow 0$ limit, one obtains the continuous-time Master equation (I.8) describing the evolution of $p_i(t)$ [86, 87]. For simplicity, especially for

the cloning part of the algorithm, we keep here a small finite dt . The reason why we use a discrete-time representation is solely for simplicity of the discussion. The main results can be derived even if we start with a continuous-time representation (see Sec. III.2.4.1). For the original dynamics described by the transition matrix (III.1), we consider an observable b_i depending on the state i and we are interested in the distribution of its time-averaged value during a time interval τ , defined as

$$B(\tau) = \frac{1}{\tau} \sum_{t=0}^{\tau/dt} dt b_{i(t)}. \quad (\text{III.2})$$

Here $i(t)$ is a trajectory of the system generated by the Markov dynamics described by $R_{j,i}$. We note that $B(\tau)$ is a path- (or history-, or realization-) dependent quantity. Since $\tau B(\tau)$ is an additive observable, the fluctuations of $B(\tau)$ depending on the realizations are small when τ is large, but one can describe the large deviations of $B(\tau)$. Those occur with a small probability, and obey a large deviation principle (I.20). We denote by $\text{Prob}(B)$ the distribution function of $B(\tau)$. The large deviation principle ensures that $\text{Prob}(B)$ takes an asymptotic form $\text{Prob}(B) \sim \exp(-\tau I(B))$ for large τ , where $I(B)$ is a large deviation function (or ‘rate function’) [1, 83]. As we mentioned in the introduction, if the rate function $I(B)$ is convex, the large deviation function is expressed as a Legendre transform of a cumulant generating function $\psi(s)$ defined as

$$\psi(s) = \lim_{\tau \rightarrow \infty} \frac{1}{\tau} \log \langle e^{-s\tau B(\tau)} \rangle,$$

so that $I(B) = -\inf_s [sB + \psi(s)]$. The large deviation function $I(B)$ and this generating function $\psi(s)$ are by definition difficult to evaluate numerically in Monte-Carlo simulations of the original system of transition rates $w(i \rightarrow j)$ (see for example Ref. [255]). To overcome this difficulty, population dynamics algorithms have been developed [7, 17–19, 33, 70]. Here, we describe this population dynamics algorithm by using a birth-death process on the occupation state n allowing us to study systematically the errors in the estimation of $\psi(s)$ within the population dynamics algorithm. We mention that, without loss of generality, we restrict our study to so-called ‘type-B’ observable (see Sec. I.5 in the Introduction) that do not depend on the transitions of the state [40], i.e., which are time integrals of the state of the system, as in Eq. (III.2). Indeed, as explained for example in Refs. [7] and [85], one can always reformulate the determination of the CGF of mixed-type observables into that of a type-B variable, by modifying the transition rates of the given system.

Note that in chapter IV, we use a continuous-time version of the algorithm [19] to study an observable of ‘type A’. This version of the algorithm differs from the one considered here, in the sense that after its selection step, a copy in the population can have strictly more than one offspring. This results in an important difference: the effective interaction between copies due to the cloning/pruning procedure is unbounded (it can *a priori* affect any proportion of the population), while in the discrete-time settings of the present chapter, this effective interaction is restricted to a maximum of one cloning/pruning event. However, we observe numerically in chapter IV that the same finite-time and finite-population size scalings are present, illustrating their universal character.

	Transition matrices
Dynamics (“mutations”)	$\mathcal{T}_{\tilde{n},n} \equiv \delta_{\tilde{n},n} + dt \sum_{i=0}^{M-1} n_i \sum_{j=0, (j \neq i)}^{M-1} w(i \rightarrow j) \left[\delta_{\tilde{n}_i, n_i-1} \delta_{\tilde{n}_j, n_j+1} \delta_{\tilde{n},n}^{i,j} - \delta_{\tilde{n},n} \right]$
Cloning (“selection”)	$\mathcal{C}_{\tilde{n},n} = \delta_{\tilde{n},n} + s dt \sum_{i=0}^{M-1} n_i \alpha_i \left[\delta_{\tilde{n}_i, n_i + \alpha_i / \alpha_i } \delta_{\tilde{n},n}^i - \delta_{\tilde{n},n} \right] + \mathcal{O}(dt^2)$
Maintaining N_c	$\mathcal{K}_{\tilde{n},n} = \delta_{\sum_{i=0}^{M-1} n_i, N_c} \delta_{\tilde{n},n} + \sum_{k=-1,1} \delta_{\sum_{i=0}^{M-1} n_i, N_c+k} \sum_{i=0}^{M-1} \delta_{\tilde{n}_i, n_i-k} \delta_{\tilde{n},n}^i \frac{n_i}{N_c+k}$
Full process	$(\mathcal{KCT})_{\tilde{n},n} = \delta_{\tilde{n},n} + dt \sum_{i=0}^{M-1} n_i \sum_{j=0, (j \neq i)}^{M-1} [w(i \rightarrow j) + s \tilde{w}_n(i \rightarrow j)] \left[\delta_{\tilde{n}_i, n_i-1} \delta_{\tilde{n}_j, n_j+1} \delta_{\tilde{n},n}^{i,j} - \delta_{\tilde{n},n} \right]$ with $\tilde{w}_n(i \rightarrow j) = \frac{n_i}{N_c} \left[\alpha_j \delta_{j \in \Omega(+)} \frac{N_c}{N_c+1} - \alpha_i \delta_{i \in \Omega(-)} \frac{N_c}{N_c-1} \right]$

Table III.2: Transition matrices (Eq. (III.3)) describing the birth-death process.

III.2.1 Transition Matrices Representing the Population Dynamics Algorithm

We denote the probability distribution of the occupation n at time t by $P_n(t)$. The time-evolution of this probability is decomposed into three parts. The first one is the original Monte-Carlo dynamics based on the transition rates $w(i \rightarrow j)$. The second one is the cloning procedure of the population dynamics algorithm, which favors or disfavors configurations according to a well-defined rule. The third one is a supplementary (but important) part which maintains the total number of clones to a constant N_c . We denote the transition matrices corresponding to these steps by \mathcal{T} , \mathcal{C} and \mathcal{K} , respectively. By using these matrices, then, the time evolution of the distribution function is given as

$$P_n(t + dt) = \sum_{\tilde{n}} (\mathcal{KCT})_{n,\tilde{n}} P_{\tilde{n}}(t). \quad (\text{III.3})$$

We derive explicit expressions of these matrices in the following sub-sections. A summary of the results obtained can be found in Table III.2.

III.2.1.1 Derivation of the Original Dynamics Part: \mathcal{T}

We first consider the transition matrix \mathcal{T} , which describes the evolution of the occupation state n solely due to the dynamics based on the rates $w(i \rightarrow j)$. During an infinitesimally small time step dt , the occupation $n = (n_0, n_1, \dots, n_{M-1})$ changes to $\tilde{n} = (n_0, n_1, \dots, n_i - 1, \dots, n_j + 1, \dots, n_{M-1})$ where $0 \leq i < M$ and $0 \leq j < M$ (for all $i \neq j$). Since there are n_i clones in the state i before the transition, the transition probability of this change is given as $n_i w(i \rightarrow j) dt$. Thus, we obtain

$$\mathcal{T}_{\tilde{n},n} \equiv \delta_{\tilde{n},n} + dt \sum_{i=0}^{M-1} n_i \sum_{j=0, (j \neq i)}^{M-1} w(i \rightarrow j) \left[\delta_{\tilde{n}_i, n_i-1} \delta_{\tilde{n}_j, n_j+1} \delta_{\tilde{n},n}^{i,j} - \delta_{\tilde{n},n} \right],$$

where $\delta_{\tilde{n},n}^{i,j}$ is a Kronecker δ for the indices except for i, j : $\delta_{\tilde{n},n}^{i,j} \equiv \prod_{k \neq i, j} \delta_{\tilde{n}_k, n_k}$. One can easily check that this matrix satisfies the conservation of the probability: $\sum_{\tilde{n}} \mathcal{T}_{\tilde{n},n} = 1$. It corresponds to the evolution of N_c independent copies of the original system with rates $w(i \rightarrow j)$.

III.2.1.2 Derivation of the Cloning Part: \mathcal{C}

In the population dynamics algorithm (for example the one described in the Appendix A of Ref. [85]), at every certain time interval Δt , one evaluates the exponential factor for all clones equal to $e^{-s \int_t^{t+\Delta t} dt' b_i(t')}$ if the clone is in state $(i(t'))_{t'=t}^{t+\Delta t}$ during a time interval $t \leq t' \leq t + \Delta t$. This cloning factor determines whether each clone is copied or eliminated after this time interval. In the continuous-time version of the algorithm this factor was given by Eq. (I.28). Although the details of how to determine this selection process can depend on the specific type of algorithms, the common idea is that each of the clones is copied or eliminated in such a way that a clone in state $i(t)$ has a number of descendant(s) proportional to the cloning factor on average after this time interval.

In order to implement this idea in our birth-death process, we assume this time step Δt to be small. For the sake of simplicity, we set this Δt to be our smallest time interval dt : $\Delta t = dt$. This condition is not mandatory whenever the $\Delta t \rightarrow 0$ limit is taken at the end (see Sec. III.2.4.1 for the case $\Delta t > dt$). Then, noticing that the time integral $\int_t^{t+\Delta t} dt' b_i(t')$ is expressed as $dt b_i(t)$ for small dt , we introduce the following quantity for each state i ($i = 0, 1, 2, \dots, M-1$):

$$\nu_i \equiv \frac{n_i e^{-s dt b_i}}{\sum_{j=0}^{M-1} n_j e^{-s dt b_j}} N_c. \quad (\text{III.4})$$

Note that there is a factor n_i in front of the exponential function $e^{-s dt b_i}$ which enumerates the number of clones that occupy the state i . The quantity ν_i is aimed at being the number of clones in state i after the cloning process, however, since ν_i is not an integer but a real number, one needs a supplementary prescription to fix the corresponding integer number of descendants. In general, in the implementation of population dynamics, this integer is generated randomly from the factor ν_i , equal either to its lower or to its upper integer part. The probability to choose either the lower or upper integer part is fixed by imposing that the number of descendants is equal to ν_i on average. For instance, if ν_i is equal to 13.2, then 13 is chosen with probability 0.8, and 14 with probability 0.2. Generically, $[\nu_i]$ and $[\nu_i] + 1$ are chosen with probability $1 + [\nu_i] - \nu_i$ and $\nu_i - [\nu_i]$, respectively. We note that we need to consider these two possibilities for all indices i . We thus arrive at the following matrix:

$$\begin{aligned} \mathcal{C}_{\tilde{n}, n} &\equiv \sum_{x_0=0}^1 \sum_{x_1=0}^1 \sum_{x_2=0}^1 \cdots \sum_{x_{M-1}=0}^1 \prod_{i=0}^{M-1} \\ &\times \delta_{\tilde{n}_i, [\nu_i] + x_i} [(\nu_i - [\nu_i]) x_i + (1 + [\nu_i] - \nu_i) (1 - x_i)]. \end{aligned} \quad (\text{III.5})$$

Now, we expand \mathcal{C} at small dt and we keep only the terms proportional to $\mathcal{O}(1)$ and $\mathcal{O}(dt)$, which do not vanish in the continuous-time limit. For this purpose, we expand ν_i as

$$\nu_i = n_i \left[1 + s dt \left(\sum_j \frac{n_j b_j}{N_c} - b_i \right) \right] + \mathcal{O}(dt^2),$$

where we have used $\sum_i n_i = N_c$. This expression indicates that $[\nu_i]$ is determined depending on the sign of $\sum_j n_j b_j / N_c - b_i$, where we assumed $s > 0$ for simplicity without loss of generality (because when $s < 0$, we can always re-define $-b$ as b to make s to be positive). By denoting this factor by α_i , i.e.,

$$\alpha_i(n) \equiv \sum_j \frac{n_j b_j}{N_c} - b_i, \quad (\text{III.6})$$

we thus define the following state-space $\Omega^{(\pm)}(n)$:

$$\Omega^{(\pm)}(n) = \{ i \mid 0 \leq i < M \text{ and } \pm \alpha_i(n) > 0 \}.$$

From this definition, for sufficiently small dt , we obtain $\lfloor \nu_i \rfloor = n_i$ for $i \in \Omega^{(+)}$, and $\lfloor \nu_i \rfloor = n_i - 1$ for $i \in \Omega^{(-)}$. Substituting these results into Eq. (III.5) and expanding in dt , we obtain (denoting here and thereafter $\alpha_i = \alpha_i(n)$):

$$\mathcal{C}_{\tilde{n},n} = \delta_{\tilde{n},n} + s dt \sum_{i=0}^{M-1} n_i |\alpha_i| [\delta_{\tilde{n}_i, n_i + \alpha_i / |\alpha_i|} \delta_{\tilde{n}, n}^i - \delta_{\tilde{n}, n}] + \mathcal{O}(dt^2), \quad (\text{III.7})$$

where $\delta_{\tilde{n}, n}^i$ is a Kronecker delta for the indices except for i : $\delta_{\tilde{n}, n}^i = \prod_{k \neq i} \delta_{\tilde{n}_k, n_k}$. One can easily check that this matrix preserves probability: $\sum_{\tilde{n}} \mathcal{C}_{\tilde{n}, n} = 1$.

III.2.1.3 Derivation of the Maintaining Part: \mathcal{K}

As directly checked, the operator \mathcal{T} preserves the total population $\sum_i n_i$. However, the operator representing the cloning \mathcal{C} , does not. In our birth-death implementation, this property originates from the rounding process $\lfloor \nu_i \rfloor$ in the definition of \mathcal{C} : even though ν_i itself satisfies $\sum_i \nu_i = N_c$, because of the rounding process of ν_i , the number of clones after multiplying by \mathcal{C} (that is designed to be proportional to ν_i on average) can change. There are several ways to keep the number N_c of copies constant without biasing the distribution of visited configurations. One of them is to choose randomly and *uniformly* δN_c clones from the ensemble, where δN_c is equal to the number of excess (resp. lacking) clones with respect to N_c , and to eliminate (resp. multiply) them.

In our birth-death description, we implement this procedure as follows. We denote by \mathcal{K} the transition matrix maintaining the total number of clones to be the constant N_c . We now use a continuous-time asymptotics $dt \rightarrow 0$. In this limit, from the expression of the transition matrix elements Eq. (III.7), we find that at each cloning step the number of copies of the cloned configuration varies by ± 1 at most. Hence, the total number of clones after multiplying by \mathcal{C} , $\sum_i n_i$, satisfies the following inequality

$$N_c - 1 \leq \sum_i n_i \leq N_c + 1.$$

Among the configurations n that satisfy this inequality, there are three possibilities, which are $\sum_i n_i = N_c$ and $\sum_i n_i = N_c \pm 1$. If n satisfies $\sum_i n_i = N_c$, we do not need to adjust n , while if n satisfies $\sum_i n_i = N_c + 1$ (resp. $\sum_i n_i = N_c - 1$), we eliminate (resp. multiply) a clone chosen randomly and uniformly. Note that, in our formulation, we do not distinguish the clones taking the same state. This means that we can choose one of the occupations n_i of a state i according to a probability proportional to the number of copies n_i in this state. In other words, the probability to choose the state i and to copy or to eliminate a clone from this state is proportional to $n_i / \sum_{j=0}^{M-1} n_j$. Therefore, we obtain the expression of the matrix \mathcal{K} as

$$\mathcal{K}_{\tilde{n}, n} = \delta_{\sum_i n_i, N_c} \delta_{\tilde{n}, n} + \sum_{k=-1, 1} \delta_{\sum_i n_i, N_c + k} \sum_{i=0}^{M-1} \delta_{\tilde{n}_i, n_i - k} \delta_{\tilde{n}, n}^i \frac{n_i}{N_c + k}$$

for \tilde{n} that satisfies $\sum_i \tilde{n}_i = N_c$, and $\mathcal{K}_{\tilde{n}, n} = 0$ otherwise.

III.2.1.4 Total Transition: \mathcal{KCT}

From the obtained expressions we calculate the matrix \mathcal{KCT} , which describes the total transition of the population dynamics (III.3)

$$(\mathcal{KCT})_{\tilde{n},n} = \delta_{\tilde{n},n} + dt \sum_{i=0}^{M-1} n_i \sum_{j=0, (j \neq i)}^{M-1} [w(i \rightarrow j) + s \tilde{w}_n(i \rightarrow j)] \times [\delta_{\tilde{n}_i, n_i - 1} \delta_{\tilde{n}_j, n_j + 1} \delta_{\tilde{n}, n}^{i,j} - \delta_{\tilde{n}, n}], \quad (\text{III.8})$$

where the population-dependent transition rate $\tilde{w}_n(i \rightarrow j)$ is given as

$$\tilde{w}_n(i \rightarrow j) = \frac{n_j}{N_c} \left[\alpha_j \delta_{j \in \Omega(+)} \frac{N_c}{N_c + 1} - \alpha_i \delta_{i \in \Omega(-)} \frac{N_c}{N_c - 1} \right].$$

The comparison of the expression (III.8) with the original part \mathcal{T} provides an insight into the result obtained. The jump ratio $w(i \rightarrow j)$ in the original dynamics is replaced by $w(i \rightarrow j) + s \tilde{w}_n(i \rightarrow j)$ in the population dynamics algorithm. We note that this transition rate depends on the population n , meaning that we cannot get a closed equation for this modified dynamics at the level of the states i in general. We finally remark that the transition matrix $\sigma(n \rightarrow \tilde{n})$ for the continuous-time limit is directly derived from Eq. (III.8) as

$$\sigma(n \rightarrow \tilde{n}) = \sum_{i=0}^{M-1} n_i \sum_{j=0, (j \neq i)}^{M-1} [w(i \rightarrow j) + s \tilde{w}_n(i \rightarrow j)] \times [\delta_{\tilde{n}_i, n_i - 1} \delta_{\tilde{n}_j, n_j + 1} \delta_{\tilde{n}, n}^{i,j}]. \quad (\text{III.9})$$

III.2.2 Derivation of the Large Deviation Results in the $N_c \rightarrow \infty$ Asymptotics

In this subsection, we study the $N_c \rightarrow \infty$ limit for the transition matrix of rates $\sigma(n \rightarrow \tilde{n})$, and derive the validity of the population dynamics algorithm.

III.2.2.1 Estimator of the Large Deviation Function

One of the ideal implementations of the population dynamics algorithm is as follows: We make copies of each clone at the end of simulation, where the number of copies for each realization is equal to the exponential weight $e^{-s\tau B(\tau)}$ in Eq. (III.2) (so that we can discuss an ensemble with this exponential weight without multiplying the probability by it, as described in the Introduction for the continuous-time case). In this implementation, the number of clones grows (or decays) exponentially proportionally as $\langle e^{-s\tau B(\tau)} \rangle$ by definition. In real implementations of the algorithm, however, since taking care of an exponentially large or small number of clones can cause numerical problems, one rather keeps the total number of clones to a constant N_c at every time step, as seen in Eq (III.4). Within this implementation, we reconstruct the exponential change of the total number of clones as follows: We compute the average of the cloning factor at each cloning step, and we store the product of these ratios along the cloning steps. At final time, this product gives the empirical estimation of total (unnormalized) population during the whole duration of the simulation [P1], i.e., an estimator of $\langle e^{-s\tau B(\tau)} \rangle$. One thus estimates the CGF $\psi(s)$ given in Eq. (III.2) [7, 17–19, 33, 70] as the logarithm of this reconstructed population, divided by the total time.

In this formulation, the average cloning ratio is given as $\sum_i n_i e^{-s \int dt b_i} / N_c$, and thus the multiplication over whole time interval reads $\prod_{t=0}^{\tau/dt} \{n_i(t) e^{-s \int dt b_i} / N_c\}$. Given that we

empirically assume that the CGF estimator converges to $\psi(s)$ in the $N_c, \tau \rightarrow \infty$ limit, the following equality is expected to hold in probability 1:

$$\psi(s) \stackrel{?}{=} \lim_{N_c \rightarrow \infty} \lim_{\tau \rightarrow \infty} \frac{1}{\tau} \sum_{t=0}^{\tau/dt} \log \sum_i \frac{n_i(t) e^{-s dt b_i}}{N_c} + O(dt). \quad (\text{III.10})$$

Since the dynamics of the population n is described by a Markov process, ergodicity is satisfied, i.e., time averages can be replaced by the expected value with respect to the stationary distribution function which applied to the right-hand side of Eq (III.10), we obtain

$$\lim_{\tau \rightarrow \infty} \frac{1}{\tau} \sum_{t=0}^{\tau/dt} \log \sum_i \frac{n_i(t) e^{-s dt b_i}}{N_c} = \frac{1}{dt} \sum_n P_n^{\text{st}} \log \sum_i \frac{n_i e^{-s dt b_i}}{N_c} + \mathcal{O}(dt),$$

where P_n^{st} is the stationary distribution function of the population n in the $dt \rightarrow 0$ limit, (namely, P_n^{st} is the stationary distribution of the dynamics of transition rates $\sigma(n \rightarrow \tilde{n})$). By expanding this right-hand side with respect to dt , we rewrite the expected equality (III.10) as

$$\psi(s) \stackrel{?}{=} -s \lim_{N_c \rightarrow \infty} \sum_n P_n^{\text{st}} \sum_i \frac{n_i b_i}{N_c} + O(dt). \quad (\text{III.11})$$

where we used that $\sum_i n_i = N_c$ is a conserved quantity. Below we demonstrate that this latter equality (III.11) is satisfied by analyzing the stationary distribution function P_n^{st} .

III.2.2.2 Connection between the Distribution Functions of the Population and of the Original System

From the definition of the stationary distribution function P_n^{st} , we have

$$\sum_{\tilde{n}} P_{\tilde{n}}^{\text{st}} \sigma(\tilde{n} \rightarrow n) - \sum_{\tilde{n}} P_n^{\text{st}} \sigma(n \rightarrow \tilde{n}) = 0, \quad (\text{III.12})$$

which is a stationary Master equation. In this equation, we use the explicit expression of σ shown in Eq. (III.9). By denoting by $n^{j \rightarrow i}$ the configuration where one clone in the state j moves to the state i : $n^{j \rightarrow i} \equiv (n_0, n_1, \dots, n_i + 1, \dots, n_j - 1, \dots, n_{M-1})$, the stationary Master equation (III.12) is rewritten as

$$\sum_{i,j(i \neq j)} \left[f_{i \rightarrow j}(n^{j \rightarrow i}) - f_{i \rightarrow j}(n) \right] = 0, \quad (\text{III.13})$$

where we defined $f_{i \rightarrow j}(n)$ as

$$f_{i \rightarrow j}(n) = P_n^{\text{st}} n_i [w(i \rightarrow j) + s \bar{w}_n(i \rightarrow j)].$$

Now we multiply expression (III.13) by n_k (k is arbitrary from $k = 0, 1, 2, \dots, M-1$), and sum it over all configurations n :

$$\sum_n \sum_{i,j(i \neq j)} n_k \left[f_{i \rightarrow j}(n^{j \rightarrow i}) - f_{i \rightarrow j}(n) \right] = 0. \quad (\text{III.14})$$

We can change the dummy summation variable n in the first term to $n^{i \rightarrow j}$, which leads to $\sum_n \sum_{i,j(i \neq j)} (n^{i \rightarrow j})_k f_{i \rightarrow j}(n)$. Since the second term has almost the same expression as the

first one except for the factor n_k , the sum in Eq. (III.14) over the indices (i, j) , where none of i nor j is equal to k , becomes 0. The remaining term in Eq. (III.14) is thus

$$0 = \sum_n \sum_{j(j \neq k)} \left((n^{k \rightarrow j})_k - n_k \right) f_{k \rightarrow j}(n) + \sum_n \sum_{i(i \neq k)} \left((n^{i \rightarrow k})_k - n_k \right) f_{i \rightarrow k}(n).$$

Using the definition of $n^{i \rightarrow j}$ in this equation, we arrive at

$$0 = \sum_n \left[\sum_{i(i \neq k)} f_{i \rightarrow k}(n) - f_{k \rightarrow i}(n) \right]. \quad (\text{III.15})$$

This equation (III.15) connects the stationary property of the population dynamics (described by the occupation states n) and the one in the original system (described by the states i).

The easiest case where we can see this connection is when $s = 0$. By defining the empirical occupation probability of the original system as $p_i \equiv \sum_n P_n^{\text{st}} n_i / N_c$, Eq. (III.15) leads to the following (stationary) master equation for $w(i \rightarrow j)$:

$$0 = \sum_j p_j w(j \rightarrow i) - \sum_j p_i w(i \rightarrow j) \quad (\text{for } s = 0) \quad (\text{III.16})$$

This is valid for any N_c , meaning that, for original Monte-Carlo simulations in $s = 0$, the empirical probability p_i is exactly equal to the steady-state probability, as being the unique solution of Eq. (III.16). It means that there are no systematic errors in the evaluation of p_i . However, in the generic case $s \neq 0$, this property is not satisfied. One thus needs to understand the $N_c \rightarrow \infty$ limit to connect the population dynamics result with the large deviation property of the original system.

III.2.2.3 Justification of the Convergence of the Large Deviation Estimator as Population Size becomes Large

We define a scaled variable x_i as n_i / N_c . While keeping this occupation fractions x_i to be $\mathcal{O}(1)$, we take the $N_c \rightarrow \infty$ limit in Eq. (III.15), which leads to

$$0 = \sum_n P_n^{\text{st}} \left[\sum_j x_j w(j \rightarrow i) - \sum_j x_i w(i \rightarrow j) - s x_i \left(b_i - \sum_k x_k b_k \right) \right] + \mathcal{O}(1/N_c). \quad (\text{III.17})$$

Inspired by this expression, we define a matrix $L_{i,j}^s$ as

$$L_{i,j}^s = w(j \rightarrow i) - \delta_{i,j} \left(\sum_k w(i \rightarrow k) + s b_i \right),$$

and a correlation function between x_i and x_j as

$$c_{i,j} = \sum_n x_i x_j P_n^{\text{st}} - p_i p_j,$$

(where we recall $p_i \equiv \sum_n x_i P_n^{\text{st}}$). From these definitions, Eq.(III.17) is rewritten as

$$\sum_j p_j L_{i,j}^s = -s p_i \sum_k p_k b_k - s \sum_k c_{i,k} b_k + \mathcal{O}\left(\frac{1}{N_c}\right).$$

Since $x_i \equiv n_i/N_c$ is an averaged quantity (an arithmetic mean) with respect to the total number of clones, we can safely assume that the correlation $c_{i,j}$ becomes 0 in $N_c \rightarrow \infty$ limit:

$$\lim_{N_c \rightarrow \infty} c_{i,k} = 0. \quad (\text{III.18})$$

For a more detailed discussion of why this is valid, see the description after Eq. (III.20). Thus, by defining $p_i^\infty \equiv \lim_{N_c \rightarrow \infty} p_i$, we obtain

$$\sum_j p_j^\infty L_{i,j}^s = -s p_i^\infty \sum_k p_k^\infty b_k.$$

From the Perron-Frobenius theory, the positive eigenvector of the matrix $L_{i,j}^s$ is unique and corresponds to its eigenvector of largest eigenvalue (in real part). This means that $-s \sum_k p_k^\infty b_k$ is the largest eigenvalue of the matrix $L_{i,j}^s$. Finally, by recalling that the largest eigenvalue of this matrix $L_{i,j}^s$ is equal to the generating function $\psi(s)$ (see Ref. [40] for example), we have justified that the CGF estimator (III.11) is valid in the large- N_c limit. This is equivalent to what we saw in Sec. I.6.1 for the continuous-time version.

III.2.3 Systematic Errors due to Finite N_c : Convergence Speed of the Large Deviation Estimator as $N_c \rightarrow \infty$

In the introduction of this chapter, we defined the **systematic errors** as the deviations of the large deviation estimator from the correct value due to a finite number of clones N_c . From Eq. (III.11), we quantitatively define this systematic error ϵ_{sys} as

$$\epsilon_{\text{sys}} \equiv \left| \psi(s) + s \sum_i p_i b_i \right|. \quad (\text{III.19})$$

From a simple argument based on a system size expansion, we show below that this ϵ_{sys} is of order $\mathcal{O}(1/N_c)$. We first show that one can perform a system size expansion (as for example in Ref. [86]) for the population dynamics. In Eq. (III.13), by recalling the definition of the vector x as $x = n/N_c$, and by denoting $\tilde{P}^{\text{st}}(x) = P_{xN_c}^{\text{st}}$, we obtain

$$0 = \sum_{i,j(i \neq j)} \sum_{r=1}^{\infty} \frac{1}{r!} \frac{1}{N_c^r} \left(\frac{\partial}{\partial x_i} - \frac{\partial}{\partial x_j} \right)^r x_i \tilde{P}^{\text{st}}(x) \times [w(i \rightarrow j) + s \tilde{w}_n(i \rightarrow j)|_{n=xN_c}]. \quad (\text{III.20})$$

This indicates that the stochastic process governing the evolution of x becomes deterministic in the $N_c \rightarrow \infty$ limit. The deterministic trajectory for x is governed by a differential equation derived from the sole term $r = 1$ in the expansion (III.20) (see Sec. 3.5.3 in ref. [87] for the detail of how to derive this property). Thus if x converges to a fixed point as N_c increases, which is normally observed in implementations of cloning algorithms, the assumption (III.18) is satisfied.

