

Reconstructing a Random Potential from its Random Walks

S. Cocco ¹, R. Monasson ²

¹ CNRS-Laboratoire de Physique Statistique de l'ENS, 24 rue Lhomond, 75005 Paris, France

² CNRS-Laboratoire de Physique Théorique de l'ENS, 24 rue Lhomond, 75005 Paris, France

The problem of how many trajectories of a random walker in a potential are needed to reconstruct the values of this potential is studied. We show that this problem can be solved by calculating the probability of survival of an abstract random walker in a partially absorbing potential. The approach is illustrated on the discrete Sinai (random force) model with a drift. We determine the parameter (temperature, duration of each trajectory, ...) values making reconstruction as fast as possible.

Introduction. Random walks (RW) in random media have been intensively studied in the past decades as a paradigm for out-of-equilibrium dynamics, and have led to the discovery and understanding of important dynamical effects as anomalous diffusion, ageing ...[1, 2]. Briefly speaking the issue is to determine the statistical properties of the walker from the ones of the energy potential. Much less attention has been devoted to the inverse problem: given one (or more) observed RW(s) can we guess the potential values? This question naturally arises in biophysics where the use of AFM, optical and magnetic tweezers make possible the mechanical separation of single protein-protein complexes [3], or the unfolding and refolding of single biomolecules[4, 5, 6]. The observed dynamics the rupture of chemical bonds, of folding/unfolding of nucleic acids, or proteins can be modeled as a RW motion affected by thermal noise, moving in a quenched potential determined by the composition of the chemical bonds, or the sequence of amino- or nucleic-acids. Reconstructing the free energy landscape of those processes is the object of current and intense efforts [3, 6, 7, 8, 9].

In this letter we show how the inverse RW problem can be practically solved within the Bayesian inference framework and address the crucial question of the accuracy of reconstruction. In practice information can be accumulated either by increasing the duration of one RW, or observing more than one RW, or combining the two. We discuss the optimal procedure minimizing the total number of data to be acquired, and show how this minimal amount of data can be calculated from the probability of survival of an abstract walker in a partially absorbing potential. The approach is illustrated in detail on the celebrated discrete random force (RF) model (Sinai model with non zero drift) [1, 2].

Inference is a key issue in information theory and statistics [10], with applications in biology [11], social science [12], finance, ... A central question is the so-called hypothesis testing problem: which one of two candidate distributions is likely to have generated a set of measured data? This question was solved in the case of independent variables by Chernoff [13], and is the core issue of the asymptotic theory of inference [10]. Chernoff showed that the probability of guessing the wrong distribution decreases exponentially with the size of the data set [13]. Large deviations techniques can be used to treat the case

of variables extracted from one recurrent realization of a finite Markov chain [14, 15]; the present work can be seen as an extension to many transient realizations of an 'infinite' chain.

Random Force model. For an illustration of the problem consider the discrete, one dimensional RF model defined on the set of sites $x = 0, 1, 2, \dots, N$ [1]. We start by choosing randomly a set of dimensionless forces $f_x = \pm 1$ on each link $(x, x + 1)$ with *a priori* probability $P_0 = \prod_x \frac{1+b f_x}{2}$ where $-1 < b < 1$ is called tilt. This defines the values of the potential \mathbf{V} on each site, $V_x = -\sum_{y < x} f_y$ (by definition $V_0 = 0$). An example of potential for $b = 0.4$ is shown on Fig. 1.

After the quenched potential has been drawn a random walker starts in $x = 0$ at time $t = 0$. The walker then jumps from one site x to one of its neighbors $x' = x \pm 1$ with rate (probability per unit of time) $r_{\mathbf{V}}(x \rightarrow x') = r_0 \times e^{(V_x - V_{x'})/(2T)}$ to satisfy detailed balance at temperature T ; the attempt rate r_0 will be set to unity in the following. Reflecting boundary conditions are imposed by setting $V_{N+1} = V_{-1} = +\infty$. We register the sequence of positions up to some time t_f : $\mathbf{X} = \{x(t), 0 \leq t \leq t_f\}$. Figure 1 shows five RWs \mathbf{X}_ρ , $\rho = 1, \dots, 5$, each starting in the origin $x(0) = 0$ and of equal duration t_f for a temperature $T = 1$. The value of the temperature strongly affects the dynamics [2], and its relevance for the inverse problem will be discussed later.

