

# Stochastic Localization Sampling For the SK Model

Mark Sellke

IAS CSDM

Ahmed El Alaoui



Andrea Montanari

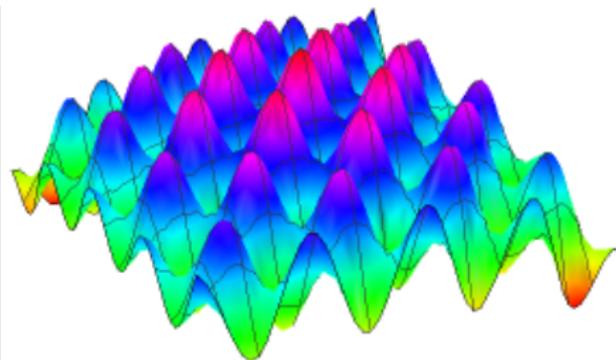
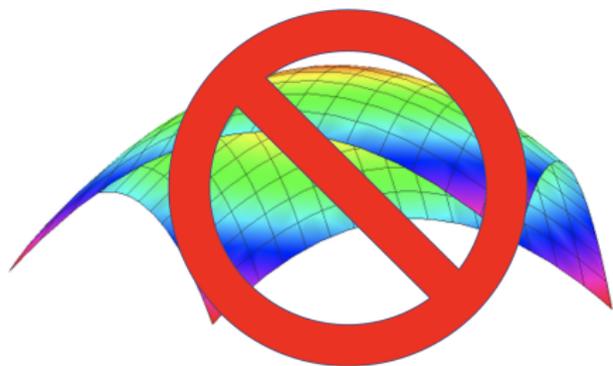


# Sampling

Goal: generate

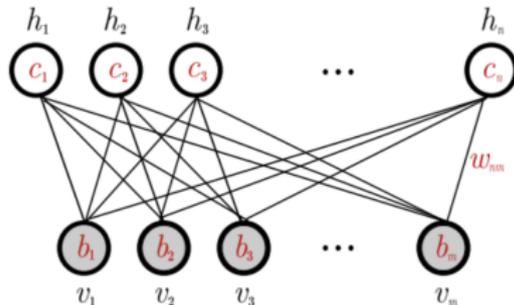
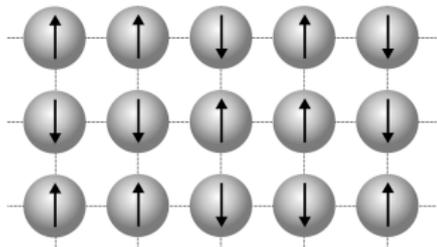
$$\mathbf{x}^* \sim \mu(d\mathbf{x}) \quad \text{given} \quad \mu \in \mathcal{P}(\mathbb{R}^n).$$

For  $\mu$  high-dimensional and NOT log-concave.



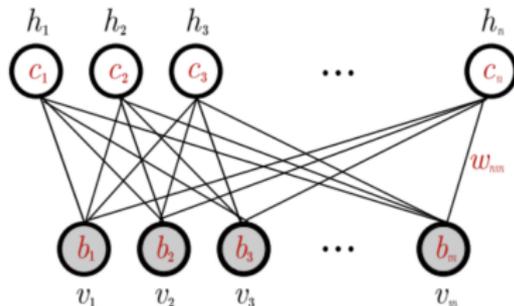
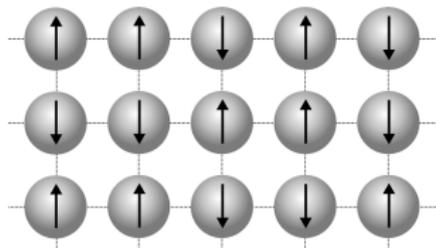
In this talk, focus on Ising models:

$$\mu_{\mathbf{A},\beta}(\mathbf{x}) = \frac{1}{Z(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle / 2}, \quad \mathbf{x} \in \{-1, +1\}^n.$$



In this talk, focus on Ising models:

$$\mu_{\mathbf{A},\beta}(\mathbf{x}) = \frac{1}{Z(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle / 2}, \quad \mathbf{x} \in \{-1, +1\}^n.$$



Glauber dynamics

- Repeatedly choose  $i \in [n]$  and resample  $x_i$  given other coordinates.
- Mixes rapidly if  $\mathbf{A}$  is small. In general, mixing can be very slow.

# Sequential Sampling

Given a distribution  $\mu \in \mathcal{P}(\{-1, +1\}^n)$ , suppose we have a conditional expectation **oracle** to evaluate

$$\mathbf{m}^t = \mathbb{E}^{\mathbf{x} \sim \mu}[\mathbf{x} \mid (x_1 = x_1^*, \dots, x_t = x_t^*)], \quad t \in \{0, 1, \dots, n-1\}.$$

# Sequential Sampling

Given a distribution  $\mu \in \mathcal{P}(\{-1, +1\}^n)$ , suppose we have a conditional expectation **oracle** to evaluate

$$\mathbf{m}^t = \mathbb{E}^{\mathbf{x} \sim \mu}[\mathbf{x} \mid (x_1 = x_1^*, \dots, x_t = x_t^*)], \quad t \in \{0, 1, \dots, n-1\}.$$

Then we can directly sample  $\mathbf{x}$ , one coordinate at a time. Namely,

$$\mathbb{P}^t[x_{t+1} = 1 \mid x_1, \dots, x_t] = \frac{\mathbf{m}_{t+1}^t + 1}{2}.$$

This is the foundation for equivalence between counting and sampling.