From the expression of ϵ_{sys} , we see that the dependence in N_c comes solely from p_i , which can be calculated from the first order correction of P_n^{st} (at large N_c). The equation to determine P_n^{st} is the stationary Master equation (III.12) or equivalently, the system-size expansion formula (III.20). We expand the jump ratio $w(i \rightarrow j) + s\tilde{w}_n(i \rightarrow j)$ in Eq. (III.20) with respect to $1/N_c$ as:

$$w(i \rightarrow j) + s\tilde{w}_n(i \rightarrow j) = w(i \rightarrow j) + s\tilde{w}_x^\infty(i \rightarrow j) + \frac{s}{N_c}\delta w_x(i \rightarrow j) + \mathcal{O}(1/N_c^2), \quad (\text{III.21})$$

where $\tilde{w}_x^\infty(i \rightarrow j)$ and $\delta w_x(i \rightarrow j)$ are defined as

$$\tilde{w}_x^\infty(i \rightarrow j) = x_j \left[\alpha_j \delta_{j \in \Omega(+)} - \alpha_i \delta_{i \in \Omega(-)} \right]$$

and

$$\delta w_x(i \rightarrow j) = -x_j \left[\alpha_j \delta_{j \in \Omega(+)} + \alpha_i \delta_{i \in \Omega(-)} \right].$$

By substituting Eq. (III.21) into the system-size expansion formula (III.20) and performing a perturbation expansion, we find that a first-order correction of p is naturally of order $\mathcal{O}(1/N_c)$, i.e., $\epsilon_{\text{sys}} = \mathcal{O}(1/N_c)$. For a practical scheme of how to implement this perturbation on a specific example, see Sec. III.4.1. In chapter IV [P3], the scaling analysis of the $1/N_c$ correction is shown to hold numerically with the continuous-time cloning algorithm. We also show that the $1/N_c$ correction behavior remains in fact valid at finite time, an open question that remains to be investigated analytically.

III.2.4 Remarks

Here, we discuss some remarks on the formulation presented in this section.

III.2.4.1 Relaxing the Condition $dt = \Delta t$

In Sec. III.2.1.2, we set the discretization time of the process dt to be equal to the time interval for cloning Δt , and we took the $dt = \Delta t \rightarrow 0$ limit at the end. We note that the condition $\Delta t = dt$ is not necessary if both limits $\Delta t \rightarrow 0$ and $dt \rightarrow 0$ (with $dt < \Delta t$) are taken at the end. This is practically important, because we can use the continuous-time process to perform the algorithm presented here by setting $dt = 0$ first, and $\Delta t \rightarrow 0$ limit afterwards. More precisely, replacing dt by Δt in the matrix \mathcal{C} and \mathcal{K} , we build a new matrix $\mathcal{K}\mathcal{C}(\mathcal{T}^{\Delta t/dt})$. Taking the $dt \rightarrow 0$ limit in this matrix while keeping Δt non-infinitesimal (but small), this matrix represents the population dynamics algorithm of a continuous-time process with a finite cloning time interval Δt . The arguments presented in this section can then be applied in the same way, replacing dt by Δt . We note that the deviation due to a non-infinitesimal Δt should thus appear as $\mathcal{O}(\Delta t)$ (see Eq. (III.10) for example).

III.2.4.2 A Continuous-Time Cloning Algorithm

The $\Delta t \rightarrow 0$ limit is the key point in the formulation developed in this section (and in this chapter). Thanks to this limit, upon each cloning step, the total number of clones $\sum_{j=0}^{M-1} n_j$ always varies only by ± 1 , which makes the expression of the matrices \mathcal{C} and \mathcal{K} simple enough to develop the arguments presented in Secs. III.2.2 and III.2.3. Furthermore, during the time interval Δt separating two cloning steps, the configuration is changing at most once. The

process between cloning steps is thus simple, which allows us to represent the corresponding time-evolution matrix as \mathcal{T} (by replacing dt by Δt as just explained above). Generalizing our analytical study to a cloning dynamics in which the limit $\Delta t \rightarrow 0$ is not taken is therefore a very challenging task, which is out of the scope of this chapter.

However, interestingly, in chapter IV [P3] we observe numerically that our predictions for the finite-time and finite-population scalings are still valid in a different version of algorithm for which $\sum_{j=0}^{M-1} n_j$ can vary by an arbitrary amount – supporting the hypothesis that the analytical arguments that we present here could be extended to more general algorithms. More precisely, in chapter IV [P3], we use a continuous-time version of the algorithm [19] to study numerically an observable of ‘type A’ [40] (See Sec. I.5). This version of the algorithm differs from that considered in this chapter, in the sense that the cloning steps are separated by non-fixed non-infinitesimal time intervals. These time intervals are distributed exponentially, in contrast to the fixed ones taken here where Δt is a constant. This results in an important difference: the effective interaction between copies due to the cloning/pruning procedure is unbounded (it can *a priori* affect any proportion of the population), while in the algorithm of the present here, this effective interaction is restricted to a maximum of one cloning/pruning event in the $\Delta t \rightarrow 0$ limit. We stress that the $dt \rightarrow 0$ limit of the cloning algorithm studied in here with a fixed Δt **does not yield the continuous-time cloning algorithm**, stressing that these two versions of the population dynamics present essential differences.

III.3 Stochastic Errors: Large Deviations of the Population Dynamics

In the previous section, we formulated the population dynamics algorithm as a birth-death process and evaluated the systematic errors (which are the deviation of the large deviation estimator from the correct value) due to a finite number of clones (Table III.3). In this section, we focus on **stochastic errors** corresponding to the run-to-run fluctuations of the large deviation estimator within the algorithm, at fixed N_c .

In order to study stochastic errors, we formulate the large deviation principle of the large deviation estimator. In the population dynamics algorithm, the CGF estimator is the time-average of the average cloning ratio of the population (see Sec. III.2.2.1):

$$\psi_{N_c, \tau}(s) \equiv -s \frac{1}{\tau} \int_0^\tau dt \sum_{i=0}^{M-1} \frac{n_i(t) b_i}{N_c}. \quad (\text{III.22})$$

As τ increases, this quantity converges to the expected value (which depends on N_c) with probability 1. However whenever we consider a finite τ , dynamical fluctuations are present, and there is a probability that this estimator deviates from its expected value. Since the population dynamics in the occupation states n is described by a Markov process, the probability of these deviations are themselves described by a large deviation principle (I.20) [1, 83]: By denoting by $\text{Prob}(\psi)$ the probability of $\psi_{N_c, \tau}(s)$, one has:

$$\text{Prob}(\psi) \sim \exp(-\tau I_{N_c, s}(\psi)),$$

where $I_{N_c, s}(\psi)$ is a large deviation ‘rate function’ (of the large deviation estimator). To study these large deviations, we can apply a standard technique using a biased evolution operator

	Magnitude of errors
Systematic errors	$\mathcal{O}(1/N_c)$
Numerical errors	$\mathcal{O}(1/(\tau N_c))$

Table III.3: Magnitudes of the numerical errors

for our population dynamics: For a given Markov system, to calculate large deviations of additive quantities such as Eq. (III.22), one biases the time-evolution matrix with an exponential factor [83]. Specifically, by defining the following matrix

$$L_{\tilde{n},n}^h = \sigma(n \rightarrow \tilde{n}) - \delta_{\tilde{n},n} \sum_{n'} \sigma(n \rightarrow n') - hs \sum_{i=0}^{M-1} \frac{n_i b_i}{N_c}. \quad (\text{III.23})$$

and by denoting the largest eigenvalue of this matrix $G(h, s)$ (corresponding, as a function of h , to a scaled cumulant generating function for the observable (III.22)), the large deviation function $I_{N_c, s}(\psi)$ is obtained as the Legendre transform $\sup_h [h\psi - G(h, s)]$. In chapter IV [P3], we show that a quadratic approximation of the rate function $I_{N_c, s}(\psi)$ (i.e., a Gaussian approximation) can be estimated directly from the cloning algorithm.

We consider the scaling properties of $I_{N_c, s}$ in the large- N_c limit. For this, we define a scaled variable $\tilde{h} \equiv h/N_c$ and a scaled function $\tilde{G}(\tilde{h}, s) \equiv G(\tilde{h}N_c, s)/N_c$. If this scaled function $\tilde{G}(\tilde{h}, s) \equiv G(\tilde{h}N_c, s)/N_c$ is well-defined in the $N_c \rightarrow \infty$ limit (which is natural as checked in the next paragraph), then we can derive that $I_{N_c, s}$ has the following scaling:

$$I_{N_c, s}(\psi) = N_c I_s(\psi) + o(N_c) \quad (\text{III.24})$$

or equivalently,

$$\text{Prob}(\psi) \sim e^{-\tau N_c I_s(\psi)}, \quad (\text{III.25})$$

where $I_s(\psi) = \max_{\tilde{h}} [\tilde{h}\psi - \tilde{G}(\tilde{h}, s)]$. The scaling form (III.24) is validated numerically in chapter IV [P3]. From this large deviation principle, we can see that the stochastic errors of the large deviation estimator is of $\mathcal{O}(1/(N_c\tau))$ as shown in Table III.3.

In the largest eigenvalue problem for the transition matrix (III.23), by performing a system size expansion (see Sec. III.2.3), we obtain

$$\begin{aligned} \tilde{G}(\tilde{h}, s) = & \sum_{i,j(i \neq j)} \left(\frac{\partial}{\partial x_i} - \frac{\partial}{\partial x_j} \right) x_i q(x) \times [w(i \rightarrow j) + s \tilde{w}_x^\infty(i \rightarrow j)] \\ & - \frac{\tilde{h}}{s} \sum_i x_i b_i q(x) + \mathcal{O}(1/N_c), \end{aligned}$$

where $q(x)$ is the right-eigenvector associated to the largest eigenvalue of $L_{\tilde{n},n}^h$ (represented as a function of $x \equiv n/N_c$). The first order of the right-hand side is of order $\mathcal{O}(N_c^0)$, so that $\tilde{G}(\tilde{h}, s)$ is also of order $\mathcal{O}(N_c^0)$ in $N_c \rightarrow \infty$. (For an analytical example of the function $\tilde{G}(\tilde{h}, s)$, see Sec. III.4.2).

III.4 Example: A Simple Two-State Model

In order to illustrate the formulation that we developed in the previous sections, here we consider a simple two state model. In this system, the dimension of the state i is two ($M = 2$) and the transition rates $w(i \rightarrow j)$ are

$$\begin{aligned} w(0 \rightarrow 1) &= c, \\ w(1 \rightarrow 0) &= d \end{aligned}$$

with $c, d > 0$ and $w(i \rightarrow i) = 0$. In this model, the quantity α_i defined in Eq. (III.6) becomes

$$\alpha_i = \delta_{i,0} \frac{n_1}{N_c} (b_1 - b_0) + \delta_{i,1} \frac{n_0}{N_c} (b_0 - b_1).$$

Hereafter, we assume that $b_1 > b_0$ without loss of generality. From this, the space $\Omega^{(\pm)}$ is determined as $\Omega^{(+)} = \{0\}$ and $\Omega^{(-)} = \{1\}$, which leads to the jump ratio $\tilde{w}_n(i \rightarrow j)$ as

$$\tilde{w}_n(i \rightarrow j) = \delta_{i,1} \delta_{j,0} \frac{n_0}{N_c} (b_1 - b_0) \left[\frac{n_1}{N_c + 1} + \frac{n_0}{N_c - 1} \right].$$

Finally, from the conservation of the total population: $n_0 + n_1 = N_c$, we find that the state of the population n can be uniquely determined by specifying only the variable n_0 . Thus the transition rate for the population dynamics is a function of n_0 (and \tilde{n}_0), $\sigma(n_0 \rightarrow \tilde{n}_0)$, which is derived as

$$\begin{aligned} \sigma(n_0 \rightarrow \tilde{n}_0) &= \delta_{\tilde{n}_0, n_0+1} \left[(N_c - n_0)d + k(n_0, N_c - n_0) \right. \\ &\quad \left. \times \left(\frac{n_0}{N_c - 1} + \frac{N_c - n_0}{N_c + 1} \right) \right] + \delta_{\tilde{n}_0, n_0-1} n_0 c, \end{aligned}$$

where we have defined

$$k(n_0, n_1) = \frac{n_0 n_1}{N_c} s [b_1 - b_0].$$

III.4.1 Systematic Errors

In order to evaluate the systematic errors (see Sec. III.2.3), we consider the distribution function P_n^{st} . Since the system is described by a one dimensional variable n_0 restricted to $0 \leq n_0 \leq N_c$, the transition rates $\sigma(n_0 \rightarrow \tilde{n}_0)$ satisfy the detailed balance condition:

$$P_{n_0}^{\text{st}} \sigma(n_0 \rightarrow n_0 + 1) = P_{n_0+1}^{\text{st}} \sigma(n_0 + 1 \rightarrow n_0).$$

We can solve this equation exactly, but to illustrate the large- N_c limit, it is in fact sufficient to study the solution in an expansion $1/N_c \ll 1$. The result is

$$P_{xN_c}^{\text{st}} = C \exp [-N_c I_{\text{conf}}(x) + \delta I(x) + \mathcal{O}(1/N_c)]$$

with $x \equiv n_0/N_c$ and explicitly

$$\begin{aligned} I_{\text{conf}}(x) &= x + \log(1-x) - \frac{d \log [d + (b_1 - b_0)sx]}{(b_1 - b_0)s} \\ &\quad - x \log \left[\frac{1}{cx} (1-x) (d + (b_1 - b_0)sx) \right] \end{aligned}$$

and

$$\begin{aligned} \delta I(x) = & -x - \frac{2dx}{(b_1 - b_0)s} + x^2 - \log x \\ & + \frac{2d^2 \log [d + (b_1 - b_0)sx]}{(b_1 - b_0)^2 s^2} + \frac{d \log [d + (b_1 - b_0)sx]}{(b_1 - b_0)s}. \end{aligned}$$

We now determine the value of x that minimizes $-N_c I_s(x) + \delta I(x)$, which leads to a finite-size correction (i.e., the systematic errors) of the population dynamics estimator. Indeed, denoting this optimal value of x by $x_{N_c}^*$, the large deviation estimator is obtained as

$$\psi_{N_c}(s) = -s [x_{N_c}^* b_0 + (1 - x_{N_c}^*) b_1]$$

(see Sec. III.2.2.1). From a straightforward calculation based on the expressions $I_{\text{conf}}(x)$ and $\delta I(x)$, we obtain the expression of $x_{N_c}^*$ as

$$x_{N_c}^* = x^* + \frac{1}{N_c} \delta x^* + \mathcal{O}((1/N_c)^2),$$

with

$$x^* = \frac{-c - d + (b_1 - b_0)s}{2(b_1 - b_0)s} + \frac{\sqrt{4d(b_1 - b_0)s + [-c - d + (b_1 - b_0)s]^2}}{2(b_1 - b_0)s}$$

and

$$\delta x^* = (2d + 2(b_1 - b_0)sx^*)^{-1} \times \frac{2c [-d - (b_1 - b_0)sx^* (1 + x^* - 2(x^*)^2)]}{\sqrt{4d(b_1 - b_0)s + [c + d - (b_1 - b_0)s]^2}}.$$

We thus arrive at

$$\psi(s) = \frac{-c - d - (b_1 + b_0)s}{2} + \frac{\sqrt{4d(b_1 - b_0)s + [-c - d + (b_1 - b_0)s]^2}}{2} \quad (\text{III.28})$$

and

$$\epsilon_{\text{sys}} = \frac{1}{N_c} \frac{1}{|d + (b_1 - b_0)sx^*|} \times \left| \frac{sc(b_0 - b_1) (d + (b_0 - b_1)s(x^* - 1)x^*(1 + 2x^*))}{\sqrt{4(b_1 - b_0)ds + [c + d + (b_0 - b_1)s]^2}} \right|$$

(see Eq. (III.19) for the definition of the systematic error ϵ_{sys}). We check easily that the expression of $\psi(s)$ is the same as the one obtained from a standard method by solving the largest eigenvalue problem of a biased time-evolution operator (as explained in Sec. I.6.1, and implemented in chapter II [P1]).

III.4.2 Stochastic Errors

We now turn our attention to the stochastic errors. The scaled cumulant generating function $N_c \tilde{G}(\tilde{h}, s)$ is the largest eigenvalue of a matrix $L_{\tilde{n}, n}^h$ (III.23). We then recall a formula to calculate this largest eigenvalue problem from the following variational principle:

$$\begin{aligned} \tilde{G}(\tilde{h}, s) = & \sup_{\phi > 0} \sum_n p_{\text{st}}(n_0) \phi(n_0)^2 \left[\frac{\sigma(n \rightarrow n+1)}{N_c} \left(\frac{\phi(n_0+1)}{\phi(n_0)} - 1 \right) \right. \\ & \left. + \frac{\sigma(n \rightarrow n-1)}{N_c} \left(\frac{\phi(n_0-1)}{\phi(n_0)} - 1 \right) - s\tilde{h} \frac{\sum_i n_i b_i}{N_c^2} \right]. \end{aligned}$$

(See e.g., the appendix G of [256] or [40] for the derivation of this variational principle). By following the usual route to solve such equations (see e.g., the Sec. 2.5 of Ref. [257]), we obtain

$$\tilde{G}(\tilde{h}, s) = \sup_x \left[- \left(\sqrt{(1-x)(d + (b_1 - b_0)sx)} - \sqrt{cx} \right)^2 - s\tilde{h} [xb_0 + (1-x)b_1] \right].$$

Thus, $\tilde{G}(\tilde{h}, s)$ is well-defined, justifying that the large deviation principle (III.25) is satisfied. Furthermore, by expanding this variational principle with respect to \tilde{h} , we obtain

$$\tilde{G}(\tilde{h}, s) = \psi(s)\tilde{h} + \frac{\kappa_s}{2}\tilde{h}^2 + \mathcal{O}(\tilde{h}^3), \quad (\text{III.29})$$

where $\psi(s)$ is given in Eq. (III.28), and the variance κ_s is given as

$$\begin{aligned} \kappa_s = c + & \frac{cs(b_1 - b_0)}{\sqrt{4(b_1 - b_0)sd + (c + d + (b_0 - b_1)s)^2}} \\ & - \frac{c(c + d)^2 + c(b_0 - b_1)(c - 3d)s}{c^2 + 2c[d + (b_0 - b_1)s] + (d + (b_1 - b_0)s)^2}. \end{aligned}$$

We note that the expansion (III.29) is equivalent to the following expansion of the large deviation function $I_s(\psi)$ (III.25) around the expected value $\psi(s)$:

$$I_s(\psi) = \frac{(\psi - \psi(s))^2}{2\kappa_s} + \mathcal{O}((\psi - \psi(s))^3).$$

The variance of the obtained large deviation estimator is thus $\kappa_s/(N_c\tau)$.

III.4.3 A Different Large Deviation Estimator

As an application of these exact expressions, we expand the systematic error ϵ_{sys} and the stochastic error (variance) κ_s with respect to s . A straightforward calculation leads to

$$\epsilon_{\text{sys}}N_c = \left| \frac{2c(b_0 - b_1)}{c + d} s \right| + \mathcal{O}(s^2)$$

and

$$\kappa_s = \frac{2(b_0 - b_1)^2 cd}{(c + d)^3} s^2 + \mathcal{O}(s^3).$$

We thus find that the first-order of the error ϵ_{sys} scales as $\mathcal{O}(s)$ at small s , but that the variance κ_s is of order $\mathcal{O}(s^2)$. From this scaling, as we explain below, one can argue that the following large deviation estimator can be better than the standard one for small s :

$$\tilde{\Psi}(s) \equiv \frac{1}{\tau} \log \overline{\prod_{t=0}^{\tau/dt} \sum_i \frac{n_i(t) e^{-sdtb_i}}{N_c}}, \quad (\text{III.30})$$

where the overline represents the averaging with respect to the realizations of the algorithm. Normally, this realization-average is taken *after* calculating the logarithm, which corresponds

to the estimator (III.10). Mathematically, this average (Eq.(III.30), before taking the logarithm) corresponds to a bias of the time-evolution matrix σ as seen in Eq. (III.23) for $h = 1$. This means that, in the limit $\tau \rightarrow \infty$ with a sufficiently large number of realizations, this averaged value behaves as $\tilde{\Psi}(s) \sim e^{\tau G(1,s)}$. By combining this result with the expansion (III.29), we thus obtain

$$\lim_{\tau \rightarrow \infty} \lim_{\substack{\text{many} \\ \text{realizations}}} \tilde{\Psi}(s) = \psi(s) + \frac{\kappa_s}{2} N_c^{-1} + \mathcal{O}(N_c^{-2}) \quad (\text{III.31})$$

(recalling $\tilde{G} = G/N_c$ and $\tilde{h} = h/N_c$). When we consider small s , by recalling $\epsilon_{\text{sys}} N_c = \mathcal{O}(s)$ and $\kappa_s = \mathcal{O}(s^2)$, we thus find that the deviations from the correct value are smaller in the estimator $\tilde{\Psi}(s)$ than in the normal estimator given in Eq. (III.10), which comes as a surprise because in Eq. (III.30) the average and the logarithm are inverted with respect to a natural definition of the CGF estimator.

To use this estimator, we need to discuss the two following points. First, since the scaled cumulant generating function $G(1, s)$ has small fluctuations, one needs a very large number of realizations in order to attain the equality (III.31). The difficulty of this measurement is the same level as the one of direct observations of a large deviation function, see for example Ref. [255]. However, we stress that this point may not be fatal in this estimator, because we do not need to attain completely this equality, i.e., our aim is the zero-th order coefficient, $\psi(s)$, in Eq. (III.31). Second, we have not proved yet the scaling properties with respect to s , which are $\epsilon_{\text{sys}} N_c = \mathcal{O}(s)$ and $\kappa_s = \mathcal{O}(s^2)$, in a general set-up aside from this simple two state model. We show in practice in the next chapter that for small values of s , the estimator (III.30) is affected by smaller systematic errors, in the numerical study of the creation-annihilation process studied in this section. This alternative way of defining the CGF estimator is studied again for the continuous-time version of the algorithm in Sec V.3.

III.5 Discussion

In this chapter, we formulated a birth-death process that describes population dynamics algorithms and evaluated numerically large deviation functions. We showed that this birth-death process leads generically to the correct large deviation results in the large limit of the number of clones $N_c \rightarrow \infty$. We also showed that the deviation of large deviation estimator from the desired value (which we called systematic errors) is small and proportional with $\mathcal{O}(N_c^{-1})$. In the next chapter, we verify and use the $1/\tau$ - and $1/N_c$ -scalings of the CGF estimator in order to interpolate its large- τ and large- N_c asymptotic value from the measured values for finite τ and N_c . We demonstrate numerically that the interpolation technique is very efficient, by a direct comparison of the resulting CGF estimation to its analytical value, which can be determined in the studied system. We also underline that this is done for a different version of the algorithm, a continuous in time population dynamics [19]. For a description of their conceptual difference refer to Secs. III.2.4.2 and IV.5.

IV – Finite-Time and Finite-Size Scalings in the Evaluation of Large-Deviation Functions: II. Numerical Approach in Continuous Time

IV.1 Introduction

In chapter III [P2], we performed an analytical study of a discrete-time version of the population dynamics algorithm. We derived the finite- N_c and finite- t scalings of the systematic errors of the LDF estimator, showing that these behave as $1/N_c$ and $1/t$ in the large- N_c and large- t asymptotics respectively. In principle, knowing the scaling *a priori* means that the asymptotic limit of the estimator in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits may be interpolated from the data at finite t and N_c . However, whether this idea is actually useful or not is a non-trivial question, as there is always a possibility that onset values of N_c - and t -scalings are too large to use these scalings. In the present chapter, we consider a continuous-time version of the population dynamics algorithms [17, 19]. We show numerically that one can indeed make use of these scaling properties in order to improve the estimation of CGF, in an application to a system with many-body interactions (a contact process). We emphasize that the two versions of the algorithm differ on a crucial point which makes that an extension of the analysis developed in chapter III [P2] cannot be done straightforwardly in order to comprehend the continuous-time case (see section IV.5). We thus stress that the observation of these scalings themselves is also non-trivial.

This chapter IV [P3] is organized as follows. In Sec. IV.3.1 we study the behavior of the CGF estimator as a function of the duration of the observation time (for a fixed population N_c) and we see how its infinite-time limit can be extracted for the numerical data. In Sec. IV.3.2 we analyze the behavior of the estimator as we increase the number of clones (for a given final simulation time) and the infinite-size limit of the LDF estimator. Based on these results, we present in Sec. IV.4 a method which allows us to extract the infinite-time, infinite-size limit of the large deviation function estimator from a finite-time, finite-size scaling analysis. In Sec. IV.5, we discuss the difficulty of an analytical approach to the continuous-time algorithm. Finally, our conclusions are made in Sec. IV.6.

IV.2 CGF Estimator: Constant-Population Approach

In practice, in order to obtain a good estimation of the CGF, it is normal to launch the simulation several times (where we denote by R the number of realizations of the same simulation), and to estimate the arithmetic mean of the obtained values of Eq. (I.30) over these R simulations. Strictly speaking (as discussed in Sec. II.2.2), as the simulation does not stop exactly at the final simulation time T but at some time $t_r^{\mathcal{F}} \leq T$ (which is different for every $r \in \{1, \dots, R\}$), the average over R realizations of $\overline{\Psi_s^{(N_c)}}$ is then correctly defined as

$$\overline{\Psi_s^{(N_c)}} = \frac{1}{R} \sum_{r=1}^R \frac{1}{t_r^{\mathcal{F}}} \log \prod_{i=1}^{\mathcal{K}_r} X_i^r. \quad (\text{IV.1})$$

However, we have observed that for not too short simulation times, $|\overline{\Psi_s^{(N_c)}}(T) - \overline{\Psi_s^{(N_c)}}(t_r^{\mathcal{F}})|$ is small. By assuming $t_r^{\mathcal{F}} \approx T$, Eq. (IV.1) can be approximated by replacing $t_r^{\mathcal{F}}$ by T (which is what we do in practice)

$$\overline{\Psi_s^{(N_c)}} \simeq \frac{1}{R} \frac{1}{T} \sum_{r=1}^R \log \prod_{i=1}^{\mathcal{K}_r} X_i^r. \quad (\text{IV.2})$$

The CGF estimator can be defined differently from Eq. (IV.1) by using an alternative way of computing the average over R realizations (as in Sec. III.4.3 in discrete-time and as in Sec. V.3 in continuous-time). Equations (IV.1) and (IV.2) allow us to estimate the CGF using the constant-population approach of the continuous-time cloning algorithm for a s -biased Markov process, given a fixed number of clones N_c , a simulation time T and R realizations of the algorithm.

IV.3 Finite-Time and Finite- N_c Behavior of CGF Estimator

In this section, we focus on the annihilation-creation process (Sec. I.8.1) for a particular value of parameter s ($s = -0.2$), which is representative of the full range of s on which we study large deviations.

IV.3.1 Finite-Time Scaling

Here, we study the large-time behavior of the CGF estimator, at fixed number of clones N_c . Fig. IV.1 presents the average over $R = 10^4$ realizations of the CGF estimator $\overline{\Psi_s^{(N_c)}}$ (IV.1) as a function of the (simulation) time for given numbers of clones $N_c = \{10, 100, 1000\}$. It is compared with the analytical value $\psi(s)$ (I.33) which is shown with a black dashed line. As can be seen in Fig. IV.1 for a small number of clones ($N_c = 10$), the CGF estimator $\overline{\Psi_s^{(N_c)}}$ highly deviates from the analytical value $\psi(s)$. However, as N_c and the simulation time t become larger, the CGF estimator gets closer to the analytical value $\psi(s)$. One can expect that in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits, $\psi(s)$ will be obtained from the estimator as

$$\lim_{N_c \rightarrow \infty} \lim_{t \rightarrow \infty} \overline{\Psi_s^{(N_c)}}(t) = \psi(s),$$

as it was derived in chapter III [P2]. However, in a practical implementation of the algorithm, this infinite-time and -size limits are not achievable and we use large but *finite* simulation

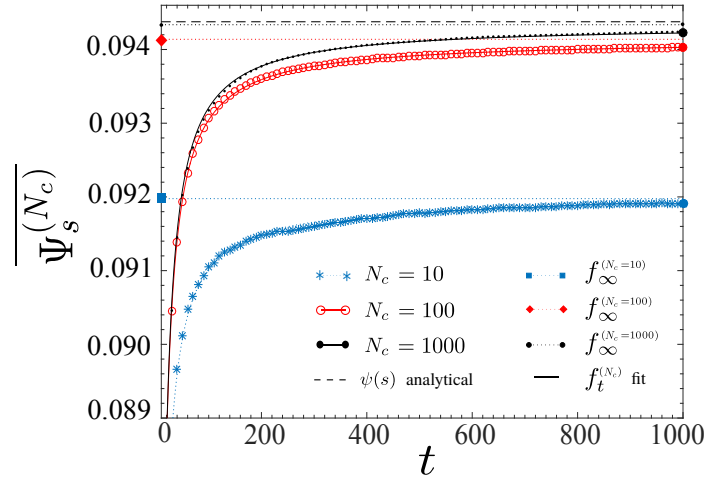


Figure IV.1: Average over $R = 10^4$ realizations of the CGF estimator $\overline{\Psi_s^{(N_c)}}$ (IV.1) as a function of duration t of the observation window, for $N_c \in \{10, 100, 1000\}$ clones, for the annihilation-creation dynamics with $c = 0.3$. The analytical expression for the large deviation function $\psi(s)$ (I.33) is shown with a black dashed line and the fitting functions $f_t^{(N_c)}$ encoding the finite- t scaling (Eq. (IV.3)) are shown with continuous curves. The (*a priori*) best estimation of the large deviation function (to which we refer as standard estimator) is given by $\overline{\Psi_s^{(N_c)}}(T)$ at the largest simulation time $T = 1000$, which are shown with solid circles (at the right end of the figure). The extracted infinite-time limits $f_\infty^{(N_c)}$ are shown as dotted lines and squares ($N_c = 10$), diamonds ($N_c = 100$) and circles ($N_c = 1000$).

time t and number of clones N_c . This fact motivates our analysis of the actual dependence of the estimator with t and N_c . The standard estimator of the large deviation function is the value of $\overline{\Psi_s^{(N_c)}}$ at the largest simulation time T and for the largest number of clones N_c , ($\overline{\Psi_s^{(N_c)}}(T)$ for $N_c = 1000$ and $T = 1000$), i.e., the black solid circle \bullet in Fig. IV.1. This value provides the (*a priori*) best estimation of the large deviation function that we can obtain from the continuous-time cloning algorithm. However encouragingly, as we detail later, this estimation can be improved by taking into account the convergence speed of the CGF estimator.