Our objective is to reconstruct the potential over a region of the lattice e.g. the value of the forces on some specific links from the observation of RWs. Within Bayes inference framework this can be done by maximizing the joint probability of the potential \mathbf{V} and of the observed RWs $\mathbf{X}_1, \dots, \mathbf{X}_R$ over \mathbf{V} [10]. P is the product of the *a priori* probability of the potential, P_0 , times the likelihood of the RWs given the potential, L . Since the RW is Markovian L depends only on the sets of total times t_x spent on every site x , and of the numbers of jumps $u(x \rightarrow x')$ from x to x' over the set of RWs:

$$L = \prod_{x, x'} e^{-t_x r_{\mathbf{V}}(x \rightarrow x')} r_{\mathbf{V}}(x \rightarrow x')^{u(x \rightarrow x')} \quad (1)$$

where the product runs over all sites x and their neighbors $x' = x \pm 1$. Expressing the rates in terms of the forces and maximizing the joint probability P we obtain the most likely values for the forces: $f_x = \text{sign}(h_x + \alpha)$ where $\alpha \equiv T \ln[(1+b)/(1-b)]$ is a global 'field' coming

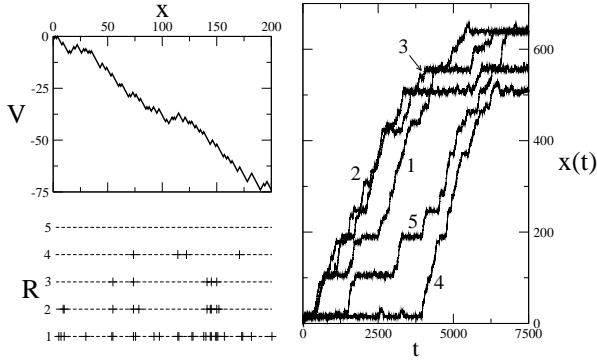


FIG. 1: Left, top: Example of potential \mathbf{V} obtained in the RF model with tilt $b = 0.4$ (size $N = 1000$, sites $x > 200$ not shown here). Right: examples of RWs, numbered from 1 to 5, in this potential at temperature $T = 1$; plateaus are in correspondence with the local minima of V . Here $\alpha \simeq 0.85$ (creep phase). Left, bottom: Predictions from the first R RWs in the right panel and (2); impulses locate incorrectly predicted forces f_x for $x \leq 200$. The number of erroneous forces decreases from 26 (for $R = 1$) to 0 ($R = 5$). Note the errors on sites $x_0 \simeq 100$ appearing when the fourth RW is taken into account; indeed this atypical RW marks no pause in the local minimum in x_0 .

from the *a priori* distribution P_0 and h_x a local contribution due to the likelihood L ,

$$h_x = 2T \sinh\left(\frac{1}{2T}\right) (t_{x+1} - t_x) + u(x \rightarrow x+1) - u(x+1 \rightarrow x). \quad (2)$$

Figure 1 (left, bottom) shows predictions made from $R = 1$ to $R = 5$ RWs for the first 200 sites. The duration t_f of the RW is chosen to be much larger than the mean first passage time in $x = 200$, and much smaller than the equilibration time $t_{eq} \sim e^{bN/T}$. In this range the quality of prediction is essentially independent of t_f as will be discussed in detail below. As expected the number of erroneous forces decreases with increasing R though atypical events may produce flaws in the prediction. The analysis of these atypical RWs, and how they lead to errors is the keystone of what follows.

Number of RWs necessary for a good reconstruction. Expression (1) for the likelihood of the RWs is true for any potential \mathbf{V} and can be geometrically interpreted as follows. Given a set of RWs we extract a signal vector \mathbf{S} whose components are: the times t_x spent on site x , the numbers $u(x \rightarrow x')$ of transitions from site x to site x' . When R is large we expect \mathbf{S} to be extensive with R and define the intensive signal $\mathbf{s} = \mathbf{S}/R$. Similarly, to each potential \mathbf{V} we associate a vector \mathbf{v} with components: minus the outgoing rate *i.e.* $-\sum_{x'(\neq x)} r_{\mathbf{V}}(x \rightarrow x')$ for each site x , the logarithm of the rate $r_{\mathbf{V}}(x \rightarrow x')$ for each pair of neighbors. Then $L = \exp(R \mathbf{s} \cdot \mathbf{v})$ from

(1) where \cdot denotes the scalar product. Maximizing the joint probability $P = P_0 \times L$ over the potential becomes equivalent, in the large R limit, to finding \mathbf{v} with the largest scalar product with the signal \mathbf{s} [20]. It is natural to partition the space of signals into ‘Voronoi cells’: $C_{\mathbf{v}}$ is the set of \mathbf{s} having a larger scalar product with \mathbf{v} than with any other potential \mathbf{v}' . Bayes rule tells us that the most likely potential given an observed signal \mathbf{s} is the one attached to the cell in which \mathbf{s} lies.