# Downsides of Sequential Sampling

Sequential sampling may be too much to hope for.

- Requires a strong oracle, especially for continuous variables.
- Maybe estimating  $m^t$  is no easier than sampling.
- Unclear how to choose a good order for the coordinates.

# Downsides of Sequential Sampling

Sequential sampling may be too much to hope for.

- Requires a strong oracle, especially for continuous variables.
- Maybe estimating  $m^t$  is no easier than sampling.
- Unclear how to choose a good order for the coordinates.

In sequential sampling, we try to reveal  $x^*$  **gradually**.  
There are other ways to do this.

## Warm-Up: Pólya's Urn

A silly way to sample  $p \sim Unif([0, 1])$ :

- Sample an infinite sequence  $(b_1, b_2, \dots)$  of i.i.d.  $Ber(p)$  bits, **without knowing  $p$** .
- Use law of large numbers to compute  $p = \lim_{t \rightarrow \infty} \frac{\sum_{s=1}^t b_s}{t}$ .

## Warm-Up: Pólya's Urn

A silly way to sample  $p \sim \text{Unif}([0, 1])$ :

- Sample an infinite sequence  $(b_1, b_2, \dots)$  of i.i.d.  $\text{Ber}(p)$  bits, **without knowing  $p$** .

- Use law of large numbers to compute  $p = \lim_{t \rightarrow \infty} \frac{\sum_{s=1}^t b_s}{t}$ .

...and this is really not so bad.

- Given  $(b_1, \dots, b_t)$ , the posterior expectation for  $p$  is given by Laplace's rule of succession:

$$\mathbb{E}^t[p] = \frac{1 + \sum_{s=1}^t b_s}{t + 2}.$$

## Warm-Up: Pólya's Urn

A silly way to sample  $p \sim Unif([0, 1])$ :

- Sample an infinite sequence  $(b_1, b_2, \dots)$  of i.i.d.  $Ber(p)$  bits, **without knowing  $p$** .

- Use law of large numbers to compute  $p = \lim_{t \rightarrow \infty} \frac{\sum_{s=1}^t b_s}{t}$ .

...and this is really not so bad.

- Given  $(b_1, \dots, b_t)$ , the posterior expectation for  $p$  is given by Laplace's rule of succession:

$$\mathbb{E}^t[p] = \frac{1 + \sum_{s=1}^t b_s}{t + 2}.$$

- Hence the sequential rule

$$\mathbb{P}^t[b_{t+1} = 1] = \frac{1 + \sum_{s=1}^t b_s}{t + 2},$$

yields an i.i.d.  $Ber(p)$  sequence for  $p \sim Unif([0, 1])$ .

# Stochastic Localization: Revealing $\mathbf{x}^*$ with Gaussian Noise

(A version of) Eldan's Stochastic localization:

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t \sim \mathcal{N}(t\mathbf{x}^*, t\mathbf{I}_n).$$

$\mathbf{x}^* \sim \mu$  is independent of Brownian motion  $\mathbf{B}_t$ .

# Stochastic Localization: Revealing $\mathbf{x}^*$ with Gaussian Noise

(A version of) Eldan's Stochastic localization:

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t \sim \mathcal{N}(t\mathbf{x}^*, t\mathbf{I}_n).$$

$\mathbf{x}^* \sim \mu$  is independent of Brownian motion  $\mathbf{B}_t$ .

Suggests a sampling algorithm:

- 1 Simulate  $\mathbf{y}_t$  for a long time  $t \in [0, T]$  without knowing  $\mathbf{x}^*$ .
- 2 Read off

$$\mathbf{x}^* \approx \frac{\mathbf{y}_T}{T}.$$

# Stochastic Localization: Revealing $\mathbf{x}^*$ with Gaussian Noise

(A version of) Eldan's Stochastic localization:

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t \quad \sim \mathcal{N}(t\mathbf{x}^*, tI_n).$$

$\mathbf{x}^* \sim \mu$  is independent of Brownian motion  $\mathbf{B}_t$ .

Suggests a sampling algorithm:

- 1 Simulate  $\mathbf{y}_t$  for a long time  $t \in [0, T]$  without knowing  $\mathbf{x}^*$ .
- 2 Read off

$$\mathbf{x}^* \approx \frac{\mathbf{y}_T}{T}.$$

Geometric motivation: decompose general  $\mu$  into posteriors

$$\mu_t(dx) \propto e^{\langle \mathbf{y}_t, \mathbf{x} \rangle - t\|\mathbf{x}\|_2^2/2} \mu(dx).$$

- If  $\mu$  log-concave, each  $\mu_t$  is **strongly** log-concave.
- KLS conjecture [Eldan 12, Lee-Vempala 17, Chen 21, Klartag-Lehec 22].

# Simulating $y_t$

Similarly to Pólya's urn, we can generate the path

$$y_t = tx^* + B_t, \quad x^* \sim \mu$$

without knowing  $x^*$ .

Similarly to Pólya's urn, we can generate the path

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t, \quad \mathbf{x}^* \sim \mu$$

without knowing  $\mathbf{x}^*$ . The **annealed** law of  $\mathbf{y}_t$  is described by

$$\begin{aligned}d\mathbf{y}_t &= \mathbf{m}_t dt + d\mathbf{W}_t; \\ \mathbf{m}_t &= \mathbb{E}[\mathbf{x}^* \mid \mathcal{F}_t] = \mathbb{E}[\mathbf{x}^* \mid \mathbf{y}_t]\end{aligned}$$

for  $W_t$  another Brownian motion.