The result of fitting $\overline{\Psi_s^{(N_c)}}(t)$ with the curve $f_t^{(N_c)}$ is shown with solid lines in Fig. IV.1. This is defined as

$$f_t^{(N_c)} \equiv f_\infty^{(N_c)} + b_t^{(N_c)} t^{-1}, \quad (\text{IV.3})$$

where the fitting parameters $f_\infty^{(N_c)}$ and $b_t^{(N_c)}$ can be determined from the least squares method by minimizing the deviation from $\overline{\Psi_s^{(N_c)}}(t)$. The clear coincidence between $\overline{\Psi_s^{(N_c)}}(t)$ and the fitting lines indicates the existence of a $1/t$ -convergence of $\overline{\Psi_s^{(N_c)}}(t)$ to $\lim_{t \rightarrow \infty} \overline{\Psi_s^{(N_c)}}(t)$ (that we call **t^{-1} -scaling**). This property can be derived from the assumption that the cloning algorithm itself is described by a Markov process: in chapter III [P2] with a different version of the algorithm, we constructed a meta-Markov process to describe the cloning algorithm by expressing the number of clones by a birth-death process. Once such meta process is

constructed, the CGF estimator (I.30) is regarded as the time-average of the observable X_i within such meta-Markov process. In other words, $t\Psi_s^{(N_c)}$ is an additive observable of the meta-process describing the cloning algorithm. We now recall that time-averaged quantities converge to their infinite-time limit with an error proportional to $1/t$ when the distribution function of the variable converges exponentially (as in Markov processes). This leads to the t^{-1} -scaling of CGF estimator (IV.3). We note that constructing such a meta-Markov process explicitly is not a trivial task, and for the algorithm discussed here, such a construction remains as an open problem.

By assuming the validity of the scaling form (IV.3), it is possible to extract the infinite-time limit of the CGF estimator from finite-time simulations. We denote this infinite-time limit as $f_\infty^{(N_c)}$ and it is expected to be a better estimator of CGF than $\overline{\Psi_s^{(N_c)}(T)}$ at finite T , provided that

$$f_\infty^{(N_c)} = \lim_{t \rightarrow \infty} \overline{\Psi_s^{(N_c)}(t)}.$$

In Fig. IV.1, we show $f_\infty^{(N_c)}$ with dotted lines and circles ($N_c = 10$), diamonds ($N_c = 100$) and squares ($N_c = 1000$). As can be seen, this parameter indeed provides a better numerical estimate of $\psi(s)$ than $\overline{\Psi_s^{(N_c)}(T)}$.

IV.3.2 Finite- N_c Scaling

Here, we study the behavior of the (standard) CGF estimator $\overline{\Psi_s^{(N_c)}(T)}$ as we increase the number of clones N_c , for a given final (simulation) time T . Similar to what we did in Sec. IV.3.1, we consider a curve in the form

$$g_{N_c}^{(T)} = g_\infty^{(T)} + \tilde{b}_{N_c}^{(T)} N_c^{-1}, \quad (\text{IV.4})$$

where $g_\infty^{(T)}$ and $\tilde{b}_{N_c}^{(T)}$ are fitting parameters which are determined by the least squares fitting to $\overline{\Psi_s^{(N_c)}(T)}$. The obtained $g_{N_c}^{(T)}$ as a function of N_c are shown in Fig. IV.2 as solid lines. We considered four values of final simulation time $T = \{200, 300, 500, 1000\}$ and population sizes in the range $10 \leq N_c \leq 1000$. As can be seen, these curves describe well the dependence in N_c of $\overline{\Psi_s^{(N_c)}(T)}$, indicating that $\overline{\Psi_s^{(N_c)}(T)}$ converges to its infinite- N_c limit with an error proportional to $1/N_c$ (that we call **N_c^{-1} -scaling**). This scaling could be proved under general assumptions in chapter III [P2], (i) however without covering the continuous-time algorithm discussed here, and (ii) for the CGF estimator $\overline{\Psi_s^{(N_c)}(T)}$ considered the $T \rightarrow \infty$ limit, instead of finite T . The generalization of the argument presented in chapter III [P2] in order to cover the general cases (i) and (ii) is an important open direction of research.

By assuming the validity of such N_c^{-1} -scaling, we can evaluate the $N_c \rightarrow \infty$ limit of $\overline{\Psi_s^{(N_c)}(T)}$ as the fitting parameter $g_\infty^{(T)}$ obtained from finite N_c simulations as

$$g_\infty^{(T)} = \lim_{N_c \rightarrow \infty} \overline{\Psi_s^{(N_c)}(T)}.$$

These parameters $g_\infty^{(T)}$ (to which we refer as infinite-size limit) are shown in Fig. IV.2 as dotted lines and provide better estimations of $\psi(s)$ than the standard estimator $\overline{\Psi_s^{(N_c)}(T)}$.

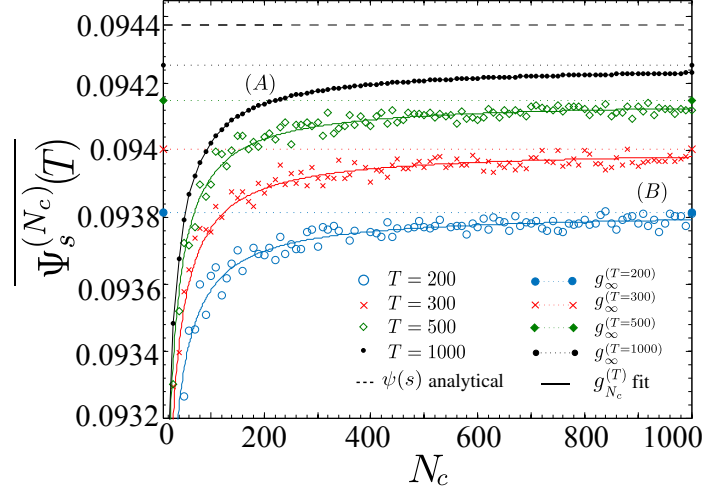


Figure IV.2: CGF estimator $\overline{\Psi_s^{(N_c)}(T)}$ (IV.1) for given final (simulation) times $T = \{200, 300, 500, 1000\}$ as a function of the number of clones N_c (on the range $10 \leq N_c \leq 1000$). The analytical value $\psi(s)$ (I.33) is shown with a dashed line and the fits $g_{N_c}^{(T)}$ (IV.4) with continuous curves. A large simulation time for a small number of clones, shown in (A), produces a better estimation compared to the one given by the largest number of clones with a relatively short simulation time, which is shown in (B). The best CGF estimation we can naively obtain would be given by $\overline{\Psi_s^{(N_c)}(T)}$ at largest simulation time T and largest number of clones N_c . However, the extracted infinite-size limits $g_{\infty}^{(T)}$ provide a better estimation in comparison. These limits are shown with dotted lines and circles ($T = 200$), crosses ($T = 300$), diamonds ($T = 500$) and dots ($T = 1000$). Additionally, $c = 0.3$ and $s = -0.2$.

IV.4 Finite-Time and Finite- N_c Scaling Method to estimate Large Deviation Functions

In the previous section, we have shown how it is possible to extract $f_{\infty}^{(N_c)}$ and $g_{\infty}^{(T)}$ from finite T - and finite N_c - simulations respectively. In this section, we combine both of these $1/t$ - and $1/N_c$ - scaling methods in order to extract the infinite-time and -size limit of the CGF estimator. This limit gives a better evaluation of the large deviation function within the cloning algorithm than the standard estimator.

We first note that either of $f_{\infty}^{(N_c)}$ or $g_{\infty}^{(T)}$ is expected to converge to $\psi(s)$ as $N_c \rightarrow \infty$ or as $T \rightarrow \infty$. We checked numerically this property by defining the distance D between $\psi(s)$ and its numerical estimator $\overline{\Psi_s^{(N_c)}}$,

$$D(\overline{\Psi_s^{(N_c)}}, \psi(s)) = |\overline{\Psi_s^{(N_c)}} - \psi(s)|. \quad (\text{IV.5})$$

This quantity is shown in Fig. IV.3 as a function of t in log-log scale. As we can see, as N_c increases, $\log D$ behaves as straight line with slope -1 on a time window which grows with N_c . In other words, when $N_c \rightarrow \infty$,

$$|\overline{\Psi_s^{(N_c)}} - \psi(s)| \sim t^{-1}. \quad (\text{IV.6})$$

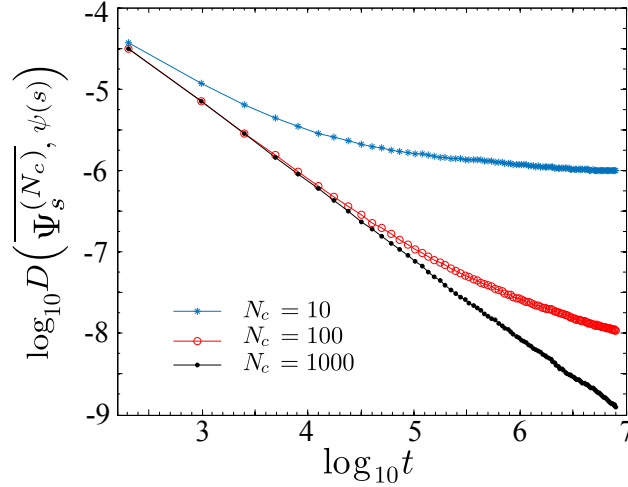


Figure IV.3: Distance D (Eq. (IV.5)) between the analytical CGF $\psi(s)$ and its numerical estimator $\overline{\Psi_s^{(N_c)}}$, as a function of time t in log-log scale. The distances are computed from the values in Fig. IV.1. This distance behaves as a power law of exponent -1 on a time window, where the size of the time window increases as N_c increases. This illustrates the scaling (IV.6). The parameters of the model are $c = 0.3$, $s = -0.2$.

Inspired by this observation, we assume the following scaling for the fitting parameter $f_\infty^{N_c}$. If we consider a set of simulations performed at population sizes $\vec{N}_c = \{N_c^{(1)}, \dots, N_c^{(j)}\}$, the obtained infinite-time limit of the CGF estimator $f_\infty^{N_c}$ behaves as a function of N_c as

$$f_\infty^{(N_c)} \simeq f_\infty^\infty + b_\infty^{(N_c)} N_c^{-1}, \quad (\text{IV.7})$$

which means that $f_\infty^{(N_c)}$ itself exhibits $1/N_c$ corrections for large but finite N_c . By using this scaling, we detail below in Sec. IV.4.1 the method to extract the infinite-time infinite- N_c limit of the CGF estimator $\overline{\Psi_s^{(N_c)}}(T)$ from finite-time and finite- N_c data. We note that this method can be used for a relatively short simulation time and a relatively small number of clones (see Fig. IV.5). In Sec. IV.4.2, we present numerical examples of the application of this method to the contact process.

IV.4.1 The Scaling Method

The procedure is summarized as follows:

1. Determine the average over R realizations $\overline{\Psi_s^{(N_c)}}(t)$ (IV.1) up to a final simulation time T for each $N_c \in \vec{N}_c$.
2. Determine the fitting parameters $f_\infty^{(N_c)}$'s defined in the form $f_t^{(N_c)} = f_\infty^{(N_c)} + b_t^{(N_c)} t^{-1}$ (IV.3) from each of the obtained $\overline{\Psi_s^{(N_c)}}(t)$'s.
3. Determine f_∞^∞ from a fit in size $f_\infty^{(N_c)} = f_\infty^\infty + b_\infty^{(N_c)} N_c^{-1}$ (IV.7) on the extracted $f_\infty^{(N_c)}$'s.

The result obtained for f_∞^∞ renders a better estimation of $\psi(s)$ than the standard estimator $\overline{\Psi_s^{(N_c)}}(t)$ evaluated for $N_c = \max \vec{N}_c$ and for $t = T$.

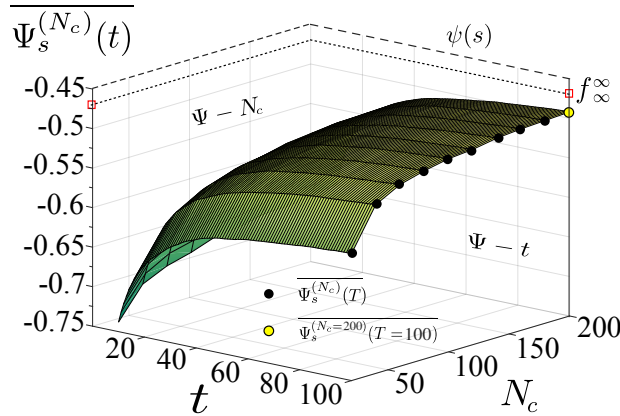


Figure IV.4: Estimator of the large deviation function $\overline{\Psi_s^{(N_c)}}(t)$ as a function of time and the number of clones. The estimator $\overline{\Psi_s^{(N_c)}}(T)$ at final simulation time $T = 100$ as a function of the number of clones (up to $N_c = 200$) is shown as black circles. The best CGF estimation under this configuration given by the standard estimator, i.e., $\overline{\Psi_s^{(N_c=200)}}(T = 100)$ is shown as a yellow circle. The analytical value of the CGF $\psi(s)$ is obtained from the largest eigenvalue of the matrix (I.24) and shown as a black dashed line. The extracted limit f_∞^∞ is shown with red squares. Additionally, $L = 6$, $s = 0.15$, $h = 0.1$, $\lambda = 1.75$ and $R = 10^3$.

IV.4.2 Application to the Contact Process

We apply the scaling method to the one-dimensional contact process (see Sec. I.8.2). We set $L = 6$, $h = 0.1$, $\lambda = 1.75$, $T = 100$ and $s = 0.15$. As we detail below, we compare the improved estimator f_∞^∞ obtained from the application of the scaling method (for $\vec{N}_c = \{20, 40, \dots, 180, 200\}$) with the standard estimator $\overline{\Psi_s^{(N_c)}}(T)$ (for $N_c = \max \vec{N}_c = 200$). Fig. IV.4 represents the behavior of the estimator $\overline{\Psi_s^{(N_c)}}(t)$ as a function of the simulation time t and of the number of clones N_c . The values of the estimator at the final simulation time T are represented with black circles for each $N_c \in \vec{N}_c$ and with a yellow circle for $N_c = \max \vec{N}_c$. The analytical expression for the large deviation function $\psi(s)$ is shown in a black dashed line.

On Fig. IV.5(a) we show the projection of the surface of Fig. IV.4 on the plane $\Psi - t$. The behavior in t of the estimator $\overline{\Psi_s^{(N_c)}}(t)$ is shown for $N_c = 20$ and $N_c = 200$, in blue dots in Fig. IV.5(a). The standard CGF estimators, $\overline{\Psi_s^{(N_c)}}(T)$, are shown in large blue dots in Fig. IV.5(a) (on the axis for $T = 100$). The fitting curves $f_t^{(N_c)}$ (Eq. (IV.3)) are shown in black continuous lines (for $N_c = 20$ and $N_c = 200$) and black dotted lines (for other intermediate values of N_c). Next, we show in Fig. IV.5(b) the projection of the surface of Fig. IV.4 on the plane $\Psi - N_c$ where the time has been set to the largest $t = T$. The standard CGF estimators, $\overline{\Psi_s^{(N_c)}}(T)$ are plotted as blue filled circles, and the fitting curve $g_{N_c}^{(T)}$ (Eq. (IV.4)) on $\overline{\Psi_s^{(N_c)}}(T)$ is shown as a blue solid line. From these curves, we determine $g_\infty^{(T)}$ (see Sec. IV.3.2), which is shown as a blue dashed line and diamonds. Finally, the parameter

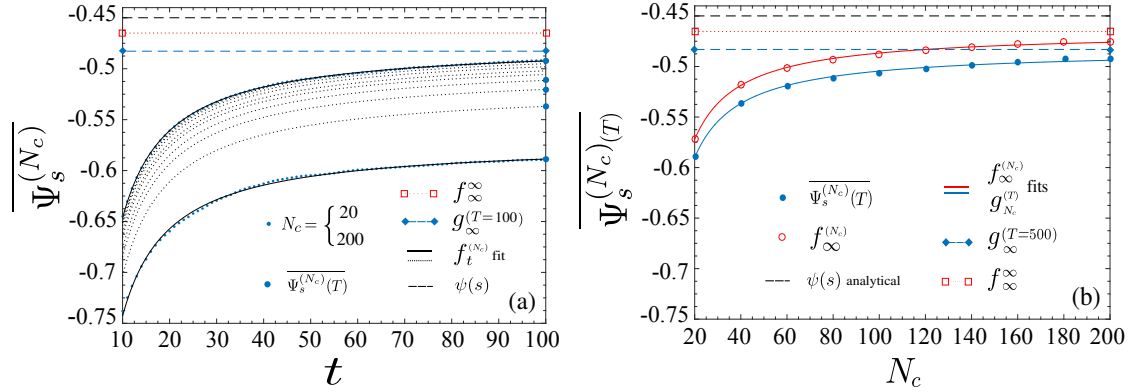


Figure IV.5: (a) Projection of the surface represented in Fig. IV.4 over the plane $\Psi - t$. $\overline{\Psi}_s^{(N_c)}(t)$ is represented for $N_c = 20$ and $N_c = 200$ with blue dots. The estimations $\overline{\Psi}_s^{(N_c)}(T)$ of the large deviation (at the final simulation time $T = 100$) are shown in large blue dots for all the values of N_c considered. The fit in time (Eq. (IV.3)) over $\overline{\Psi}_s^{(N_c)}(t)$ is shown as black solid lines (for $N_c = 20$ and $N_c = 200$) and dotted lines (for other values of N_c). (b) Projection at the final simulation time $T = 100$ on the plane $\Psi - N_c$, $\overline{\Psi}_s^{(N_c)}(T)$ is shown in large blue dots. The infinite-time limit $f_\infty^{(N_c)}$ as a function of N_c (see Eq. (IV.3)) is represented in red circles. The results of fitting $\overline{\Psi}_s^{(N_c)}(T)$ (Eq. (IV.4)) and $f_\infty^{(N_c)}$ (Eq. (IV.7)) are shown with blue and red solid curves respectively. The infinite- N_c limit $g_\infty^{(T)}$ is shown with blue dashed line and diamonds meanwhile the infinite-size and time limit f_∞ is shown with a red dotted line in both of (a) and (b). The extracted limit f_∞ renders a better estimation of the large deviation function than $\overline{\Psi}_s^{(N_c=200)}(T = 100)$ (and also than $g_\infty^{(T)}$) demonstrating the efficacy of the method proposed.

$f_\infty^{(N_c)}$ extracted from the fitting on $\overline{\Psi}_s^{(N_c)}(t)$ (for each value of N_c) is shown as red circles in Fig. IV.5(b). These values also scale as $1/N_c$ (Eq. (IV.7)) and their fit is shown as a red solid curve. The scaling parameter f_∞ obtained from this last step provides a better estimation of the large deviation function than the standard estimator $\overline{\Psi}_s^{(N_c=200)}(T = 100)$ that is widely used in the application of cloning algorithms. This improvement is valid on a wide range of values of the parameter s as can be visualized in Fig IV.6, where we represented the relative systematic error $[\Psi(s) - \psi(s)]/\psi(s)$ between the standard and improved estimators $\Psi(s)$ and the analytical CGF $\psi(s)$.

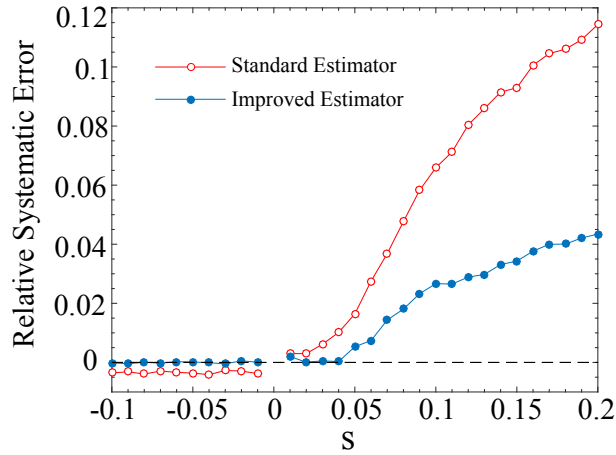


Figure IV.6: Relative systematic error $[\Psi(s) - \psi(s)]/\psi(s)$ between the numerical estimators $\Psi(s)$ and the analytical CGF $\psi(s)$. The error for the standard estimator $\overline{\Psi}_s^{(N_c)}(T)$ is shown in blue and for the improved one, f_∞^∞ (Eq. (IV.7)) in red. The scaling method proposed in this chapter was tested on the contact process (with $L = 6$, $h = 0.1$, and $\lambda = 1.75$) for a set of populations $\vec{N}_c = \{20, \dots, 200\}$, a simulation time $T = 100$, and $R = 1000$ realizations. As can be seen, the errors due to finite-size and -time effects can be reduced through the improved estimator.

IV.5 Issues on an Analytical Approach

In chapter III [P2], we considered a *discrete-time* version of the population dynamics algorithm, where a cloning procedure is performed every small time interval Δt . We have proved the convergence of the algorithm in the large- N_c , $-t$ limits, and we also derived that the systematic error of the LDF estimator (i.e., the deviation of the estimator from the desired LDF) decayed proportionally to $1/t$ and $1/N_c$. From a practical point of view, however, the formulation used there had one problem. In order to prove the result, we took the large frequency limit of cloning procedure or, in other words, we took the $\Delta t \rightarrow 0$ limit. A rough estimate of the error due to non-infinitesimal Δt proves to be $O(\Delta t)$. For a faster algorithm, it is better to take this value to be larger, and indeed empirically, we expect that this error to be very small (or rather disappearing in the large t, N_c limits). However, the detailed analytical estimation of this error is still an open problem.

In the main part of this chapter, from a different point of view, we consider the *continuous-time* version of the population dynamics algorithm [17, 19]. Here, the cloning is performed at each change of state of a copy. The time intervals Δt which separate those changes of state are non-infinitesimal, which means that the formulation used in chapter III [P2] cannot be applied to understand its convergence. Furthermore, because these time intervals are of non-constant duration and stochastically distributed, the continuous-time algorithm is more difficult to handle analytically than the discrete-time version. Instead of pursuing the analytical study within the continuous-time algorithm, we performed a numerical study, and we have shown that the $1/N_c$ and $1/t$ scalings are also observed for the continuous-time

algorithm. Although the proof of these scalings are beyond the scope of this chapter, these numerical observations support a conjecture that such scaling in large t and in large N_c limits are generally valid in cloning algorithms to calculate large deviation functions.

IV.6 Conclusions

Direct sampling of the distribution of rare trajectories is a rather difficult numerical issue (see for instance Ref. [255]) because of the scarcity of the non-typical trajectories. We have shown how to increase the efficiency of a commonly used numerical method (the so-called cloning algorithm) in order to improve the evaluation of large deviation functions which quantify the distribution of such rare trajectories, in the large time limit. We used the finite-size and finite-time scaling behavior of CGF estimators in order to propose an improved version of the continuous-time cloning algorithm which provides more reliable results, less affected by finite-time and $-N_c$ effects. We verified the results observed for the discrete-time version of the cloning algorithm in chapter III [P2] and we showed their validity also for the continuous case [P3]. Importantly, we showed how these results can be applied to more complex systems.

We note that the scalings which rule the convergence to the infinite-size infinite-time limits (with corrections in $1/N_c$ and in $1/t$) have to be taken into account properly: indeed, as power laws, they present no characteristic size and time above which the corrections would be negligible. The situation is very similar to the study of the critical depinning force in driven random manifolds: the critical force presents a corrections in one over the system size [92] which has to be considered properly in order to extract its actual value. Generically, such scalings also provide a convergence criterion to the asymptotic regimes of the algorithm: one has to confirm that the CGF estimator does present corrections (first) in $1/t$ and (second) in $1/N_c$ with respect to an asymptotic value in order to ensure that such value does represent a correct evaluation of the CGF.

It would be interesting to extend our study of these scalings to systems presenting dynamical phase transitions (in the form of a non-analyticity of the CGF), where it is known that the finite-time and the finite-size scalings of the CGF estimator can be very hard to overcome [19]. In particular, in this context, it would be useful to understand how the dynamical phase transition of the original system translates into anomalous features of the distribution of the CGF estimator in the cloning algorithm. These phase transitions are normally accompanied with an infinite system-size limit (although there was a report of dynamical phase transitions without taking a such limit [258]). To overcome these difficulties (caused by a large system size and/or by the presence of a phase transition), it may be useful to use the adaptive version of the cloning algorithm [259], which has been recently developed to study such phase transitions, with the scaling method presented in this chapter.

V – Fluctuations of CGF Estimator

In order to complement the main discussion done in the previous chapter, here we study the fluctuations of the CGF estimator [P3] as defined in Eq. (IV.1). This is done by studying its distribution and its dependence with the simulation time and the number of clones. Compatible with the central limit theorem, we show how a proper rescaling of the CGF estimator produces a collapse of the distributions into a normal standard distribution for different values of N_c and simulation times. Additionally, we discuss in Sec. V.3 an alternative way of defining it which was already introduced in Sec. III.4.3 for the discrete-time version.

V.1 Central Limit Theorem

From relation (IV.1), one can infer that the dispersion of the distribution of $\Psi_s^{(N_c)}$ depends on the simulation time t . This determines whether or not a large number of realizations R is required in order to minimize the statistical error. In fact, as seen in Fig. V.1, the dispersion of $\Psi_s^{(N_c)}$ is concentrated around its mean value, which approaches the analytical value $\psi(s)$ as the simulation time and the number of clones increase. We numerically confirm that these distributions are well-approximated by a Gaussian distribution

$$P\left(\Psi_s^{(N_c)}\right) \sim A e^{-\frac{1}{C^2}\left(\Psi_s^{(N_c)}-B\right)^2} \quad (\text{V.1})$$

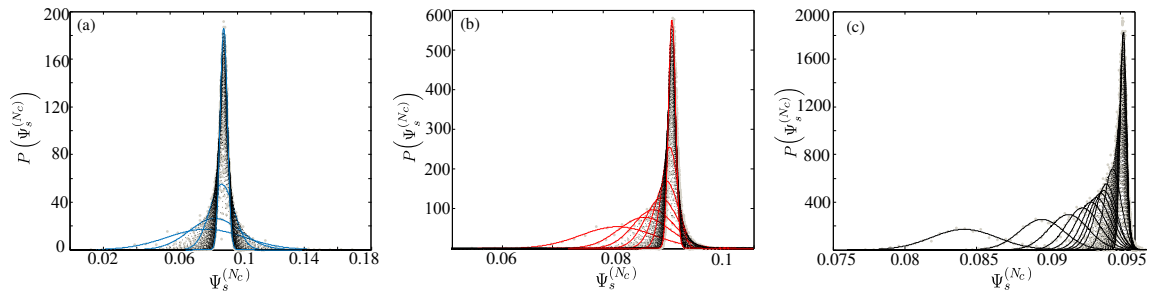


Figure V.1: Distribution $P\left(\Psi_s^{(N_c)}\right)$ of the CGF estimator $\Psi_s^{(N_c)}$ for (a) $N_c = 10$, (b) $N_c = 100$ and (c) $N_c = 1000$ and for simulation times $t \in [10, 1000]$. Each realization ($R = 10^4$ for each simulation time) is shown with gray dots meanwhile its respective Gaussian fit (Eq. (V.1)) is shown with a dotted or a continuous curve. The dispersion of $\Psi_s^{(N_c)}$ is wider for shorter simulation times and small N_c . The mean value of the distribution converges to the theoretical value as the simulation time and the number of clones increase.

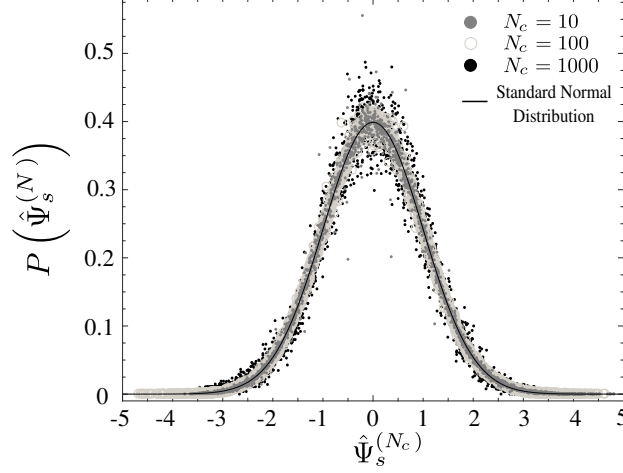


Figure V.2: The distribution function of the rescaled variable $\hat{\Psi}_s^{(N_c)}$ (Eq. (V.2)). Compatible with the central limit theorem, a collapse of the distribution function into a standard normal distribution for different number of clones is observed.

where the parameter B is equal to $\overline{\Psi_s^{(N_c)}}(T)$ and the parameters A and $1/C^2$ are respectively of the order of $N_c^{1/2}$ and N_c . A mathematical argument to explain this obtained Gaussian distribution is given as follows: At any given time (not necessarily at T), let us perform the following rescaling

$$\hat{\Psi}_s^{(N_c)} = \frac{\Psi_s^{(N_c)} - \overline{\Psi_s^{(N_c)}}}{\sigma_{\Psi_s^{(N_c)}}}, \quad (\text{V.2})$$

where

$$\sigma_{\Psi_s^{(N_c)}}^2 = \frac{1}{R-1} \sum_{r=1}^R \left| (\Psi_s^{(N_c)})_r - \overline{\Psi_s^{(N_c)}} \right|^2$$

is the variance of the R realizations of $\Psi_s^{(N_c)}$. Then, this rescaling produces a collapse of the distributions $P(\hat{\Psi}_s^{(N_c)})$, for any t and any N_c (Fig. V.2). We remark then that the CGF estimator (IV.1) is an additive observable of the history of the population, which follows a Markov dynamics. Hence, the rescaled estimator $\hat{\Psi}_s^{(N_c)}$ follows a standard normal distribution in the large time limit, according to the central limit theorem (CLT):

$$P(\hat{\Psi}_s^{(N_c)}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\hat{\Psi}_s^{(N_c)})^2}.$$

The verification of the CLT allows us to ensure if the steady-state of the population dynamics has been reached. Note that in general the typical convergence time to the steady state is larger than the inverse of the spectral gap of the biased evolution operator as discussed in Secs. I.8.1 and II.4.1.

