

6.S896 - Algorithmic Statistics

Lecture 7: MRF Distribution Learning (from finite samples)

- Motivation: estimating parameters/structure requires upper/lower bounds on the edge strengths, o.w. it cannot be done from a finite sample

philosophical Question: if it takes a gazillion number of samples to distinguish whether or not an edge is present, in what applications does it really matter to find it?

- An alternative Goal: given finite # samples learn the target distribution in some distance

↳ TV, KL, ...

↳ could mean proper learning: return a distribution from the same family that is close to target
e.g. proper learning an Ising model must return an Ising model

or improper learning: return any distribution that is close to target

From Parameter to Distribution Learning?

In settings where we can estimate parameters, we commonly can also do dist'n learning

e.g. suppose in Ising Model

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

we have estimated all θ_{ij} 's & θ_i 's to high enough accuracy:

$$\forall i \neq j: |\hat{\theta}_{ij} - \theta_{ij}| \leq \frac{\varepsilon^2}{n^2}$$

↑ we can do this when $\lambda(\theta) = \max_i (\sum_{j \neq i} |\theta_{ij}| + |\theta_i|) < \infty$

using $\sqrt{2} \left(\frac{\lambda^2 \cdot e^{12\lambda} \cdot \text{poly}(n)}{\varepsilon^8} \right)$

samples using main theorem of lecture 6

$$\forall i: |\hat{\theta}_i - \theta_i| \leq \frac{\varepsilon^2}{n}$$

↑ can also be done from aforementioned number of samples using techniques of lecture 6 (small extension of lemma 4)

recall:
 $\text{SKL}(p, q) = \text{KL}(p \parallel q) + \text{KL}(q \parallel p)$

Now: claim 1: $\text{SKL}(P_\theta, P_{\hat{\theta}}) \leq 4\varepsilon^2$

$$\Rightarrow \text{TV}(P_\theta, P_{\hat{\theta}}) \leq \varepsilon$$

b.c. by Pinsker:

$$\text{TV}(p, q) \leq \sqrt{\frac{1}{2} \text{KL}(p \parallel q)}$$

proof of claim 1: from lecture 3/4

$$\text{SKL}(P_\theta, P_0) = \sum_{i \neq j} (\theta_{ij} - \hat{\theta}_{ij}) \cdot \left(\mathbb{E}_\theta [x_i x_j] - \mathbb{E}_{\hat{\theta}} [x_i x_j] \right)$$
$$+ \sum_i (\theta_i - \hat{\theta}_i) \cdot \left(\mathbb{E}_\theta [x_i] - \mathbb{E}_{\hat{\theta}} [x_i] \right)$$

b.c.

$$|\mathbb{E}_\theta x_i x_j - \mathbb{E}_{\hat{\theta}} x_i x_j| \leq 2$$
$$|\mathbb{E}_\theta x_i - \mathbb{E}_{\hat{\theta}} x_i| \leq 2$$

$$\leq n^2 \cdot \frac{\epsilon^2}{n^2} \cdot 2 + n \cdot \frac{\epsilon^2}{n} \cdot 2 \leq 4\epsilon^2$$

• Can we do distribution learning without paying for low temperature?

& more generally, without estimating parameters?

enters ... Tournament-Based Approach

↳ very general, quick & dirty method to understand sample complexity of distribution learning

↳ also provides algorithm to select among a set of hypotheses distributions one that is close to some target

- Tournament-Based Approach: 3 steps at a high level

↳ ① For the distance of interest d , the accuracy of interest ϵ , and the family of dist'ns of interest \mathcal{H} find a small ϵ -cover $\mathcal{H}_{(d, \epsilon)} \subseteq \mathcal{H}$

↳ for all $p \in \mathcal{H}$ $\exists q \in \mathcal{H}_{(d, \epsilon)}$ s.t. $d(p, q) < \epsilon$

- ② Define a pairwise comparison procedure

$\text{comp}(P; H_1, H_2)$: output¹ H_1 as winner
 want this to
 only use sample access to P
 which is unknown dist'n