Similarly to Pólya's urn, we can generate the path

$$y_t = tx^* + B_t, \quad x^* \sim \mu$$

without knowing  $x^*$ . The **annealed** law of  $y_t$  is described by

$$\begin{aligned} dy_t &= m_t dt + dW_t; \\ m_t &= \mathbb{E}[x^* \mid \mathcal{F}_t] = \mathbb{E}[x^* \mid y_t] \end{aligned}$$

for  $W_t$  another Brownian motion.

Equivalence:

- Quadratic variation is Brownian in either case.
- $y_t - \int_0^t m_t dt$  is a martingale in either case since  $m_t = \mathbb{E}[x_* \mid \mathcal{F}_t]$ .
- Now use Lévy's characterization of Brownian motion.

$$dy_t = m_t dt + dW_t,$$

A continuous-time stochastic process is not really an algorithm.

Of course, we should discretize.

# The Resulting Algorithm

$$dy_t = m_t dt + dW_t,$$

---

**Input:** Data: Probability measure  $\mu$

**Input:** Result: Sample  $x^* \sim \mu$

**for**  $t \in [0, \delta, \dots, T - \delta]$  **do**

    | Sample  $g_t \sim \mathcal{N}(0, I_n)$

    | Set  $y_{t+\delta} = y_t + \hat{m}_t(y_t)\delta + \sqrt{\delta}g_t$

**end**

Set  $x^* = \text{Round}(y_T/T) \in \{-1, +1\}^n$

**return**  $x^*$

---

# The Resulting Algorithm

$$dy_t = m_t dt + dW_t,$$

---

**Input:** Data: Probability measure  $\mu$

**Input:** Result: Sample  $x^* \sim \mu$

**for**  $t \in [0, \delta, \dots, T - \delta]$  **do**

    | Sample  $g_t \sim \mathcal{N}(0, I_n)$

    | Set  $y_{t+\delta} = y_t + \hat{m}_t(y_t)\delta + \sqrt{\delta}g_t$

**end**

Set  $x^* = \text{Round}(y_T/T) \in \{-1, +1\}^n$

**return**  $x^*$

---

Main requirement: a good approximation  $\hat{m}_t(y_t) \approx \mathbb{E}[x^* \mid y_t]$ .

# Where Do We Stand?

So far:

- General sampling procedure.
- Requires estimating  $\mathbf{m}_t(\mathbf{y}_t) \approx \mathbb{E}[x^* \mid \mathbf{y}_t]$ .

We have replaced the need for one oracle with another...is it any better?

# Where Do We Stand?

So far:

- General sampling procedure.
- Requires estimating  $\mathbf{m}_t(\mathbf{y}_t) \approx \mathbb{E}[x^* \mid \mathbf{y}_t]$ .

We have replaced the need for one oracle with another...is it any better?

Remainder of the talk: example where the answer is **yes**.

- SK model: coupling matrix  $\mathbf{A}$  is GOE.
- Computing  $\mathbf{m}_t(\mathbf{y}_t)$  falls into the wheelhouse of high-dimensional statistics/optimization.

# Sherrington-Kirkpatrick Model

Ising model with random couplings:

$$\mu_{\mathbf{G},\beta}(\mathbf{x}) = \frac{1}{Z_n(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{G} \mathbf{x} \rangle / 2}.$$

Random symmetric matrix  $\mathbf{G} \sim GOE(n)$ :

- $\mathbf{G} = \mathbf{G}^\top$ . Entries otherwise independent.
- $G_{i,j} \sim \mathcal{N}(0, 1/n)$  for  $i < j$ .

Ising model with random couplings:

$$\mu_{\mathbf{G},\beta}(\mathbf{x}) = \frac{1}{Z_n(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{G} \mathbf{x} \rangle / 2}.$$

Random symmetric matrix  $\mathbf{G} \sim GOE(n)$ :

- $\mathbf{G} = \mathbf{G}^\top$ . Entries otherwise independent.
- $G_{i,j} \sim \mathcal{N}(0, 1/n)$  for  $i < j$ .

Goal: given  $\mathbf{G} \sim GOE(n)$ , generate a sample from  $\mu_{\mathbf{G},\beta}$ .

Dobrushin's condition for fast mixing of Glauber works if  $\beta \leq cn^{-1/2}$ .  
But we would like  $\beta$  to be constant size.

# Brief History of the SK Model

[[Ising 1925](#)]: Ising model for ferromagnets.

[[Sherrington-Kirkpatrick 1975](#)]: model for **disordered** magnets.

[[Parisi 1982](#)]: non-rigorous solution via replica symmetry breaking.

[Ising 1925]: Ising model for ferromagnets.

[Sherrington-Kirkpatrick 1975]: model for **disordered** magnets.

[Parisi 1982]: non-rigorous solution via replica symmetry breaking.

[Talagrand 2005] proves the Parisi formula.

- Huge amount of other important work including [Aizenman-Ruelle-Lebowitz 82, Ruelle 87, Chatterjee 09, Panchenko 14, Ding-Sly-Sun 15, Auffinger-Chen 17, . . .].

SK model is a prototype for disordered, random probability measures.

- Random MaxCut and  $K$ -SAT.
- Coloring random graphs.
- Posteriors in high-dimensional statistics.

SK model is a prototype for disordered, random probability measures.

- Random MaxCut and  $K$ -SAT.
- Coloring random graphs.
- Posteriors in high-dimensional statistics.

