
Learning Discrete Distributions from Metastable Data via Pseudo-Likelihood

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Physically motivated stochastic dynamics are standard tools for sampling high-
2 dimensional distributions but often mix slowly due to metastability. We show that
3 for multivariate discrete distributions, *the true stationary model can be learned*
4 *from i.i.d. samples drawn from a metastable distribution*. The key observation is
5 that for strongly metastable states of a reversible chain with stationary distribution
6 μ , the single-variable conditionals of the metastable law are, on average, close to
7 those of μ , even when the two distributions are far in global metrics. This enables
8 accurate parameter recovery with conditional-likelihood estimators such as pseudo-
9 likelihood (PL). We formalize these guarantees and illustrate them numerically on
10 the Curie-Weiss model, where PL succeeds while maximum likelihood fails.

11 Motivation and Overview

12 Markov chains used for sampling may become trapped for long periods in regions of state space,
13 a hallmark of metastability and a principal cause of poor mixing. From a learning standpoint
14 this is troubling: data collected from a natural dynamics or a sampling algorithm may reflect a
15 metastable law ν that can be *globally* far from the true equilibrium μ in total variation or KL
16 divergence. Global-metric-based estimators such as maximum likelihood are therefore ill-suited: they
17 attempt to match sufficient statistics of ν , not μ . A more fruitful perspective is local: many efficient
18 estimators for undirected graphical models learn by *matching single-variable conditionals* (e.g.,
19 pseudo-likelihood/logistic regressions) [2, 4]. Our central finding is that if ν is “strongly metastable”
20 for a reversible chain with equilibrium μ , then the average discrepancy between the single-variable
21 conditionals of ν and those of μ is small. Consequently, PL trained on metastable samples is nearly
22 optimal for μ , and parameters of pairwise Ising models can be recovered (up to a small bias controlled
23 by the metastability level). The effect is especially pronounced in slow-mixing regimes, where the
24 bias decays with system size and is dominated by statistical error.

25 Setup and Main Definitions

26 We consider n discrete variables $\sigma = (\sigma_1, \dots, \sigma_n) \in \mathcal{Q}^n$ and a reversible Markov chain P with
27 stationary distribution μ . We adopt two notions of metastability.

28 **Definition 1 (Metastability)** *A distribution ν is η -metastable (w.r.t. P) if $\|\nu - \nu P\|_{\text{TV}} \leq \eta$.*

29 **Definition 2 (Strong metastability)** *A distribution ν is η -strongly metastable (w.r.t. P) if*

$$\frac{1}{2} \sum_{i,j} |P(i|j)\nu(j) - P(j|i)\nu(i)| \leq \eta.$$

30 Strong metastability implies metastability and captures small violations of detailed balance. It arises
 31 naturally from bottlenecks in probability flow. For a set A with conductance $\Gamma(A)$, the truncated
 32 distribution $\mu_A(\sigma) \propto \mu(\sigma)\mathbf{1}\{\sigma \in A\}$ is η -strongly metastable with $\eta = \Gamma(A)$, hence slow-mixing
 33 chains admit exponentially small η on suitable A . In these cases, we should think of η as scaling as
 34 the inverse of the mixing time of the chain and hence as a quantity that decays fast as function of n .

35 We assume the usual single-site-update setting and a bounded spin-flip probability:

$$\frac{P(\sigma_u \rightarrow q | \sigma)}{\mu(q | \sigma_{\setminus u})} \geq \omega_P > 0, \quad \text{e.g., } \omega_P = \frac{1}{n} \text{ for Glauber dynamics.}$$

36 **Closeness of conditionals.** The central technical statement of our work is the following:

37 **Theorem 1 (Conditionals of strongly metastable states)** *Let ν be η -strongly metastable for a re-*
 38 *versible P with equilibrium μ . If the bounded spin-flip condition holds with constant ω_P , then*

$$\sum_{u=1}^n \mathbb{E}_{\sigma \sim \nu} \|\nu(\cdot | \sigma_{\setminus u}) - \mu(\cdot | \sigma_{\setminus u})\|_{\text{TV}} \leq \frac{\eta}{\omega_P}. \quad (1)$$

39 Thus, even if ν and μ are far in global TV, their *local* conditionals are close on average. Via reverse
 40 Pinsker, the same holds in average conditional KL, which is precisely the PL loss.