² H_2 as winner
³ a tie

- ③ Set-up a tournament among all distributions in $\mathcal{H}_{(d, \epsilon)}$ which runs a sequence of pairwise comparisons of distributions in $\mathcal{H}_{(d, \epsilon)}$ to select a winner

WANT: winner of tournament is close to unknown P

- Today: tournament for TV distance

[Derroyé - Lugosi '96] : Scheffé Estimate
 [Yatracos '85], [Daskalakis et al '12], [Daskalakis - Kamath '14]
 [Acharya et al '14]

Can do other distances too [e.g. Feldman - O'Donnell - Srivastava '16]

• Main Primitive: $\text{Comp}(P; H_1, H_2)$ pairwise comparison procedure

(A) { uses: sample access to $P, H_1, H_2 \in \mathcal{H}$
PDF comparator: given x , compares $H_1(x) \& H_2(x)$
will assume PDFs exist }

behavior that I want: Comp draws $m = m(\epsilon, \delta)$ samples from P, H_1, H_2 and has the following properties:

For $(i,j) = (1,2)$ or $(2,1)$:

1. If $d(P, H_i) \leq \epsilon$ & $d(P, H_j) > 8\epsilon$:

H_i is declared winner, w.pr. $> 1 - \delta$

2. If $d(P, H_i) \leq \epsilon$ & $d(P, H_j) > 4\epsilon$:

H_i wins or it's a tie, w.pr. $> 1 - \delta$

3. If $d(H_1, H_2) \leq 5\epsilon$:

it's a tie, w.pr. $> 1 - \delta$

Lemma 1: For $d = TV$ distance & assuming (A) there exists a poly-time algorithm satisfying behavior (B) where:

$$m = \Theta\left(\frac{\log(1/\delta)}{\epsilon^2}\right).$$

- Suppose we have the following

(A') {

- we have a comparator which assumes (A) and satisfies (B)
- we have a cover $\mathcal{H}_{(d,\epsilon)}$ of our set of dist'ns \mathcal{H}
- we have sample access to some unknown $P \in \mathcal{H}$.

How do we set-up a tournament to select some $H \in \mathcal{H}$ that is close in d to P ?

Lemma 2: Assume (A') . There exists an algorithm that:

- draws $m(\epsilon, \frac{\delta}{2 \cdot N})$ samples from P
- makes $\binom{N}{2}$ comparisons
- outputs H s.t. $d(P, H) < 8\epsilon$ w.p. $> 1 - \delta$

Daskalakis-Kamath
 Acharya et al
 improve this
 to $N \cdot \log N$

Lemma 1 + Lemma 2 \Rightarrow

Corollary: For $d = \text{TV}$ distance, suppose we have:

- a set of dist'ns \mathcal{H} & ε -cover \mathcal{H}_ε of cardinality N ;
- sample access to unknown $P \in \mathcal{H}$;
- pdf comparator for all $H_1, H_2 \in \mathcal{H}_\varepsilon$;
- sample access to all $H \in \mathcal{H}_\varepsilon$;

Then there is an algorithm that

- draws $O\left(\frac{\log(N/\delta)}{\varepsilon^2}\right)$ samples from P (as well as all $H \in \mathcal{H}_\varepsilon$)
- runs in time $O\left(N^2 \cdot \frac{\log(N/\delta)}{\varepsilon^2}\right) \rightsquigarrow$ can be improved to $N \cdot \log N \cdot \frac{\log(N/\delta)}{\varepsilon^2}$
- and outputs H s.t. $d(P, H) \leq 8\varepsilon$, wpr $\geq 1 - \delta$.