E.g. optimal MaxCut in a random sparse graph ([Dembo-Montanari-Sen 17]).

For  $G \sim G\left(n, \frac{\lambda}{n}\right)$ :

$$\text{MaxCut}(G) = n \left( \frac{\lambda}{4} + C_* \sqrt{\frac{\lambda}{4}} + o(\sqrt{\lambda}) \right) + o(n).$$

$$\mu_{\mathbf{G},\beta}(\mathbf{x}) = \frac{1}{Z_n(\beta)} e^{\beta\langle \mathbf{x}, \mathbf{G}\mathbf{x} \rangle/2}.$$

Expect: efficient sampling possible for  $\beta < 1$ , impossible for  $\beta > 1$ .

- Replica symmetric iff  $\beta \leq 1$ .

$$\mu_{\mathbf{G},\beta}(\mathbf{x}) = \frac{1}{Z_n(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{G} \mathbf{x} \rangle / 2}.$$

Expect: efficient sampling possible for  $\beta < 1$ , impossible for  $\beta > 1$ .

- Replica symmetric iff  $\beta \leq 1$ .

Recent progress: Glauber mixes in  $O(n \log n)$  steps for  $\beta < 1/4$ .

[Bodineau-Bauerschmidt 20, Eldan-Koehler-Zeitouni 21, Anari-Jain-Koehler-Pham-Vuong 21].

A different method for tensor analogs: [Adhikari-Brennecke-Xu-Yau 22]

$$\mu_{\mathbf{G},\beta}(\mathbf{x}) = \frac{1}{Z_n(\beta)} e^{\beta \langle \mathbf{x}, \mathbf{G} \mathbf{x} \rangle / 2}.$$

Expect: efficient sampling possible for  $\beta < 1$ , impossible for  $\beta > 1$ .

- Replica symmetric iff  $\beta \leq 1$ .

Recent progress: Glauber mixes in  $O(n \log n)$  steps for  $\beta < 1/4$ .

[Bodineau-Bauerschmidt 20, Eldan-Koehler-Zeitouni 21, Anari-Jain-Koehler-Pham-Vuong 21].

A different method for tensor analogs: [Adhikari-Brennecke-Xu-Yau 22]

Our result: stochastic localization succeeds (in a weaker sense) for  $\beta < 1$ .  
(Originally  $\beta < 1/2$ , improvement by [Celentano 22].)

Given  $\mu_1, \mu_2 \in \mathcal{P}(\{-1, 1\}^n)$ , define the normalized Wasserstein metric

$$W_{1,n}(\mu_1, \mu_2) = \inf_{(\mathbf{x}_1, \mathbf{x}_2) \sim \text{Coupling}(\mu_1, \mu_2)} \frac{\mathbb{E}[\|\mathbf{x}_1 - \mathbf{x}_2\|_{\ell^1}]}{n}.$$

$W_{1,n}(\mu_1, \mu_2) \leq o(1)$  means that  $\mathbf{x}_1, \mathbf{x}_2$  differ by  $o(n)$  coordinates under an optimal coupling. We will consider such pairs of points to be close.

Theorem (El Alaoui-Montanari-S 22, Celentano 22)

For any  $\beta < 1$  and  $\varepsilon > 0$ , there exists a randomized algorithm with complexity  $O(n^2)$  which given  $\mathbf{G}$  outputs  $\mathbf{x} \sim \mu_{\mathbf{G},\beta}^{\text{alg}}$  such that

$$\mathbb{E}[W_{1,n}(\mu_{\mathbf{G},\beta}^{\text{alg}}, \mu_{\mathbf{G},\beta})] \leq \varepsilon.$$

# Algorithmic Stability

Our algorithm is **stable** with respect to  $(\mathbf{G}, \beta)$ : just uses  $O_{\beta, \epsilon}(1)$  matrix-vector products, and some 1-dimensional non-linearities.

Concretely, from i.i.d.  $\mathbf{G} = \mathbf{G}_0$  and  $\mathbf{G}_1$ , consider perturbation path

$$\mathbf{G}_s = \sqrt{1-s^2} \mathbf{G}_0 + s \mathbf{G}_1.$$

Stability of the algorithm means:

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_0, \beta}^{\text{alg}}, \mu_{\mathbf{G}_s, \beta}^{\text{alg}})] = 0.$$

# Algorithmic Stability

Our algorithm is **stable** with respect to  $(\mathbf{G}, \beta)$ : just uses  $O_{\beta, \epsilon}(1)$  matrix-vector products, and some 1-dimensional non-linearities.

Concretely, from i.i.d.  $\mathbf{G} = \mathbf{G}_0$  and  $\mathbf{G}_1$ , consider perturbation path

$$\mathbf{G}_s = \sqrt{1-s^2} \mathbf{G}_0 + s \mathbf{G}_1.$$

Stability of the algorithm means:

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_0, \beta}^{\text{alg}}, \mu_{\mathbf{G}_s, \beta}^{\text{alg}})] = 0.$$

A purely structural consequence with an algorithmic proof:

**Theorem (El Alaoui-Montanari-S 22; Celentano 22)**

*The **true** SK Gibbs measures are stable when  $\beta < 1$ :*

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_0, \beta}, \mu_{\mathbf{G}_s, \beta})] = 0.$$

*Similar stability holds for small perturbations in  $\beta$ .*

The stability property

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_{0,\beta}}, \mu_{\mathbf{G}_s,\beta})] = 0.$$

for the true Gibbs measure is **false** for  $\beta > 1$ . Combination of:

The stability property

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_{0,\beta}}, \mu_{\mathbf{G}_{s,\beta}})] = 0.$$

for the true Gibbs measure is **false** for  $\beta > 1$ . Combination of:

Theorem (Chatterjee 09; Disorder Chaos)

Let  $(\mathbf{x}_0, \mathbf{x}_s) \sim \mu_{\mathbf{G}_{0,\beta}} \times \mu_{\mathbf{G}_{s,\beta}}$ . For all  $\beta \in \mathbb{R}$  and  $s > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{E}[|\langle \mathbf{x}_0, \mathbf{x}_s \rangle|/n] = 0.$$

The stability property

$$\lim_{s \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}_{0,\beta}}, \mu_{\mathbf{G}_{s,\beta}})] = 0.$$

for the true Gibbs measure is **false** for  $\beta > 1$ . Combination of:

**Theorem (Chatterjee 09; Disorder Chaos)**

Let  $(\mathbf{x}_0, \mathbf{x}_s) \sim \mu_{\mathbf{G}_{0,\beta}} \times \mu_{\mathbf{G}_{s,\beta}}$ . For all  $\beta \in \mathbb{R}$  and  $s > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{E}[|\langle \mathbf{x}_0, \mathbf{x}_s \rangle|/n] = 0.$$

**Theorem (Replica Symmetry Breaking)**

Let  $\mathbf{x}_0, \mathbf{x}'_0 \sim \mu_{\mathbf{G}_{0,\beta}}$  be independent. For all  $\beta > 1$ ,

$$\liminf_{n \rightarrow \infty} \mathbb{E}[|\langle \mathbf{x}_0, \mathbf{x}'_0 \rangle|/n] \geq c(\beta) > 0.$$

The previous results show that  $\mu_{\mathbf{G}_0, \beta}$  and  $\mu_{\mathbf{G}_s, \beta}$  must be significantly different. Therefore:

## Theorem (El Alaoui-Montanari-S 22)

Let  $\mu_{\mathbf{G}, \beta}^{\text{alg}}$  be the law of  $\text{ALG}_n(\mathbf{G}, \beta, \omega)$  conditional on  $\mathbf{G}$ . If  $\text{ALG}_n$  is *stable*, then for all  $\beta > 1$ ,

$$\liminf_{n \rightarrow \infty} \mathbb{E}[W_{1,n}(\mu_{\mathbf{G}, \beta}^{\text{alg}}, \mu_{\mathbf{G}, \beta})] > c(\beta) > 0.$$

Stability holds for gradient-based methods such as Langevin dynamics and AMP, at least on **dimension-independent** time-scales.

[Back to the Main Story...](#)

To sample for  $\beta < 1$ , our main requirement is to estimate  $\mathbf{m}_t = \mathbb{E}[\mathbf{x}^* \mid \mathbf{y}_t]$  for

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t.$$

The solution goes through several ideas in high-dimensional statistics and optimization.

To sample for  $\beta < 1$ , our main requirement is to estimate  $\mathbf{m}_t = \mathbb{E}[\mathbf{x}^* \mid \mathbf{y}_t]$  for

$$\mathbf{y}_t = t\mathbf{x}^* + \mathbf{B}_t.$$

The solution goes through several ideas in high-dimensional statistics and optimization.

Two phase procedure:

- Rough estimate for  $\mathbf{m}_t$  using approximate message passing.
- High-accuracy estimate for  $\mathbf{m}_t$  using gradient descent on a well-chosen potential.

## Step 1: Rough Estimate of $\mathbf{m}_t$

Self-consistent “naive mean-field” equation for  $\mathbf{m}_t = \mathbb{E}[\mathbf{x} \mid \mathbf{y}_t]$ :

$$\mathbf{m}_t \approx \tanh(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)$$

- Intuitively,  $(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)_i$  is the effective field on  $x_i$ .
- $\tanh(\cdot)$  converts from field on  $\{-1, +1\}$  to probabilities

## Step 1: Rough Estimate of $\mathbf{m}_t$

Self-consistent “naive mean-field” equation for  $\mathbf{m}_t = \mathbb{E}[\mathbf{x} \mid \mathbf{y}_t]$ :

$$\mathbf{m}_t \approx \tanh(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)$$

- Intuitively,  $(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)_i$  is the effective field on  $x_i$ .
- $\tanh(\cdot)$  converts from field on  $\{-1, +1\}$  to probabilities
- Not quite right. It actually should be

$$\mathbf{m}_t = \mathbb{E}^t[\tanh(\beta \mathbf{G} \mathbf{x} + \mathbf{y}_t)].$$

$\tanh(\cdot)$  is non-linear and although  $\mathbb{E}^t[\mathbf{G} \mathbf{x}] = \mathbf{G} \mathbf{m}_t$  there is nontrivial conditional randomness left.

## Step 1: Rough Estimate of $\mathbf{m}_t$

Self-consistent “naive mean-field” equation for  $\mathbf{m}_t = \mathbb{E}[\mathbf{x} \mid \mathbf{y}_t]$ :

$$\mathbf{m}_t \approx \tanh(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)$$

- Intuitively,  $(\beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t)_i$  is the effective field on  $x_i$ .
- $\tanh(\cdot)$  converts from field on  $\{-1, +1\}$  to probabilities
- Not quite right. It actually should be

$$\mathbf{m}_t = \mathbb{E}^t[\tanh(\beta \mathbf{G} \mathbf{x} + \mathbf{y}_t)].$$

$\tanh(\cdot)$  is non-linear and although  $\mathbb{E}^t[\mathbf{G} \mathbf{x}] = \mathbf{G} \mathbf{m}_t$  there is nontrivial conditional randomness left.