By considering the scaling (V.2) we focus only on the small fluctuations of $\Psi_s^{(N_c)}$ around $\overline{\Psi_s^{(N_c)}}$. But in general, the distribution function is not Gaussian, and in that case we need to consider a large deviation principle as below.

V.2 Logarithmic Distribution of CGF Estimator

Since $\Psi_s^{(N_c)}$ is itself an additive observable of the dynamics of the ensemble of clones (chapter III [P2]), the distribution of the CGF estimator $\Psi_s^{(N_c)}$ satisfies itself a large deviation principle

$$P(\Psi_s^{(N_c)}) \sim e^{-t I_{N_c}(\Psi_s^{(N_c)})}, \quad (\text{V.3})$$

where $I_{N_c}(\Psi_s^{(N_c)})$ is the rate function. This rate function could be evaluated in principle from the empirical distribution $P(\Psi_s^{(N_c)})$ as

$$I_{N_c}(\Psi_s^{(N_c)}) \approx -\frac{1}{t} \log P(\Psi_s^{(N_c)})$$

for a large t . Here we try to estimate the rate function from this equation. The numerical estimation of the right-hand side of the last expression at final simulation time T is shown in Fig. V.3(a), where we have defined

$$\hat{I}_{N_c}(\Psi_s^{(N_c)}) \equiv -\frac{1}{t} \log P(\Psi_s^{(N_c)}) + \frac{1}{t} \log P(\overline{\Psi_s^{(N_c)}}) \quad (\text{V.4})$$

so that $\hat{I}_{N_c}(\overline{\Psi_s^{(N_c)}}) = 0$. In the same figure, we also show $\overline{\Psi_s^{(N_c)}}(T)$ as vertical dotted lines which correspond to the minima of the logarithmic distribution $\hat{I}_{N_c}(\Psi_s^{(N_c)})$. As can be seen, these minima are displaced towards the analytical value $\psi(s)$ (shown with a dashed line) as $N_c \rightarrow \infty$. The logarithmic distribution \hat{I}_{N_c} also becomes more concentrated as N_c increases.

Next, in order to study this decreasing of the width, we show a rescaled logarithmic distribution function $(1/N_c)\hat{I}_{N_c}(\Psi_s^{(N_c)})$ in Fig. V.3(b). The minimum converges to the analytical value $\psi(s)$ (black dashed line) as $N_c \rightarrow \infty$. In the infinite-time infinite-size limit of $\Psi_s^{(N_c)}$, it would be thus compatible with a logarithmic distribution function given by

$$I(\Psi_s^{(N_c)}) = -\lim_{N_c \rightarrow \infty} \frac{1}{N_c} \lim_{t \rightarrow \infty} \frac{1}{t} \log P(\Psi_s^{(N_c)}(t))$$

which is shown (rescaled) with black dots in Fig. V.3(b). By performing the shift $\check{\Psi}_s^{(N_c)} = (\Psi_s^{(N_c)} - \overline{\Psi_s^{(N_c)}})$ we can see in the inset of Fig. V.3(b) the superposition of quadratic deviations of the numerical estimator $\Psi_s^{(N_c)}$ around the minimum of \hat{I}_{N_c} (especially for $N_c = 100, 1000$). This indicates the decreasing of the fluctuation of CGF estimator proportional with both of T and N_c (see chapter III [P2] for more detailed explanation).

The obtained logarithmic distribution is well-approximated by a quadratic form, although these large deviations are in general not quadratic (chapter III [P2]). This means that the direct observation discussed here cannot capture the large deviations of the CGF estimator (see also Ref. [255] for more detailed study of the direct estimation of rate functions). However we note that, for practical usage of the algorithm, we only consider small fluctuations described by central limit theorem, although these large fluctuations might play an important role in more complicated systems, such as the ones presenting dynamical phase transitions.

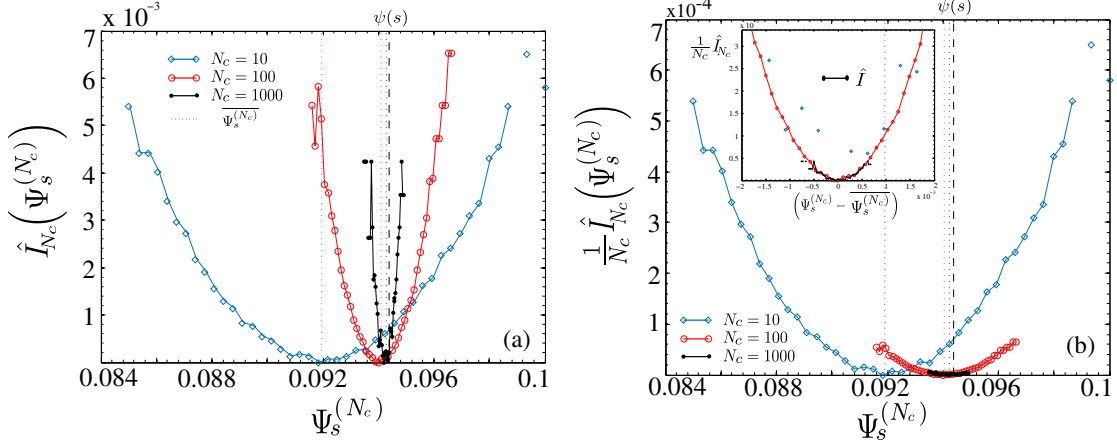


Figure V.3: (a) Logarithmic distribution $\hat{I}_{N_c}(\Psi_s^{(N_c)})$ (Eq. (V.4)). Numerical evaluations were made for three fixed population sizes $N_c \in \{10, 100, 1000\}$ with a fixed simulation time $T = 1000$. The logarithmic distribution presents a smaller width as N_c increases. The average over R realizations of the CGF estimator $\overline{\Psi_s^{(N_c)}}(T)$ corresponds to the minimum of $\hat{I}_{N_c}(\Psi_s^{(N_c)})$ (dotted lines) and converges to the analytical value $\psi(s)$ (dashed lines) as $N_c \rightarrow \infty$. (b) Rescaled logarithmic distribution $\frac{1}{N_c} \hat{I}_{N_c}(\Psi_s^{(N_c)})$ as a function of $\Psi_s^{(N_c)}$ and as a function of $\check{\Psi}_s^{(N_c)} = (\Psi_s^{(N_c)} - \overline{\Psi_s^{(N_c)}})$ (inset) for a final simulation time $T = 1000$.

V.3 A Different CGF Estimator

Normally, the CGF estimator is defined as an arithmetic mean over many realizations, as seen in Eq. (IV.1). Here we show that another definition of the CGF estimator can be used, which indeed provides better results than the ones from the standard estimator (in some parameter ranges). We define a new estimator as

$$\Phi_s^{(N_c)} = \frac{1}{T} \log \overline{\prod_{i=1}^{K_r} X_i^r}, \quad (\text{V.5})$$

where we note that the average with respect to realizations are taken *inside* the logarithm. As we discussed in Sec. III.4.3, this estimator provides a correct value of CGF $\psi(s)$ in the infinite-time infinite- N_c limits. This is thanks to the fact that the distribution of $\Psi_s^{(N_c)}$ concentrates around $\psi(s)$ in those limits (the so-called “self-averaging” property). At any finite population, one can rewrite $\Phi_s^{(N_c)}$ using the large-time LDF principle (V.3) as follows:

$$\begin{aligned} \Phi_s^{(N_c)} &= \frac{1}{T} \log \overline{e^{T\Psi_s^{(N_c)}}} \\ &= \frac{1}{T} \log \int d\Psi e^{-T[I_{N_c}(\Psi) + \Psi]} \end{aligned}$$

which proves that in the large- T limit,

$$\Phi_s^{(N_c)} = \min_{\Psi} [I_{N_c}(\Psi) + \Psi],$$

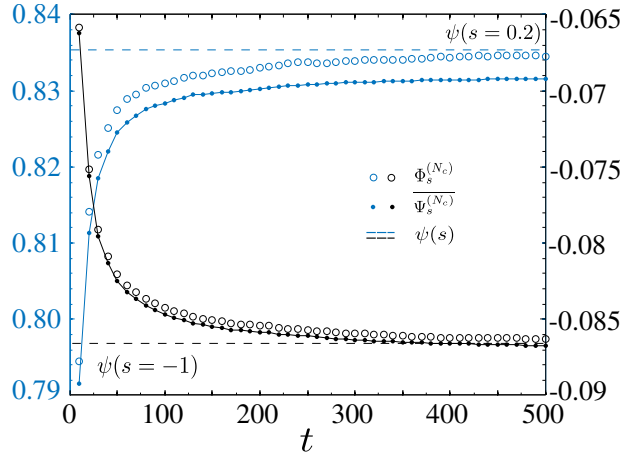


Figure V.4: Comparison between two different estimators of the large deviation function, $\overline{\Psi}_s^{(N_c)}$ (Eq. (IV.1)) shown in dots and $\Phi_s^{(N_c)}$ (Eq. (V.5)) in circles, for the annihilation-creation dynamics (Sec. I.8.1). The analytical value $\psi(s)$ (Eq (I.33)) is shown with a dashed line. Here we have also compared two different values of parameter $s = 0.2$ (blue) and $s = -1$ (black). Additionally, $N_c = 100$, $c = 0.4$, $T = 500$ and $R = 500$. As discussed in the text, $\Phi_s^{(N_c)}$ provides a better numerical evaluation of the CGF at small s .

to be compared to

$$\overline{\Psi}_s^{(N_c)} = \underset{\Psi}{\operatorname{argmin}} I_{N_c}(\Psi).$$

On one hand, the definition (V.5) amounts to estimate ψ from the exponential growth rate of the average of the final- T population of many small (non-interacting) “islands”, where the cloning algorithm would be operated. On the other hand, the estimator (IV.1) amounts to estimate ψ from growth rate of a large “island” gathering the full set of the R populations. The latter is thus expected to be a better estimator of $\psi(s)$ than the former because it corresponds to a large population, where finite-size effects are less important. As a consequence, the estimator $\Phi_s^{(N_c)}$ appears *a priori* to be worse estimator than $\overline{\Psi}_s^{(N_c)}$ of $\psi(s)$. However, as shown in Sec. III.4.3, at small $|s|$ and finite- N_c , a supplementary bias introduced by taking Eq. (V.5) in fact *compensates* the finite- N_c systematic error presented by Eq. (IV.1), for a simple two state model. Namely, the error is $O(sN_c^{-1})$ for Eq. (IV.1) while it is $O(s^2N_c^{-1})$ for Eq. (V.5). This fact is illustrated on Fig. V.4, where we show that at small $s = 0.2$, $\Phi_s^{(N_c)}$ provides a better estimation of $\psi(s)$ than $\overline{\Psi}_s^{(N_c)}$, while at larger $|s|$ ($s = -1$) the two estimators yield a comparable error.

VI – Breakdown of the Finite-Time and Finite- N_c Scalings in the Large- L Limit

VI.1 Introduction

The analysis of the finite- t and finite- N_c scalings in the evaluation of the LDF was performed following two different approaches: an analytical one, in chapter III [P2], using a discrete-time version of the population dynamics algorithm [18], and a numerical one, in chapter IV [P3], using a continuous-time version [17, 19]. In both cases, the systematic errors of these scalings were found to behave as $1/t$ and $1/N_c$ in the large- t and large- N_c asymptotics respectively. Moreover, it was shown how these scaling properties can be used in order to improve the LDF estimation by the implementation of a scaling method (Sec. IV.4.1). This was done considering that the asymptotic behavior of the estimator in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits may be interpolated from the data obtained from simulations at **finite and relative small** simulation time and number of clones.

However, the validity of these scalings and the method efficiency were proved only in cases for which the number of sites L (where the dynamics occurs) was small: a simple two-states annihilation-creation dynamics (Sec. I.8.1) (in one site) and a one-dimensional contact process (Sec. I.8.2) (with $L = 6$ sites). Here, we complement the results presented in chapters III [P2] and IV [P3] by extending the analysis to a large- L contact process. This is done by introducing the exponents γ_t and γ_{N_c} such that the generalized $t^{-\gamma_t}$ - and $N_c^{-\gamma_{N_c}}$ -scalings allow to characterize the scaling behavior in the large- L limit where we verify that t^{-1} and N_c^{-1} -scalings are no longer valid.

This chapter VI [P4] is organized as follows: The generalization to large- L systems of the finite-time and finite- N_c scalings of the LDF is done in Sec. VI.2.2. We make use of these results in Sec. VI.3 where we check the validity of the t^{-1} - and N_c^{-1} -scalings (Sec. VI.3.1), their behavior in the s -modified dynamics (Sec. VI.3.2) as well as the applicability of the scaling method (Sec. VI.3.3) for a contact process with $L = 100$ sites. This analysis is generalized in Sec. VI.4 where we characterize the finite- t and finite- N_c scalings of the LDF in the plane $s - L$. Before presenting our conclusions in Sec. VI.6, we discuss about the effects of the dynamical phase transition in the scalings in Sec. VI.5.

VI.2 Finite Scalings of the Large Deviation Function Estimator

Below, we summarize the finite-time and finite- N_c scalings of the CGF estimator and its generalization to large- L systems.

VI.2.1 Large-Time and Large- N_c Limit

When we analyze the time behavior of the CGF estimator (IV.2) for a fixed number of clones N_c , we observe this can be well described by a curve $f_t^{(N_c)}$ (IV.3) indicating the existence of a t^{-1} -convergence to the value $f_\infty^{(N_c)}$. We call this **t^{-1} -scaling** and is valid independently if N_c is small or large. The curve $f_t^{(N_c)}$ is determined from a fit in time over $\overline{\Psi_s^{(N_c)}(t)}$ up to (the final simulation) time T and allows the extraction of the infinite-time limit of the CGF estimator $f_\infty^{(N_c)} = \lim_{t \rightarrow \infty} \overline{\Psi_s^{(N_c)}(t)}$ which provides a better CGF estimation¹.

When we repeat this procedure for different values of population size $N_c \in \vec{N}_c = \{N_c^{(1)}, \dots, N_c^{(j)}\}$, extracting in each case the corresponding $f_\infty^{(N_c)}$, we observe they exhibit $1/N_c$ corrections in N_c (**N_c^{-1} -scaling**). In other words, the $f_\infty^{(N_c)}$'s satisfy a equation of the form (IV.7) which can be obtained from a fit in N_c over the extracted $f_\infty^{(N_c)}$'s. Thus, the t^{-1} - and N_c^{-1} -scalings of the CGF estimator are given by Eqs. (IV.3) and (IV.7), i.e.,

$$\begin{aligned} f_t^{(N_c)} &= f_\infty^{(N_c)} + b_t^{(N_c)} t^{-1}, \\ f_\infty^{(N_c)} &= f_\infty + b_\infty^{(N_c)} N_c^{-1}. \end{aligned}$$

These equations imply that $\overline{\Psi_s^{(N_c)}(t)}$ converges to its infinite- t and infinite- N_c limit, $f_\infty^\infty = \lim_{N_c \rightarrow \infty} f_\infty^{(N_c)}$, proportionally to $1/t$ and $1/N_c$. Importantly, this limit can be obtained using a small number of clones and simulation time by making use of the **scaling method** (Sec. IV.4.1 [P3]). The results obtained for f_∞^∞ rendered a better estimation of $\psi(s)$ than the standard estimator $\overline{\Psi_s^{(N_c)}(t)}$ evaluated for $N_c = \max \vec{N}_c$ and for $t = T$.

VI.2.2 Scalings in the Large- L Limit

In order to verify whether the scalings observed in small systems are also valid in the large- L limit, we assume that the CGF estimator can be described by equations of the form

$$\chi_t^{(N_c)} \equiv \chi_\infty^{(N_c)} + \kappa_t^{(N_c)} t^{-\gamma t}, \quad (\text{VI.1})$$

$$\chi_\infty^{(N_c)} \equiv \chi_\infty + \kappa_\infty^{(N_c)} N_c^{-\gamma N_c}, \quad (\text{VI.2})$$

redefining in a more general way the scalings (IV.3) and (IV.7). We will refer to Eq. (VI.1) as **$t^{-\gamma t}$ -scaling** whereas Eq. (VI.2) as **$N_c^{-\gamma N_c}$ -scaling**. The problem reduces in determining

¹Additionally, the behavior of the standard CGF estimator $\overline{\Psi_s^{(N_c)}(T)}$ as a function of the population size N_c is well described by a behavior of the form (IV.4)

$$g_{N_c}^{(T)} = g_\infty^{(T)} + \tilde{b}_{N_c}^{(T)} N_c^{-1}$$

indicating that $\overline{\Psi_s^{(N_c)}(T)}$ also converges to its infinite- N_c limit $g_\infty^{(T)} = \lim_{N_c \rightarrow \infty} \overline{\Psi_s^{(N_c)}(T)}$ with an error proportional to $1/N_c$.

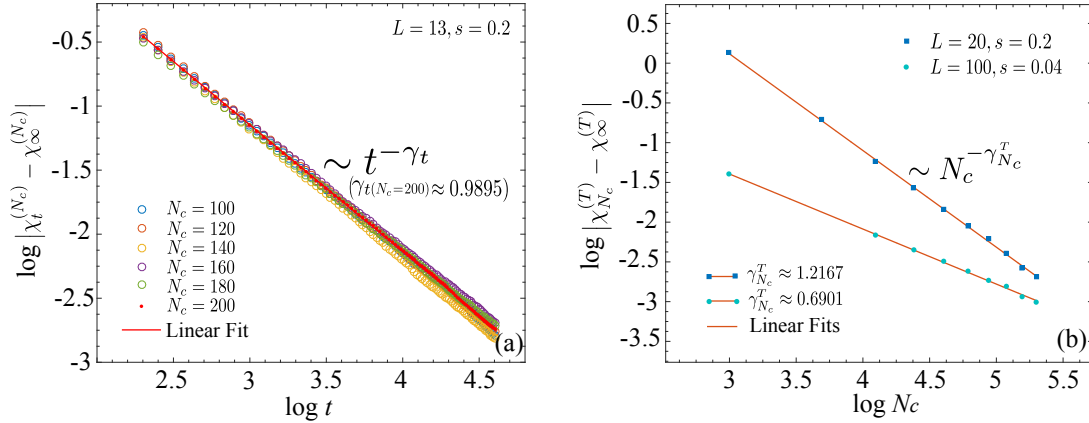


Figure VI.1: Exponents **(a)** γ_t & **(b)** γ_{N_c} which characterize the finite-time and -size scalings of the CGF estimator of the activity for a contact process with $\lambda = 1.75$ and $h = 0.1$. These exponents were determined from the slope of a linear fit in log-log scale over Eqs. (VI.4) and (VI.3), respectively. For γ_t , we used $L = 13$, $s = 0.2$ and $N_c \in \vec{N}_c = \{100, 120, \dots, 200\}$. Meanwhile, γ_{N_c} was computed for $L = 20$ and $s = 0.2$, and for $L = 100$ and $s = 0.04$. Additionally, in all the cases $T = 100$ and $R = 500$.

the exponents γ_t and γ_{N_c} in order to verify if effectively $\gamma_t \approx 1$ and $\gamma_{N_c} \approx 1$ and whether the terms $\chi_\infty^{(N_c)}$ and χ_∞ represent the limits in $t \rightarrow \infty$ and $N_c \rightarrow \infty$ of the CGF estimator. Thus, a value of the exponent $\gamma_t \approx 1$, verifies $\chi_\infty^{(N_c)} \approx f_\infty^{(N_c)}$ and $\gamma_{N_c} \approx 1$, verifies $\chi_\infty \approx f_\infty$. This is done in Sec. VI.3 on a contact process with $L = 100$ sites. Below we describe the procedure followed in order to obtain these exponents ².

VI.2.2.1 Determination of the Exponents γ_t & γ_{N_c}

From Eqs. (VI.1) and (VI.2) we expect that, independently of N_c , T , L or s , a power law behavior of the form

$$|\chi_t^{(N_c)} - \chi_\infty^{(N_c)}| \sim t^{-\gamma_t}, \quad (\text{VI.4})$$

$$|\chi_\infty^{(N_c)} - \chi_\infty| \sim N_c^{-\gamma_{N_c}}, \quad (\text{VI.5})$$

be observed. Thus, the exponents γ_t & γ_{N_c} can be obtained from the slope of a straight curve in log-log scale of Eqs. (VI.4) and (VI.5). This can be seen in Fig. VI.1. Despite only some representative configurations have been considered, we confirm this power law behavior independently of the parameters chosen.

² Additionally to Eqs. (VI.1) and (VI.2), the N_c -behavior of $\overline{\Psi_s^{(N_c)}(T)}$ can be described by the equation

$$\chi_{N_c}^{(T)} = \chi_\infty^{(T)} + \tilde{\kappa}_{N_c}^{(T)} N_c^{-\gamma_{N_c}^T}, \quad (\text{VI.3})$$

where $\chi_\infty^{(T)} = \lim_{N_c \rightarrow \infty} \overline{\Psi_s^{(N_c)}(T)}$. Here it is important to remark that both $\chi_\infty^{(N_c)}$ and $\overline{\Psi_s^{(N_c)}(T)}$ scale in the same way in N_c . In other words, $\gamma_{N_c} \approx \gamma_{N_c}^T$.

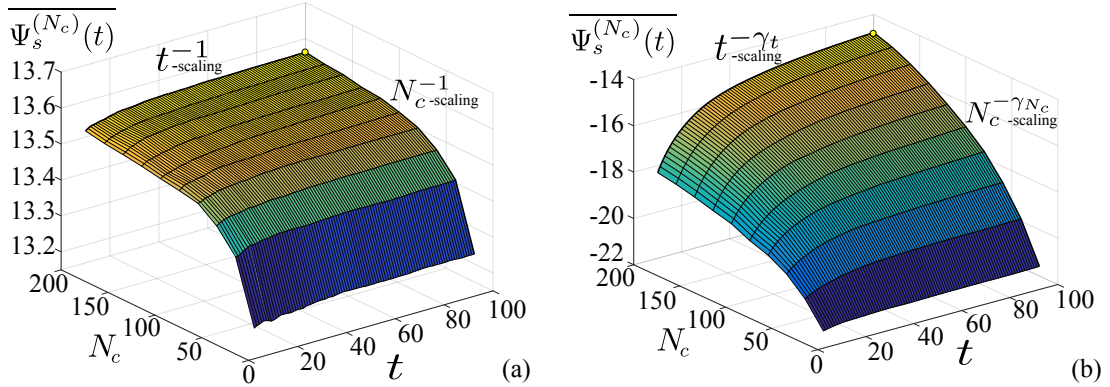


Figure VI.2: Large deviations of the activity for a contact process with $L = 100$, $\lambda = 1.75$ and $h = 0.1$. The CGF estimator $\overline{\Psi_s^{(N_c)}(t)}$ (Eq. (IV.1)) is presented as a function of time t and the number of clones N_c for (a) $s = -0.1$ and (b) $s = 0.2$. These surfaces were computed using the continuous-time cloning algorithm up a final simulation time $T = 100$, $\vec{N}_c = \{20, \dots, 200\}$ and $R = 500$ realizations. The N_c^{-1} -scaling observed in small- L systems holds only for $s = -0.1$ whereas for $s = 0.2$ a $N_c^{-\gamma_{N_c}}$ -scaling is observed ($\gamma_{N_c}(s = 0.2) \approx -0.16$). Similarly, for the time-scaling for which $\gamma_t(s = 0.2) \approx 0.7$.

VI.3 Finite Scalings for a Large- L Contact Process

In Fig. VI.2, we compare the behavior of $\overline{\Psi_s^{(N_c)}(t)}$ as function of t and N_c , for two representative values of the parameter s , $s = -0.1$ (left) and $s = 0.2$ (right). The size of the system is $L = 100$ sites. Each point of these surfaces was obtained using the cloning algorithm (Eq. (IV.1)) up to time $T = 100$, for $\vec{N}_c = \{20, 40, \dots, 180, 200\}$ and for $R = 500$ realizations. The best possible CGF estimation (i.e., at largest T and N_c) in both cases is shown with solid circles which, according to the results presented in chapter IV [P3], could be improved by using the t^{-1} and N_c^{-1} -scalings (if still valid for large- L).

VI.3.1 Finite-Time and Finite- N_c Scalings

Although the exponents γ_t and γ_{N_c} can be computed in principle for any value of $N_c \in \vec{N}_c$ and for any $t \leq T$, as we saw above, from now on, we will consider these exponents defined at the highest number of clones and at final simulation time, i.e.,

$$\begin{aligned}\gamma_t &:= \gamma_t(N_c = \max \vec{N}_c), \\ \gamma_{N_c} &:= \gamma_{N_c}(t = T).\end{aligned}$$

Thus, the exponent γ_t is obtained as described in Sec. VI.2.2.1 after adjusting Eq. (VI.1) to $\overline{\Psi_s^{(N_c)}(t)}$ for $N_c = \max \vec{N}_c = 200$. On the other hand, γ_{N_c} is determined after fitting $\chi_\infty^{(N_c)}$ with Eq. (VI.2) at $T = 100$ or, as $\gamma_{N_c} \approx \gamma_{N_c}^T$, after fitting $\overline{\Psi_s^{(N_c)}(t = T)}$ using Eq. (VI.3). In simple words, these exponents can be obtained from an adequate fit over the thick curves in Fig. (VI.2). They characterize the finite- t and finite- N_c scalings of the large deviations of the dynamical activity K .

Following this approach, we found that the t^{-1} -scaling (IV.3) is satisfied only for $s = -0.1$, meaning that the exponent γ_t was found to be $\gamma_t \approx 1$. As a consequence, the parameter $\chi_\infty^{(N_c)}$ obtained from Eq. (VI.1) effectively represents the limit in $t \rightarrow \infty$ of the CGF estimator, i.e., $\chi_\infty^{(N_c)} \approx f_\infty^{(N_c)}$. This is not the case for $s = 0.2$ for which $\gamma_t(s = 0.2) \approx 0.7$. Similarly, a $N_c^{-\gamma_{N_c}}$ -scaling is observed for $s = 0.2$, whereas for $s = -0.1$, the N_c^{-1} -scaling (IV.7) holds. It is important to remark that a value of exponent $\gamma_{N_c} > 0$ could still guaranty the convergence of the CGF estimator in the infinite- N_c limit. However, even though $\gamma_{N_c}(t = 10) > 0$ at initial times, at final time T , the exponent is negative ($\gamma_{N_c}(t = T) \approx -0.16$), which would imply that χ_∞^∞ does not corresponds to the $t, N_c \rightarrow \infty$ limit of the CGF estimator. This fact will be addressed later in Sec. VI.5. Below, we present how the change in the scalings is produced depending on s .

VI.3.2 Exponents Characterization & s -Dependence

Here, we consider values of s ranging in the interval $s \in [-0.1, 0.2]$. For $s < 0$, the exponent $\gamma_t(s)$ varies around 1. However for $s > 0$, γ_t deviates slightly from 1 decreasing with s up to $\gamma_t \approx 0.7$ at $s = 0.2$. In order to describe the behavior of this exponent, results convenient to define s' as the value of the parameter $s \in [s_a, s_b]$ such that $\gamma_t(s < s') \approx 1$, i.e., until which the t^{-1} -scaling holds. Thus,

$$\gamma_t(L = 100) : \begin{cases} \gamma_t(s) \approx 1, & \text{for } s < s' \\ 0 < \gamma_t(s) < 1, & \text{otherwise.} \end{cases}$$

If the scaling holds $\forall s \in [s_a, s_b]$ (given some system size L), then $s' = s_b$.

On the other hand, the value of $s \in [s_a, s_b]$ which signals the validity of the N_c^{-1} -scaling is denoted by s^* . From this point, γ_{N_c} decreases until eventually it becomes negative, as can be seen in Fig. VI.3. Here, we introduce s^{**} such that $\gamma_{N_c}(s = s^{**}) = 0$. Thus, $\gamma_{N_c} < 0$ for $s > s^{**}$. This behavior was not observed in chapter IV [P3] for $L = 6$ for which the N_c^{-1} -scaling was valid independently of s . In those cases, $s^* = s_b$ and $\nexists s^{**}$. Instead of confirming for $L = 100$ the N_c^{-1} -scaling of the CGF estimator presented in chapter IV [P3], here we have been able to distinguish clearly three stages for the exponent $\gamma_{N_c}(s)$:

$$\gamma_{N_c}(L = 100) : \begin{cases} \gamma_{N_c}(s) \approx 1, & \text{for } s < s^* \\ 0 < \gamma_{N_c}(s) < 1, & \text{for } s^* < s < s^{**} \\ \gamma_{N_c}(s) < 0, & \text{for } s > s^{**}. \end{cases}$$

The possibility of extracting the infinite- t and infinite- N_c limit of the CGF estimator relied on the validity of the t^{-1} - and N_c^{-1} -scalings. How the results obtained from the application of the scaling method are affected by γ_t and γ_{N_c} are presented below.