41 Learning via Pseudo-Likelihood

42 We model $\mu(\sigma) \propto \exp(E^*(\sigma))$ and estimate within a parametric class $\{E(\sigma; \theta) : \theta \in \Theta\}$ using PL.
 43 Now we assume a finite interaction strength for the models in the parametric class. Now if the data
 44 comes from a metastable state of a Markov chain attempting to sample from $\mu(\sigma)$, then we can show
 45 that the true parameters of the model (θ^*) are almost optimal for the pseudo-likelihood loss function,

46 **Near-optimality of PL on metastable data.** With M' i.i.d. samples from ν , the PL estimator $\hat{\theta}$
 47 satisfies

$$\frac{1}{M'} \sum_{t,u} \left(L_u(\theta^*; \sigma^{(t)}) - L_u(\hat{\theta}; \sigma^{(t)}) \right) \lesssim \underbrace{O\left(\frac{\eta}{\omega_P}\right)}_{\text{bias from metastability}} + O\left(\underbrace{M'^{-\frac{1}{2}}}_{\text{statistics}}\right),$$

48 so in slow-mixing regimes (small η), the bias is dominated by sampling noise. (Constants depend on
 49 the size of the discrete alphabet and the interaction strength)

50 **Ising models: couplings and structure.** For the specific case of Ising models, we can use this
 51 result to show that PL learning from metastable data correctly recovers the parameters of the
 52 true distribution. In this case, $E^*(\sigma) = \sum_{i<j} \theta_{ij}^* \sigma_i \sigma_j + \sum_i \theta_i^* \sigma_i$, ℓ_1 -constrained PL yields high-
 53 dimensional guarantees analogous to the i.i.d. case. With $M = \tilde{O}(e^{8\gamma} \gamma^4 \varepsilon^{-4} \log(n))$ samples from
 54 ν ,

$$\max_{v \neq u} |\hat{\theta}_{uv} - \theta_{uv}^*| \leq \varepsilon + 4e^{2\gamma} \sqrt{\frac{(1+\gamma)\eta}{\omega_P}},$$

55 and edges are recovered by simple thresholding provided nonzero couplings exceed a (metastability-
 56 dependent) level. In particular, for Glauber dynamics $\omega_P = \frac{1}{n}$, any $\eta = o(n^{-1})$ yields vanishing bias
 57 with growing n . The proof here uses a perturbative generalization of the techniques in [3].

58 Illustration: Curie–Weiss Numerics

59 We illustrate the theory on the Curie–Weiss ferromagnet with spins $\sigma_i \in \{\pm 1\}$ as it is the canonical
 60 example of a system that exhibits metastability [1].

$$E(\sigma) = \frac{J}{n} \sum_{i<j} \sigma_i \sigma_j - h \sum_i \sigma_i, \quad m(\sigma) = \frac{1}{n} \sum_i \sigma_i.$$

61 The equilibrium distribution over magnetization concentrates near minima of the free energy $\Psi(m) =$
 62 $-\frac{J}{2}m^2 + hm - S(m)$, where S is the entropy. In the low-temperature, weak-field regime ($J > 1$,
 63 small h), there are two well-separated minima, and Glauber dynamics started at all-plus becomes
 64 trapped near the positive-magnetization well for exponentially long times.