Revised Thouless-Anderson-Palmer (TAP) equation:

$$\mathbf{m}_t \approx \tanh \left( \beta \mathbf{G} \mathbf{m}_t + \mathbf{y}_t - \beta^2 \left( 1 - \frac{\|\mathbf{m}_t\|_2^2}{n} \right) \mathbf{m}_t \right).$$

## Step 1: Rough Estimate of $\mathbf{m}_t$

Turn the TAP equation into a **recursion** and repeat until convergence to an approximate **fixed point**:

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right),$$
$$b_k = \beta^2 \left( 1 - \frac{\|\mathbf{m}_t^{(k)}\|_2^2}{n} \right).$$

## Step 1: Rough Estimate of $\mathbf{m}_t$

Turn the TAP equation into a **recursion** and repeat until convergence to an approximate **fixed point**:

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right),$$
$$b_k = \beta^2 \left( 1 - \frac{\|\mathbf{m}_t^{(k)}\|_2^2}{n} \right).$$

- This is an **approximate message passing** algorithm. Generalizes belief propagation to dense matrices  $\mathbf{G}$ .
  - Onsager term  $b_k \hat{\mathbf{m}}_t^{(k-1)}$  cancels “backtracking” paths.

## Step 1: Rough Estimate of $\mathbf{m}_t$

Turn the TAP equation into a **recursion** and repeat until convergence to an approximate **fixed point**:

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right),$$
$$b_k = \beta^2 \left( 1 - \frac{\|\mathbf{m}_t^{(k)}\|_2^2}{n} \right).$$

- This is an **approximate message passing** algorithm. Generalizes belief propagation to dense matrices  $\mathbf{G}$ .
  - Onsager term  $b_k \hat{\mathbf{m}}_t^{(k-1)}$  cancels “backtracking” paths.
  - By now, a major tool in high-dimensional statistics.

[Bolthausen 14, Donoho-Maleki-Montanari 09, Bayati-Montanari 11, Javanmard-Montanari 12, Rush-Venkataramanan 18, Chen-Lam 20, Fan 20, Dudeja-Lu-Sen 22]
- In our case, the AMP state evolution is unclear.  $\mathbf{y}_t = t\mathbf{x}^* + B_t$  for  $\mathbf{x}^* \sim \mu_{\mathbf{G},\beta}$  has a complicated distribution.

# Contiguity with a Simpler Spiked Model

To analyze the AMP recursion, we consider a **spiked** joint distribution  $\mathbb{Q}$  over  $(\mathbf{G}, \mathbf{x}^*, \mathbf{y}_t)$ . Under  $\mathbb{Q}$ :

$$\mathbf{x}^* \sim \text{Unif}(\{-1, 1\}^n), \quad \mathbf{y}_t = t\mathbf{x}^* + B_t,$$

$$\mathbf{G} \sim \text{GOE}(n) + \frac{\beta \mathbf{x} \mathbf{x}^\top}{n}.$$

# Contiguity with a Simpler Spiked Model

To analyze the AMP recursion, we consider a **spiked** joint distribution  $\mathbb{Q}$  over  $(\mathbf{G}, \mathbf{x}^*, \mathbf{y}_t)$ . Under  $\mathbb{Q}$ :

$$\mathbf{x}^* \sim \text{Unif}(\{-1, 1\}^n), \quad \mathbf{y}_t = t\mathbf{x}^* + B_t,$$

$$\mathbf{G} \sim \text{GOE}(n) + \frac{\beta \mathbf{x} \mathbf{x}^\top}{n}.$$

The resulting conditional law  $\mathbb{Q}[\mathbf{G} \mid \mathbf{x}^*]$  looks similar to  $\mathbb{P}[\mathbf{x}^* \mid \mathbf{G}]$  for the SK model:

$$\mathbb{Q}[\mathbf{G} \mid \mathbf{x}^*] \propto e^{\beta \langle \mathbf{x}^*, \mathbf{G} \mathbf{x}^* \rangle / 2}.$$

# Contiguity with a Simpler Spiked Model

To analyze the AMP recursion, we consider a **spiked** joint distribution  $\mathbb{Q}$  over  $(\mathbf{G}, \mathbf{x}^*, \mathbf{y}_t)$ . Under  $\mathbb{Q}$ :

$$\begin{aligned}\mathbf{x}^* &\sim \text{Unif}(\{-1, 1\}^n), & \mathbf{y}_t &= t\mathbf{x}^* + B_t, \\ \mathbf{G} &\sim \text{GOE}(n) + \frac{\beta \mathbf{x} \mathbf{x}^\top}{n}.\end{aligned}$$

The resulting conditional law  $\mathbb{Q}[\mathbf{G} \mid \mathbf{x}^*]$  looks similar to  $\mathbb{P}[\mathbf{x}^* \mid \mathbf{G}]$  for the SK model:

$$\mathbb{Q}[\mathbf{G} \mid \mathbf{x}^*] \propto e^{\beta \langle \mathbf{x}^*, \mathbf{G} \mathbf{x}^* \rangle / 2}.$$

Swapping the order distorts probabilities by a partition function factor

$$Z_{SK}(\mathbf{G}) = \sum_{\mathbf{v} \in \{-1, +1\}^n} e^{\beta \langle \mathbf{v}, \mathbf{G} \mathbf{v} \rangle / 2}.$$

- $Z_{SK}(\mathbf{G})$  fluctuates **mildly** for  $\beta < 1$  [Aizenman-Ruelle-Lebowitz 82]. The spiked model is **contiguous** with the original.