Consider now RWs taking place in a given potential \mathbf{V} . From the law of large number the signal \mathbf{s} is equal, in the infinite R limit, to $\mathbf{s}_{\mathbf{v}}^* = \{t_x^*, u^*(x \rightarrow x') = t_x^* r_{\mathbf{V}}(x \rightarrow x')\}$ where t_x^* is the average sojourn time on site x over RWs of duration t_f . As $\mathbf{s}_{\mathbf{v}}^* \in C_{\mathbf{v}}$ [21] reconstruction becomes flawless in the limit of an infinite number of data as expected. For large albeit finite R , \mathbf{s} typically deviates from $\mathbf{s}_{\mathbf{v}}^*$ by $O(R^{-\frac{1}{2}})$; finite deviations have exponentially small-in- R probabilities, $e^{-R\omega_{\mathbf{V}}(\mathbf{s})}$, controlled by a rate function $\omega_{\mathbf{V}}(\mathbf{s})$ [14]. The probability to predict an erroneous potential is the probability that the stochastic signal \mathbf{s} does not belong to cell $C_{\mathbf{v}}$. This probability of error thus decays exponentially with R over a typical number of RWs

$$R_c(\mathbf{V}) = \left[\min_{\mathbf{s} \notin C_{\mathbf{v}}} \omega_{\mathbf{V}}(\mathbf{s}) \right]^{-1}, \quad (3)$$

where the minimum is taken over signals outside the ‘true’ cell. It depends on the temperature, the duration of the RW, ...

As the RWs are independently drawn $\omega_{\mathbf{V}}$ is a convex function of \mathbf{s} [14]. The minimum in (3) is thus reached on the boundary between the true cell and another, bad cell, say, $C_{\bar{\mathbf{v}}}$. The attached potential, $\bar{\mathbf{V}}$, is the most ‘dangerous’ one from the inference point of view. RWs generated from \mathbf{V} and $\bar{\mathbf{V}}$ are hardly told from each other unless more than $R_c(\mathbf{V})$ of them are observed.

Assume $\bar{\mathbf{V}}$ is known. Then the boundary between $C_{\mathbf{v}}$ and $C_{\bar{\mathbf{v}}}$ is the set of signals $\mathbf{s} \perp \mathbf{v} - \bar{\mathbf{v}}$. We deduce

$$R_c(\mathbf{V}) = \left[\max_{\mu} \min_{\mathbf{s}} (\omega_{\mathbf{V}}(\mathbf{s}) + \mu \mathbf{s} \cdot (\bar{\mathbf{v}} - \mathbf{v})) \right]^{-1} \quad (4)$$

where the Lagrange multiplier $\mu \in [0; 1]$ ensures that \mathbf{s} is confined to the boundary. The Legendre transform of $\omega_{\mathbf{V}}$ appearing in (4) is intimately related to the evolution operator of an abstract random walk process, denoted by $\text{RW}(\mu)$ to distinguish from the original RW [16]. This $\text{RW}(\mu)$ -er moves with the rates $r_{(1-\mu)\mathbf{V} + \mu\bar{\mathbf{V}}}(x \rightarrow x')$ and may die on every site x with positive rate

$$d_{\mathbf{V}, \bar{\mathbf{V}}, \mu}(x) = \sum_{x'(\neq x)} \left[(1-\mu) r_{\mathbf{V}}(x \rightarrow x') + \mu r_{\bar{\mathbf{V}}}(\bar{\mathbf{v}}(x \rightarrow x')) - r_{(1-\mu)\mathbf{V} + \mu\bar{\mathbf{V}}}(x \rightarrow x') \right]. \quad (5)$$

Consider now the probability $\pi(\mu)$ that $\text{RW}(\mu)$ -er, initially at the origin, has survived up to time t_f (the duration of the original RW). Then $R_c(\mathbf{V}) = \min_{\mu \in [0; 1]} 1/|\ln \pi(\mu)|$.

