When MCMC Meets Variational Methods

Frederic Koehler (Stanford => UChicago)

Holden Lee (JHU)



Andrej Risteski (CMU)









Reminder:

Ising models,

sampling,

and all that



Ising Model

 $p_{J,h}(x) \propto \exp\left(\frac{1}{2}\langle x, Jx \rangle\right)$

$\Box J: n \times n$ is an arbitrary symmetric interaction matrix, and $h \in \mathbb{R}^n$ is a vector of bias/external field

General class including models of magnetism, spin glasses, neurons, social networks, Bayesian statistics...

□ Maximum entropy distribution. Analogue of Gaussian.

$$x \rangle + \langle h, x \rangle \bigg), \qquad x \in \{\pm 1\}^n$$





Sampling and Statistical Inference

 \Box Sampling: given parameters J, h generate $X \sim p_{Jh}$ Many connections to statistical inference. **Bayesian posteriors often Ising.** Ising is the posterior on $X \sim Uni\{\pm 1\}^n$ given $Y = X + N(0,\Sigma)$ for some Σ, Y Sampling used to optimally estimate J,h from samples (via MLE.)





stochastic block model (community detection) posterior is Ising



Alternatives to MLE (maximum likelihood estimation) are usually not as accurate...



How to sample? Typically MCMC. Especially nice is the Gibbs sampler/Glauber dynamics: \Box Let initial $X \in \{\pm 1\}^n$ be arbitrary. For t = 1 to T : \square Pick $i \in [n]$ randomly and resample $X_i \mid X_{\sim i}$ But: it is hard to know if it is working! Sometimes it doesn't work! How many steps to mix to correct distribution? Can be exponential time... Some Ising models are NP-hard to sample from.





What can we say theoretically?



Useful example: Curie-Weiss model

 $p(x) = \frac{1}{Z} \exp(\frac{1}{Z})$

\Box Rapid mixing for $\beta < 1$ (Dobrushin).

n

When $\beta > 1$: bottleneck emerges between two clusters of spins. Gibbs sampler is exponentially slow to mix (runtime c^n)...

$$xp\left(\frac{\beta}{2n}\left(\sum_{i}x_{i}\right)^{2}\right)$$

 $x_i \approx m^*$

n



Gibbs Sampler at High Temperature

Rapid mixing of Gibbs/Glauber for "small" J is a generic phenomena. Lots of work for many years...

sampler mixes rapidly.

Not true for "larger" J due to bottleneck (previous slide). What are the alternatives to Gibbs?

Theorem [..., **BB** '19, **EKZ** '21, **AJKPV** '22]: if $\lambda_{max}(J) - \lambda_{min}(J) < 1$, then Gibbs



Variational Inference Picture

□ Variational inference: approximate by a simpler distribution. Popular alternative to MCMC, comes from statistical physics.



□ When $\beta > 1$, Curie-Weiss model is close to a mixture of two product measures centered at fixed points of "mean-field equation" $m^* = \tanh(\beta m^* + h)$





Structure of low-rank Ising models

Rigorous Naive Mean-Field Approximation [..., Eldan-Gross '18, Austin '19,...] shows approximate mixture of product decomposition for all low-rank J.

- \Box If rank J = o(n) then Ising \approx mixture of $2^{o(n)}$ product measures (with means m satisfying $m \approx \tanh(Jm + h)$).
- \Box Only a rough approximation (o(n) in Wasserstein).
- □ Not constructive (unless you use our results!)

 $m_5 \approx \tanh(Jm_5 + h)$

 $m_3 \approx \tanh(Jm_3 +$

 $p_{J,h} \approx$ $m_1 \approx \tanh(Jm_1 + h)$

 $m_4 \approx \tanh(Jm_4 + h)$

 $m_2 \approx \tanh(Jm_2 + h)$

Mixture of products



MCMC versus Variational Inference

□ Variational inference world

GOOD: makes sense in multimodal/low-temperature settings.

□ BAD: only approximates the true distribution, structural results are not algorithmic.

□ MCMC world

GOOD: GOOD: when it works, it really samples from the true distribution.

BAD: Gibbs sampler fails in multimodal case. Unclear how to fix...



MCMC meets Variational Inference

Theorem (this work): new sampler for approximately low-rank Ising models! Runtime parameterized by (# of spectral outliers/"threshold rank" of J).

Size one interval

$$J = J_{LR} + J_{small} \qquad \lambda$$



Runtime $n^{O(\#outliers)}$. Close to lower bound from ETH [JKR '19]. Negative eigenvalues must be O(1) in size. Large negative eigenvalue problem is NP-hard.

spectral outliers

 $\lambda_{max}(J_{\text{small}}) - \lambda_{min}(J_{\text{small}}) < 1$





Proof Ideas



Key: a new structural result

□ We decompose the Ising model as a mixture of high-temperature Ising models: $(J_{eff}, h^j)_{i=1}^M$ with $||J_{eff}||_{OP}$ small. □ Unlike mixture of product measures: decomposition is (1) constructive, and (2) very accurate.

Enables polynomial-time sampling of the real distribution!

 $p_{J,h} \approx (J_{eff},h^2)$

lsing (J_{eff}, h^1)

lsing (J_{eff}, h^4)

Mixture of hightemperature models!



Two key principles



O(1) negative outliers

Size one interval (containing zero)

II. Eliminate negative eigenvalues: replace by an "effective field".

O(1) positive outliers

I. Eliminate positive eigenvalues: low-dimensional decomposition



Step I: positive eigenvalues

Suppose $J = J_{LR} + J'$ (LR = low rank) with J_{LR} PSD and J' small.

