# 6.S896 - Algorithmic Statistics

## Lecture 6: MRF Structure Learning (from finite samples)

So far: Introduced MRFs, GRFs, Ising Models
 └ capture cond. independence structure
    expressible via undirected graph

Showed that MRF is exploitable assumption
  to test hypotheses about high-dim dist'ns
  from polynomial in the dimension samples
    e.g. for uniformity testing (lectures 3,4)

In terms of learning:

- can identify <u>structure</u> of tree-structured MRF
  using infinitely-many samples (Chow-Liu)

- can learn tree-structured Ising model
  in <u>total variation distance</u> computationally
  efficiently using a tight $\Theta(\frac{n\log n}{\varepsilon^2})$ -samples

[Open question: continuous tree-structured MRFs]

 └ upper bound: Chow-Liu
   lower bound: Fano

How about learning general MRFs?
  └ structure learning : today
    TV learning : Thursday

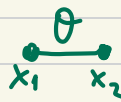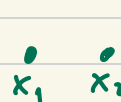# Structure Learning of MRFs from finite samples

Focus: Ising models
[Similar ideas: Gibbs Random Fields]

$$p(x) \propto \exp\left( \sum_{i \neq j} \theta_{ij} x_i x_j + \sum_i \theta_i x_i \right)$$

Infeasible goal: identify support of $(\theta_{ij})_{ij}$ matrix

i.e. all pairs $(i,j)$ s.t. $\theta_{ij} \neq 0$

why infeasible? B.c. for any finite $N$ number of
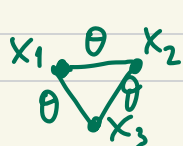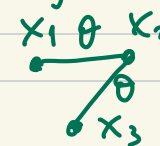samples exists small enough $\theta > 0$
can't distinguish w pr of error $\leq 0.49$

between $\underset{x_1 \quad x_2}{\bullet\!-\!\overset{\theta}{}\!-\!\bullet}$ and $\underset{x_1 \quad x_2}{\bullet \quad \bullet}$

More Realistic goal: identify $(i,j)$ pairs s.t. $\theta_{ij}$ "large enough"

↳ depends on #samples

Another caveat though: can still get confused
if $\theta_{ij}$ are too large

for any finite $N$ number of samples exists large
enough $\theta > 0$ s.t. can't distinguish wpr of error $\leq 0.49$

between $\underset{x_3}{\overset{x_1 \;\; \theta \;\; x_2}{\bigtriangledown}}$ and $\underset{x_3}{\overset{x_1 \; \theta \; x_2}{\diagup}}$

thus  #samples should depend on the "edge strengths"

<u>Theorem</u> [Klivans-Meka'17, Rigollet&Hutter'17, Wu-Sanghavi-Dimakis'19]
                improving on Bresler'15, Vuffray-Misra-Lokhov-Chertkov'16

Exist polynomial-time algorithm which,

given  $N \geq \dfrac{\lambda^2 \exp(12 \cdot \lambda) \log(n/\delta)}{\varepsilon^4}$  samples

ignores
constant
factors   from an Ising model  $p(x) \propto \exp\left(\sum\limits_{i \neq j} \theta_{ij} x_i x_j + \sum\limits_{j} \theta_i x_i\right)$

where $\lambda = \lambda(\theta) \triangleq \max\limits_{i} \left\{\sum\limits_{j \neq i} |\theta_{ij}| + |\theta_i|\right\}$ , outputs

$(\hat{\theta}_{ij})_{ij}$   s.t. with prob. $\geq 1-\delta$ satisfy:

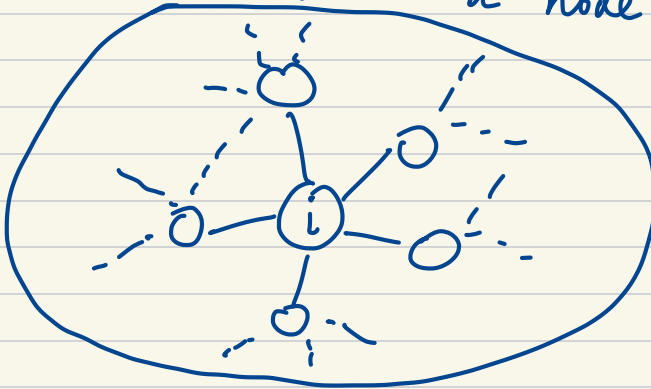$$\forall i,j : \left|\hat{\theta}_{ij} - \theta_{ij}\right| < \varepsilon$$

<u>Corollary</u>: If all $\theta_{ij} \neq 0$ satisfy  $|\theta_{ij}| > \eta > 0$ we
can identify support of $(\theta_{ij})_{ij}$ matrix using

$$N \geq \dfrac{\lambda^2}{\eta^4} \exp(12 \cdot \lambda) \cdot \log(n/\delta) \text{ samples.}$$

with success prob. $\geq 1-\delta$.