Optimal Working Point for the RF model. We apply the above theory to the discrete RF model, and want to predict the value of the force f_y on the link $(y, y+1)$ for some specific y . The dangerous potential is \mathbf{V} obtained from \mathbf{V} upon reversal of the force $f_y \rightarrow -f_y$. We aim at calculating the probability $\pi(\mu)$ of survival of RW(μ)-er moving with rate $r(x \rightarrow x') = r_{\mathbf{V}}(x \rightarrow x')$ and dying on site x with rate $d(x) = 0$ except: $r(y \rightarrow y+1) = 1/r(y+1 \rightarrow y) = e^{(1-2\mu)f_y/(2T)}$, $d(y) = D(f_y)$, $d(y+1) = D(-f_y)$ where $D(f) \equiv (1-\mu)e^{f/(2T)} + \mu e^{-f/(2T)} - e^{(1-2\mu)f/(2T)}$ from (5). From the previous section the number of RWs required for a reliable prediction of f_y is $R_c(y; \mathbf{V}) = \min_{\mu} 1/|\ln \pi(\mu)|$.

Let $\pi_x(\mu, t)$ be the probability that RW(μ), initially on site x , is still alive at time t . The time-evolution of π_x is described by

$$\frac{\partial \pi_x}{\partial t} = \sum_{x' (\neq x)} r(x \rightarrow x') (\pi_{x'} - \pi_x) - d(x) \pi_x, \quad (6)$$

with initial condition $\pi_y(\mu, 0) = 1$ (by convention $\pi_{-1} = \pi_{N+1} = 0$). After Laplace transform over time, eqns (6) are turned into recurrence equations for the ratios π_x/π_{x+1} and solved with great numerical accuracy. We obtain this way the probability of survival, $\pi(\mu) = \pi_0(\mu, t_f)$, and optimize over μ . Though R_c depends on the potential \mathbf{V} its general behavior for tilt $b > 0$ as a function of the duration t_f is sketched in Fig. 2. Three regimes are observed:

- for $t_f \ll \tau_y$ (mean first passage time in y) RW(μ) has a low probability to visit y and is almost surely alive, hence R_c is very large;
- for $\tau_y \ll t_f \ll t_{eq}$ RW(μ) has visited the region surrounding y and escaped from this region (transient regime), hence its probability of survival remains constant, and so does R_c ;
- for $t_f \gg t_{eq}$ RW(μ) visits again and again the region surrounding y , hence the probability of survival decreases exponentially with the duration: $R_c \propto 1/t_f$.

The total time $R_c \times t_f$ for a good reconstruction is minimal when we choose $t_f \gtrsim \tau_y$. This marginally transient regime corresponds to the plateau of Fig. 2: RWs are long enough to visit site y but short enough not to wander much away from y . To calculate the corresponding value of R_c we take the limits, in order, $N \rightarrow \infty$, $t_f \rightarrow \infty$, and look for the stationary solution of (6) with boundary condition $\pi_{x \rightarrow \infty} = 1$. The result for the probability of survival is

$$\pi(\mu) = \frac{e^{-\frac{\mu}{T}}}{1 - \mu + \mu e^{-\frac{1}{T}} + \mu(1-\mu)t_{y+1}^* (e^{\frac{1}{4T}} - e^{-\frac{3}{4T}})^2}, \quad (7)$$

where the mean sojourn time on site $y+1$ in \mathbf{V} is [2]

$$t_{y+1}^* = \sum_{z \geq 0} \exp \left[\frac{1}{T} \left(\frac{V_{y+z+2} + V_{y+z+1}}{2} - V_{y+1} \right) \right]. \quad (8)$$

Distribution of R_c over potentials. The number $R_c(y; \mathbf{V})$ of RWs necessary to predict the value of f_y de-

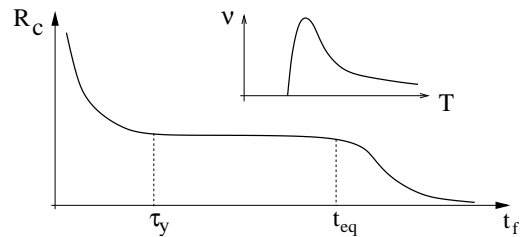


FIG. 2: Sketch of the number $R_c(y; \mathbf{V})$ of RWs necessary for a good inference of the force f_y as a function of the RW duration t_f . τ_y is the typical first-passage time in y from the origin, t_{eq} the equilibration time (comparable to the first-passage time from the extremity N when $y \ll N$). Inset: rate of reconstruction (9) as a function of temperature at fixed tilt.

pends on the potential \mathbf{V} through the sojourn time t_{y+1}^* (8). By randomly drawing potentials (or varying site y) we obtain the distribution of R_c shown in Fig. 3. Main features are:

- Small R_c correspond to sites where the RW spends long time t^* (traps)[22]: $R_c \sim \frac{1}{|\ln \pi|} \sim \frac{1}{\ln t^*}$ from (7). The power law tail of the distribution of sojourn times, $P(t^*) \sim (t^*)^{-(\alpha+1)}$ [2], gives rise to an essential singularity at the origin in the cumulative distribution, $\mathcal{Q}(R_c) \sim e^{-\alpha/R_c}$. The potential is easy to predict over trapping regions since RWer spends a long time there, and accumulates information about the energy landscape.