VI.3.3 Implementation of the Scaling Method

The scaling method allows to determine the asymptotic limit to which the CGF estimator (IV.1) converges in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits. This limit, that we have denoted f_∞^∞ (Eq. (IV.7)), was proved to render a better estimation of the analytical CGF $\psi(s)$ than the standard estimator $\Psi_s^{(\max N_c)}(T)$, at least for the cases analyzed in chapter IV [P3]. However, the results we just presented would suggest that the determination of f_∞^∞ could be affected

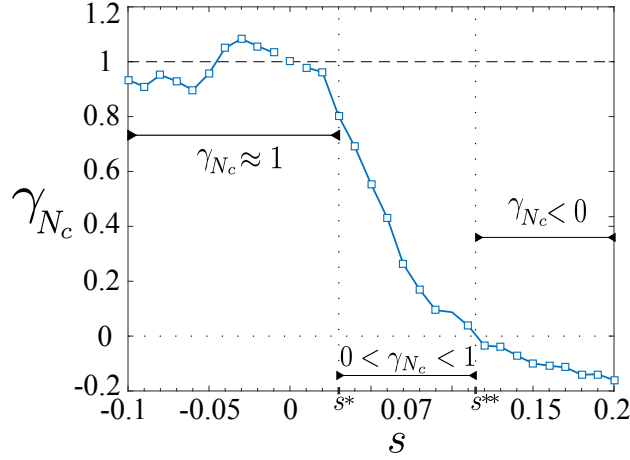


Figure VI.3: Dependence of the $N_c^{-\gamma_{N_c}}$ -scaling with the parameter $s \in [-0.1, 0.2]$. The exponent $\gamma_{N_c}(s)$ is obtained by fitting $\overline{\Psi_s^{(N_c)}}(T=100)$ as function of $N_c \in \vec{N}_c = \{20, 40, \dots, 180, 200\}$ by Eq. (VI.3) in log-log scale as described in Sec. VI.2.2.1. Three stages of $\gamma_{N_c}(s)$ can be clearly distinguish for $L = 100$: (i) $\gamma_{N_c}(s < s^*) \approx 1$, (ii) $0 < \gamma_{N_c}(s^* < s < s^{**}) < 1$, and (iii) $\gamma_{N_c}(s > s^{**}) < 0$. The exponent γ_{N_c} for $s = 0$ was set to $\gamma_{N_c}(s = 0) = 1$.

depending whether $\gamma_t \approx \gamma_{N_c} \approx 1$ or not. If this holds, the scaling method could render valid results in our example only for $s < 0$. Solely in this region the extracted χ_∞^∞ (obtained from Eq. (VI.2)) would represent the infinite- t and $-N_c$ limit of the CGF estimator. Indeed, this can be observed in Fig. VI.4 where we have applied the scaling method to our example.

The method can be performed following two different approaches: (i) $(t^{-1}, N_c^{-\gamma_{N_c}})$: First, imposing a t^{-1} -scaling for $\overline{\Psi_s^{(N_c)}}(t)$ (setting $\gamma_t = 1$ in Eq. (VI.1)) and then, considering a $N_c^{-\gamma_{N_c}}$ -scaling (VI.2) for the extracted $\chi_\infty^{(N_c)}$'s. Alternatively, (ii) $(t^{-\gamma_t}, N_c^{-\gamma_{N_c}})$: Leaving γ_t and γ_{N_c} as free parameters in Eqs. (VI.1) and (VI.2). Both resulting estimators $\chi_\infty^\infty(i)$ and $\chi_\infty^\infty(ii)$ are shown in Fig. VI.4 with squares and circles, respectively. Additionally, the infinite- N_c limit χ_∞^T (VI.3) is also presented with diamonds. The standard CGF estimator $\overline{\Psi_s^{(\max \vec{N}_c)}}(T)$ (in dots) serves as reference.

As can be seen in Fig. VI.4, the different estimators correspond to each other up to $s = s^*$. From this point, their distance with respect to $\overline{\Psi_s^{(\max \vec{N}_c)}}(T)$ increases rapidly with s up to $s = s^{**}$ where a discontinuity occurs. In fact, the behavior observed in Fig. VI.4 keeps correspondence with the $N_c^{-\gamma_{N_c}}$ -scaling of the CGF estimator. Specifically, with the stages of the exponent γ_{N_c} that were presented in Sec. VI.3.2 and Fig. VI.3. Thus, the discontinuity in χ_∞^∞ is related precisely with the change in sign of γ_{N_c} in $s = s^{**}$ in the same way as the divergence of the estimators from the standard one at $s = s^*$ is related with the fact that from this point, $\gamma_{N_c} \neq 1$.

The example presented through this section related the effectiveness of the scaling method (as proposed in chapter IV [P3]) with the actual scaling of the CGF estimator in large- L systems. Depending on the value of the exponents (γ_t, γ_{N_c}) it is possible to extract or not the infinite- t and infinite- N_c limit of the CGF estimator. Moreover, we also presented

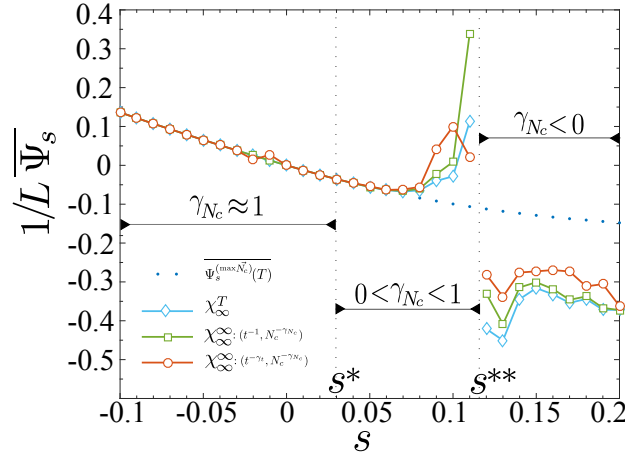


Figure VI.4: Different estimators of the large deviations of the activity as function of the parameter $s \in [-0.1, 0.2]$ for a contact process with $L = 100$ sites. The standard CGF estimator $\overline{\Psi_s^{(N_c)}(t)}$ (evaluated at $N_c = \max \vec{N}_c = 200$, $t = T = 100$ and for $R = 500$ realizations) is shown with dots meanwhile the ones obtained from the scaling method χ_∞^∞ are presented in squares and circles, and χ_∞^T in diamonds. The legend $(t^{-1}, N_c^{-\gamma_{N_c}})$ refers to the assumption of a t^{-1} -scaling for $\overline{\Psi_s^{(N_c)}(t)}$ (setting $\gamma_t = 1$ in Eq. (VI.1)) and a $N_c^{-\gamma_{N_c}}$ -scaling (VI.2) for the $\chi_\infty^{(N_c)}$'s. On the other hand, $(t^{-\gamma_t}, N_c^{-\gamma_{N_c}})$ refers to the fact that we have left γ_t and γ_{N_c} as free parameters. The different estimators correspond to each others up to $s = s^*$ from which they diverge up to $s = s^{**}$ where there is a discontinuity and they become negative for $s > s^{**}$. This is directly related with the behavior of the exponent γ_{N_c} observed in Fig. VI.3.

the dependence of these exponents with the parameter s . Below we extend our analysis by considering the scaling behavior on a wider range of values of L . This will provide a complete overview of how the CGF estimator behaves and how the change in scaling is given.

VI.4 L -Dependence of the Finite Scalings

In this section, we detail the behavior of the finite- t and $-N_c$ scalings of the CGF estimator for $s > 0$ and L ranging in the interval $L \in [3, 100]$. For each pair (s, L) , the exponents γ_t and γ_{N_c} were computed as described in Sec. VI.2.2.1 for $T = 100$ and $\vec{N}_c = \{20, 40, \dots, 180, 200\}$.

VI.4.1 Characterization of the Exponent $\gamma_t(s, L)$

The contour plot in Fig. VI.5 shows the value of the exponent γ_t as it changes depending on the parameters s and L . We have focused in the region for $s \in [0.02, 0.2]$ as for $s < 0$, $\gamma_t \approx 1$ and thus, the t^{-1} -scaling (IV.3) holds. The values closest to 1 are presented with the darkest tone while smaller values are shown with clearer tones. As can be seen, the exponent γ_t decreases gradually as L and s increase.

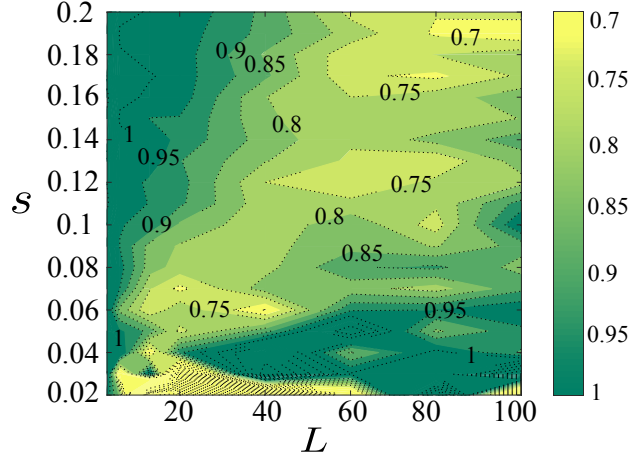


Figure VI.5: Behavior of the exponent γ_t in a region of the plane $s - L$ given by the parameters $s \in [0.02, 0.2]$ and $L \in [3, 100]$. This exponent characterizes the finite-time scaling of the large deviations of the activity in the contact process ($t^{-\gamma_t}$ -scaling). The values of γ_t closest to 1 are presented with the darkest tones whereas smaller values are shown with clearer ones. Although for small values of L the t^{-1} -scaling holds independently of s , in general the exponent γ_t decreases gradually as s and L increase.

For a given system size L , we can describe qualitatively the behavior of γ_t with respect to s in a similar way as we did for $L = 100$ in Sec. VI.3.2. In order to extend that description into the plane $s - L$, we introduce a number of sites dependency of the bound s' . We denote by $s'(L)$ the value of s until which the t^{-1} -scaling is valid given a particular L . Similarly, $\gamma_t^\circ(L)$ is the lower bound of $\gamma_t^{(L)}(s)$. Thus, the exponent γ_t which characterizes the $t^{-\gamma_t}$ -scaling (VI.1) of the CGF estimator is given by

$$\gamma_t : \begin{cases} \gamma_t^{(L)}(s) \approx 1, & \text{for } s < s'(L) \\ \gamma_t^\circ(L) \leq \gamma_t^{(L)}(s) \lesssim 1, & \text{otherwise,} \end{cases}$$

where $s'(L) > 0$, $\gamma_t^\circ(L) > 0$ and L is large. In fact, for this case, $\gamma_t^\circ(L) > 1/2$, for all L .

VI.4.2 Characterization of the Exponent $\gamma_{N_c}(s, L)$

Similarly as above, in Fig. VI.6 we present the exponent γ_{N_c} as it changes depending of some particular choice of the parameters (s, L) within the intervals considered. The surface in Fig. VI.6(a) illustrates clearly the change in the N_c -scaling of the CGF estimator. For every value of L considered, the exponent γ_{N_c} is approximately 1 up to some value of s , denoted as $s^*(L)$ (Sec. VI.3.2). However, from this point, its value decreases as s and L increases, becoming, in some cases, negative. This change in the $N_c^{-\gamma_{N_c}}$ -scaling is also shown in the contour plot in Fig. VI.6(b) where we have focus in the region for $s > 0$. The values of γ_{N_c} closer to 1 are shown in dark tones.

In Sec. VI.3.2, we also defined s^{**} such that $\gamma_{N_c}^{(L)}(s^{**}) = 0$. This value of course depends on L and in some cases it does not even exist. However, for some particular values of L (large), the exponent $\gamma_{N_c}^{(L)}$ changes sign twice (as can be seen in Fig. VI.6(b)). We will use

this fact in order to characterize the $N_c^{-\gamma_{N_c}}$ -scaling depending on the number of zeros of the exponent $\gamma_{N_c}^{(L)}(s)$ for a given L .

We define \mathcal{L}_I as the set of values of L , for which the exponent $\gamma_{N_c}^{(L)}(s)$ has no zeros, \mathcal{L}_{II} : if has two zeros ($s_1^{**}(L)$ and $s_2^{**}(L)$, with $s_2^{**}(L) > s_1^{**}(L)$) and \mathcal{L}_{III} : if has one zero ($s^{**}(L)$). These regions are bounded by L_{inf} and/or by L_{sup} , where L_{inf} is the smallest value of L such that the curve $L = L_{inf}$ is tangent to $\gamma_{N_c}(s, L) = 0$ in one single point. On the other hand, L_{sup} is the largest L such that the curve $L = L_{sup}$ cuts $\gamma_{N_c}(s, L) = 0$ in two points. Thus, the region \mathcal{L}_I groups the values of L such that $L < L_{inf}$, \mathcal{L}_{II} the values of L within the interval $L_{inf} < L < L_{sup}$ and \mathcal{L}_{III} , the values of L such that $L > L_{sup}$. Thus, the exponent γ_{N_c} which characterizes the $N_c^{-\gamma_{N_c}}$ -scaling (VI.2) of the CGF estimator is given by

$$\gamma_{N_c} : \begin{cases} \mathcal{L}_I : \begin{cases} \gamma_{N_c}^{(L)}(s) \approx 1, & \text{for } s < s^*(L) \\ 0 < \gamma_{N_c}^{(L)}(s) \lesssim 1, & \text{otherwise.} \end{cases} \\ \mathcal{L}_{II} : \begin{cases} \gamma_{N_c}^{(L)}(s) \approx 1, & \text{for } s < s^*(L) \\ 0 < \gamma_{N_c}^{(L)}(s) < 1, & \text{for } s^*(L) < s < s_1^{**}(L) \\ & \text{and } s > s_2^{**}(L) \\ \gamma_{N_c}^{(L)}(s) < 0, & \text{for } s_1^{**}(L) < s < s_2^{**}(L) \end{cases} \\ \mathcal{L}_{III} : \begin{cases} \gamma_{N_c}^{(L)}(s) \approx 1, & \text{for } s < s^*(L) \\ 0 < \gamma_{N_c}^{(L)}(s) < 1, & \text{for } s^*(L) < s < s^{**}(L) \\ \gamma_{N_c}^{(L)}(s) < 0, & \text{for } s > s^{**}(L) \end{cases} \end{cases}$$

VI.5 Dynamical Phase Transition and Finite-Scalings

In Sec. I.5, we introduced the biasing parameter (or field) s (conjugated to an observable \mathcal{O}) to characterize a non-equilibrium ensemble of trajectories. Within this “ s -ensemble”, space-time or **dynamical phase transitions** manifest themselves as singularities in the CGF and, in our case, express a dynamical coexistence of histories with high and low activity K [35].

The contact process [38, 107, 108] is well know to exhibit a dynamical phase transition in the $L \rightarrow \infty$ limit [19, 35, 74, 252] even in one-dimension [252]. However in Ref. [19] evidence of the presence of a phase transition (in the active phase of λ) was reported to occur at $s_c \approx 0.057$ for finite- L . There, the authors used the same version of the contact process and the same approach we used throughout this thesis (i.e., the cloning algorithm). On the other hand, in Ref. [74], using a density matrix re-normalization group approach (DMRG) [260–264], it was showed that for every value of infection rate λ , either if this belong to the absorbing or to the active phase, there exists a phase transition as a function of s . For the case of the active phase, this transition was found to occur at $s_c = 0$. It is important to remark that even if the versions of the contact process used in Refs. [19] and [74] are different, both present a dynamical phase transition. Meanwhile in the later case the particles are created just at the boundaries, in Ref. [19] (and here) they are created at every site and also, the spontaneous rate of creation h is considered different from 0 (in order to circumvent the absorbing state in finite size [35]).

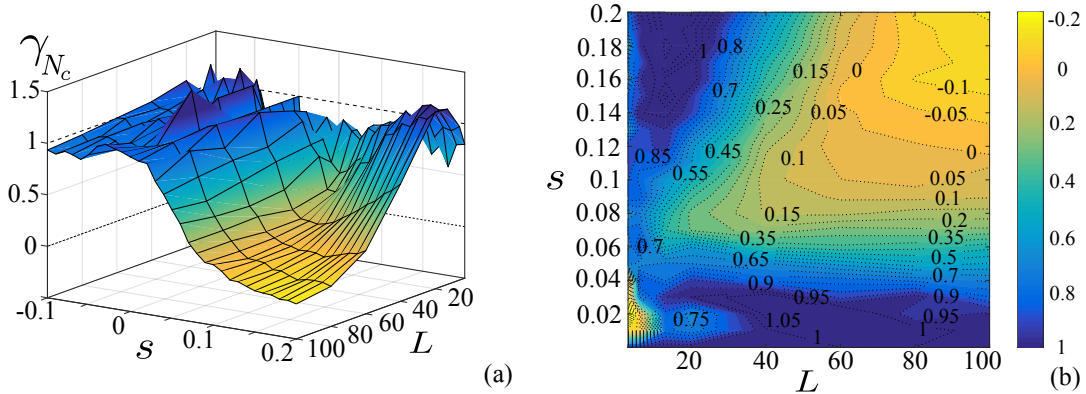


Figure VI.6: Behavior of the exponent γ_{N_c} in a region of the plane $s - L$ given by the parameters $s \in [-0.1, 0.2]$ and $L \in [3, 100]$. This exponent characterizes the finite- N_c scaling of the large deviations of the activity in the contact process (**$N_c^{-\gamma_{N_c}}$ -scaling**). (a) The surface $\gamma_{N_c}(s, L)$ illustrates the change in the scaling of the CGF estimator. The exponent $\gamma_{N_c}(s, L) \approx 1$ up to some value of s (which depends of L), from which it decreases as s and L increases, even becoming negative for some values of L . (b) Projection of the surface in (a) on the plane $s - L$. The values of $\gamma_{N_c}(s, L)$ closer to 1 are shown in dark tones while the smaller tones with clearer ones. The $N_c^{-\gamma_{N_c}}$ -scaling can be characterized depending on the number of zeros of the exponent $\gamma_{N_c}(s)$ for a given L in three regions: \mathcal{L}_I if has no zeros, \mathcal{L}_{II} : if have two zeros and \mathcal{L}_{III} : if have one zero.

Despite our main interest is not the study of the dynamical phase transition in the contact process, what does concern us is how this could affect the finite scalings of the CGF. Importantly, the relation that s^* (but also s' or s^{**}) could have with s_c where this transition occurs. In Sec. VI.2 we showed how the scaling behavior given by Eqs. (VI.1) and (VI.2) was robust independently of T , N_c , s or L (Fig. VI.1), but not the exponents γ_t and γ_{N_c} whose behavior change depending on s and L specially in $s \geq 0$ as L becomes larger. We remark that even if the infinite- L limit is not achievable numerically, the effects induced by a dynamical phase transition should become more evident as L increases (which could explain many of the behavior observed throughout this chapter). This was clearly illustrated for $L = 100$ for which γ_{N_c} has an abrupt change for $s \geq 0$, where we know the dynamical phase transition occurs, even taking negative values and inducing a divergence of the infinite- t and infinite- N_c limit of the CGF estimator (Fig. VI.4).

We recall here that our purpose was to verify the validity of the scalings (and the scaling method) presented in chapters III [P2] and IV [P3] (for small size systems) in the large- L limit where a main result was the possibility of extracting the infinite- N_c infinite- t limit of the CGF estimator from finite and small number of clones and time. An analysis of the dynamical phase transition, on the other hand, would require a large- N_c and $-t$ configuration which under our approach is a task difficult to fulfill. This however does not represent any surprise given that is well know that the existing methods [12, 13, 18, 265, 266] perform poorly in the vicinity of a dynamical phase transition, or they are numerically expensive in order to obtain accurate estimations [266–268] developing if not important finite-size effects [22]. However, recently has been proposed a promising method [85, 259] which combines the existing cloning

algorithm [7, 12, 13, 17–19, 265, 266, P2, P3] with a modification of the dynamics [88–91] resulting in a significant improvement of its computational efficiency. The method was successfully applied to the study of the dynamical phase transition of 1D FA model [42] using a relatively small N_c and L . The implementation of this method will provide in a next stage a clear contrast between the results obtained following the two different approaches and a correct relation between s_c and s^* .

VI.6 Conclusion

The dependence of the CGF estimator (and of its accuracy) with the simulation time T and number of clones N_c was studied in chapters III [P2] and IV [P3] where the finite- t and finite- N_c scalings of the systematic errors of the LDF were found to behave as $1/N_c$ and $1/t$ in the large- N_c and large- t asymptotics, respectively. By making use of these convergence-speeds, a scaling method was proposed which allowed to extract the asymptotic behavior of the CGF estimator in the $t \rightarrow \infty$ and $N_c \rightarrow \infty$ limits. At least for the cases analyzed in chapters III [P2] and IV [P3], this infinite-time and infinite- N_c limit resulted to render a better LDF estimation in comparison with the standard estimator. However, the validity of the method and of these scalings were proved only for a simple one-site annihilation-creation dynamics and for a contact process with $L = 6$ sites, leaving an analysis of the dependence with the number of sites L pending. In order to do so, in this chapter, we redefined these scalings in a more general way. We assume the behavior of the CGF estimator described by a $t^{-\gamma_t}$ -scaling (Eq. (VI.1)) and a $N_c^{-\gamma_{N_c}}$ -scaling (Eq. (VI.2)). This redefinition allowed us to verify in large- L systems if effectively $\gamma_t \approx 1$ and $\gamma_{N_c} \approx 1$ and whether the terms $\chi_\infty^{(N_c)}$ and χ_∞^∞ represent the limits in $t \rightarrow \infty$ and $N_c \rightarrow \infty$ of the CGF estimator.

This was done at first in Sec. VI.3.1 where we considered a contact process with $L = 100$ sites and two representative values of the parameter s . Although the t^{-1} -scaling and N_c^{-1} -scaling were proved to hold for $s = 0.1$, this was not the case for $s = 0.2$. How this change in the scaling is produced depending on the parameter s was presented in detail in Sec. VI.3.2 where the exponents $\gamma_t(s)$ and $\gamma_{N_c}(s)$ were characterized. Particularly, for $\gamma_{N_c}(s)$, we were able to distinguish three stages in its behavior, where, the N_c^{-1} -scaling was valid up to $s = s^*$, then γ_{N_c} decreases to 0 at $s = s^{**}$ and finally, it becomes negative for $s > s^{**}$. In Sec. VI.3.3 we showed how these scalings affect the determination of the infinite- t and infinite- N_c limit of the CGF estimator. This occurs because the scaling method relied on the validity of the t^{-1} - and N_c^{-1} -scalings. As for $L = 100$ this was not the case, it was possible to see how the different estimators corresponded to each others up to $s = s^*$ from which they diverge up to $s = s^{**}$ where there is a discontinuity.

This analysis was extended to the plane $s-L$ in Sec. VI.4 where the exponents γ_t and γ_{N_c} were computed for a grid of values of the parameters (s, L) . Their characterization was done introducing a number-of-sites dependency of the bounds s' , s^* and s^{**} previously defined in Sec. VI.3 as well as the use of the number of zeros of the exponent $\gamma_{N_c}^{(L)}(s)$ in order to characterize the different groups of L . Whether the results presented through this chapter are restricted only to the contact process or not is left as a pending problem and a possible direction for future research.

VII – Intra-day Seasonalities in High Frequency Financial Time Series

VII.1 Introduction

From the statistical study of financial time series have arisen a set of properties or empirical laws sometimes called “stylized facts” or seasonalities. These properties have the characteristic of being common and persistent across different markets, time periods and assets [93–99]. As it has been suggested in Ref. [99], the reason why these “patterns” appear could be because markets operate in synchronization with human activities which leave a trace in the financial time series. However using the “right clock” might be of primary importance when dealing with statistical properties and the patterns could vary depending if we use daily data or intra-day data and event time, trade time or arbitrary intervals of time (e.g. $T = 1, 5, 15$ minutes, etc). For example, it is a well-known fact that empirical distributions of financial returns and log-returns are fat tailed [102, 103], however as one increases the time scale the fat-tail property becomes less pronounced and the distribution approach the Gaussian form [104]. As was stated in Ref. [96], the fact that the shape of the distribution changes with time makes it clear that the random process underlying prices must have a non-trivial temporal structure. In a previous work, Allez et al. [99] established several new stylized facts concerning the intra-day seasonalities of single and cross-sectional stock dynamics. This dynamics is characterized by the evolution of the moments of its returns during a typical day. Following the same approach, we show the bin size dependence of these patterns for the case of returns and, motivated by the work of Kaisoji [100], we extend the analysis to relative prices and show how in this case, these patterns are independent of the size of the bin, also independent of the index we consider but characteristic for each index. These facts could be used in order to detect an anomalous behavior during the day, like market crashes or intra-day bubbles [100, 101]. The work presented in this chapter VII [P0] is completely empirical but it could offer signs of the underlying stochastic process that governs the financial time series.

VII.2 Definitions

The data consists in two sets of intra-day high frequency time series, the CAC 40 and the S&P 500. For each of the $D = 22$ days of our period of analysis (March 2011), we dispose with the evolution of the prices of each of the stocks that composes our indexes during a specific day from 10 : 00 a.m. to 16 : 00 p.m. The main reasons why we chose to work with

these two indexes are: The number of stocks that compose them ($N_1 = 40$ and $N_2 = 500$), the time gap between their respective markets and the different range of stock prices (between 5 and 600 USD for the S&P 500 and between 5 and 145 EU for the CAC 40).

As the changes in prices are not synchronous between different stocks (Fig. VII.1), we manipulated our original data in order to construct a new homogeneous matrix $P_D^{(j)}$ of bin prices. In order to do this, we divided our daily time interval $[10 : 00, 16 : 00]$ in K bins of size T (minutes), i.e., $B_1 = [10 : 00, 10 : 00 + T]$, $B_2 = [10 : 00 + T, 10 : 00 + 2T]$, \dots , $B_K = [16 : 00 - T, 16 : 00]$, where the right endpoints of these intervals are called bin limits. For a particular day j , the prices that conform the matrix $P_D^{(j)}$ are given by the last prices that reaches that stock i just before a specific bin limit.

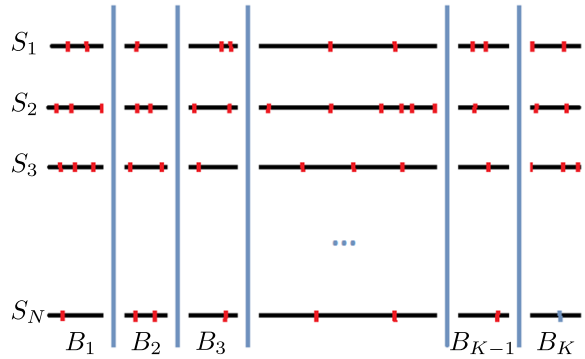


Figure VII.1: Intra-day asynchronous financial time series. S_i are the stocks and B_k are bins. The asynchronous prices are show in red and the bin limit in blue.

Each row in the matrix below represents the evolution of the prices of a particular stock as function of the bins. For example, the element $(P_D^{(j)})_{ik}$, represents the price for a particular day j of the stock i and just before the bin limit of the bin B_k .

$$P_D^{(j)} = \begin{pmatrix} P_{11}^{(j)} & P_{12}^{(j)} & \dots & \dots & \dots & P_{1K}^{(j)} \\ P_{21}^{(j)} & P_{22}^{(j)} & \dots & \dots & \dots & P_{2K}^{(j)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & P_{ik}^{(j)} & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ P_{N1}^{(j)} & P_{N2}^{(j)} & \dots & \dots & \dots & P_{NK}^{(j)} \end{pmatrix}$$

In a similar way, we can construct the matrix $P_S^{(i)}$ for each of the $i = 1, \dots, N_{1,2}$ stocks. $(P_S^{(i)})_{jk}$ is the price of the stock (i) in the day j and just before the bin limit of the bin B_k .

$$P_S^{(i)} = \begin{pmatrix} P_{11}^{(i)} & P_{12}^{(i)} & \dots & \dots & \dots & P_{1K}^{(i)} \\ P_{21}^{(i)} & P_{22}^{(i)} & \dots & \dots & \dots & P_{2K}^{(i)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & P_{jk}^{(i)} & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ P_{D1}^{(i)} & P_{D2}^{(i)} & \dots & \dots & \dots & P_{DK}^{(i)} \end{pmatrix}$$

In the following and for simplicity, we will refer to the price P of a particular stock $i = \alpha$ during a particular day $j = t$ and just before the bin limit of the bin B_k as $P_\alpha(k, t)$ where $P_\alpha(k, t) = P_{tk}^{(\alpha)} = P_{\alpha k}^{(t)}$. We will perform our statistical analysis over the variable $x_\alpha(k, t)$ that can be computed from the matrices above. For our interests we will be working with returns

$$x_\alpha^{(1)}(k, t) = \frac{P_\alpha(k+1, t) - P_\alpha(k, t)}{P_\alpha(k, t)},$$

and relative prices [100, 101]

$$x_\alpha^{(2)}(k, t) = \frac{P_\alpha(k, t) - P_\alpha(1, t)}{P_\alpha(1, t)}.$$

The single or collective stock dynamics is characterized by the evolution of the moments of the returns (or relative prices). Below, we show how we computed these moments [99].

VII.2.1 Single Stock Properties

The distribution of the stock α in bin k is characterized by its four first moments: mean $\mu_\alpha(k)$, standard deviation (volatility) $\sigma_\alpha(k)$, skewness $\zeta_\alpha(k)$ and kurtosis $\kappa_\alpha(k)$ defined as

$$\begin{aligned} \mu_\alpha(k) &= \langle x_\alpha(k, t) \rangle, \\ \sigma_\alpha^2(k) &= \langle x_\alpha^2(k, t) \rangle - \mu_\alpha^2(k), \\ \zeta_\alpha(k) &= \frac{6}{\sigma_\alpha(k)} (\mu_\alpha(k) - m_\alpha(k)), \\ \kappa_\alpha(k) &= 24 \left(1 - \sqrt{\frac{\pi}{2}} \frac{\langle |x_\alpha(k, t) - \mu_\alpha(k)| \rangle}{\sigma_\alpha(k)} \right) + \zeta_\alpha^2(k), \end{aligned}$$

where $m_\alpha(k)$ is the median of all values of $x_\alpha(k, t)$ and time averages for a given stock in a given bin are expressed with angled brackets $\langle \dots \rangle$.

