

Lecture 18 - scribbles

Friday, November 9, 2018 13:26

today: finish exp. family
sampling

(cts. exp family)

example 3: UGM ?

Let $p \in \mathcal{L}(G)$, G is undirected

with $\psi_c(x_c) > 0 \quad \forall c, x_c$

$$\begin{aligned} p(x) &= \frac{1}{Z} \prod_{c \in \beta} \psi_c(x_c) = \exp\left(\sum_c \ln \psi_c(x_c) - \log Z\right) \\ &= \exp\left(\sum_{c \in \beta} \sum_{y \in X_c} \underbrace{\mathbb{1}_{\{y_c = x_c\}}}_{T_{c,y}(x)} \underbrace{\log \psi_c(x_c)}_{n_{c,x_c}} - \log Z\right) \end{aligned}$$

$$T(x) = \begin{pmatrix} \vdots \\ \mathbb{1}_{\{x_c = y_c\}} \\ \vdots \end{pmatrix} \quad y \in X_c$$

$$n(\theta) = \begin{pmatrix} \vdots \\ \log \psi_c(y_c) \\ \vdots \end{pmatrix} \quad c \in \beta$$

notes: a) Mult(π) is special case with complete graph (1 big clique)

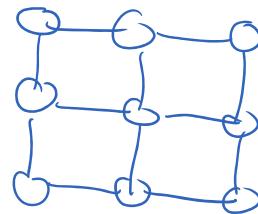
..

- b) feature perspective: instead of using all indicators $\{x_c = y_c\}$
 you could choose relevant subset for a task
 for example; if x_i is a word
 feature on $x_i \{x_{i+1} \mid \{x_i \text{ is a verb} \wedge x_{i+1} \text{ is a noun}\}$

c) binary Ising model

$$x_i \in \{0, 1\} \quad |C| \leq 2$$

Suppose use nodes $\{\text{pairs}\}$ as cliques
 (edges)

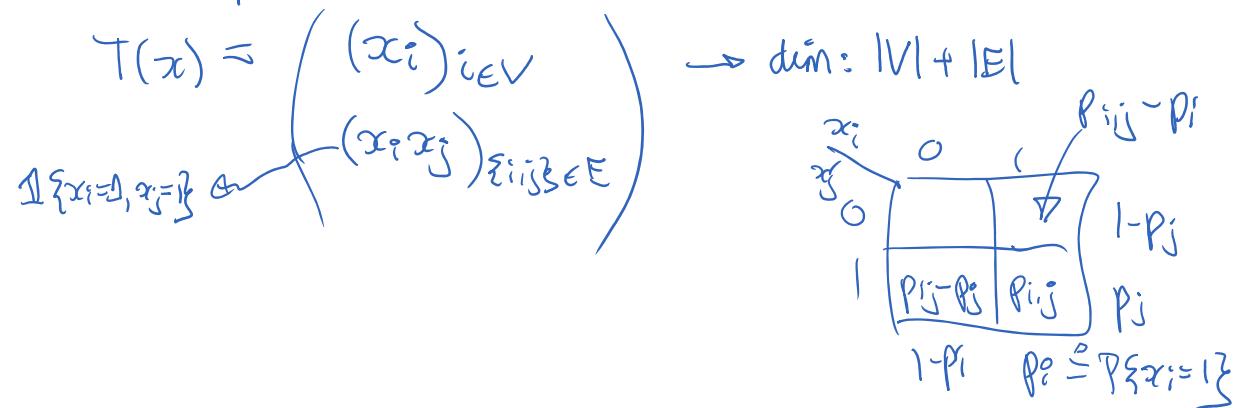


\Rightarrow dimension of $T(x)$ is $2|V| + 4|E|$ \rightsquigarrow "overparametrized" exp. family

$$\sum_{y_C} T_{C,y_C}(x) = 1 \Rightarrow \text{not a minimal exp. family}$$

* a minimal representation:

$$T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{\{i, j \in E\}} \end{pmatrix} \rightarrow \dim: |V| + |E|$$



properties of A :

- $A(\cdot)$ is C^∞ for $m \in \text{int}(\Omega)$
 $\mu(\cdot)$ is convex (Ω is a convex set)
- $\nabla_m A(n) = \mathbb{E}_{p(x|n)}[T(x)] \triangleq \mu(n)$ "moment vector" (for $n \in \text{int}(\Omega)$)
- $\left(\frac{\partial^2}{\partial n_i \partial n_j} A(n) \right)_{ij} = \mathbb{E}_{p(x|n)}[(T(x) - \mu(n))(T(x) - \mu(n))^T] = \text{cov}(T(x))$
 (prop as exercise?)
- $Z(n)$ is the mgf for $p(x; n)$ $t \mapsto Z(m+t)$ for $t \in \mathbb{R}^p$
 $\mathbb{E}_x \exp(t^T T(x))$