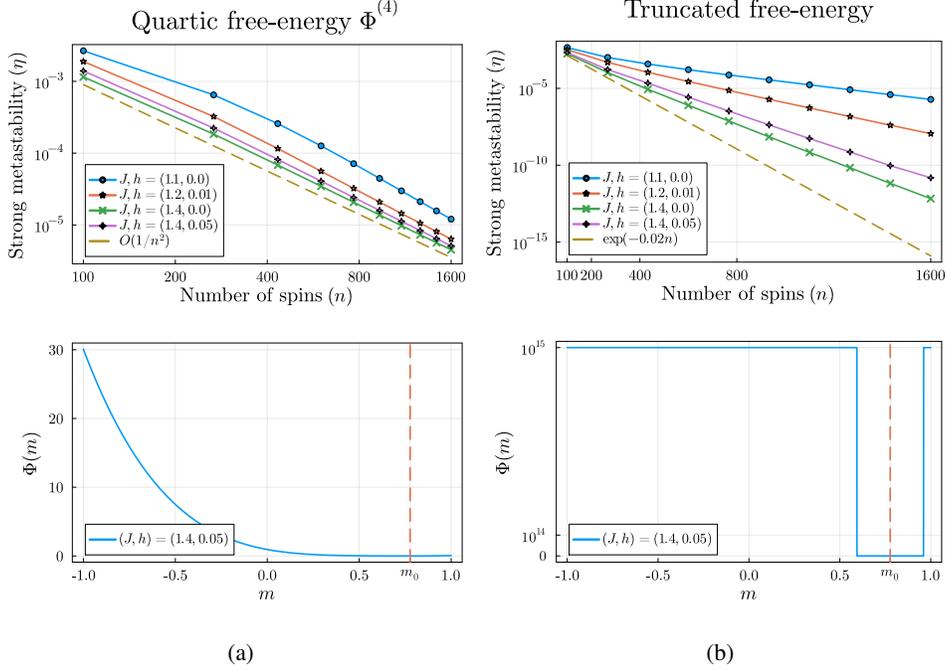


Figure 1: Strongly metastable states in the CW model. Here we plot the violation in detailed balance condition (η in the definition of strong metastability) computed by projecting to the magnetization space. (a) Fourth-order approximation to the free energy at m_0 (b) truncated free energy around m_0 .

65 **Metastable states near free-energy minima.** Thanks to permutation symmetry, the metastability
 66 of this model can be entirely studied in the magnetization space ($m(\sigma) = \frac{\sum_i \sigma_i}{n}$). One can construct
 67 strongly metastable laws by expanding the free energy Ψ around its positive minimizer $m_0 > 0$:

$$\Phi^{(K)}(m) = \sum_{k=2}^K \frac{\Psi^{(k)}(m_0)}{k!} (m - m_0)^k, \quad \nu^{(K)}(\sigma) \propto \exp\{-n\Phi^{(K)}(m(\sigma))\} \binom{n}{\frac{n(1+m(\sigma))}{2}}.$$

68 Numerically we confirm that $\nu^{(4)}$ is $O(n^{-2})$ -strongly metastable for $J > 1$ and small h , and that
 69 truncating Ψ in a narrow window around m_0 yields η decaying exponentially in n . Numerical results
 70 are given in Figure 1.

71 **Pseudo-likelihood succeeds, MLE fails.** We simulate Glauber dynamics in a regime where the
 72 true Gibbs measure strongly favors *negative* magnetization, but the chain started at all-plus remains
 73 stuck near the *positive* well. Training PL on these metastable samples accurately recovers the model
 74 parameters, whereas MLE (which matches global statistics of the data) yields a spurious solution,
 75 even flipping the sign preferred by the true magnetization. Concretely, for $(J, h) = (1.2, 0.04)$ on
 76 large n , PL trained on metastable data recovers parameters with small error, while the negative
 77 log-likelihood landscape attains its minimum far from the ground truth. This numerical experiment
 78 mirrors the theory: local conditionals are close, but global sufficient statistics are not. Details of these
 79 experiments are in Figure 2.

80 Discussion

81 The picture that emerges is simple: metastability hurts global-metric estimators but not conditional
 82 ones. For reversible chains, the *local* information carried by a strongly metastable law is already
 83 almost that of the equilibrium model, and PL exploits exactly this information. In slow-mixing
 84 regimes (small conductance), the resulting bias is tiny compared to sampling noise, so learning from
 85 “wrong” data is still possible, provided we learn from a metastable distribution that has the right local
 86 statistics.