VII.2.2 Cross-Sectional Stock Properties

The cross-sectional distributions (i.e., the dispersion of the values of the variable x of the N stocks for a given bin k in a given day t) are also characterized by the four first moments

$$\begin{aligned} \mu_d(k, t) &= [x_\alpha(k, t)], \\ \sigma_d^2(k, t) &= [x_\alpha^2(k, t)] - \mu_d^2(k, t), \\ \zeta_d(k, t) &= \frac{6}{\sigma_d(k, t)} (\mu_d(k, t) - m_d(k, t)), \\ \kappa_d(k) &= 24 \left(1 - \sqrt{\frac{\pi}{2}} \frac{[|x_\alpha(k, t) - \mu_\alpha(k)|]}{\sigma_d(k)} \right), \end{aligned}$$

where $m_d(k, t)$ is the median of all the N values of the variable x for a given (k, t) and the square brackets [...] represent averages over the ensemble of stocks in a given bin and day. If $x_\alpha(k, t)$ are the returns, $\mu_d(k, t)$ can be seen as the return of an index equi-weighted on all stocks.

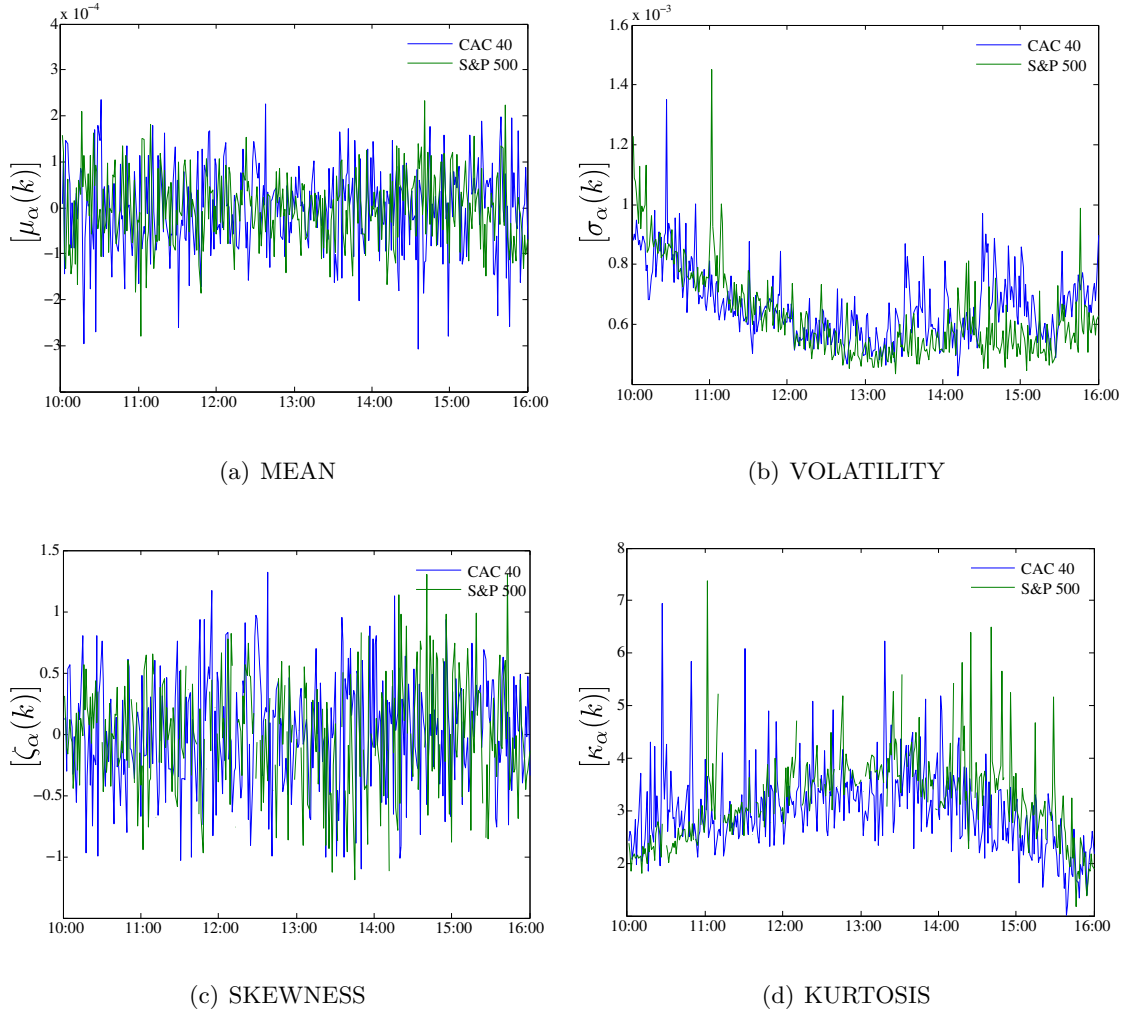


Figure VII.2: Single Stock Intra-day Seasonalities: Stock average of the single stock mean, volatility, skewness and kurtosis for the CAC 40 (blue) and the S&P 500 (green). $T = 1$.

VII.3 Intra-day Seasonalities for Returns

The following results are in complete agreement with the ones presented in [97–99].

VII.3.1 Single Stock Intra-day Seasonalities

Figure VII.2 shows the stock average of the single stock mean $[\mu_\alpha(k)]$, volatility $[\sigma_\alpha(k)]$, skewness $[\zeta_\alpha(k)]$ and kurtosis $[\kappa_\alpha(k)]$ for the CAC 40 (blue) and the S&P 500 (green), and $T = 1$ minute bin. As can be seen in Fig. VII.2(a), the mean tends to be small (in the order of 10^{-4}) and noisy around zero. The average volatility reveals the well known U-shaped pattern (Fig. VII.2(b)), high at the opening of the day, decreases during the day and increases again at the end of the day. The average skewness (Fig. VII.2(c)) is also noisy around zero. The

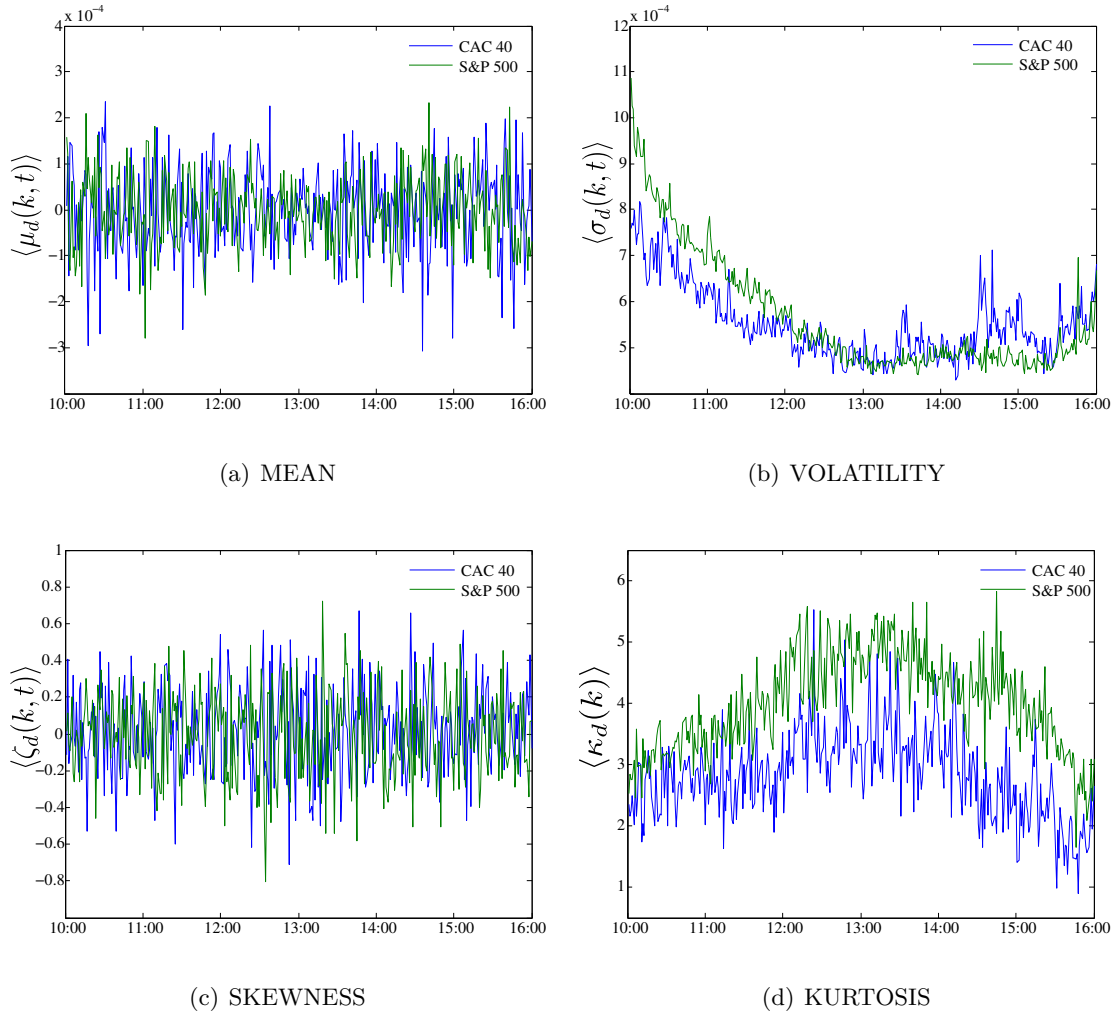


Figure VII.3: Cross-Sectional Intra-day Seasonalities: Time average of the cross-sectional mean, volatility, skewness and kurtosis for the CAC 40 (blue) and the S&P 500 (green), and $T = 1$ minute bin.

average kurtosis exhibits an inverted U-pattern (Fig. VII.2(d)), it increases from around 2 at the beginning of the day to around 4 at mid day, and decreases again during the rest of the day.

VII.3.2 Cross-Sectional Intra-day Seasonalities

As the time average of the cross sectional mean is equal to the stock average of the single stock mean, the result we show in Fig. VII.3(a) is exactly the same as the one shown in Fig. VII.2(a). The time average of the cross sectional volatility $\langle \sigma_d(k, t) \rangle$ (Fig. VII.3(b)) reveals a U-shaped pattern very similar to the stock average volatility, but less noisy (less pronounced peaks). The dispersion of stocks is stronger at the beginning of the day and

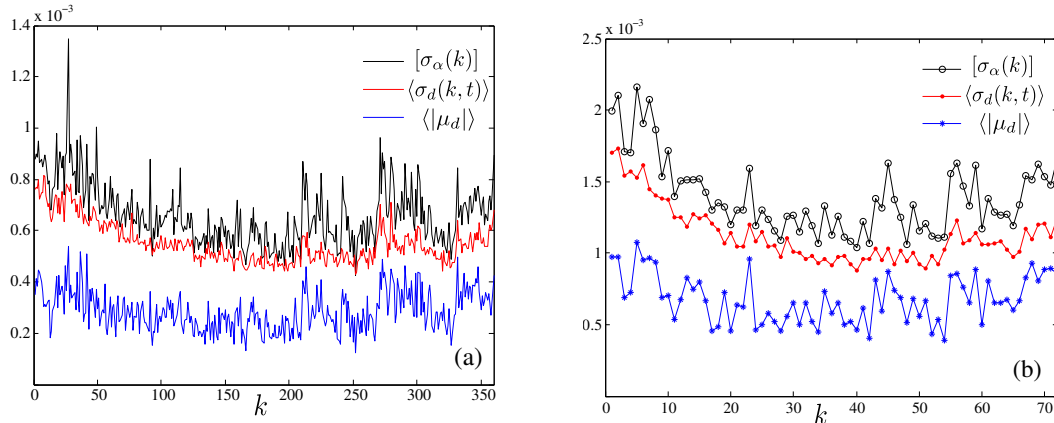


Figure VII.4: U-Pattern Volatilities: Stock average of single stock volatility $[\sigma_\alpha(k)]$ (black), time average of the cross-sectional volatility $\langle\sigma_d(k,t)\rangle$ (red) and the average absolute value of the equi-weighted index return $\langle|\mu_d|\rangle$ (blue) for the CAC 40, for (a) $T = 1$ minute bin and (b) $T = 5$ minute bin. Similar results were obtained for the S&P 500.

decreases as the day proceeds. The average skewness $\langle\zeta_d(k,t)\rangle$ is noisy around zero without any particular pattern (Fig. VII.3(c)). The cross sectional kurtosis $\langle\kappa_d(k)\rangle$ (Fig. VII.3(d)) also exhibits an inverted U-pattern as in the case of the single stock kurtosis. It increases from around 2.5 at the beginning of the day to around 4.5 at mid day, and decreases again during the rest of the day. This means that at the beginning of the day the cross-sectional distribution of returns is on average closer to Gaussian.

VII.3.3 U-Pattern Volatilities

In Fig. VII.4, we compare the stock average of single stock volatility $[\sigma_\alpha(k)]$ (black), the time average of the cross-sectional volatility $\langle\sigma_d(k,t)\rangle$ (red) and the average absolute value of the equi-weighted index return $\langle|\mu_d|\rangle$ (blue) for the CAC 40, and for $T = 1$ (left) and $T = 5$ minute bin (right). Similar results were obtained for the S&P 500. As can be seen, the average absolute value of the equi-weighted index return also exhibits a U-shaped pattern and it is a proxy for the index volatility. One thing that results interesting to observe is that the values of these volatilities actually depends of the size of the bin that we consider. For $T = 5$ minute bin, the volatilities double the values found for $T = 1$ minute bin (we will discuss this result in the next sections).

VII.3.4 Intra-day Seasonalities in the Stock Correlation

In order to compute the correlation between stocks, we first normalize the returns by the dispersion of the corresponding bin [99] i.e.,

$$\hat{x}_\alpha(k,t) = x_\alpha^{(1)}(k,t)/\sigma_d(k,t)$$

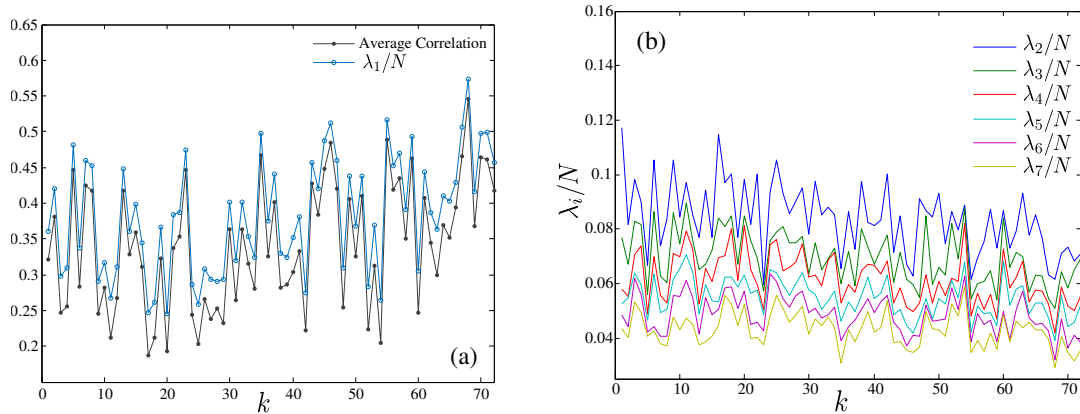


Figure VII.5: Largest eigenvalues structure for the CAC 40, $T = 5$ minute bin. **(a)** Average correlation between stocks (black) and top eigenvalue λ_1/N (blue) of the correlation matrix $C_{\alpha\beta}(k)$. **(b)** Smaller eigenvalues.

The $N \times N$ correlation matrix for a given bin k would be given by

$$C_{\alpha\beta}(k) = \frac{\langle \hat{x}_\alpha(k, t) \hat{x}_\beta(k, t) \rangle - \langle \hat{x}_\alpha(k, t) \rangle \langle \hat{x}_\beta(k, t) \rangle}{\sigma_\alpha(k) \sigma_\beta(k)}.$$

In Fig. VII.5(a) we show the average correlation between stocks (blue) and top eigenvalue λ_1/N (green) for the CAC 40. As can be seen the largest eigenvalue is a measure of the average correlation between stocks [99, 269–272]. This average correlation increases during the day from a value around 0.35 to a value around 0.45 when the market closes. For the case of smaller eigenvalues, what we can see is that the amplitude of risk factors decreases during the day (Fig. VII.5(b)), as more and more risk is carried by the market factor (Fig. VII.5(a)) [99].

In order to simplify the computation of the N^2 correlation matrices for each bin k in the case of the S&P 500, we computed the correlation matrix $C_{\alpha\beta}$ for 4 different sets of stocks: r_0 : composed by the 100 first stocks of the S&P 500; $r_{1,2}$: composed by 100 stocks randomly picked; and r_3 : composed by 200 stocks randomly picked. Figure VII.6(a) shows $\frac{\lambda_1}{N}$ as function of the bins. Although the values of the eigenvalues seem to be out of scale, it can be seen clearly that the average correlation increases during the day. This scale conflict is solved by normalizing the value of the top eigenvalue not by N but by the sample size N_0 (i.e., 100 or 200) (Fig. VII.6(b)). As can be seen the average correlation of the index can be computed by taking a subset of it which means that actually just the more capitalized stocks in the index drive the rest of stocks.

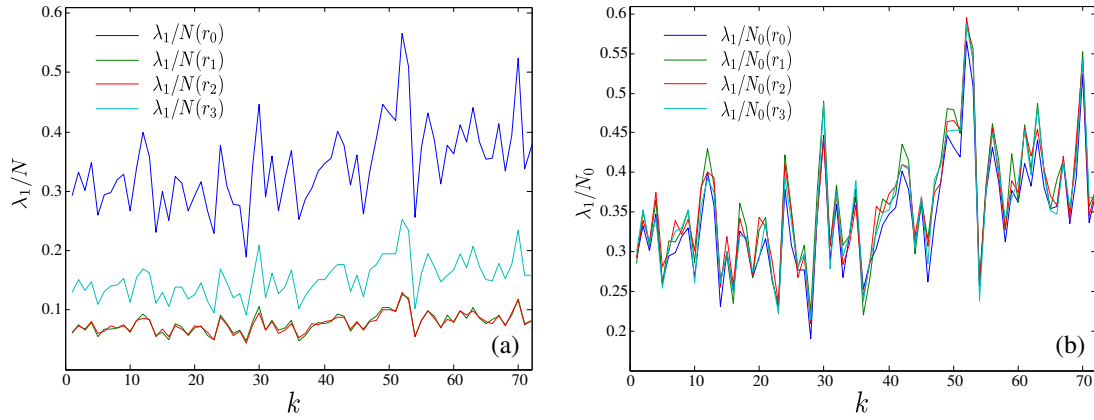


Figure VII.6: **(a)** Top eigenvalue λ_1/N and **(b)** λ_1/N_0 for the S&P 500 for 4 different sets of stocks: r_0 (blue), r_1 (green), r_2 (red) and r_3 (clear blue). $T = 5$ minute bin.

VII.4 Intra-day Seasonalities for Relative Prices

In this section, we will report the results we found for the S&P 500. Similar results were found also for the CAC 40. We will see how in the case of the relative prices these intra-day seasonalities are independent of the size of the bin, also independent of the index we consider (but characteristic for each index) however this is not the case for the returns.

VII.4.1 Single Stock Intra-day Seasonalities

Each path in Fig. VII.7 represents the evolution of a particular moment of one of the stocks that compose the S&P 500 (i.e., one path, one stock moment). The stock average of the single stock mean $[\mu_\alpha(k)]$, volatility $[\sigma_\alpha(k)]$, skewness $[\zeta_\alpha(k)]$ and kurtosis $[\kappa_\alpha(k)]$ of the S&P 500 are shown in black. The stock average of the single stock mean varies around zero. The average volatility increases logarithmically with time. The skewness varies between $[-3, 3]$ with an average value of zero. The single stock kurtosis takes values between $[-2, 6]$ with an average value of one and its stock average starts from a value around 2 in the very beginning of the day and decreases quickly to the mean value 1 in the first minutes of the day.

VII.4.2 Cross-Sectional Intra-day Seasonalities

Each path in Fig. VII.8 represents the evolution of a particular index moment during a particular day (i.e., one path, one day moment). As in the case of the single stock volatility, the cross-sectional dispersion $\langle\sigma_d(k)\rangle$ increases logarithmically with respect to the time (Fig. VII.8(b)). The cross-sectional skewness $\langle\zeta_d(k)\rangle$ takes values in the interval $[-1, 1]$ with an average value of zero (Fig. VII.8(c)). The average kurtosis $\langle\kappa_d(k)\rangle$ starts from a value around 2.5 in the very beginning of the day and decreases quickly to the mean value 2 in the first minutes of the day (Fig. VII.8(d)).

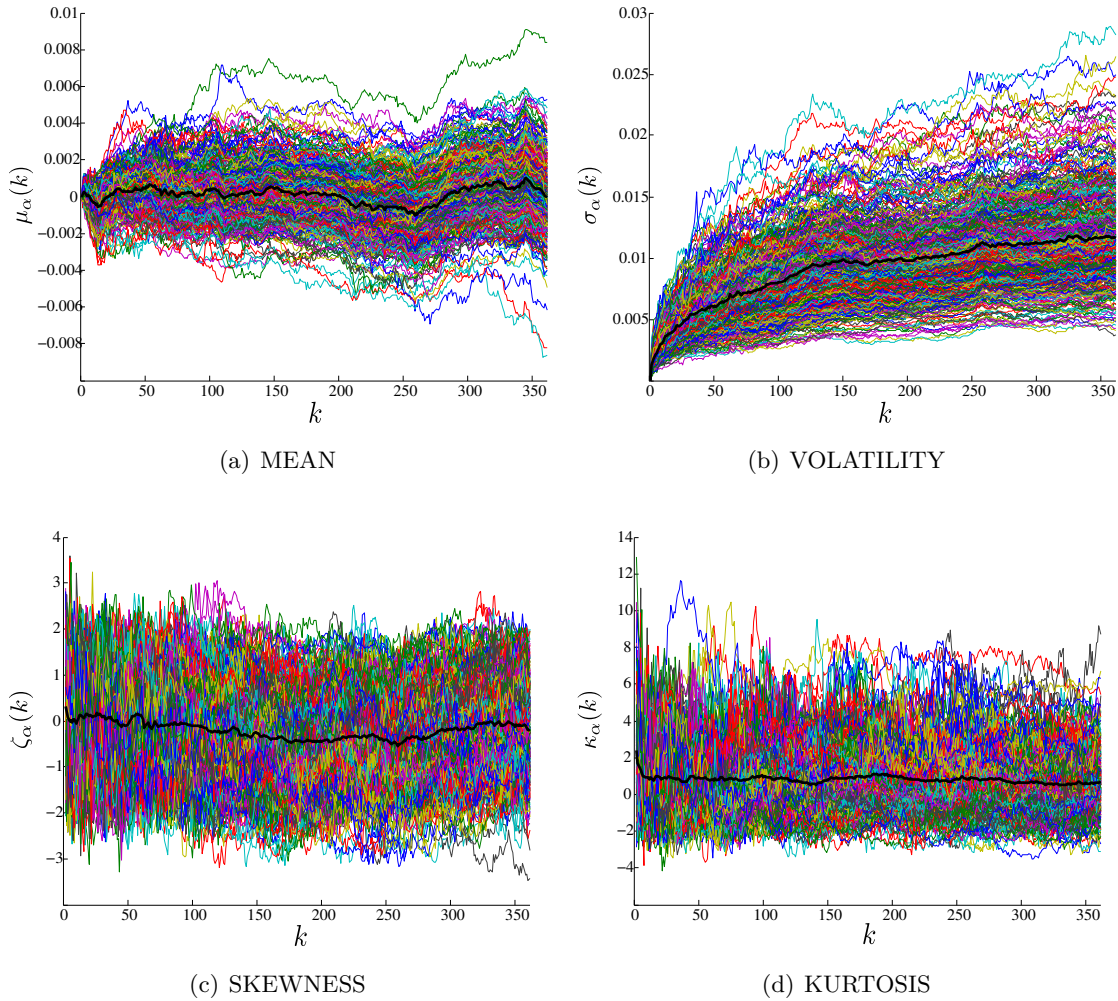


Figure VII.7: Single Stock Intra-day Seasonalities: Stock average of the single stock mean, volatility, skewness and kurtosis for the S&P 500 (black). $T = 1$ minute bin.

VII.4.3 C-Pattern Volatilities

Similarly as we did in Sec. VII.3.3 for returns, in Fig. VII.9 we show a comparative plot between the stock average of the single stock volatility $[\sigma_\alpha(k)]$, the time average of the cross-sectional volatility $\langle \sigma_d(k, t) \rangle$ and the average absolute value of the cross-sectional mean $\langle |\mu_d| \rangle$ for the relative prices of the S&P 500, and for $T = 1$ and $T = 5$ minute bin. As can be seen, these three measures exhibit the same kind of intra-day pattern (as it did in the case of the returns). But the most important fact is to notice that this intra-day seasonality is independent of the size of the bin, also independent of the index we consider, but characteristic for each index (see inset Fig. VII.9).

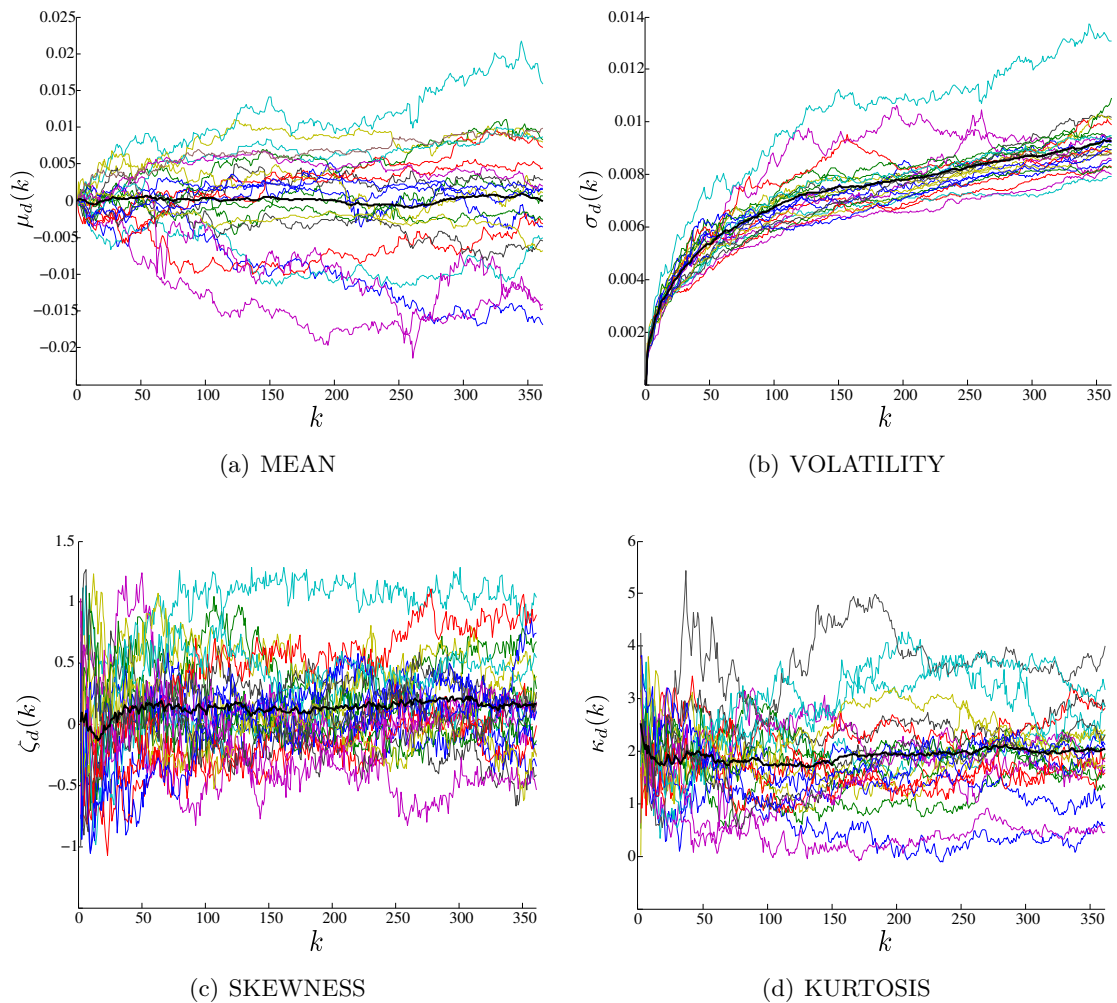


Figure VII.8: Cross-Sectional Intra-day Seasonalities: Time average of the cross-sectional mean, volatility, skewness and kurtosis for the S&P 500 (black). $T = 1$ minute bin.