$A(n) \approx \log \text{mgf} \rightarrow$ cumulant generating functions

→ why derivative of A gives cumulants

1st → mean
 2nd → covariance
 3rd → 3rd cumulant

14h08

Approximate inference → sampling

example: cannot do exact inference in Ising model (tree width is too big)
 ↳ need approximations [NP-hard]

why sampling?

$$X = (X_1, \dots, X_p)$$

a) simulation: $X^{(i)} \sim p$

b) approximate $p(x_i)$ "marginal inference"

u) -----

b) approximate $p(x)$ "marginal inference"

→ special case of expectations

consider: $f: \mathbb{R}^P \rightarrow \mathbb{R}^Q$,

We want to approximate $\mu = \mathbb{E}_p[f(x)]$

e.g. if $f(x) = \mathbf{1}_{\{X_A = x_A\}}$, $\mathbb{E}_p[f(x)] = p(X_A = x_A)$

Monte-Carlo integration / estimation → appears in physics, applied math, ML, etc...

to approximate $\mu = \mathbb{E}_p[f(x)]$

[MC estimate]

alg: • n samples $X^{(i)} \stackrel{\text{iid.}}{\sim} P$

• estimate: $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) = \mathbb{E}_{\hat{P}_n}[f(x)]$

this is true even if
 $X^{(i)}$ were dependent

properties: 1) unbiased $\mathbb{E}[\hat{\mu}] = \underset{\substack{\text{expectation over } (X^{(i)})_{i=1}^n}}{\mathbb{E}} \left[\frac{1}{n} \sum_i f(X^{(i)}) \right] = \frac{n}{n} \mu = \mu$

$$X^{(i)} \sim P \Rightarrow X^{(i)} \stackrel{\text{d}}{\equiv} X \sim P$$

2) expected l_2 error

$$\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu}))$$

$$\begin{aligned} & \mathbb{E}\left[\left\| \hat{\mu} - \mu \right\|_2^2\right] = \mathbb{E}\left[\frac{1}{n^2} \sum_{i,j} \langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle\right] \\ & \xrightarrow{\text{by independence}} \text{off diagonal terms are zero} \\ & \frac{1}{n^2} \sum_{i=j} \underbrace{\mathbb{E}\left[\langle f(X^{(i)}) - \mu, f(X^{(i)}) - \mu \rangle\right]}_{\mathbb{E}\left[\|f(X^{(i)}) - \mu\|^2\right] \triangleq \sigma^2} \end{aligned}$$

$$\begin{aligned} & \mathbb{E}\left[\|f(X^{(i)}) - \mu\|^2\right] \triangleq \sigma^2 \\ & = \text{tr}(\text{cov}(f(x), f(x))) \end{aligned}$$

$$\mathbb{E} L \left[\|f(x) - u\|^2 \right] \stackrel{\triangle}{=} \sigma^2 \\ = \text{tr}(\text{cov}(f(x), f(x)))$$

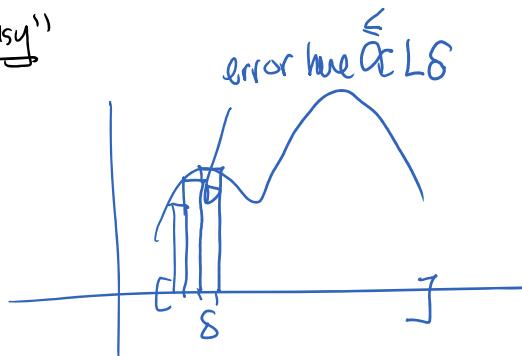
$$\boxed{\mathbb{E} [\|\hat{\mu} - \mu\|^2] = \frac{\sigma^2}{n}}$$

note that no dimension in rate (apart σ^2)

14h33

aside on numerical computing: 1d is "easy"

- numerical integration in 1d



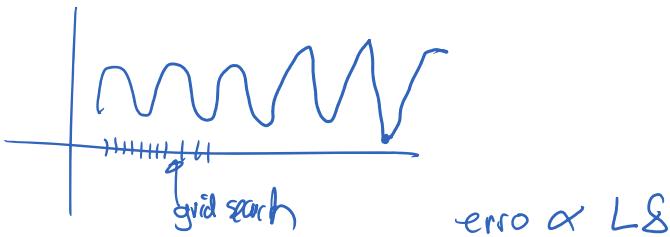
for function L-Lipschitz

$$|f(x) - f(y)| \leq L \|x - y\|$$

error $\leq \varepsilon$ \Rightarrow complexity $O(\frac{1}{\varepsilon})$
 \sim use $\delta \approx \varepsilon$ of approximating integral within ε

but grid in dimension d $O(\frac{1}{\varepsilon^d})$ \rightarrow curse of dimensionality?

