

# Lecture 18 - scribbles

Friday, November 15, 2019 13:57

- today:
- finish properties of exp. family
  - estimation of graph model parameters
  - sampling (approximate inference)

## Exp. family (ct.ed)

example 2: 1D Gaussian

$$X \sim N(\mu, \sigma^2) \quad X = \mathbb{R} \quad \Theta = (\mu, \sigma^2) \quad \text{"moment parameterization"}$$

$$p(x; (\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x^2}{2} \left[\frac{1}{\sigma^2}\right] + x \left[\frac{\mu}{\sigma^2}\right] - \left[\frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right]\right)$$

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix} \quad \eta(\theta) = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}$$

$$\eta_2 = \frac{1}{\sigma^2} = \text{precision} > 0$$

$$\eta_1 = \eta_2 \cdot \mu$$

$$\Omega = \left\{ (\eta_1, \eta_2) : \eta_2 > 0, \eta_1 \in \mathbb{R} \right\}$$

[we'll see later: multivariate Gaussian  $\Lambda = \Sigma^{-1}$   
 $\eta = \Lambda \mu = \Sigma^{-1} \mu$ ]

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix}$$

Example 3: discrete UGM?

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

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

$$\begin{aligned} p(x) &= \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(x_c) = \exp\left(\sum_c \ln \psi_c(x_c) - \log Z\right) \\ &= \exp\left(\sum_{c \in \mathcal{C}} \sum_{y_c \in \mathcal{X}_c} \underbrace{1_{\{x_c = y_c\}}}_{\text{indicator}} \underbrace{\log \psi_c(y_c)}_{\text{potential}} - \log Z\right) \end{aligned}$$

$$T(x) = \begin{pmatrix} \vdots \\ \mathbb{1}\{x_c = y_c\} \\ \vdots \end{pmatrix} \left\{ \begin{array}{l} y_c \in X_c \\ c \in E \end{array} \right.$$

$$n(\theta) = \begin{pmatrix} \vdots \\ \log \pi_c(y_c) \\ \vdots \end{pmatrix}$$

$$X_C = \{ (y_i)_{i \in C} : \exists x. y \in X \}$$

$$= \prod_{i \in C} X_i$$

notes: a) Mult( $\pi$ ) is special case where complete graph (1 big clique)

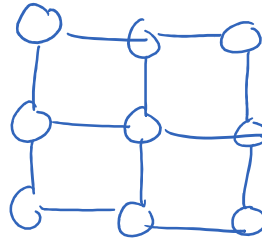
b) feature perspective: instead of using all indicators  $\mathbb{1}\{x_c = y_c\}$ , you could choose relevant subset for a task

for example  $x$  is a sentence and  $x_i$  is a word

feature on  $x_i$  &  $x_{i+1}$  e.g.  $\mathbb{1}\{x_i \text{ is a verb} \& x_{i+1} \text{ is a noun}\}$

### c) binary Ising model

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



suppose use nodes & pairs (edges) as cliques

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

$$\sum_{y_c} T_{c, y_c}(x) = 1 \text{ for any } c$$

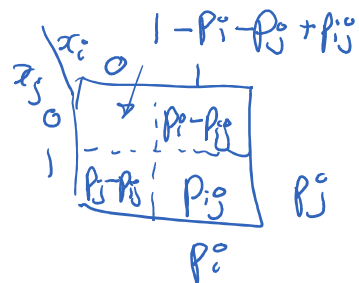
$\Rightarrow$  not a min. exp. family

\* a minimal representation:

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

$\mathbb{1}\{x_i = 1, x_j = 1\}$

$\rightarrow$  dim:  $|V| + |E|$



properties of  $A$ :

## properties of $A$ :

$\mu$ :

- $A(\cdot)$  is  $C^\infty$  for  $n \in \text{int}(\Omega)$
- $A(\cdot)$  is convex ( $\Omega$  is convex set)
- $\nabla_n A(n) = \mathbb{E}_{p(x|n)} [T(x)] \triangleq \mu(n)$  "moment vector" (for  $n \in \text{int}(\Omega)$ )
- $\left( \frac{\partial^2}{\partial n_i \partial n_j} A(n) \right)_{i,j} = \mathbb{E}_{p(x|n)} [ [T(x) - \mu(n)] [T(x) - \mu(n)]^T ] = \text{cov}(T(x))$   
(Hessian) (proof as exercise)
- $Z(n)$  is the mgf for  $p(x; n)$   $t \mapsto Z(n+it)$  for  $t \in \mathbb{R}^p$   
 $\mathbb{E}_x \exp(t^T T(x))$