VII.5 Intra-day Patterns and Bin Size

As we saw in the last section, the volatilities for the relative prices exhibit the same kind of intra-day pattern (Fig. VII.9). This intra-day seasonality is independent of the size of the bin, and the index we consider, but characteristic for each index. Actually, this is not true in the case of the returns as we already suggested in Sec. VII.3.3 from Fig. VII.4. If we consider the odd moments (mean and skewness) of the returns, the behavior is basically the same (noisy around zero) and without any particular pattern, independently of the bin size (as can be seen in Figs. VII.2, VII.3 and VII.10). But for the case of the even moments of returns, although they exhibit the well known U and inverted U-patterns, these patterns depend on the bin size. This fact is well illustrated through Figs. VII.11 and VII.12 where we have chosen 5 different values of bin size from $T = 0.5$ to $T = 10$ minutes. In these figures

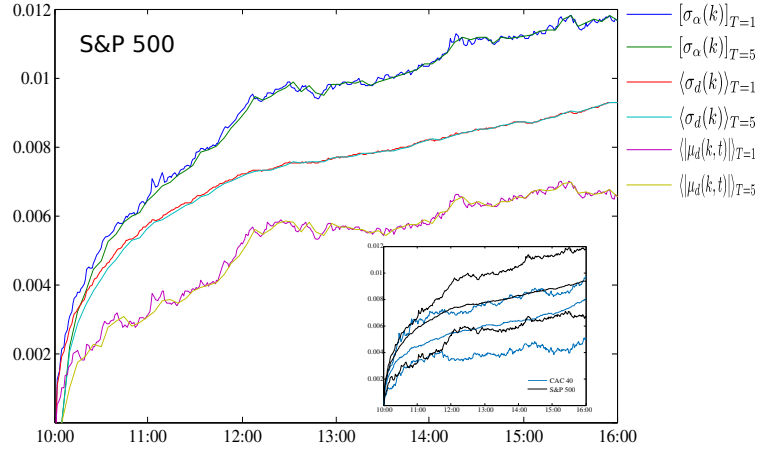


Figure VII.9: C-Pattern Volatilities: Stock average of the single stock volatility $[\sigma_\alpha(k)]$, time average of the cross-sectional volatility $\langle\sigma_d(k,t)\rangle$ and the average absolute value of the cross-sectional mean $\langle|\mu_d|\rangle$ for the relative prices of the S&P 500. $T = 1$ and $T = 5$ minute bin. Inset: CAC 40 (blue) and S&P 500 (black).

we show the time average of the cross-sectional volatility and kurtosis for the S&P 500 but a similar bin size dependence can be shown for the CAC 40 or any other index and also for the time average of the single stock volatility and kurtosis. By other hand the kurtosis is a decreasing function of the size of the bin and the inverted U-pattern is evident just when we consider “small” bin sizes, in our case this occurs for $T = 1$ and $T = 0.5$ minute bin (Fig. VII.12). This represents a confirmation that on small scales the returns have heavier tails, and on long time scales they are more Gaussian [96, 102–104].

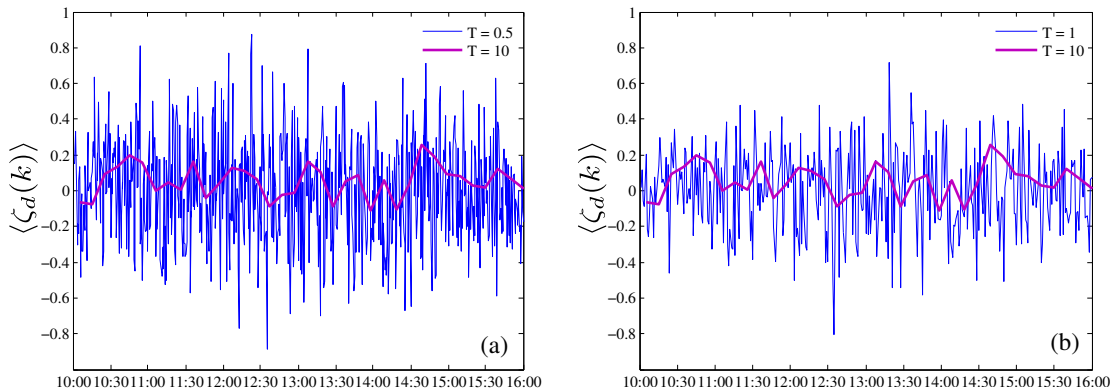


Figure VII.10: Time average of the cross-sectional skewness: Comparison of the intra-day patterns for (a) $T = 0.5$ and (b) $T = 1$ against $T = 10$ minute bin for the S&P 500.

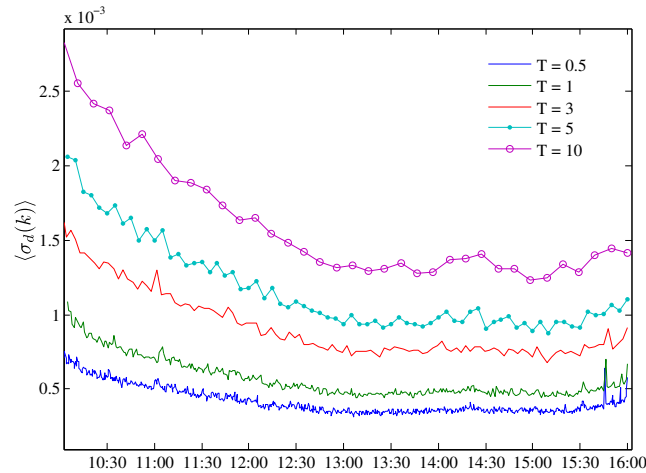


Figure VII.11: Bin size dependence in the U-pattern volatilities: Time average of the cross-sectional volatility for the S&P 500 for 5 different values of bin size T .

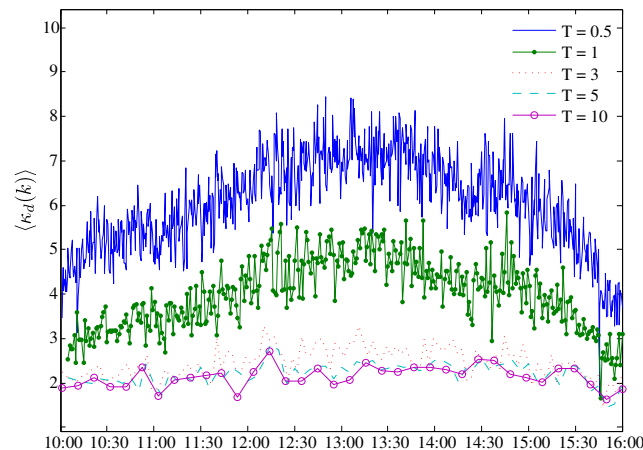


Figure VII.12: Bin size dependence in the inverted U-pattern kurtosis: Time average of the cross-sectional kurtosis for the S&P 500 for 5 different values of bin size T .

VII.6 Intra-day Abnormal Patterns

One of the motivations to explore into the intra-day seasonalities for relative prices was due to Kaisoji's previous work [101]. In his work he found that the upper tail of the complementary cumulative distribution function of the ensemble of the relative prices in the high value of the price is well described by a power-law distribution which when its exponent approached two, the Japan's internet bubble burst. Taking into consideration our recent findings we suggest the use of the bin size independence for intra-day patterns in relative prices in order to characterize "atypical days" for indexes and "anomalous behaviors" for stocks.

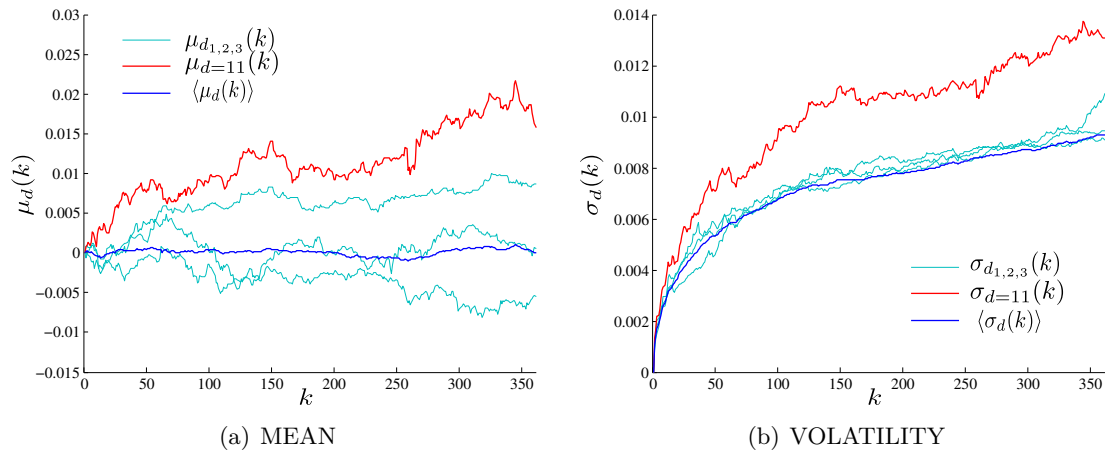


Figure VII.13: S&P 500 Atypical Day: Time average of the cross-sectional mean and volatility (blue), cross-sectional mean and volatility of the S&P 500 during day 11 (red) and during three days chosen at random (clear blue).

The time average of the cross-sectional moments represents the average behavior of a particular index moment during an average day. In Fig. VII.8 each path represents the evolution of a particular index moment for one of the days of the period under analysis (i.e., one path, one day moment). If we look directly into the prices of the CAC 40 and S&P 500, we can observe during day 11 a fall of the prices of the stocks that compose both indexes. During the days before and following day 11, the (index) moments move along our intra-day pattern. Moreover, if we pick randomly one day from our period of analysis, in most of the cases our index during that day will behave as our intra-day seasonality (as in Fig. VII.13), but the one for day 11 will not. In Fig. VII.13 we show the (cross-sectional) intra-day seasonalities for the (a) mean and (b) volatility in blue and in clear blue the respective cross-sectional stock moments for 3 days randomly picked. The average behavior (of the moments) of our index during these days moves along with our intra-day pattern. This is not the case of the curve corresponding to the day 11 shown in red which clearly diverges from the expected behavior. This is what could be called as an “atypical day” for the S&P 500.

We could use the same reasoning as before in order to characterize “anomalous behaviors” in stocks. Each path in Fig. VII.7 represents the average evolution of a particular moment of one of the stocks that compose the S&P 500. The stock average of those single stock moments represents the average behavior of that moment for an average stock during an average day of our period of analysis. Meaning that if we pick randomly one stock from our set of stocks, in most of the cases (its moments) will behave as the intra-day seasonality. This is clearly illustrated in Fig. VII.14 where we present the intra-day seasonalities for the (a) mean and (b) volatility in blue and the respective single stock moments for 3 stocks randomly picked in clear blue. As can be seen, the average behavior of the moments of these stocks move along with our intra-day patterns. However this is not the case for the curves shown in red which have been chosen on purpose to illustrate how in this case the stock 228 behaves in an anomalous way with respect to what is expected.

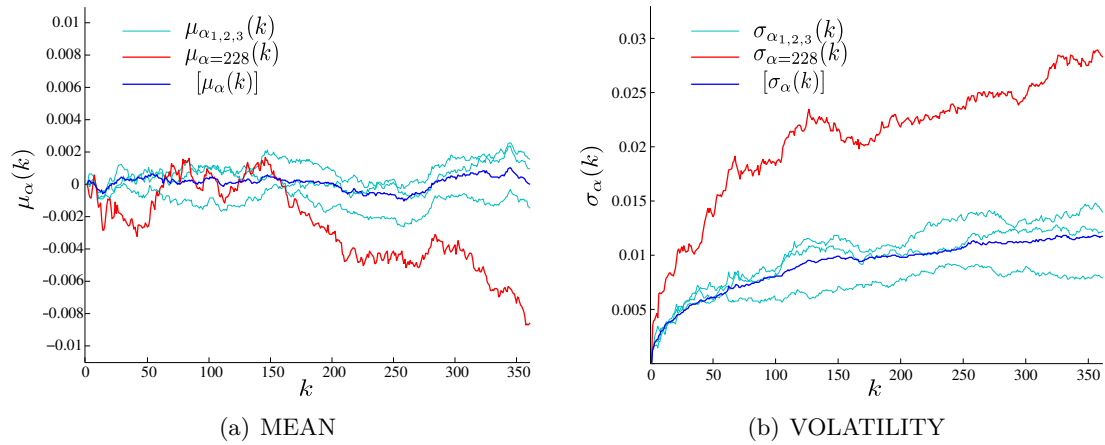


Figure VII.14: Anomalous Stock Behavior: Stock average of the single stock mean and volatility (blue), single stock mean and volatility of stock 228 (red) and from three stocks chosen at random from the S&P 500 (clear blue).

VII.7 Discussion

In this chapter, we have analyzed the intra-day seasonalities of the single and cross-sectional (or collective) stock dynamics by the evolution of the moments of its returns (and relative prices) during a typical day. What we have called “single stock intra-day seasonalities” is the average behavior of the moments of the returns (and relative prices) of an average stock in an average day. In the same way, the cross-sectional intra-day seasonality is not more than the average day behavior of an index moment. We presented these intra-day seasonalities for returns (Figs. VII.2 and VII.3) and relative prices (Figs. VII.7 and VII.8) and compared the stock average of single stock volatility $[\sigma_\alpha(k)]$, the time average of the cross-sectional volatility $\langle\sigma_d(k,t)\rangle$ and the average absolute value of the equi-weighted index $\langle|\mu_d|\rangle$ (Figs. VII.4 and VII.9).

Notably, in the case of the returns, is that these “patterns” actually depend on the size of the bin. This fact was well illustrated with 5 different values of bin size through Fig. VII.11 for volatilities and Fig. VII.12 for kurtosis in which its inverted U-pattern is evident just when we consider “small” bin sizes.

In the case of relative prices, the volatilities also exhibit the same kind of intra-day pattern (Fig. VII.9), but contrary with the returns, it is independent of the size of the bin, and the index we consider, but characteristic for each index. We suggested in Sec. VII.6 how this bin size independence of intra-day patterns in relative prices could be used in order to characterize “atypical days” for indexes and “anomalous behaviors” in stocks. This was presented in Figs. VII.13 and VII.14 where we presented our intra-day seasonalities for the (a) mean and (b) volatility in blue and the respective the cross-sectional moments for 3 days (and the single stock moments for 3 stocks) randomly picked in clear blue and we saw how the average behavior of their moments move along with our intra-day patterns which was not the case for the day 11 and the stock 228.

VIII – Conclusion

In this thesis, we studied issues that arise from the evaluation of the large deviation function (LDF) from population dynamics algorithms. Different versions of the cloning algorithm were used which differ among them in crucial aspects as the way in which the selection mechanism is performed or on the restriction in the growth of the total population of copies of the system. The LDF behaviour for these different versions and its related features were analyzed. We gave particular attention to the dependence of the estimator with number of clones N_c and the simulation time t (the two parameters introduced by the method) by studying the finite- t and finite- N_c effects, its convergence in the infinite- t and infinite- N_c limit as well as its behavior in the large system size L limit. Moreover, different ways and methods to improve the LDF estimation were proposed.

In chapter II [P1] using a non-constant population approach of the cloning algorithm, we analyzed the discreteness effects at initial times in population dynamics. We show how these effects play an important role in the determination of the large deviation function which may be obtained from the growth rate of an average log-population. Fluctuations at initial times produce that some populations remain in their initial states for much longer than others, producing a gap in their individual evolution. This induces a relative shift between populations that lasts forever which supplemented with a short-time evolution affect strongly the average population and thus the LDF estimation. We argue in Sec. II.4.1 that these lags between populations could be compensated by performing a time translation (Eq. (II.2)) over populations in order to emphasize the effects of the exponential growth regime. This along with a discarding of initial regimes in the evolution of the population surpasses the influence of initial discreteness effects.

The finite- t and finite- N_c scalings in the evaluation of large deviation functions were studied in chapters III [P2], IV [P3] and V [P3] following two different approaches: an analytical one, in chapter III, using a discrete-time version of the population dynamics algorithm [18], and a numerical one, in chapters IV and V, using a continuous-time version [17, 19]. In both cases, we derived that the deviations of large deviation estimator from the desired value (which we called systematic errors) were small and behaved as $1/t$ and $1/N_c$ in the large- t and large- N_c asymptotics respectively. Importantly, in chapter IV [P3] we showed the validity of these results in more complex systems. Such scalings also provided a convergence criterion to the asymptotic regimes of the algorithm: In order to ensure a correct LDF evaluation, one has to confirm that the LDF estimator does present corrections (first) in $1/t$ and (second) in $1/N_c$ with respect to an asymptotic value. We discussed in Secs. III.2.4.2 and IV.5 how these two versions differ on a crucial point which makes that an extension of the analysis developed in chapter III cannot be done straightforwardly in order to comprehend the continuous-time case in chapter IV and thus the observation of these scalings themselves

is also non-trivial. This finite- t and finite- N_c scaling behavior was used in chapter IV [P3] in order to interpolate the large- t and large- N_c asymptotic value of the LDF estimator from the measured values for finite and small t and N_c . This allowed us to propose an improved version of the continuous-time cloning algorithm in Sec. IV.4.1 providing more reliable results, less affected by finite- t and $-N_c$ effects. We demonstrated numerically that the interpolation technique is very efficient, by a direct comparison of the resulting LDF estimation with the analytical value, which can be determined in the studied system. However, the validity of the method and of these scalings were proved only for a simple one-site annihilation-creation dynamics and for a contact process with $L = 6$ sites, leaving an analysis of the dependence with the system size (number of sites) L pending.

In order to prove whether the finite- t and $-N_c$ scalings observed in small (number of sites L) systems are also valid in the large- L limit, we redefined these scalings in a more general way in chapter VI [P4]. We assumed a $t^{-\gamma_t}$ (VI.1) and a $N_c^{-\gamma_{N_c}}$ (VI.2) -scaling behavior for the LDF estimator. This redefinition allowed us to verify in large- L systems if effectively $\gamma_t \approx 1$ and $\gamma_{N_c} \approx 1$ and whether the extracted quantities from the application of the scaling method represented the limits in $t \rightarrow \infty$ and $N_c \rightarrow \infty$. First, we considered a contact process with $L = 100$ sites and two representative values of the parameter s . Although the t^{-1} -scaling and N_c^{-1} -scaling were proved to hold for $s < 0$, this was not the case for $s > 0$, being this fact valid in general for large- L systems. As the scaling method relied on the validity of the t^{-1} - and N_c^{-1} -scalings, in Sec. VI.3.3 we showed how the determination of the infinite- t and infinite- N_c limit of the LDF estimator is affected. In order to have a clear picture of the change in the scalings of the LDF estimator, the analysis was extended to the plane $s - L$ where the exponents γ_t and γ_{N_c} were computed and characterized for a grid of values of the parameters (s, L) . Moreover, we discussed how this breakdown in the scalings in the large- L limit could be related to the dynamical phase transition of the contact process.

Although our study on the cloning algorithm is closed in chapter VI [P4], the study of rare events is complemented with chapter VII [P0] using a completely different approach. This is the empirical study of the patterns that hide behind financial time series, known as stylized facts. We analyzed the intra-day seasonalities of the single and cross-sectional (or collective) stock dynamics by characterizing the dynamics of a stock (or a set of stocks) by the evolution of the moments of its returns (and relative prices) during a typical day. We showed how these patterns actually depend on the size of the bin in the case of the returns. However, in the case of relative prices, these patterns are independent of the size of the bin, and the index we consider, but characteristic for each index. We suggested in Sec. VII.6 how this bin size independence of intra-day patterns in relative prices could be used in order to characterize “atypical days” for indexes and “anomalous behaviours” in stocks.

IX – Perspectives

Below we mention some questions that arose from our study which remain open and may constitute possible directions for future research.

The analysis of the discreteness effects at initial times in population dynamics developed in chapter II [P1] (using a non-constant population approach of the cloning algorithm), was performed only on a simple system: a one-site annihilation-creation dynamics (Sec. I.8.1). However we hope it can be extended to more complex phenomena so that our results can be verified or else, more interesting features can be found. Nevertheless even for that simple system there remain pending issues. Some of them related to the fact that the duration of the initial discrete-population regime could be understood from an analytical study of the population dynamics itself. On the other hand, the results presented support a power-law behaviour in time of the variance of the delays. Additionally, the distribution of the delays was found to take an universal form, after rescaling the variance. Both of which could be explored deeply.

From a constant population approach, as the one used in chapters III, IV, V and VI, is still possible to reconstruct the evolution in time of the population of clones. Thus, it would result interesting to compare both approaches but importantly, the properties of the reconstructed populations in contrast with actual populations (obtained from a non-constant population as in chapter II).

From the analytical study of the finite- t and $-N_c$ scalings of the LDF estimator developed in chapter III [P2], we mention two open questions. The first is related to the precise estimate of the error due to a non-infinitesimal time interval Δt between cloning steps: As explained in Sec. III.2.4.1 and Sec. III.2.4.2, taking the $\Delta t \rightarrow 0$ limit is important in our analysis, in order to make the estimator converge to the correct LDF. From a practical point of view, taking this limit can however be problematic, since it requires infinitely many cloning procedures per unit time (as $\Delta t \rightarrow 0$). Interestingly, most of existing algorithms do not take such a limit (for instance the original version of the algorithm, Ref. [18]). Empirically, one thus expects that the error goes to zero as $N_c \rightarrow \infty$ while keeping Δt finite. Within the method developed in chapter III [P2] the analytical estimation of this error is challenging (see Sec. III.2.4.2) and remains an open problem.

The second question is related with possible extensions of the formulation developed in chapter III [P2]. As the cloning procedure is performed for a fixed time interval, the formulation cannot cover the case of algorithms where Δt itself is statistically distributed, as in continuous-time cloning algorithms [19]. Moreover, the formulation is limited to Markov systems, although population dynamics algorithms are applied to chaotic deterministic dynamics [33, 70] or to non-Markovian evolutions [273]. Once one removes the Markov condition in the dynamics, developing analytical approaches becomes more challenging. However, as the

physics of those systems are important scientifically and industrially [274], the understanding of such dynamics cannot be avoided for the further development of population algorithms.

Results evident to question whether the numerical study developed in (specially) chapter IV [P3] can be extended to systems presenting dynamical phase transition (DPT) in the form of a non-analyticity of the LDF. In particular, in this context, it would be useful to understand how the dynamical phase transition of the original system translates into anomalous features of the distribution of the LDF estimator in the cloning algorithm. Although the system used in chapter IV [P3] (the contact process, Sec. I.8.2) is known to exhibit a DPT in the $L \rightarrow \infty$ limit [19, 35, 74, 252], the finite- t and $-N_c$ scalings of the LDF estimator were studied on a small system with $L = 6$ sites, for which the effects of the DPT cannot be observed. In chapter VI [P4] we extended our analysis to a large- L contact process ($L = 100$ sites, where DPT manifests [19, 258]) showing evidence of a changing in the LDF scalings with the size of the system L which could be related to a DPT. However, a study of the DPT effects would require a large- N_c and $-t$ configuration, which under our approach was a task not possible to fulfill (as the main objective in chapters IV [P3] and VI [P4] was the possibility of extracting the infinite- N_c infinite- t limit of the LDF estimator from data for a small number of clones N_c and time t).

Taking in consideration that it is well known that the existing methods [12, 13, 18, 265, 266] perform poorly in the vicinity of a dynamical phase transition, or they are numerically expensive in order to obtain accurate estimations [266–268] developing if not important finite-size effects [22], the analysis of this problem in the large- t and $-N_c$ limits is not necessarily the best option. Recently has been proposed a promising method [85, 259] which combines the existing cloning algorithm [7, 12, 13, 17–19, 265, 266] with a modification of the dynamics [88–91, 275–277] resulting in a significant improvement of its computational efficiency. The method was successfully applied to the study of the dynamical phase transition of 1D FA model [42] using a relatively small N_c and L . The implementation of this method will provide in a next stage a clear contrast between the results obtained following the two different approaches and a better understanding of their limitations and advantages.

Publications

- [P0] E. Guevara Hidalgo, Bin Size Independence in Intra-day Seasonalities for Relative Prices, *Physica A* **468** 722–732 (2017).
- [P1] E. Guevara Hidalgo, Vivien Lecomte, Discreteness Effects in Population Dynamics, *J. Phys. A: Math. Theor.* **49** 205002 (2016).
- [P2] Takahiro Nemoto, E. Guevara Hidalgo, Vivien Lecomte, Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: Analytical Study using a Birth-Death Process, *Phys. Rev. E* **95** 012102 (2017).
- [P3] E. Guevara Hidalgo, Takahiro Nemoto, Vivien Lecomte, Finite-Time and -Size Scalings in the Evaluation of Large Deviation Functions: Numerical Analysis in Continuous Time, *Phys. Rev. E* **95** 062134 (2017).
- [P4] E. Guevara Hidalgo, Breakdown of the Finite-Time and -Population Scalings of the Large-Deviation Function in the Large-Size Limit of a Contact Process, under review, [arXiv:1709.09322](https://arxiv.org/abs/1709.09322) (2017).

8

Bibliography

- [1] H. Touchette, *Physics Reports* **478**, 1 (2009).
- [2] H. Touchette and R. J. Harris, “Large deviation approach to nonequilibrium systems,” in *Nonequilibrium Statistical Physics of Small Systems* (Wiley-VCH Verlag GmbH Co. KGaA, 2013) pp. 335–360.
- [3] H. Touchette, [arxiv:1106.4146](https://arxiv.org/abs/1106.4146) (2011).
- [4] S. R. S. Varadhan, *Science* **247**, 1351 (1990).
- [5] T. Bodineau and B. Derrida, *Phys. Rev. E* **72**, 066110 (2005).
- [6] J. Mehl, T. Speck, and U. Seifert, *Phys. Rev. E* **78**, 011123 (2008).
- [7] C. Giardinà, J. Kurchan, V. Lecomte, and J. Tailleur, *J. Stat. Phys.* **145**, 787 (2011).
- [8] J. Bucklew, *Introduction to Rare Event Simulation* (Springer Science & Business Media, 2013).
- [9] H. Kahn and T. E. Harris, National Bureau of Standards applied mathematics series **12**, 27 (1951).
- [10] F. Cérou and A. Guyader, *Stochastic Analysis and Applications* **25**, 417 (2007).
- [11] W. G. Cochran, *Sampling Techniques* (John Wiley, 1977).
- [12] L. O. Hedges, R. L. Jack, J. P. Garrahan, and D. Chandler, *Science* **323**, 1309 (2009).
- [13] T. Speck, A. Malins, and C. P. Royall, *Phys. Rev. Lett.* **109**, 195703 (2012).
- [14] P. G. Bolhuis, D. Chandler, C. Dellago, and P. L. Geissler, *Annual Review of Physical Chemistry* **53**, 291 (2002).
- [15] D. Aldous and U. Vazirani, in *Foundations of Computer Science, 35th* (IEEE, 1994) pp. 492–501.
- [16] P. Grassberger, *Computer Physics Communications* **147**, 64 (2002), proceedings of the Europhysics Conference on Computational Physics Computational Modeling and Simulation of Complex Systems.
- [17] J. Tailleur and V. Lecomte, *AIP Conf. Proc.* **1091**, 212 (2009).

- [18] C. Giardinà, J. Kurchan, and L. Peliti, *Phys. Rev. Lett.* **96**, 120603 (2006).
- [19] V. Lecomte and J. Tailleur, *J. Stat. Mech.* **2007**, P03004 (2007).
- [20] P. Del Moral, A. Doucet, and A. Jasra, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**, 411 (2006).
- [21] P. I. Hurtado and P. L. Garrido, *Phys. Rev. Lett.* **102**, 250601 (2009).
- [22] P. I. Hurtado and P. L. Garrido, *J. Stat. Mech.* **2009**, P02032 (2009).
- [23] F. Cérou, A. Guyader, T. Lelièvre, and D. Pommier, *The Journal of Chemical Physics* **134**, 054108 (2011).
- [24] W. E and E. Vanden-Eijnden, *Journal of Statistical Physics* **123**, 503 (2006).
- [25] B. Derrida and J. L. Lebowitz, *Phys. Rev. Lett.* **80**, 209 (1998).
- [26] B. Derrida, *Journal of Statistical Mechanics: Theory and Experiment* **2011**, P01030 (2011).
- [27] L. Bertini, A. De Sole, D. Gabrielli, G. Jona-Lasinio, and C. Landim, *Rev. Mod. Phys.* **87**, 593 (2015).
- [28] J. K. Weber, R. L. Jack, and V. S. Pande, *Journal of the American Chemical Society* **135**, 5501 (2013).
- [29] S. Vaikuntanathan, T. R. Gingrich, and P. L. Geissler, *Phys. Rev. E* **89**, 062108 (2014).
- [30] J. K. Weber, D. Shukla, and V. S. Pande, *Proceedings of the National Academy of Sciences* **112**, 10377 (2015).
- [31] R. Chetrite and H. Touchette, *Annales Henri Poincaré* **16**, 2005 (2015).
- [32] J. B. Anderson, *The Journal of Chemical Physics* **63**, 1499 (1975).
- [33] J. Tailleur and J. Kurchan, *Nature Phys* **3**, 203 (2007).
- [34] M. Baiesi, C. Maes, and K. Netočný, *Journal of Statistical Physics* **135**, 57 (2009).
- [35] V. Lecomte, C. Appert-Rolland, and F. van Wijland, *Journal of Statistical Physics* **127**, 51 (2007).
- [36] H. Spohn, *Large Scale Dynamics of Interacting Particles* (Springer Verlag, Heidelberg, 1991).
- [37] H. Spohn, *Journal of Physics A: Mathematical and General* **16**, 4275 (1983).
- [38] T. E. Harris, *Ann. Probability* **2**, 969 (1974).
- [39] J. P. Garrahan, R. L. Jack, V. Lecomte, E. Pitard, K. van Duijvendijk, and F. van Wijland, *Phys. Rev. Lett.* **98**, 195702 (2007).