[Wainwright-Santhanam] : lower bound on $N \geq \dfrac{2^{\lambda/4} \log n}{\eta \, 2^{3\eta}}$
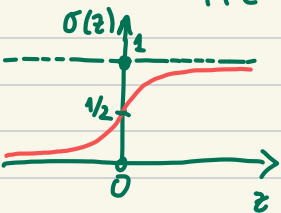
# Idea for algorithm: look at the neighborhood of a node



$$pr(X_i = s \mid X_{-i}) = \frac{\exp\left(\sum_{j \neq i} \theta_{ij} x_j s + \theta_i s\right)}{\exp\left(\sum_{j \neq i} \theta_{ij} x_j s + \theta_i s\right) + \exp\left(-\sum_{j \neq i} \theta_{ij} x_j s - \theta_i s\right)}$$

$$= \frac{1}{1 + \exp\left(-2 \cdot s \cdot \left(\sum_{j \neq i} \theta_{ij} x_j + \theta_i\right)\right)}$$

$$\simeq \frac{1}{1 + \exp\left(-2 \cdot s \cdot \left(\langle \theta_{i \cdot}, x_{-i} \rangle + \theta_i\right)\right)}$$

vector $(\theta_{ij})_{j \neq i}$

sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$= \sigma\left(2 \cdot s \cdot \left(\langle \theta_{i \cdot}, x_{-i} \rangle + \theta_i\right)\right)$$



**Idea:** think of $X_i$ as $\{\pm 1\}$ - outcome: $Y$
& $X_{-i}$ as feature vector: $\vec{z}$
in a linear logistic model

$$Pr[Y \mid \vec{z}] = \sigma\left(s \cdot (2\langle \theta_{i \cdot}, \vec{z} \rangle + 2\theta_i)\right)$$

want to estimate

- **Plan:** estimate neighborhood-by-neighborhood
  do logistic regression to estimate $\theta_{i\cdot}$ and $\theta_i$

- From now on focus on node $i$

  - given samples $X^{(1)}, X^{(2)}, \dots, X^{(N)}$ from $P_\theta$
    create dataset for logistic regression

$$Y^{(\ell)} = X_i^{(\ell)} \;;\; Z^{(\ell)} = X_{-i}^{(\ell)} \quad,\quad \ell = 1,\dots,N$$

  - do $\overset{\text{empirical}}{\text{MLE}}$ to estimate logistic model

$$\underset{\vec{w},c}{\min} \; \frac{1}{N} \sum_{\ell=1}^{N} \log\left(1 + e^{-Y^{(\ell)}\left(\langle w, Z^{(\ell)}\rangle + c\right)}\right)$$

suppose $\rightarrow$ $(\hat{w}, \hat{c})$
is argmin $|w|_1 + c \leq \lambda$

$\hat{L}(w,c)$ : empirical negative
log likelihood

  - consider also <u>population negative log likelihood</u>
    $\llcorner$ means w/ infinite many samples

$$L(w,c) = \underset{\substack{(Y,Z)\sim P_\theta \\ (X_i, X_{-i})}}{\mathbb{E}}\left[\log\left(1 + e^{-Y(\langle w,Z\rangle + c)}\right)\right]$$

known fact!
$$(w^*, c^*) = (\theta_{i\cdot}, \theta_i) \text{ is an optimal solution}$$
$$\text{to} \quad \underset{\substack{w,c \\ |w|_1 + c \leq \lambda}}{\min} L(w,c)$$

goal: compare $(\hat{w}, \hat{c})$ with $(w^*, c^*)$

want to show they are close w.pr $\geq 1-\delta$
over the randomness in $X^{(1)}, \ldots, X^{(N)} \sim P_\theta$
if $N$ is large enough

step 1: If $N \geq \Omega\left(\lambda^2 \log \frac{n}{\delta} / \gamma^2\right)$ then:

(will show)
later

$$L(\hat{w}, \hat{c}) - L(w^*, c^*) \leq \gamma \ , \ \text{w pr} \geq 1-\delta$$

step 2: for any $w, c$:

(will show)
later

$$L(w,c) - L(w^*, \hat{c}) \geq 2 \cdot \underset{z = X_{-i} \sim P_\theta}{\mathbb{E}} \left[ \left( \sigma(\langle w, z \rangle + c) - \sigma(\langle w^*, z \rangle + c^*) \right)^2 \right]$$

step 3:
(will show)
later

$\forall \ w, c, w', c' :$

$$\underset{z = X_{-i} \sim P_\theta}{\mathbb{E}} \left[ \left( \sigma(\langle w, z \rangle + c) - \sigma(\langle w', z \rangle + c') \right)^2 \right] \leq \gamma$$

$$\Rightarrow \|w - w'\|_\infty \leq e^{\|w\|_1 + |c| + \|w'\|_1 + |c'|} \cdot \sqrt{16 \gamma \cdot e^{2\lambda}}$$

Putting everything together: choose $\gamma \leq O\left(\varepsilon^2 \cdot e^{-6\lambda}\right)$

use step 1
+
step 2  setting $(w, c) = (\hat{w}, \hat{c})$
+
step 3  setting $(w, c) = (\hat{w}, \hat{c})$
           $(w', c') = (w^*, c^*)$

$\Rightarrow N \geq \Omega\left(\lambda^2 e^{12\lambda} \log \frac{n}{\delta} / \varepsilon^4\right)$
as promised

$\Rightarrow \|\hat{w} - w^*\|_\infty \leq \varepsilon$
w pr $\geq 1-\delta$

Made with Goodnotes

## Step 1

Lemma 1 (see eg. Shalev-Shwartz & Ben-David book)

Suppose $(z, Y) \sim D$ s.t. wpr 1 under $D$: $|z|_\infty \le 1$
$\rightarrow \in \{\pm 1\}$
$\hookrightarrow \in \mathbb{R}^{q-1}$

Take $L(w,c) = \underset{(z,Y) \sim D}{\mathbb{E}} \left[ \log \left( 1 + e^{-Y \cdot (\langle w, z \rangle + c)} \right) \right]$

$$\hat{L}(w,c) = \frac{1}{N} \sum_{\ell=1}^{N} \left[ \log \left( 1 + e^{-Y_i^{(\ell)}(\langle w, z_i^{(\ell)} \rangle + c)} \right) \right]$$

where $(z^{(1)}, Y^{(1)}), \ldots, (z^{(N)}, Y^{(N)}) \overset{iid}{\sim} D$

Then wpr $\ge 1 - \delta$, for all $w, c$ s.t. $|w|_1 + c \le \lambda$:

$$L(w, c) \le \hat{L}(w,c) + 2 \cdot \lambda \cdot \sqrt{\frac{2 \log(2n)}{N}} + \lambda \sqrt{\frac{2 \log(2/\delta)}{N}}$$

Proof: via Rademacher complexity analysis. $\boxtimes$

if $N \ge \Omega\left( \lambda^2 \log(n/\delta)/\gamma^2 \right)$

Lemma 1 $\Rightarrow$ $L(\hat{w}, \hat{c}) - L(w^*, c^*) \le O(\gamma)$, wpr $\ge 1 - \delta$

why? B.c. $L(\hat{w}, \hat{c}) \le \hat{L}(\hat{w}, \hat{c}) + O(\gamma)$   (By lemma 1 & choice of N)

$\le \hat{L}(w^*, c^*) + O(\gamma)$   (optimality of $(\hat{w}, \hat{c})$ for $\hat{L}$)

$\le L(w^*, c^*) + O(\gamma)$

$\underset{}{\llcorner}$ by Chernoff: $\hat{L}(w^*, c^*) - L(w^*, c^*)$

are $\sqrt{\frac{\log(1/\delta)}{N}}$

close w pr $\ge 1 - \delta$

## Step 2

**Lemma 2.** In same setting as Lemma 1, suppose
$$\Pr_D[Y=1\mid z] = \sigma(\langle w^*, z\rangle + c^*) \quad \text{for some } (w^*, c^*)$$