Application of Corollary to learn Ising Models in TV

For $\lambda > 0$ define:

$$\mathcal{H}(\lambda) = \{ \text{all Ising models s.t. } \max_i \left(\sum_{j \neq i} |\theta_{ij}| + \sum_j |\theta_{ij}| \right) \leq \lambda \}$$

Theorem: Given $\Omega \left(\frac{n^2 \log \frac{n\lambda}{\varepsilon}}{\varepsilon^2} \right)$ samples from Ising model $p_\theta \in \mathcal{H}(\lambda)$, can find Ising model $p_{\theta'} \in \mathcal{H}(\lambda)$ such that $TV(p_\theta, p_{\theta'}) \leq \varepsilon$, w.p.r.t. $\geq 1 - \delta$

Proof: Recall: $SKL(p_\theta, p_{\theta'}) \leq \sum_{i \neq j} |\theta_{ij} - \theta'_{ij}| + \sum_i |\theta_i - \theta'_i|$,

$$\Rightarrow \begin{cases} SKL(p_\theta, p_{\theta'}) \leq \frac{\varepsilon^2}{16} \\ TV(p_\theta, p_{\theta'}) \leq \frac{\varepsilon}{8} \end{cases} \quad \begin{array}{l} \text{if } |\theta_{ij} - \theta'_{ij}| \leq \frac{\varepsilon^2}{32n^2} \forall ij \\ \text{and } |\theta_i - \theta'_i| \leq \varepsilon/32n, \forall i \end{array}$$

\Rightarrow gridding over θ_{ij} 's & θ_i 's exists $\frac{\varepsilon}{8}$ -cover of \mathcal{H}_1 in TV that has size $N \leq \left(\frac{32\lambda n^2}{\varepsilon^2} \right)^{n^2} \left(\frac{32\lambda n}{\varepsilon^2} \right)^n$ (using that all Ising models in $\mathcal{H}(\lambda)$ satisfy that $|\theta_{ij}| \leq \lambda$ $|\theta_i| \leq \lambda$)

tournament corollary

$$\Rightarrow \frac{\log(N/\delta)}{\varepsilon^2} \text{ samples suffice to learn within } \varepsilon \text{ in TV} \quad \checkmark$$

Remark: Can strengthen Theorem to remove dependence on λ entirely! [see Theorem 9 of Braverman-Cai-Daskalakis '20]

Theorem 9 from Brustle - Cai - Daskalakis '20:

Theorem 9 (Learnability of MRFs in Total Variation and Prokhorov Distance). Suppose we are given sample access to an MRF p , as in Definition 7, defined on an unknown graph with hyper-edges of size at most d .

- **Finite alphabet Σ :** Given $\frac{\text{poly}(|V|^d, |\Sigma|^d, \log(\frac{1}{\epsilon}))}{\epsilon^2}$ samples from p we can learn some MRF q whose hyper-edges also have size at most d such that $\|p - q\|_{TV} \leq \epsilon$. If the graph on which p is defined is known, then $\frac{\text{poly}(|V|, |E|, |\Sigma|^d, \log(\frac{1}{\epsilon}))}{\epsilon^2}$ -many samples suffice. Moreover, the polynomial dependence of the sample complexity on $|\Sigma|^d$ cannot be improved, and the dependence on ϵ is tight up to $\text{poly}(\log \frac{1}{\epsilon})$ factors.
- **Alphabet $\Sigma = [0, H]$:** If the log potentials $\phi_v(\cdot) \equiv \log(\psi_v(\cdot))$ and $\phi_e(\cdot) \equiv \log(\psi_e(\cdot))$ for every node v and every edge e are C -Lipschitz w.r.t. the ℓ_1 -norm, then given $\text{poly}\left(|V|^{d^2}, \left(\frac{H}{\epsilon}\right)^d, C^d\right)$ samples from p we can learn some MRF q whose hyper-edges also have size at most d such that $\|p - q\|_P \leq \epsilon$. If the graph on which p is defined is known, then $\text{poly}\left(|V|, |E|^d, \left(\frac{H}{\epsilon}\right)^d, C^d\right)$ -many samples suffice.