# State Evolution for AMP

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right)$$

Idea of AMP: for fixed  $\mathbf{v}, \mathbf{w}$ , the vectors

$$(\mathbf{G}\mathbf{v}, \mathbf{G}\mathbf{w})$$

each have i.i.d. Gaussian coordinates. Covariance between  $(\mathbf{G}\mathbf{v})_i$  and  $(\mathbf{G}\mathbf{w})_j$  equals  $\langle \mathbf{v}, \mathbf{w} \rangle$ .

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right)$$

Idea of AMP: for fixed  $\mathbf{v}, \mathbf{w}$ , the vectors

$$(\mathbf{G}\mathbf{v}, \mathbf{G}\mathbf{w})$$

each have i.i.d. Gaussian coordinates. Covariance between  $(\mathbf{G}\mathbf{v})_i$  and  $(\mathbf{G}\mathbf{w})_j$  equals  $\langle \mathbf{v}, \mathbf{w} \rangle$ .

- **Onsager term** lets us apply this recursively to each  $\hat{\mathbf{m}}_t^{(k+1)}$ , despite accumulating dependence on  $\mathbf{G}$ .
- In spiked model, correlation with  $x_i$  also enters the recursion.

# State Evolution for AMP

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right)$$

Idea of AMP: for fixed  $\mathbf{v}, \mathbf{w}$ , the vectors

$$(\mathbf{G}\mathbf{v}, \mathbf{G}\mathbf{w})$$

each have i.i.d. Gaussian coordinates. Covariance between  $(\mathbf{G}\mathbf{v})_i$  and  $(\mathbf{G}\mathbf{w})_j$  equals  $\langle \mathbf{v}, \mathbf{w} \rangle$ .

- **Onsager term** lets us apply this recursively to each  $\hat{\mathbf{m}}_t^{(k+1)}$ , despite accumulating dependence on  $\mathbf{G}$ .
- In spiked model, correlation with  $x_i$  also enters the recursion.

State evolution:  $i$ -th coordinate of  $\hat{\mathbf{m}}_t^{(k)}$  behaves like

$$\tanh(a_t^{(k)} x_i + b_t^{(k)} Z), \quad Z \sim \mathcal{N}(0, 1).$$

- $(a_t^{(k)}, b_t^{(k)})$  determined recursively, converge to  $(a_t^\infty, b_t^\infty)$ .

# State Evolution for AMP

From  $(a_t^\infty, b_t^\infty)$ , one can read off the asymptotic MSE

$$E_* = \lim_{k \rightarrow \infty} \text{p-lim}_{n \rightarrow \infty} \mathbb{E} \|\hat{\mathbf{m}}_t^{(k)} - \mathbf{x}\|_2^2.$$

# State Evolution for AMP

From  $(a_t^\infty, b_t^\infty)$ , one can read off the asymptotic MSE

$$E_* = \lim_{k \rightarrow \infty} \text{p-lim}_{n \rightarrow \infty} \mathbb{E} \|\hat{\mathbf{m}}_t^{(k)} - \mathbf{x}\|_2^2.$$

If we can show

$$E_* \approx \text{MMSE}(t) \equiv \mathbb{E} \|\mathbf{m}_t - \mathbf{x}\|_2^2,$$

then we conclude  $\hat{\mathbf{m}}_t^{(k)} \approx \mathbf{m}_t$ .

# State Evolution for AMP

From  $(a_t^\infty, b_t^\infty)$ , one can read off the asymptotic MSE

$$E_* = \lim_{k \rightarrow \infty} \text{p-lim}_{n \rightarrow \infty} \mathbb{E} \|\hat{\mathbf{m}}_t^{(k)} - \mathbf{x}\|_2^2.$$

If we can show

$$E_* \approx \text{MMSE}(t) \equiv \mathbb{E} \|\mathbf{m}_t - \mathbf{x}\|_2^2,$$

then we conclude  $\hat{\mathbf{m}}_t^{(k)} \approx \mathbf{m}_t$ .

I-MMSE Area Law [Guo-Shamai-Verdu 04, Deshpande-Abbe-Montanari 15]:

$$\int_0^\infty \text{MMSE}(t) dt = 2 \cdot \text{Ent}(\mathbf{x}^*).$$

- Verify explicitly that  $\int_0^\infty E_*(t)$  asymptotically matches  $\text{Ent}(\mathbf{x}^*)$ .

## Conclusion of Step 1: Rough Estimate for $\mathbf{m}_t$

$$\hat{\mathbf{m}}_t^{(k+1)} = \tanh \left( \beta \mathbf{G} \hat{\mathbf{m}}_t^{(k)} + \mathbf{y}_t - b_k \hat{\mathbf{m}}_t^{(k-1)} \right),$$

Proposition (El Alaoui-Montanari-S 22)

For  $\beta < 1$  and any  $\varepsilon, t \geq 0$  there exists  $k_0(t, \varepsilon)$  such that for all  $k \geq k_0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ \|\hat{\mathbf{m}}_t^{(k)} - \mathbf{m}_t\| \leq \varepsilon \sqrt{n} \right] = 1.$$

## Step 2: Refined Estimate of $\mathbf{m}_t$

Surprisingly, this is not quite enough.

- Two types of error: SDE  $\delta$ -discretization and  $\hat{\mathbf{m}}_t^{(k)} \approx \mathbf{m}_t$ .
- Simply sending  $(\delta, k) \rightarrow (0, \infty)$  doesn't work. Not Lipschitz enough.