Then:
$$L(w,c) - L(w^*, c^*) \ge 2\cdot \mathbb{E}_z\left[\left(\sigma(\langle w,z\rangle + c) - \sigma(\langle w^*, z\rangle + c^*)\right)^2\right]$$

**Proof:**

$$L(w,c) - L(w^*, c^*) = \underset{(z,Y)\sim D}{\mathbb{E}}\left[-\frac{Y+1}{2}\log\left(\sigma(\langle w,z\rangle + c)\right) - \frac{1-Y}{2}\log\left(1-\sigma(\langle w,z\rangle + c)\right)\right.$$

$$\left.+ \frac{Y+1}{2}\log\left(\sigma(\langle w^*, z\rangle + c^*)\right) + \frac{1-Y}{2}\log\left(1-\sigma(\langle w^*, z\rangle + c^*)\right)\right]$$

← claim 1 (next page)

$$\Pr[Y=1\mid z] = \sigma(\langle w^*, z\rangle + c) \qquad = \underset{z}{\mathbb{E}}\underset{Y\mid z}{\mathbb{E}}\left[\frac{Y+1}{2}\log\frac{\sigma(\langle w^*, z\rangle + c^*)}{\sigma(\langle w, z\rangle + c)} + \frac{1-Y}{2}\log\frac{1-\sigma(\langle w^*, z\rangle + c^*)}{1-\sigma(\langle w, z\rangle + c)}\right]$$

$$\downarrow$$

$$= \underset{z}{\mathbb{E}}\left[\sigma(\langle w^*, z\rangle + c^*)\log\frac{\sigma(\langle w^*, z\rangle + c^*)}{\sigma(\langle w, z\rangle + c)} + \left(1-\sigma(\langle w^*, z\rangle + c^*)\right)\cdot\log\frac{1-\sigma(\langle w^*, z\rangle + c^*)}{1-\sigma(\langle w, z\rangle + c)}\right]$$

$$= \underset{z}{\mathbb{E}}\left[KL\left(\text{Bernoulli}\left(\sigma(\langle w^*, z\rangle + c^*)\right)\,\|\,\text{Bernoulli}\left(\sigma(\langle w, z\rangle + c)\right)\right)\right]$$

$$\ge \underset{z}{\mathbb{E}}\left[2\cdot\left(\sigma(\langle w^*, z\rangle + c^*) - \sigma(\langle w, z\rangle + c)\right)^2\right]$$

$$\underset{\frac{1}{2}KL(\text{Bernoulli}(p)\,\|\,\text{Bernoulli}(q)) \ge (p-q)^2 \text{ by Pinsker}}{}$$

**Claim 1:** $L(w,c) = \underset{(z,Y) \sim D}{\mathbb{E}} \left[ - \frac{Y+1}{2} \log \left( \sigma(\langle w,z \rangle + c) \right) \right.$

$$\left. - \frac{1-Y}{2} \log \left( 1 - \sigma(\langle w,z \rangle + c) \right) \right]$$

**Proof:** $L(w,c) = \underset{(z,Y) \sim D}{\mathbb{E}} \left[ -\log \left( \sigma \left( Y \cdot (\langle w,z \rangle + c) \right) \right) \right]$

$$= \underset{(z,Y) \sim D}{\mathbb{E}} \left[ - \frac{1+Y}{2} \cdot \log \left( \sigma(\langle w,z \rangle + c) \right) \right.$$

$$\left. - \frac{1-Y}{2} \log \left( \underbrace{\sigma(-\langle w,z \rangle - c)}_{1 - \sigma(\langle w,z \rangle + c)} \right) \right] \boxtimes$$

$\boxed{\text{Step 3}}$

<u>Lemma 3</u> : Suppose $\underline{X} \sim P_\Theta$ $\leftarrow$ Ising & $\lambda(\theta) = \max_i \left( \sum_j |\theta_{ij}| + \theta_i \right)$

Then $\min_i \min_{s \in \{\pm 1\}} \min_{S_{-i}} Pr[X_i = s \mid X_{-i} = s_{-i}] \geq \frac{1}{2} e^{-2\lambda(\theta)}$