- Conversely the largest value of R_c , denoted by R_c^H , correspond to the homogeneous potential $V_x^H = -x$ in which the walker is never trapped and is quickly driven to $+\infty$. R_c^H can be calculated from (7) by setting $f_x = +1$ for all sites in (8). The singularity in \mathcal{Q} when $R_c \rightarrow R_c^H$ corresponds to quasi-homogeneous potentials, where one force, say, on site ℓ , is -1 . Such potentials have exponential-in- ℓ small probabilities, but give values of R_c on site $y = 0$ exponentially close to R_c^H . On the overall we find $1 - \mathcal{Q}(R_c^H - \epsilon) \sim \epsilon^\beta$ where the exponent is $\beta = T \ln \frac{1+b}{2}$.

- In between \mathcal{Q} shows marked steps at well defined and b -independent values of R_c , which correspond to specific local force patterns beyond site y . A ℓ -pattern is defined as a sequence of forces on sites $y+1$ to $y+\ell+1$, followed by all $+$ forces; the corresponding R_c can be exactly calculated from (7,8), and is shown for 7 among the 16 $\ell = 4$ -patterns in Fig. 3. The histogram of R_c can be accurately approximated for any tilt $b > 0$ based on the above local pattern description. Given a length ℓ we enumerate all the 2^ℓ patterns, calculate the corresponding R_c , and weight them with probability $(\frac{1+b}{2})^{\#f_x=+} \times (\frac{1-b}{2})^{\#f_x=-}$. In practice we choose $\ell \sim 10/\ln[2/(1-b)]$, to ensure that patterns with more than ℓ negative forces have negligible weights ($< e^{-10}$). The resulting histograms are in excellent agreement with \mathcal{Q} for intermediate values of R_c (dashed lines in Fig. 3).

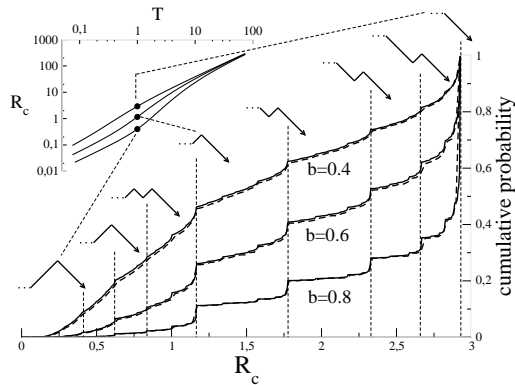


FIG. 3: Cumulative probability distribution \mathcal{Q} of $R_c(y; \mathbf{V})$ at temperature $T = 1$ and for three tilt values b . Full lines are numerical results from 10^6 samples, and dashed lines are the outcomes from the ℓ -pattern approximation. Inset: R_c vs. T for the 3-patterns $+++$, $+-+$, $---$ (from top to down).

Tuning temperature for fast reconstruction. The dependence of R_c upon temperature is shown for three patterns in the Inset of Fig. 3. We have $R_c \sim 4T$ as $T \rightarrow \infty$ independently of the pattern, and $R_c \sim 2T/(h+3)$ when $T \rightarrow 0$ where h is the highest barrier to the right of y in the potential defined by the pattern (Fig. 3). When the temperature exceeds the temperature T_b such that $\alpha = 1$ the velocity of the RWer is finite $\frac{y}{\tau_y} \sim v(T) > 0$ [2]. The reconstruction rate (number of correctly predicted forces per unit of time) is equal to the velocity $v(T)$ divided by R_c ,

$$v(T) = \frac{1 - \cosh \frac{1}{T} + b \sinh \frac{1}{T}}{\cosh \frac{1}{2T} - b \sinh \frac{1}{2T}} \times \int_0^{R_c^H} dR_c \frac{\mathcal{Q}'(R_c)}{R_c} \quad (9)$$

after averaging over the quenched potential. The depen-

dence of ν upon temperature is sketched in the Inset of Fig. 2; it is maximal and equal to ν^M for some temperature T^M realizing a trade-off between fast motion (large velocity) and accurate reading-out (small R_c). Even in the small b limit the optimal reconstruction rate is finite, $\nu^M \sim b^2$, by working at high temperature $T^M \sim \frac{1}{b}$, while in the absence of optimization procedure the number of predicted forces scales only as the squared logarithm of the time [19].