- [40] J. P. Garrahan, R. L. Jack, V. Lecomte, E. Pitard, K. van Duijvendijk, and F. van Wijland, *J. Phys. A* **42**, 075007 (2009).
- [41] F. Ritort and P. Sollich, *Advances in Physics* **52**, 219 (2003).
- [42] G. H. Fredrickson and H. C. Andersen, *Phys. Rev. Lett.* **53**, 1244 (1984).
- [43] J. Jäckle and S. Eisinger, *Zeitschrift für Physik B Condensed Matter* **84**, 115 (1991).
- [44] J. Jäckle and A. Kronig, *Journal of Physics: Condensed Matter* **6**, 7633 (1994).
- [45] W. Kob and H. C. Andersen, *Phys. Rev. E* **48**, 4364 (1993).
- [46] A. Kronig and J. Jäckle, *Journal of Physics: Condensed Matter* **6**, 7655 (1994).
- [47] J. Kurchan, L. Peliti, and M. Sellitto, *EPL (Europhysics Letters)* **39**, 365 (1997).
- [48] P. Sollich and M. Evans, *Phys. Rev. Lett.* **83**, 3238 (1999).
- [49] M. Einax and M. Schulz, *The Journal of Chemical Physics* **115**, 2282 (2001).
- [50] J. P. Garrahan and D. Chandler, *Phys. Rev. Lett.* **89**, 035704 (2002).
- [51] D. Aldous and P. Diaconis, *Journal of Statistical Physics* **107**, 945 (2002).
- [52] Y. Jung, J. P. Garrahan, and D. Chandler, *Phys. Rev. E* **69**, 061205 (2004).
- [53] C. Toninelli, G. Biroli, and D. S. Fisher, *Journal of Statistical Physics* **120**, 167 (2005).
- [54] A. C. Pan, J. P. Garrahan, and D. Chandler, *Phys. Rev. E* **72**, 041106 (2005).
- [55] P. L. Geissler and D. R. Reichman, *Phys. Rev. E* **71**, 031206 (2005).
- [56] M. D. Ediger, C. A. Angell, and S. R. Nagel, *The Journal of Physical Chemistry* **100**, 13200 (1996).
- [57] C. A. Angell, *Science* **267**, 1924 (1995).
- [58] K. Binder and W. Kob, *Glassy materials and disordered solids: An introduction to their statistical mechanics* (World Scientific, 2011).
- [59] P. I. Hurtado, C. Pérez-Espigares, J. J. del Pozo, and P. L. Garrido, *Proceedings of the National Academy of Sciences* **108**, 7704 (2011).
- [60] P. I. Hurtado and P. L. Garrido, *Phys. Rev. E* **81**, 041102 (2010).
- [61] P. I. Hurtado and P. L. Garrido, *Phys. Rev. Lett.* **107**, 180601 (2011).
- [62] P. Glasserman, P. Heidelberger, P. Shahabuddin, and T. Zajic, *Proceedings of the 28th Conference on Winter Simulation*, WSC '96 (IEEE Computer Society, Washington, DC, USA, 1996) pp. 302–308.
- [63] Y. Iba, *Transactions of the Japanese Society for Artificial Intelligence* **16**, 279 (2001).

- [64] P. L'Ecuyer, V. Demers, and B. Tuffin, in *Proceedings of the 2006 Winter Simulation Conference* (2006) pp. 137–148.
- [65] O. Cappé, A. Guillin, J. M. Marin, and C. P. Robert, *Journal of Computational and Graphical Statistics* **13**, 907 (2004).
- [66] T. S. van Erp, D. Moroni, and P. G. Bolhuis, *The Journal of Chemical Physics* **118**, 7762 (2003).
- [67] R. J. Allen, P. B. Warren, and P. R. ten Wolde, *Phys. Rev. Lett.* **94**, 018104 (2005).
- [68] P. Del Moral and J. Garnier, *Ann. Appl. Probab.* **15**, 2496 (2005).
- [69] T. Dean and P. Dupuis, *Stochastic Processes and their Applications* **119**, 562 (2009).
- [70] T. Laffargue, K.-D. N. T. Lam, J. Kurchan, and J. Tailleur, *Journal of Physics A: Mathematical and Theoretical* **46**, 254002 (2013).
- [71] M. Merolle, J. P. Garrahan, and D. Chandler, *Proceedings of the National Academy of Sciences* **102**, 10837 (2005).
- [72] R. L. Jack, J. P. Garrahan, and D. Chandler, *The Journal of Chemical Physics* **125**, 184509 (2006).
- [73] P. Visco, F. van Wijland, and E. Trizac, *Phys. Rev. E* **77**, 041117 (2008).
- [74] J. Hooyberghs and C. Vanderzande, *Journal of Statistical Mechanics: Theory and Experiment* **2010**, P02017 (2010).
- [75] R. L. Jack and J. P. Garrahan, *Phys. Rev. E* **81**, 011111 (2010).
- [76] D. Chandler and J. P. Garrahan, *Annual Review of Physical Chemistry* **61**, 191 (2010).
- [77] J. P. Garrahan and I. Lesanovsky, *Phys. Rev. Lett.* **104**, 160601 (2010).
- [78] J. P. Garrahan, A. D. Armour, and I. Lesanovsky, *Phys. Rev. E* **84**, 021115 (2011).
- [79] S. Genway, J. P. Garrahan, I. Lesanovsky, and A. D. Armour, *Phys. Rev. E* **85**, 051122 (2012).
- [80] C. Ates, B. Olmos, J. P. Garrahan, and I. Lesanovsky, *Phys. Rev. A* **85**, 043620 (2012).
- [81] J. M. Hickey, C. Flindt, and J. P. Garrahan, *Phys. Rev. E* **90**, 062128 (2014).
- [82] R. L. Jack and P. Sollich, *Journal of Physics A: Mathematical and Theoretical* **47**, 015003 (2014).
- [83] A. Dembo and O. Zeitouni, *Large deviations techniques and applications*, Applications of mathematics (Springer, New York, Berlin, Heidelberg, 1998).
- [84] T. Bodineau and B. Derrida, *Phys. Rev. Lett.* **92**, 180601 (2004).
- [85] T. Nemoto, F. Bouchet, R. L. Jack, and V. Lecomte, *Phys. Rev. E* **93**, 062123 (2016).

- [86] N. V. Kampen, ed., *Stochastic Processes in Physics and Chemistry (Third Edition)*, third edition ed., North-Holland Personal Library (Elsevier, Amsterdam, 2007).
- [87] C. W. Gardiner, *Handbook of stochastic methods : for physics, chemistry and the natural sciences*, Springer series in synergetics (Springer, Berlin, 1983).
- [88] R. L. Jack and P. Sollich, *Prog. Theor. Phys. Supplement* **184**, 304 (2010).
- [89] V. Popkov, G. M. Schütz, and D. Simon, *Journal of Statistical Mechanics: Theory and Experiment* **2010**, P10007 (2010).
- [90] R. Chetrite and H. Touchette, *Phys. Rev. Lett.* **111**, 120601 (2013).
- [91] T. Nemoto and S.-i. Sasa, *Phys. Rev. Lett.* **112**, 090602 (2014).
- [92] A. B. Kolton, S. Bustingorry, E. E. Ferrero, and A. Rosso, *J. Stat. Mech.* **2013**, P12004 (2013).
- [93] A. R. Admati and P. Pfleiderer, *The Review of Financial Studies* **1**, 3 (1988).
- [94] T. Andersen and T. Bollerslev, *Journal of Empirical Finance* **4**, 115 (1997).
- [95] R. Cont, *Quantitative Finance* **1**, 223 (2001).
- [96] A. Chakraborti, I. M. Toke, M. Patriarca, and F. Abergel, *Quantitative Finance* **11**, 991 (2011).
- [97] L. Borland, *Quantitative Finance* **12**, 1367 (2012).
- [98] L. Borland and Y. Hassid, *arxiv:1010.4917* (2010).
- [99] R. Allez and J.-P. Bouchaud, *New Journal of Physics* **13**, 025010 (2011).
- [100] T. Kaizoji, *Eur. Phys. J. B* **50**, 123 (2006).
- [101] F. Abergel, N. Huth, and I. Muni Toke, *SSRN.1474612* (2009).
- [102] P. Gopikrishnan, M. Meyer, L. A. N. Amaral, and H. E. Stanley, *Eur. Phys. J. B* **3**, 139 (1998).
- [103] P. Gopikrishnan, V. Plerou, L. A. N. Amaral, M. Meyer, and H. E. Stanley, *Phys. Rev. E* **60**, 5305 (1999).
- [104] L. Kullmann, J. Töyli, J. Kertesz, A. Kanto, and K. Kaski, *Physica A: Statistical Mechanics and its Applications* **269**, 98 (1999).
- [105] M. Tchernookov and A. R. Dinner, *J. Stat. Mech.* **2010**, P02006 (2010).
- [106] A. Kundu, S. Sabhapandit, and A. Dhar, *Phys. Rev. E* **83**, 031119 (2011).
- [107] P. Grassberger and A. de la Torre, *Annals of Physics* **122**, 373 (1979).
- [108] T. Liggett, *Interacting particle systems* (Springer-Verlag Berlin Heidelberg, 2005).

- [109] L. Boltzmann, On the Relation between the Second Law of the Mechanical Theory of Heat and Probability, and the Theorems Concerning Thermal Equilibrium, Respectively **2**, 373 (1877).
- [110] R. S. Ellis, *Physica D: Nonlinear Phenomena* **133**, 106 (1999).
- [111] J. W. Gibbs, *Elementary principles in statistical mechanics* (Charles Scribner's Sons, 1902).
- [112] H. Cramér, "Sur un nouveau théorème-limite de la théorie des probabilités." Actual. sci. industr. 736, 5-23. (Confér. internat. Sci. math. Univ. Genève. Théorie des probabilités. III: Les sommes et les fonctions de variables aléatoires.) (1938).
- [113] I. N. Sanov, *On the probability of large deviations of random variables*, Tech. Rep. (North Carolina State University. Dept. of Statistics, 1958).
- [114] M. Donsker and S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **28**, 1 (1975).
- [115] M. D. Donsker and S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **28**, 279 (1975).
- [116] M. D. Donsker and S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **29**, 389 (1976).
- [117] M. D. Donsker and S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **36**, 183 (1983).
- [118] S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **19**, 261 (1966).
- [119] M. I. Freidlin and A. D. Wentzell, *Random Perturbations of Dynamical Systems* (Springer, 1998) pp. 15–43.
- [120] J. Gärtner, *Theory of Probability and Its Applications* **22**, 24 (1977).
- [121] R. S. Ellis, *Ann. Probab.* **12**, 1 (1984).
- [122] R. S. Ellis, *Scandinavian Actuarial Journal* **1995**, 97 (1995).
- [123] R. S. Ellis, in *Lectures for the International Seminar on Extreme Events in Complex Dynamics* (2006).
- [124] R. S. Ellis, *Entropy, large deviations, and statistical mechanics* (Springer, 2007).
- [125] A. Martin-Lüf, *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete* **59**, 101 (1982).
- [126] W. Bryc, *Statistics and Probability Letters* **18**, 253 (1993).
- [127] O. E. Lanford, in *Statistical Mechanics and Mathematical Problems*, edited by A. Lenard (Springer Berlin Heidelberg, Berlin, Heidelberg, 1973) pp. 1–113.

- [128] N. O’Connell, *From laws of large numbers to large deviation principles* (Hewlett-Packard Laboratories, 1997).
- [129] N. O’Connell, *Mathematical Proceedings of the Cambridge Philosophical Society* **128**, 561–569 (2000).
- [130] A. Einstein, *Annalen der Physik* **14**, 280 (1907).
- [131] A. Einstein, *Annalen der Physik* **14**, 368 (1910).
- [132] E. T. Jaynes, *Phys. Rev.* **106**, 620 (1957).
- [133] E. T. Jaynes, *Phys. Rev.* **108**, 171 (1957).
- [134] E. T. Jaynes, *Probability Theory: The Logic of Science* (Cambridge University Press, 2003).
- [135] L. Onsager and S. Machlup, *Phys. Rev.* **91**, 1505 (1953).
- [136] S. Machlup and L. Onsager, *Phys. Rev.* **91**, 1512 (1953).
- [137] D. Falkoff, *Progress of Theoretical Physics* **16**, 530 (1956).
- [138] D. Falkoff, *Annals of Physics* **4**, 325 (1958).
- [139] D. Ruelle, *Statistical Mechanics: Rigorous Results* (World Scientific, 1999).
- [140] J. T. Lewis, in *On Stochastic Mechanics and Stochastic Processes* (Springer-Verlag New York, Inc., New York, NY, USA, 1986) pp. 141–155.
- [141] J. T. Lewis, in *Mark Kac Seminar on Probability and Physics, vol. 17, Math. Centrum Centrum Wisk. Inform.* (1988) pp. 85–102.
- [142] J. T. Lewis, in *Mathematical Methods in Statistical Mechanics, in: Leuven Notes Math. Theoret. Phys. Ser.A Math. Phys., vol. 1* (Leuven Univ. Press, 1989) pp. 77–90.
- [143] J. T. Lewis, C.-E. Pfister, and W. G. Sullivan, “Large deviations and the thermodynamic formalism: A new proof of the equivalence of ensembles,” in *On Three Levels: Micro-, Meso-, and Macro-Approaches in Physics*, edited by M. Fannes, C. Maes, and A. Verbeure (Springer US, Boston, MA, 1994) pp. 183–192.
- [144] J. Lewis, C.-E. Pfister, and W. Sullivan, *Markov Process. Related Fields* **1**, 319 (1995).
- [145] J. T. Lewis and C.-E. Pfister, *Russian Mathematical Surveys* **50**, 279 (1995).
- [146] Y. Oono, *Progress of Theoretical Physics Supplement* **99**, 165 (1989).
- [147] A. Amann and H. Atmanspacher, *J. Sci. Exploration* **13**, 639 (1999).
- [148] T. Eisele and R. S. Ellis, *Journal of Statistical Physics* **52**, 161 (1988).
- [149] S. Orey, *Stochastics* **25**, 3 (1988).
- [150] R. S. Ellis and K. Wang, *Stochastic Processes and their Applications* **35**, 59 (1990).

- [151] M. Costeniuc, R. S. Ellis, and H. Touchette, *Journal of Mathematical Physics* **46**, 063301 (2005).
- [152] J. Barré, D. Mukamel, and S. Ruffo, *Phys. Rev. Lett.* **87**, 030601 (2001).
- [153] R. S. Ellis, H. Touchette, and B. Turkington, *Physica A: Statistical Mechanics and its Applications* **335**, 518 (2004).
- [154] R. S. Ellis, P. T. Otto, and H. Touchette, *Ann. Appl. Probab.* **15**, 2203 (2005).
- [155] J. Barré, F. Bouchet, T. Dauxois, and S. Ruffo, *Journal of Statistical Physics* **119**, 677 (2005).
- [156] M. Kastner and O. Schnetz, *Journal of Statistical Physics* **122**, 1195 (2006).
- [157] L. Casetti and M. Kastner, *Physica A: Statistical Mechanics and its Applications* **384**, 318 (2007).
- [158] I. Hahn and M. Kastner, *Phys. Rev. E* **72**, 056134 (2005).
- [159] I. Hahn and M. Kastner, *The European Physical Journal B - Condensed Matter and Complex Systems* **50**, 311 (2006).
- [160] A. Campa, S. Ruffo, and H. Touchette, *Physica A: Statistical Mechanics and its Applications* **385**, 233 (2007).
- [161] L. Landau and E. Lifshitz, *Statistical Physics, 3rd ed., in: Landau and Lifshitz Course of Theoretical Physics, vol. 5* (Butterworth Heinemann, Oxford, 1991).
- [162] I. Ispolatov and E. Cohen, *Physica A: Statistical Mechanics and its Applications* **295**, 475 (2001).
- [163] M. K.-H. Kiessling and T. Neukirch, *Proceedings of the National Academy of Sciences* **100**, 1510 (2003).
- [164] R. S. Ellis, K. Haven, and B. Turkington, *Nonlinearity* **15**, 239 (2002).
- [165] R. S. Ellis, K. Haven, and B. Turkington, *Journal of Statistical Physics* **101**, 999 (2000).
- [166] M. K.-H. Kiessling and J. L. Lebowitz, *Letters in Mathematical Physics* **42**, 43 (1997).
- [167] P. H. Chavanis, *International Journal of Modern Physics B* **20**, 3113 (2006).
- [168] D. Lynden-Bell, *Physica A: Statistical Mechanics and its Applications* **263**, 293 (1999), proceedings of the 20th IUPAP International Conference on Statistical Physics.
- [169] J. Lebowitz, *Physics Today* **46**, 32 (1993), cited By 216.
- [170] B. Derrida, *Journal of Statistical Mechanics: Theory and Experiment* **2007**, P07023 (2007).

- [171] C. Gardiner, *Handbook of Stochastic Methods for Physics, Chemistry and the Natural Sciences* (Springer, New York, 1985).
- [172] G. L. Eyink, *Journal of Statistical Physics* **61**, 533 (1990).
- [173] Bertini, A. D. Sole, D. Gabrielli, G. Jona-Lasinio, and C. Landim, *Journal of Statistical Mechanics: Theory and Experiment* **2007**, P07014 (2007).
- [174] B. Derrida, *Physics Reports* **301**, 65 (1998).
- [175] O. Benois, C. Kipnis, and C. Landim, *Stochastic Processes and their Applications* **55**, 65 (1995).
- [176] C. Landim and M. Mourragui, *Annales de l'Institut Henri Poincare (B) Probability and Statistics* **33**, 65 (1997).
- [177] T. Bodineau and B. Derrida, *Journal of Statistical Physics* **123**, 277 (2006).
- [178] C. Kipnis and S. Olla, *Stochastics and Stochastic Reports* **33**, 17 (1990).
- [179] B. Derrida, J. L. Lebowitz, and E. R. Speer, *Phys. Rev. Lett.* **87**, 150601 (2001).
- [180] B. Derrida, J. L. Lebowitz, and E. R. Speer, *Journal of Statistical Physics* **107**, 599 (2002).
- [181] B. Derrida, J. L. Lebowitz, and E. R. Speer, *Phys. Rev. Lett.* **89**, 030601 (2002).
- [182] B. Derrida, J. L. Lebowitz, and E. R. Speer, *Journal of Statistical Physics* **110**, 775 (2003).
- [183] D. J. Evans and D. J. Searles, *Advances in Physics* **51**, 1529 (2002).
- [184] J. Kurchan, *Journal of Statistical Mechanics: Theory and Experiment* **2007**, P07005 (2007).
- [185] J. Kurchan, [arxiv:0901.1271](https://arxiv.org/abs/0901.1271) (2009).
- [186] D. J. Evans, E. G. D. Cohen, and G. P. Morriss, *Phys. Rev. Lett.* **71**, 2401 (1993).
- [187] D. J. Evans and D. J. Searles, *Phys. Rev. E* **50**, 1645 (1994).
- [188] G. M. Wang, E. M. Sevick, E. Mittag, D. J. Searles, and D. J. Evans, *Phys. Rev. Lett.* **89**, 050601 (2002).
- [189] G. Gallavotti and E. G. D. Cohen, *Phys. Rev. Lett.* **74**, 2694 (1995).
- [190] G. Gallavotti and E. G. D. Cohen, *Journal of Statistical Physics* **80**, 931 (1995).
- [191] J. Kurchan, *Journal of Physics A: Mathematical and General* **31**, 3719 (1998).
- [192] J. L. Lebowitz and H. Spohn, *Journal of Statistical Physics* **95**, 333 (1999).
- [193] C. Maes, *Journal of Statistical Physics* **95**, 367 (1999).

- [194] D. Andrieux, P. Gaspard, S. Ciliberto, N. Garnier, S. Joubaud, and A. Petrosyan, *Phys. Rev. Lett.* **98**, 150601 (2007).
- [195] R. van Zon, S. Ciliberto, and E. G. D. Cohen, *Phys. Rev. Lett.* **92**, 130601 (2004).
- [196] N. Garnier and S. Ciliberto, *Phys. Rev. E* **71**, 060101 (2005).
- [197] S. Aumaître, S. Fauve, S. McNamara, and P. Poggi, *The European Physical Journal B - Condensed Matter and Complex Systems* **19**, 449 (2001).
- [198] K. Feitosa and N. Menon, *Phys. Rev. Lett.* **92**, 164301 (2004).
- [199] A. Puglisi, P. Visco, A. Barrat, E. Trizac, and F. van Wijland, *Phys. Rev. Lett.* **95**, 110202 (2005).
- [200] P. Visco, A. Puglisi, A. Barrat, E. Trizac, and F. van Wijland, *EPL (Europhysics Letters)* **72**, 55 (2005).
- [201] P. Visco, A. Puglisi, A. Barrat, E. Trizac, and F. van Wijland, *Journal of Statistical Physics* **125**, 533 (2006).
- [202] S. Ciliberto and C. Laroche, *Journal De Physique. IV : JP* **8**, 215 (1998).
- [203] S. Ciliberto, N. Garnier, S. Hernandez, C. Lacpatia, J.-F. Pinton, and G. R. Chavarria, *Physica A: Statistical Mechanics and its Applications* **340**, 240 (2004).
- [204] B. Cleuren, C. Van den Broeck, and R. Kawai, *Phys. Rev. E* **74**, 021117 (2006).
- [205] C. Beck and F. Schlögl, *Thermodynamics of Chaotic Systems: An Introduction* (Cambridge University Press, Cambridge, 1993).
- [206] G. Paladin and A. Vulpiani, *Physics Reports* **156**, 147 (1987).
- [207] J. L. McCauley, *Physics Reports* **189**, 225 (1990).
- [208] K. Falconer, *Techniques in Fractal Geometry* (Wiley, New York, 1997).
- [209] P. Gaspard, *Chaos, Scattering and Statistical Mechanics* (Cambridge University Press, Cambridge, 1998).
- [210] V. Alekseev and M. Yakobson, *Physics Reports* **75**, 290 (1981).
- [211] J. P. Eckmann and D. Ruelle, *Rev. Mod. Phys.* **57**, 617 (1985).
- [212] M. M. A. Lasota, *Chaos, Fractals, and Noise: Stochastic Aspects of Dynamics, Applied Mathematical Sciences, vol. 97* (Springer, New York, 1994).
- [213] G. Zohar, *Stochastic Processes and their Applications* **79**, 229 (1999).
- [214] D. Veneziano, *Fractals* **10**, 117 (2002).
- [215] D. Harte, *Multifractals: Theory and Applications* (CRC Press, New York, 2001).

- [216] G. Keller, *Equilibrium States in Ergodic Theory*, *London Math. Soc. Student Texts*, vol. 42 (Cambridge University Press, Cambridge, 1998).
- [217] D. Ruelle, *Thermodynamic Formalism* (Cambridge University Press, Cambridge, 2004).
- [218] D. Ruelle, *Communications in Mathematical Physics* **125**, 239 (1989).
- [219] Y. Sinai, *Russian Mathematical Surveys* **27**, 21 (1972).
- [220] Y. Sinai, *Topics in Ergodic Theory* (Princeton University Press, Princeton, 1994).
- [221] L.-S. Young, *Trans. Amer. Math. Soc.* **318**, 525 (1990).
- [222] A. O. Lopes, *Nonlinearity* **3**, 527 (1990).
- [223] S. Waddington, *Annales de l'Institut Henri Poincaré (C) Non Linear Analysis* **13**, 445 (1996).
- [224] M. Pollicott, R. Sharp, and M. Yuri, *Nonlinearity* **11**, 1173 (1998).
- [225] O. Z. N. Gantert, in *Proceedings of the Conference on Random Walks, Bolyai Society Mathematical Studies*, vol.9 (1999) pp. 127–165.
- [226] F. Comets, N. Gantert, and O. Zeitouni, *Probability Theory and Related Fields* **118**, 65 (2000).
- [227] S. R. S. Varadhan, *Communications on Pure and Applied Mathematics* **56**, 1222 (2003).
- [228] O. Zeitouni, *Journal of Physics A: Mathematical and General* **39**, R433 (2006).
- [229] T. C. Dorlas and J. R. Wedagedera, *International Journal of Modern Physics B* **15**, 1 (2001).
- [230] T. C. Dorlas and W. M. B. Dukes, *Journal of Physics A: Mathematical and General* **35**, 4385 (2002).
- [231] M. Talagrand, *Journal of Statistical Physics* **126**, 837 (2007).
- [232] W. Cegła, J. T. Lewis, and G. A. Raggio, *Communications in Mathematical Physics* **118**, 337 (1988).
- [233] M. van den Berg, J. T. Lewis, and J. V. Pulé, *Communications in Mathematical Physics* **118**, 61 (1988).
- [234] T. C. Dorlas, P. A. Martin, and J. V. Pule, *Journal of Statistical Physics* **121**, 433 (2005).
- [235] J. L. Lebowitz, M. Lenci, and H. Spohn, *Journal of Mathematical Physics* **41**, 1224 (2000).
- [236] G. Gallavotti, J. L. Lebowitz, and V. Mastropietro, *Journal of Statistical Physics* **108**, 831 (2002).

- [237] F. Hiai, M. Mosonyi, and T. Ogawa, *Journal of Mathematical Physics* **48**, 123301 (2007).
- [238] M. Lenci and L. Rey-Bellet, *Journal of Statistical Physics* **119**, 715 (2005).
- [239] K. Netočný and F. Redig, *Journal of Statistical Physics* **117**, 521 (2004).
- [240] S. Tănase-Nicola and J. Kurchan, *Phys. Rev. Lett.* **91**, 188302 (2003).
- [241] S. Tănase-Nicola and J. Kurchan, *Journal of Statistical Physics* **116**, 1201 (2004).
- [242] J. Tailleur, S. Tănase-Nicola, and J. Kurchan, *Journal of Statistical Physics* **122**, 557 (2006).
- [243] P. A. M. Dirac, *Mathematical Proceedings of the Cambridge Philosophical Society* **35**, 416 (1939).
- [244] A. Rákos and R. J. Harris, *Journal of Statistical Mechanics: Theory and Experiment* **2008**, P05005 (2008).
- [245] M. Gorissen, J. Hooyberghs, and C. Vanderzande, *Phys. Rev. E* **79**, 020101 (2009).
- [246] M. Gorissen and C. Vanderzande, *Journal of Physics A: Mathematical and Theoretical* **44**, 115005 (2011).
- [247] M. Gorissen, A. Lazarescu, K. Mallick, and C. Vanderzande, *Phys. Rev. Lett.* **109**, 170601 (2012).
- [248] C. P. Espigares, P. L. Garrido, and P. I. Hurtado, *Phys. Rev. E* **87**, 032115 (2013).
- [249] P. I. Hurtado, C. P. Espigares, J. J. del Pozo, and P. L. Garrido, *Journal of Statistical Physics* **154**, 214 (2014).
- [250] P. Tsobgni Nyawo and H. Touchette, *Phys. Rev. E* **94**, 032101 (2016).
- [251] C. Bezuidenhout and G. Grimmett, *The Annals of Probability* **18**, 1462 (1990).
- [252] J. Marro and R. Dickman, *Nonequilibrium Phase Transitions in Lattice Models*, Collection Alea-Saclay: Monographs and Texts in Statistical Physics (Cambridge University Press, 1999).
- [253] H. Hinrichsen, *Advances in Physics* **49**, 815 (2000).
- [254] G. Ódor, *Rev. Mod. Phys.* **76**, 663 (2004).
- [255] C. M. Rohwer, F. Angeletti, and H. Touchette, *Phys. Rev. E* **92**, 052104 (2015).
- [256] T. Nemoto and S.-i. Sasa, *Phys. Rev. E* **84**, 061113 (2011).
- [257] T. Nemoto, V. Lecomte, S. Sasa, and F. van Wijland, *Journal of Statistical Mechanics: Theory and Experiment* **2014**, P10001 (2014).
- [258] P. Tsobgni Nyawo and H. Touchette, *EPL (Europhysics Letters)* **116**, 50009 (2016).

- [259] T. Nemoto, R. L. Jack, and V. Lecomte, *Phys. Rev. Lett.* **118**, 115702 (2017).
- [260] S. R. White, *Phys. Rev. Lett.* **69**, 2863 (1992).
- [261] S. R. White, *Phys. Rev. B* **48**, 10345 (1993).
- [262] U. Schollwöck, *Rev. Mod. Phys.* **77**, 259 (2005).
- [263] M. Kaulke and I. Peschel, *The European Physical Journal B - Condensed Matter and Complex Systems* **5**, 727 (1998).
- [264] E. Carlon, M. Henkel, and U. Schollwöck, *The European Physical Journal B - Condensed Matter and Complex Systems* **12**, 99 (1999).
- [265] E. Pitard, V. Lecomte, and F. van Wijland, *EPL (Europhysics Letters)* **96**, 56002 (2011).
- [266] T. Speck and D. Chandler, *The Journal of Chemical Physics* **136**, 184509 (2012).
- [267] D. T. Limmer and D. Chandler, *Proceedings of the National Academy of Sciences* **111**, 9413 (2014).
- [268] T. R. Gingrich and P. L. Geissler, *The Journal of Chemical Physics* **142**, 234104 (2015).
- [269] J.-P. Bouchaud and M. Potters, *Theory of financial risk and derivative pricing: from statistical physics to risk management* (Cambridge University Press, 2003).
- [270] J.-P. Bouchaud and M. Potters, *arxiv: 0910.1205* (2009).
- [271] V. Plerou, P. Gopikrishnan, B. Rosenow, L. A. N. Amaral, T. Guhr, and H. E. Stanley, *Phys. Rev. E* **65**, 066126 (2002).
- [272] A. Utsugi, K. Ino, and M. Oshikawa, *Phys. Rev. E* **70**, 026110 (2004).
- [273] M. Cavallaro and R. J. Harris, *Journal of Physics A: Mathematical and Theoretical* **49**, 47LT02 (2016).
- [274] J. S. Wettlaufer, *Phys. Rev. Lett.* **116**, 150002 (2016).
- [275] R. L. Jack and P. Sollich, *The European Physical Journal Special Topics* **224**, 2351 (2015).
- [276] U. Ray, G. Kin-Lic Chan, and D. Limmer, *arxiv: 1708.00459* (2017).
- [277] U. Ray, G. Kin-Lic Chan, and D. Limmer, *arxiv: 1708.09482* (2017).
- [278] J.-D. Deuschel and D. W. Stroock, *Large deviations*, Vol. 342 (American Mathematical Soc., 2001).
- [279] R. J. Harris and G. M. Schütz, *Journal of Statistical Mechanics: Theory and Experiment* **2007**, P07020 (2007).
- [280] R. J. Harris, A. Rákos, and G. M. Schütz, *EPL (Europhysics Letters)* **75**, 227 (2006).