- global optimization in 1d



How to sample?

Complexity $O(\frac{1}{\epsilon})$
(# of evaluations)
to find global optimum

- 1) $X \sim \text{Unif}([0, 1]) \rightarrow$ pseudo random generator "rand"
- 2) $X \sim \text{Bernoulli}(p) \quad X = \mathbb{1}\{U \leq p\}$ where $U \sim \text{Unif}([0, 1])$
- 3) inverse transform sampling: cumulative dist. fn.

Let F be target cdf of distribution p for X $F(x) \triangleq P\{X \leq x\} \quad (x \in \mathbb{R})$

(first, suppose F is invertible)

$$\boxed{\text{Let } X \triangleq F^{-1}(U) \text{ with } U \sim \text{Unif}([0, 1])}$$

claim that X has cdf $F(x)$

F is invertible

$$\text{proof: } P\{X \leq y\} = P\{F^{-1}(U) \leq y\} \stackrel{F \text{ is invertible}}{=} P\{U \leq F(y)\} = F(y)$$

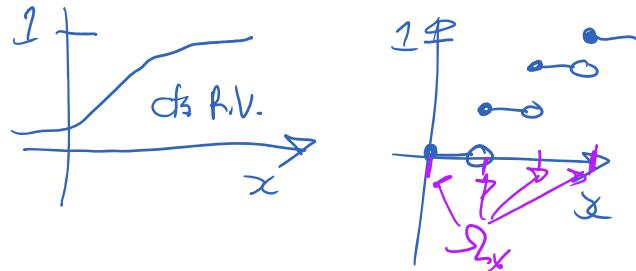
[if F is not invertible, define $\boxed{X \triangleq \min\{x \in \mathbb{R}: F(x) \geq u\}}$] (recall that F is cts from right)

example:

want $X \sim \text{Exp}(\lambda)$ density $p(x) = \lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}^+}(x)$
 $F(x) = 1 - e^{-\lambda x}$

inverse $F^{-1}(u) = -\frac{1}{\lambda} \ln(1-u)$

~



$$\text{universe } F^{-1}(u) = \frac{1}{\lambda} \ln(1-u)$$

Multivariate distribution?

generalize above trick using "chain rule"

$$X_{1:p} \quad (\dim p) \quad \text{cdf } F(x_{1:p}) \triangleq P\{X_1 \leq x_1, \dots, X_p \leq x_p\}$$

$$F_{X_{1:p}}(x_{1:p}) = F_{X_1}(x_1) \underbrace{F_{X_2|X_1}(x_2|x_1)}_{\rightarrow} \cdots F_{X_p|X_{1:p-1}}(x_p|x_1, \dots, x_{p-1})$$

$$F_{X_2|X_1}(x_2|x_1) \triangleq P\{X_2 \leq x_2 | X_1 \leq x_1\}$$

Could use U_1, \dots, U_p ^{i.i.d.} Uniform

$$X_1 = F_{X_1}^{-1}(U_1)$$

$$\vdots$$

$$X_p = F_{X_p|X_{1:p-1}}^{-1}(U_p | X_{1:p-1})$$

is very complicated set.
(cause of dimensionality?)

[aside: "copulas" \rightarrow model for multivariate dependence with uniform marginals]

Exception is multivariate Gaussian:

$$N(\mu, \Sigma) \quad \Sigma = U D U^T$$

(where $U^T U = I_p$)
and D is diagonal
(Cholesky decomposition)

generate $V \sim N(0, I_p)$

$\{v_p \stackrel{\text{i.i.d.}}{\sim} N(0, 1)\}$

$$X \triangleq \underbrace{U D^{1/2}}_{\text{Cholesky}} V + \mu$$

and L is diagonal
(Cholesky decomposition)

$$L = U L^{1/2}$$

$$\Sigma = L L^T$$

$$X \triangleq \underbrace{U L^{1/2}}_L V + \mu$$

$$\mathbb{E} X = \mu$$

$$\text{cov}(X) = \Sigma$$

$$\text{so } X \sim N(\mu, \Sigma)$$

Box-Muller transform to sample (2d) Gaussian

$$\begin{aligned} x = r \cos \theta &\rightsquigarrow r^2 \sim \text{Exp}(1) \\ y = r \sin \theta & \quad \theta \text{ is uniform } [0, 2\pi] \end{aligned} \Rightarrow \begin{pmatrix} x \\ y \end{pmatrix} \sim N(0, I)$$