Conclusion. We have shown how the number of RWs required for a good reconstruction of the potential can be deduced from the probability of survival of an absorbing RW process. This result is of practical interest since the survival probability can be estimated through numerical simulations e.g. in dimension ≥ 2 . Furthermore we have determined, for the special case of the RF model, the optimal ‘experimental’ protocol for reconstruction (number of RWs, duration, temperature).

Our formalism applies to continuously parametrized potentials e.g. RF model with forces taking continuous instead of binary values. The aim is now to predict the true potential values up to some accuracy on each site; this in turn determines an acceptable neighborhood around \mathbf{s}_v^* in the space of signals. The rate function ω_v is generically parabolic around \mathbf{s}_v^* , with a curvature matrix called Fisher information matrix [10]. Finding R_c amounts to minimize this (positive) quadratic form on the boundary of the neighborhood, a task which can be carried out efficiently [17]. Our approach can be easily extended to the case of a finite delay between two measures of the positions, and Chernoff’s result is recovered in the finite N , infinite delay limits [8, 13].

Acknowledgments. We are grateful to D. Thirumalai for his suggestion of illustrating our formalism on the RF model. This work was partially funded by ANR under contract 06-JCJC-051.

-
- [1] B. D. Hughes, *Random walks and random environments*, Oxford University Press (1996).
 - [2] J-P. Bouchaud, A. Georges, *Physics Reports* **195**, 127 (1990).
 - [3] R. Merkel *et al.* *Nature* **397**,50-53 (1999).
 - [4] J.M. Fernandez, H. Li *Science* **303**, 1674 (2004).
 - [5] B. Essevaz-Roulet, U. Bockelmann, F. Heslot, *Proc. Natl. Acad. Sci. (USA)* **94**, 11935 (1997).
 - [6] M.T. Woodside *et al.* *Science* **314**, 1001 (2006).
 - [7] C. Hyeon, D. Thirumalai *Proc Natl Acad Sci USA* **100**,10249-53 (2003).
 - [8] V. Baldazzi *et al.* *Phys. Rev. Lett.* **96** 128102 (2006); *Phys. Rev. E* **75**, 011904 (2007)
 - [9] M. Manosas, D. Collin, F. Ritort *Phys. Rev. Lett.* **96**, 218301 (2006).
 - [10] T.M. Cover, J.A. Thomas, *Elements of Information Theory*, Wiley (1991)
 - [11] L. Ein-Dor, O. Zuk, E. Domany, *Proc. Nat. Acad. Sci. (USA)* **103**, 5923-5928 (2006).
 - [12] T. W. Anderson, L. Goodman *The Annals of Mathematical Statistics* **28**, 89 (1957)
 - [13] H. Chernoff, *Ann. Math. Statis.* **23**, 493 (1952)
 - [14] A. Dembo, O. Zeitouni, *Large deviations Techniques and Applications* Springer-Verlag (1998)
 - [15] L.B. Boza, *Ann. Math. Statis.* **42**, 1992 (1971)
 - [16] S. Cocco, R. Monasson, *in preparation*
 - [17] M.R. Garey, D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W.H. Freeman (1979).
 - [18] O. Adelman, N. Enriquez, *Israel J. Math.* **142**, 205-220 (2004).
 - [19] P. Androletti, *preprint arxiv:math.PR/0612208* (2006).
 - [20] The irrelevance of the *a priori* distribution in the asymptotic case of large data set is well-known [10] and can be

checked for the RF model: the local field (2) is extensive in R , while the global field α remains finite.

[21] Let $\mathbf{v}' \neq \mathbf{v}$; $\mathbf{s}_v^* \cdot (\mathbf{v} - \mathbf{v}') = \sum_{x \neq x'} u^*(x \rightarrow x') G(r_{\mathbf{v}'}(x \rightarrow x')/r_{\mathbf{v}}(x \rightarrow x'))$ where $G(z) = z - \ln z - 1 > 0$ for $z \neq 1$.

[22] $\text{RW}(\mu)$, due to conditioning to survival, is likely to stay for $\sim 1/d(y) \ll t^*$ in the trap only.