↑ Prokhorov Distance

Definition 7. A Markov Random Field (MRF) is a distribution defined by a hypergraph $G = (V, E)$. Associated with every vertex $v \in V$ is a random variable X_v taking values in some alphabet Σ , as well as a potential function $\psi_v : \Sigma \rightarrow [0, 1]$. Associated with every hyperedge $e \subseteq V$ is a potential function $\psi_e : \Sigma^e \rightarrow [0, 1]$. In terms of these potentials, we define a probability distribution p associating to each vector $x \in \Sigma^V$ probability $p(x)$ satisfying:

$$p(x) = \frac{1}{Z} \prod_{v \in V} \psi_v(x_v) \prod_{e \in E} \psi_e(x_e), \quad (1)$$

where for a set of nodes e and a vector x we denote by x_e the restriction of x to the nodes in e , and Z is a normalization constant making sure that p , as defined above, is a distribution. In the degenerate case where the products on the RHS of (1) always evaluate to 0, we assume that p is the uniform distribution over Σ^V . In that case, we get the same distribution by assuming that all potential functions are identically 1. Hence, we can in fact assume that the products on the RHS of (1) cannot always evaluate to 0.



→ this is general no matter what d is and
 compl., ...)
 is as long
 as it satisfies
 desiderata
 stated above

Proof of Lemma 2:

- For all pairs $H_i, H_j \in \mathcal{H}_\varepsilon$ run

$\text{comp}(P, H_i, H_j)$

using $m(\varepsilon, \frac{\delta}{2N})$ samples

- Output any H_k that never lost (win or tied in every comparison it participated in)
 If no such H_k exists output "failure"

- Claim 2.1: Suppose $H^* \in \mathcal{H}_\varepsilon$ satisfies $d(P, H^*) \leq \varepsilon$.
 { exists b.c. \mathcal{H}_ε is a cover}

With prob. $\geq 1 - \frac{\delta}{2}$, H^* will never lose
 thus the algorithm will not output "failure".

Proof: Take any $H' \in \mathcal{H}_\varepsilon$.

- If $d(P, H') > 4\varepsilon$, by property B.2 of the comparator, H' either won or tied against H^*

- If $d(P, H') \leq 4\varepsilon \Rightarrow d(H, H') \leq 5\varepsilon$
 so by property B.3 of the comparator,
 H' tied w/ H' wpr $\geq 1 - \frac{\delta}{2N}$

doing a union bound over all $H' \in \mathcal{H}_\varepsilon$ proves
 the claim \square

Claim 2.2: If $H \in \mathcal{H}_\epsilon$ never lost, then $d(P, H) \leq 8\epsilon$,
w.pr. $\geq 1 - \frac{\delta}{2}$.

Proof: Suppose $d(P, H) > 8\epsilon$, and take H^* s.t. $d(P, H^*) \leq \epsilon$.

By property B.1 of the comparator H loses to H^*
w.pr. $\geq 1 - \frac{\delta}{2N}$

By union bound, all $H \in \mathcal{H}_\epsilon$ s.t. $d(P, H) > 8\epsilon$,
lose to H^* with prob. $\geq 1 - \frac{\delta}{2}$. \square

Claim 2.1 + Claim 2.2 \Rightarrow w.pr. $\geq 1 - \delta$, algorithm does not
fail & its output
dist'n is 8ϵ -close
to P . \square

→ this is special for TV

Proof of Lemma 1: Set up a competition between H_1, H_2 .

• Define set $\mathcal{W}_1 = \mathcal{W}(H_1 \parallel H_2) = \{x \text{ s.t. } H_1(x) > H_2(x)\}$

↳ to determine for some x whether $x \in \mathcal{W}_1$ will use PDF comparator

• Define $p_1 = H_1(\mathcal{W}_1)$, $p_2 = H_2(\mathcal{W}_1)$

clearly $p_1 > p_2$ and $TV(H_1, H_2) = p_1 - p_2$

• Here is how $\text{Comp}(P, H_1, H_2)$ works

1a. draw $m = O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$ samples from P & call $\hat{\tau}$ the fraction of them that fall in \mathcal{W}_1

1b. draw m samples from H_1 , call \hat{p}_1 the fraction in \mathcal{W}_1

1c. $-/-/-/-$ $-/-H_2$, $-/-\hat{p}_2$ $-/-$

2. if $\hat{p}_1 - \hat{p}_2 \leq 6\epsilon$, declare a tie

3. If $\hat{\tau} > \hat{p}_1 - 2\epsilon$, declare H_1 winner

4. If $\hat{\tau} < \hat{p}_2 + 2\epsilon$, declare H_2 winner

5. declare a tie.