## Step 2: Refined Estimate of $\mathbf{m}_t$

Surprisingly, this is not quite enough.

- Two types of error: SDE  $\delta$ -discretization and  $\hat{\mathbf{m}}_t^{(k)} \approx \mathbf{m}_t$ .
- Simply sending  $(\delta, k) \rightarrow (0, \infty)$  doesn't work. Not Lipschitz enough.

Second step: by construction,  $\hat{\mathbf{m}}_t^{(k)}$  is an approximate stationary point for the TAP free energy:

$$F_{TAP}(\mathbf{m}, \mathbf{y}_t) = -\frac{\beta}{2} \langle \mathbf{m}, \mathbf{G} \mathbf{m} \rangle - \langle \mathbf{y}_t, \mathbf{m} \rangle - \sum_{i=1}^n h(m_i).$$

- Refine  $\hat{\mathbf{m}}_t^{(k)}$  to  $\hat{\mathbf{m}}_t = \arg \min_{\mathbf{m}} F_{TAP}(\mathbf{m}, \mathbf{y}_t)$  via gradient descent.

## Step 2: Refined Estimate of $\mathbf{m}_t$

Surprisingly, this is not quite enough.

- Two types of error: SDE  $\delta$ -discretization and  $\hat{\mathbf{m}}_t^{(k)} \approx \mathbf{m}_t$ .
- Simply sending  $(\delta, k) \rightarrow (0, \infty)$  doesn't work. Not Lipschitz enough.

Second step: by construction,  $\hat{\mathbf{m}}_t^{(k)}$  is an approximate stationary point for the TAP free energy:

$$F_{TAP}(\mathbf{m}, \mathbf{y}_t) = -\frac{\beta}{2} \langle \mathbf{m}, \mathbf{G} \mathbf{m} \rangle - \langle \mathbf{y}_t, \mathbf{m} \rangle - \sum_{i=1}^n h(m_i).$$

- Refine  $\hat{\mathbf{m}}_t^{(k)}$  to  $\hat{\mathbf{m}}_t = \arg \min_{\mathbf{m}} F_{TAP}(\mathbf{m}, \mathbf{y}_t)$  via gradient descent.
- [Celentano 22]:  $F_{TAP}$  is strongly convex near  $\mathbf{m}_t$  for  $\beta < 1$ .  
Implies  $\mathbf{y}_t \mapsto \hat{\mathbf{m}}_t$  is  $C_\beta$ -Lipschitz. ( $\mathbf{y}_t \mapsto \hat{\mathbf{m}}_t^{(k)}$  is  $C_\beta^k$ -Lipschitz.)

This type of algorithm must be completely impractical, right?

Not quite...

Recall:

$$\begin{aligned} \mathbf{m}_t(\mathbf{y}_t) &= \mathbb{E}[\mathbf{x} \mid \mathbf{y}_t], \quad \mathbf{y}_t = t\mathbf{x}_t + \sqrt{t}\mathbf{g}, \quad \mathbf{g} \sim \mathcal{N}(0, I_n), \\ \mathbf{m}_t &= \arg \min_{\phi: \mathbb{R}^n \rightarrow \mathbb{R}^n} \mathbb{E}[\|\phi(\mathbf{y}_t) - \mathbf{x}\|_2^2]. \end{aligned}$$

i.e.:

**Bayes-optimal inversion of Gaussian noise suffices to sample.**

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be i.i.d. natural images. Generate noisy versions  $\mathbf{y}_i$ .

Choose  $\hat{\mathbf{m}}_t = \phi(\mathbf{y}_i)$  minimizing empirical loss

$$\frac{1}{n} \sum_{i=1}^n \|\phi(\mathbf{y}_i) - \mathbf{x}_i\|_2^2$$

...for  $\phi \in \mathcal{F}$  constrained inside some **function class** such as **convolutional neural networks**.

# Image Generation

These are diffusion models! [Song-Ermon 19], DALL-E 2, Imagen.



# Image Generation

These are diffusion models! [Song-Ermon 19], DALL-E 2, Imagen.



- Equivalent setup: turn  $x \sim \mu$  into Gaussian noise with OU flow. Then simulate the time-reversal (corresponds to  $y_t/t$ ).
- Mean-estimation is done using “forward” sample paths.

# Image Generation

These are diffusion models! [Song-Ermon 19], DALL-E 2, Imagen.



- Equivalent setup: turn  $x \sim \mu$  into Gaussian noise with OU flow. Then simulate the time-reversal (corresponds to  $y_t/t$ ).
- Mean-estimation is done using “forward” sample paths.

[Chen-Chewi-Li-Li-Salim-Zhang 22, Lee-Lu-Tan 22a,22b,22c]: estimating  $m_t$  in  $L^2$  suffices for sampling if  $y_t \mapsto m_t$  is **globally Lipschitz**.

- For us: proxy  $\hat{m}_t$  is **typically locally** Lipschitz near the sample path.

Stochastic localization for the SK model: interaction with high-dimensional probability enables a rigorous, end-to-end analysis.

Our algorithm produces Wasserstein-approximate samples for  $\beta < 1$ . For  $\beta > 1$ , disorder chaos is a natural barrier for stable algorithms.

- What other distributions are stochastic localization sampleable?
- Sharp thresholds in related models.
  - **Shattering** may obstruct efficient sampling even when replica symmetric. Absent in SK model, expected for pure spherical  $p$ -spin.