Goal: Goal: find a mixing distribution q over a low-dimensional space so that approximately,

$p_{J,h} \approx \int p_{J',h+b} q(b) db$

 $p_{J,h+b}$

Mixture of high-temperature Ising with weights q(b)



Hubbard-Stratonovich transform [Hubbard '58]. Quadratic to linear reduction:

 $e^{\langle x,J_{LR}x\rangle/2} = \frac{1}{(2\pi)^{d/2}} \int_{span(J_{LR})} e^{\langle x,J_{LR}x\rangle/2} e^{\langle x,J_{LR}x\rangle/2} d^{\prime} d$

 \Box Using $J = J' + J_{IR}$ yields mixture decomposition: $e^{\langle x,Jx\rangle/2} \propto \int_{span(J_{LR})} e^{\langle x,Jx\rangle/2} \propto e^{\langle x,Jx\rangle/2} \propto e^{\langle x,Jx\rangle/2} e^{\langle x,Jx\rangle/2} \propto e^{\langle x,Jx\rangle/2} e^{$ $J_{span}(J_{LR})$

$$e^{\langle J_{LR}^{1/2}x,z\rangle}e^{-\|z\|^{2}/2}dz$$

$$e^{\langle x,J'x\rangle/2+\langle J_{LR}^{1/2}x,z\rangle}e^{-\|z\|^2/2}dz$$

 $\propto \qquad p_{J',J_{LR}^{1/2}z}(x) \ Z_{J',J_{LR}^{1/2}z} e^{-\|z\|^2/2} \, dz$

q(z)



Two key principles



O(1) negative outliers

Size one interval

II. Eliminate negative eigenvalues: replace by an "effective field".

O(1) positive outliers

I. Eliminate positive eigenvalues: low-dimensional decomposition (sketch done!)



Step II: handling negative eigenvalues

 \Box Trick: write $J = J_+ - J_-$ with $J_+, J_- \geq 0$ and run SGD on $G(\mu) := \log \mathbb{E}_{P_{J+h}}[e$ \Box Why? Postulate that J_x concentrates near deterministic quantity J_μ . \Box Then $\langle x, Jx \rangle \approx \langle x, J_+x \rangle - \langle x, J_-\mu \rangle$ so $P_{J,h} \approx P_{J_+,h-J_-\mu}$ (kill negative eigenvalues) $\Box \text{ In particular } J_{\mu} \approx \mathbb{E}_{P_{J,h}}[J_x] \approx \mathbb{E}_{P_{J+,h-J_{\mu}}}[J_x]$ (so $\nabla G(\mu) \approx 0$) FACTS: (1) G has a critical point, (2) at any critical point $P_{J,h} \approx P_{J_+,h-J_-\mu}$ for rejection sampling

$$\left[\left\langle \mu, -J_{x} \right\rangle \right] + \left\langle \mu, J_{\mu} \right\rangle / 2$$



The Simplified Algorithm

□ Now we know there exists an approximate mixture decomposition into "hot" models: $p_{J,h}(x) \approx p_{J',h+b-J_{\mu}(b)}(x) q(b) db$

□ Yields a natural sampling algorithm:

 \Box (1) Riemann integration for integral, (2) SGD+Glauber to compute critical point $\mu(b)$, (3) Glauber dynamics for $p_{I'}$, (4) rejection sampling.

Using this gives suboptimal $poly(1/\epsilon)$ runtime in error ϵ . We can get $\log(1/\epsilon)$ runtime by designing a "simulated tempering" chain. (see paper)



 $p_{J',h'}$

 $||J'||_{OP} \le 1$

Example Application



Task: Dense MAX-CUT

 \Box Input : adjacency matrix A of a graph on n vertices with $\Theta(n^2)$ edges

Goal: find a set $S \subseteq [n]$ to maximize size of cut $\#\{(u, v) \in E : u \in S, v \in S^C\}$



 \square

 $\mathbf{\mathcal{O}}$





"Statistical Physicsy" Approach

 $\Box \text{ Let OPT} := \frac{|E|}{2n} + \max_{x \in \{\pm 1\}^n} \frac{-1}{4n} \langle x, Ax \rangle.$

 \Box Gibbs measure at inverse temperature $\beta \geq 0$:

Exercise (typical sample is a $(1 + 1/\beta)$ -apx to OPT): $\mathbf{OPT} \ge \frac{|E|}{2\beta n} + \mathbb{E}_{p_{\beta}} \left| \frac{-1}{4n} \langle x, Ax \rangle \right| \ge \mathbf{OPT} - \frac{n \log 2}{\beta}$

Exercise (few large eigenvalues): $\frac{1}{n^2} \sum_{i=1}^n \lambda_i(A)^2 = O(1)$

((1/n) * Dense MAX-CUT so OPT is $\Theta(n)$)

 $p_{\beta}(x) \propto \exp\left(\frac{-\beta}{4n}\langle x, Ax \rangle\right)$



Approximation from Sampling

 \Box Yields a $(1 + 1/\beta)$ approximation to Dense MAX-CUT. **D** Matches $n^{O(1/\epsilon^2)}$ runtime [AKK '92] for Dense MAX-CUT but now get all near-optimal cuts! \Box No computational phase transition in β ! Sampling easier than optimization ! **Q:** What other PTAS's can be found by "just" sampling the Gibbs measure?

 \Box CORR: for any $\beta \ge 0$, can sample Gibbs measure $p_{\beta}(x) \propto \exp\left(\frac{-\beta}{4n}\langle x, Ax \rangle\right)$ in time $n^{O(\beta^2)}$



Conclusion

□ We developed a new algorithm which provably samples from all "approximately low rank" Ising models.

□ Many other interesting applications: Hopfield networks, Ferromagnetic SK, Contextual SBM, mixture models, Ising models on expanders, ...

□ Variational inference ideas help us predict which problems are tractable.

□ What aspects of this story extend beyond Ising ?





Thanks!