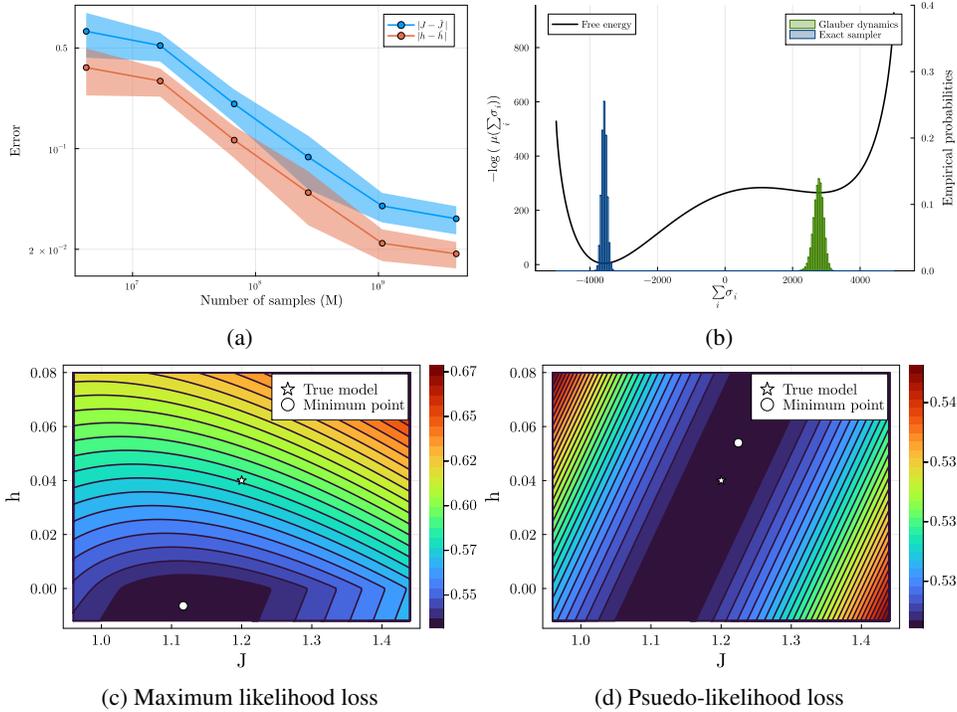


Figure 2: **(a)** Error in learning the Curie-Weiss model on 5000 spins. Samples here are produced by Glauber dynamics “stuck” at the positive minima of the free energy. True parameters here are $J = 1.2, h = 0.04$. **(b)** The true distribution is highly biased towards towards negative magnetization as seen by free energy curve. There is a metastable distribution with positive magnetization that is highly suppressed in terms of probability. The empirical distributions of samples ($M = 4 \times 10^9$) drawn by an exact sampler and Glauber dynamics is overlaid on top of the free energy. This shows that the Markov chain is effectively stuck around the positive minima. **(c) and (d)** Comparison of the loss function landscape for the CW model with true parameters $J = 1.2, h = 0.04$. These are plotted with $M = 2^{32}$ samples produced by Glauber dynamics “stuck” at the positive minima of the free energy. Negative log-likelihood computed from this data clearly has it’s minimum far from the true model. The sign of the magnetization is opposite of the true model. This is expected as maximum likelihood tries to match the sufficient statistics of the data to the model. PL loss function has the minima close to the true model and learns the magnetic field with the right sign.

87 References

- 88 [1] David A Levin, Malwina J Luczak, and Yuval Peres. Glauber dynamics for the mean-field Ising
89 model: cut-off, critical power law, and metastability. *Probability Theory and Related Fields*,
90 146:223–265, 2010.
- 91 [2] Andrey Y Lokhov, Marc Vuffray, Sidhant Misra, and Michael Chertkov. Optimal structure and
92 parameter learning of Ising models. *Science advances*, 4(3):e1700791, 2018.
- 93 [3] Marc Vuffray, Sidhant Misra, Andrey Lokhov, and Michael Chertkov. Interaction screening:
94 Efficient and sample-optimal learning of Ising models. In *Advances in Neural Information
95 Processing Systems*, pages 2595–2603, 2016.
- 96 [4] Shanshan Wu, Sujay Sanghavi, and Alexandros G Dimakis. Sparse logistic regression learns
97 all discrete pairwise graphical models. In *Advances in Neural Information Processing Systems*,
98 pages 8071–8081, 2019.