<u>Proof:</u> $Pr[X_i = s \mid X_{-i} = s_{-i}] = \dfrac{1}{1 + exp\left( -2 \cdot \left( \sum_j \theta_{ij} s_j + \theta_i \right) \cdot s \right)}$

$$\geq \dfrac{1}{1 + exp\left( 2\left( \sum_j |\theta_{ij}| + |\theta_i| \right) \right)}$$

$$\geq \dfrac{1}{1 + exp(2\lambda(\theta))} \geq \frac{1}{2} e^{-2\lambda(\theta)} \qquad . \ \boxtimes$$

<u>Def:</u> A dist'n $D$ over boolean vectors $X$ is $J$-unbiased iff for all $i$, $\forall s \in \{\pm 1\}$, $\forall s_{-i}$ : $Pr[X_i = s \mid X_{-i} = s_{-i}] \geq J$.

Ising model is $\left( \frac{1}{2} e^{-2\lambda(\theta)} \right)$-unbiased.

**Lemma 4:** Suppose dist'n $D$ over $\{\pm 1\}$-vectors is $J$-unbiased. Then $\forall w, w', c, c'$:

$$\underset{z \sim D}{\mathbb{E}}\left[\left(\sigma(\langle w, z\rangle + c) - \sigma(\langle w', z\rangle + c')\right)^2\right] \leq \gamma$$

$$\Rightarrow \|w - w'\|_\infty \leq \overbrace{e^{\|w\|_1 + |c| + \|w'\|_1 + |c'|}}^{=: \Lambda} \cdot \sqrt{\frac{8\gamma}{J}}$$

**Proof:** Pick arbitrary coordinate, say coordinate $k$:

$$\gamma \geq \underset{z_{-k}}{\mathbb{E}}\left[\underset{z_k}{\mathbb{E}}\left[\left(\sigma(\langle w, z\rangle + c) - \sigma(\langle w', z\rangle + c')\right)^2 \mid z_{-i}\right]\right]$$

$$\geq \underset{z_{-k}}{\mathbb{E}}\left[\Pr[z_k = 1 \mid z_{-k}] \cdot \left(\sigma(w_k + \underbrace{\langle w_{-k}, z_{-k}\rangle + c}_{A(z_{-k})}) - \sigma(w_k' + \underbrace{\langle w'_{-k}, z_{-k}\rangle + c'}_{B(z_{-k})})\right)^2\right.$$

$$\left. + \Pr[z_k = -1 \mid z_{-k}] \cdot \left(\sigma(-w_k + A(z_{-k})) - \sigma(-w_k' + B(z_{-k}))\right)^2\right]$$

<span style="color:green">claim 2<br>next<br>page</span>

$$\geq \underset{z_{-k}}{\mathbb{E}}\left[\frac{J \cdot e^{-2\Lambda}}{16} \cdot |w_k - w_k' + A(z_{-k}) - B(z_{-k})|^2 \right.$$

$$\left. + \frac{J \cdot e^{-2\Lambda}}{16} \cdot |w_k' - w_k + A(z_{-k}) - B(z_{-k})|^2\right]$$

$$\geq J \cdot e^{-2\Lambda} \cdot \frac{1}{8}|w_k - w_k'|^2 \quad \Rightarrow \quad |w_k - w_k'| \leq e^{\Lambda}\sqrt{\frac{8\gamma}{J}} \quad \boxtimes$$

<span style="color:gray">Made with Goodnotes</span>

**Claim 2:** $\forall\ x, y \in \mathbb{R}:\ |\sigma(x) - \sigma(y)| \geq \frac{1}{4}\,e^{-|x|}\cdot e^{-|y|}\cdot|y-x|$

**Proof:** $|\sigma(x) - \sigma(y)| = \left|\dfrac{1}{1+e^{-x}} - \dfrac{1}{1+e^{-y}}\right|$

$\left(\begin{array}{c}\text{suppose} \\ \text{wlog } y \geq x\end{array}\right)$

$$= \frac{|e^{-y} - e^{-x}|}{(1+e^{-x})(1+e^{-y})} = \frac{e^{-y}|1 - e^{y-x}|}{(1+e^{-x})(1+e^{-y})}$$

$$\geq \frac{|y-x|}{(1+e^{-x})\cdot(1+e^{y})}$$

$$\geq \frac{e^{-|x|-|y|}}{4}\cdot|y-x| \qquad \boxtimes$$