Claim: For $(i, j) = (1, 2)$ or $(2, 1)$, suppose $TV(P, H_i) \leq \varepsilon$. Then:

1. If $TV(P, H_j) > 8\varepsilon$, then

H_i is declared winner w.p. $\geq 1 - 6e^{-m\varepsilon^2/2}$ (B.1)

2. If $TV(P, H_j) > 4\varepsilon$, then

H_i wins or it's a tie w.p. $\geq 1 - 6e^{-m\varepsilon^2/2}$ (B.2)

3. If $TV(H_1, H_2) \leq 5\varepsilon$, then

it's a tie w.p. $\geq 1 - 6e^{-m\varepsilon^2/2}$

Proof: • By Chernoff, w.p. $\geq 1 - 6e^{-m\varepsilon^2/2}$ the following are true:

$$|P_1 - \hat{P}_1| < \frac{\varepsilon}{2}, \quad |P_2 - \hat{P}_2| < \frac{\varepsilon}{2}, \quad |\bar{z} - \hat{z}| < \frac{\varepsilon}{2} \quad (*)$$

• Assuming (*) we have:

case 1: $TV(P, H_i) \leq \varepsilon$, $TV(P, H_j) > 8\varepsilon \Rightarrow \underbrace{TV(H_i, H_j)}_{\substack{\text{triangle} \\ || \\ P_1 - P_2}} > 7\varepsilon$

$$\Rightarrow \hat{P}_1 - \hat{P}_2 \geq P_1 - P_2 - \varepsilon > 6\varepsilon$$

\Rightarrow algorithm will go beyond step 2

Now $TV(P, H_i) \leq \varepsilon \Rightarrow |\bar{z} - P_i| \leq \varepsilon \Rightarrow |\hat{z} - \hat{P}_i| \leq 2\varepsilon$
 \Rightarrow steps 3, 4 will choose H_i as winner no matter if $i=1$ or $i=2$

indeed: suppose $i=1$; then b.c. $|\hat{z} - \hat{P}_1| \leq 2\varepsilon \Rightarrow$ H_1 declared winner at step 3
if $i=2$, then $\hat{P}_1 > \hat{P}_2 + 6\varepsilon > \hat{z} + 4\varepsilon \Rightarrow$ algorithm goes to step 4
at step 4, b.c. $\hat{z} < \hat{P}_2 + 2\varepsilon \Rightarrow H_2$ declared winner

case 2: $TV(P, H_i) \leq \varepsilon$, $TV(P, H_j) > 4\varepsilon$

subcase 2.1: If $\hat{p}_i - \hat{p}_j \leq 6\varepsilon$, then algorithm stops
at step 2 declaring a tie

subcase 2.2: If $\hat{p}_i - \hat{p}_j > 6\varepsilon$, algorithm goes to step 3

{ - if $i=1$, then $TV(P, H_1) \leq \varepsilon \Rightarrow |\hat{z} - p_1| \leq \varepsilon \Rightarrow |\hat{z} - \hat{p}_1| \leq 2\varepsilon \rightarrow$ algorithm declares H_1 as winner
- if $i=2$, then $TV(P, H_2) \leq \varepsilon \Rightarrow |\hat{z} - p_2| \leq 2\varepsilon$
thus $\hat{p}_1 > \hat{p}_2 + 6\varepsilon > \hat{z} + 4\varepsilon$
so algorithm does not stop at step 3
& at step 4, declares H_2 as winner
so in this case H_1 is declared the winner
regardless of whether $i=1$ or 2

so in case 2, either H_i wins or it's a tie

case 3: $TV(H_1, H_2) \leq 5\varepsilon \Rightarrow \hat{p}_1 - \hat{p}_2 \leq 6\varepsilon \Rightarrow$ step 2
declares a tie
□