$A(n) = \log \text{mgf} \rightarrow$  cumulant generating function

$\leadsto$  why derivatives of  $A$  gives cumulants

1st  $\rightarrow$  mean

2nd  $\rightarrow$  covariance

3rd  $\rightarrow$  3rd cumulant

14h56

## Estimation of parameters for AGM

DGM: parametric family  $\mathcal{P}_{\Theta} = \left\{ p(x) = \prod_i p(x_i | x_{\pi_i}, \Theta_i) \right\}$   
 $\Theta = (\Theta_1, \dots, \Theta_{|\mathcal{V}|})$   
 $\Theta \in \Theta = \Theta_1 \times \dots \times \Theta_{|\mathcal{V}|}$

independent parameterization

i.e. no tying of parameters

$\Rightarrow$  MLE decouples in  $|\mathcal{V}|$  independent problems

$$\{x^{(i)}\}_{i=1}^n \quad p(\text{data} | \Theta) = \prod_{i=1}^n p(x^{(i)} | \Theta) = \prod_{i=1}^n \prod_{j=1}^{|\mathcal{V}|} p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j)$$

$$\log [ \quad ] = \sum_{j=1}^{|\mathcal{V}|} \underbrace{\left( \sum_{i=1}^n \log p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j) \right)}_{\mathcal{J}_j(\Theta_j)}$$

$$j=1 \underbrace{\{i=1, \dots, n\}}_{f_j(\theta_j)}$$

Example: for discrete R.V.  $\Rightarrow \Theta_j^{ML} =$  proportion of observations

$$\frac{\#(x_j = k, x_{\pi_j} = \text{something})}{\# x_{\pi_j} = \text{something}}$$

gives  $\hat{p}(x_j = k | x_{\pi_j})$

⊗ if have latent variable (ie. unobserved variable)

$\Rightarrow$  use EM. (like we did for HMM)

UGM:

example for exp family

$$p(x|n) = \exp\left(\sum_c n_c^T T_c(x_c) - A(n)\right)$$

$\rightarrow$  unlike in a DGM,  $\log p(x|n)$  does not separate as just  $\sum_c f_c(n_c)$

gradient ascent on log-likelihood

$$\frac{\partial}{\partial \mu} \log p(x^{(i)}|n) = \sum_c n_c^T \left( \frac{1}{n} \sum_{i=1}^n T_c(x_c^{(i)}) \right) - \frac{\partial A(n)}{\partial \mu}$$

$\mu_c$

$$\nabla_n [J] = \hat{\mu}_c - \mu_c(n)$$

$$\hookrightarrow \mathbb{E}_{p(x|n)} [T_c(x_c)]$$

to compute this, need inference

e.g. Ising model  $T_{ij}(x_i, x_j) = x_i \cdot x_j$

$$\mathbb{E}[T_{ij}] = \mu_{ij} = p(x_i = 1, x_j = 1 | n)$$

here need to use approximate inference  $\leftarrow$  sampling variational method

Approximate inference

# sampling

example: Not hard to do exact inference in Ising model  
→ need approximations

## Why sampling?

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

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

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

→ special case of expectations

consider:  $f: \mathbb{R}^p \rightarrow \mathbb{R}^d$

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

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

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

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

MC estimation: alg: • n samples  $X^{(i)} \sim p$   
• estimate  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) = \mathbb{E}_n[f(x)]$

properties: 1) unbiased  $\mathbb{E}_p[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f(X^{(i)})] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$   
expectation over  $(X^{(i)})_{i=1}^n$  ↖ this is true even if  $X^{(i)}$  are dependent

2) expected error (L2 error)  $\mathbb{E}[\|\mu - \hat{\mu}\|_2^2] = \mathbb{E}\left[\frac{1}{n^2} \sum_{i,j} \langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle\right]$   
=  $\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu}))$  by independence → off diagonal terms are zero  
=  $\frac{1}{n^2} \sum_{i=j} \mathbb{E}[\langle f(X^{(i)}) - \mu, f(X^{(i)}) - \mu \rangle]$

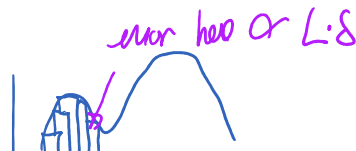
$$\mathbb{E}[\|\hat{\mu} - \mu\|_2^2] = \frac{\sigma^2}{n}$$

$$\mathbb{E}[\|f(x) - \mu\|_2^2] \triangleq \sigma^2 = \text{tr}(\text{cov}(f(x), f(x)))$$

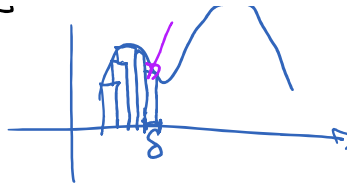
note that no dimension in rate (apart  $\sigma^2$ )

aside on numerical computing: 1d is "easy"

• numerical integration in 1D



• numerical integration in 1D



for function  $f$   $L$ -Lipschitz

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

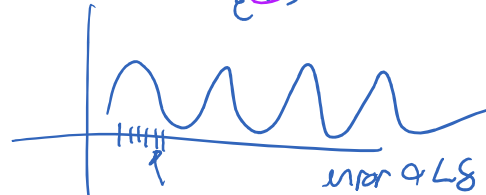
$$\text{error} \leq \epsilon$$

$\leadsto$  use  $\delta \approx \epsilon$

$\Rightarrow$  complexity  $O(\frac{1}{\epsilon})$   
of approximating  
integral within  $\epsilon$

but grid in dimension  $d$   $O(\frac{1}{\epsilon^d})$  cause of dimensionality!

• global optimization in 1d



complexity  $O(\frac{1}{\epsilon})$

to find global optimum