

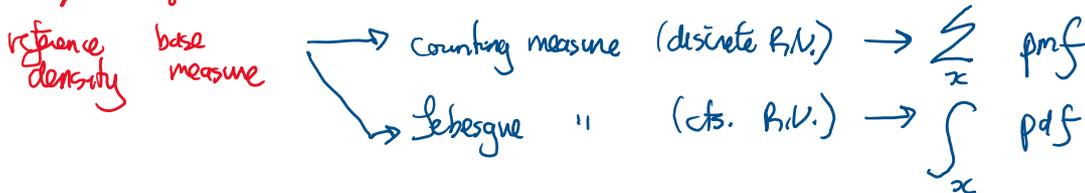
today: exp. family  
estimation for PGM

## Exponential family

a (flat/canonical) exponential family on  $X$

is a parametric family of dist. on  $X$  defined by two quantities

I)  $h(x) d\mu(x) \rightarrow$  reference measure



II)  $T: X \rightarrow \mathbb{R}^p$  called "sufficient statistics" vector aka. feature vector  
members of the family will have pmf/pdf

$$p(x; \eta) d\mu(x) = \exp(\eta^T T(x) - A(\eta)) h(x) d\mu(x)$$

"canonical parameter"  $\rightarrow$  log-normalization or cumulant generating function  
log partition fct.

if  $\Omega_X$  is discrete; then  $p(x; \eta)$  is a pmf

" " cts. ; " " " pdf

$$\mathbb{E}[f] = \int_x f(x) d\mu(x)$$

$$\begin{aligned} \textcircled{*} \text{ want } 1 &= \int_X p(x; \eta) d\mu(x) \\ &= \int_X \exp(\eta^T T(x)) e^{-A(\eta)} h(x) d\mu(x) \end{aligned}$$

or  $\sum_x f(x) \mu(x)$  if  $\mu$  is discrete

[discrete:  $\sum_x p(x; \eta)$ ]  $\Rightarrow A(\eta) \triangleq \log \left( \int_X \exp(\eta^T T(x)) h(x) d\mu(x) \right)$   
 $Z(\eta)$

domain  $\Omega \triangleq \{ \eta \in \mathbb{R}^p \mid A(\eta) < \infty \}$

set of valid canonical parameters

class of notation:

$$\Omega \neq \Omega_X \subseteq X \subseteq \mathbb{R}^p$$

note:  $A(\eta)$  is convex in  $\eta \Rightarrow \Omega$  is convex

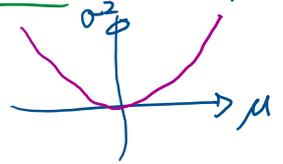
$\textcircled{*}$  more generally, consider a reparametrization of a subset of the flat exponential family  
by defining  $\eta: \mathcal{P} \rightarrow \Omega$   
new set

by defining  $\eta: \Theta \rightarrow \mathcal{L}$   
 $\eta$   
 new set  
 of parameters

$$p(x; \theta) \triangleq p(x; \eta(\theta)) \text{ for } \theta \in \Theta$$

(get a "curved exponential family" if  $\eta(\Theta)$  is curved manifold in  $\Omega$ )

↳ e.g. could consider Gaussians where  $N(\mu, \mu^2)$



\* note: any single dist  $p(x)$  can be part of an exponential family

by using  $h(x) = p(x)$

two examples of family not an exp-family: •  $\text{unif}(0, \theta)$

• mixture of Gaussians  
 (latent variable model)

Example 1: (multinomial)

$$X \sim \text{Mult}(\pi) \quad X = \{0, 1\}^k$$

$$\Omega_X = \Delta_k \cap X \text{ (one hot encodings)}$$

parameter  $\pi \in \Delta_k$ ; suppose  $\pi_i > 0 \forall i$

$$p(x; \pi) = \prod_{j=1}^k \pi_j^{x_j} = \exp\left(\sum_{j=1}^k x_j \log \pi_j\right)$$

$x \in X$  "think as '0'"

$$= \exp\left(\sum_j x_j \log \pi_j - 0\right)$$

$$\text{we have } \eta_j(\pi) = \log \pi_j$$

$$T(x) = x \quad \Omega_X \in \mathbb{R}^k$$

$d(x) = \text{counting measure on } X$

$$h(x) = \mathbb{1}_{\{x \in \Omega_X\}} = \mathbb{1}_{\{x \in \Delta_k \cap X\}}$$

$$A(\pi) = 0??$$

$$\Theta = \text{int}(\Delta_k) \quad A(\eta(\pi)) = 0 \quad \forall \pi \in \Theta$$

$$\Theta \rightarrow \dim k-1$$

$$\eta(\Theta) \rightarrow \text{" "}$$

$$\Omega \in \mathbb{R}^k \rightarrow \dim k$$

we do not have a "minimal exp family"

we do not have a "minimal exp family"

note: here, for any  $x$  s.t.  $h(x) \neq 0$

$$\sum_{j=1}^k T_j(x) = \sum_{j=1}^k x_j = 1$$

affine linear dep between components of  $T$

$\Rightarrow$  multiple  $\eta$ 's give rise to same distribution  $\rightarrow$  "overparameterization"

e.g.  $(n + \mathbb{1}_\alpha)^T T(x) = n^T T(x) + \alpha \underbrace{\mathbb{1}^T T(x)}_{=1}$

vector (!)  
 $\downarrow$

$\downarrow$  not a minimal exp. family

$\otimes$  for a multinomial; a min. exp. family

$$T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix} \quad \left[ "x_k = 1 - \sum_{j=1}^{k-1} x_j \right]$$

$$Z(\eta) = \sum_{x \in \Omega} \exp(\eta^T T(x)) = \sum_{j=1}^{k-1} e^{\eta_j} + 1$$

$$p(x; \eta) = \exp\left( \sum_{j=1}^{k-1} \eta_j x_j - \underbrace{\log\left(1 + \sum_{j=1}^{k-1} e^{\eta_j}\right)}_{A(\eta)} \right)$$

recall:  $\nabla_{\eta} A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)]$  (valid  $\eta \in \text{int}(\Omega)$ )

for multinomial;  $\frac{\partial A}{\partial \eta_j} = \frac{1}{Z(\eta)} e^{\eta_j} = p("x=j" | \eta)$

$= \mathbb{E}_{p(x; \eta)} [T(x)]$   
as required //

moment matching can be different than MLE in exp. family  
(with wrong moments?)

e.g. gamma dist  $T(\alpha, \beta)$  has  $T(x) = \begin{bmatrix} \log x \\ x \end{bmatrix}$

so moment matching with  $T(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$  will yield a different estimate than MLE

15h35

example 2: 1d Gaussian

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

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

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

A(n)

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

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

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

$$\Omega = \{(\eta_1, \eta_2) : \eta_2 > 0, \eta_1 \in \mathbb{R}\} \quad h(x) = 1 \quad (\text{but some people use } h(x) = \frac{1}{\sqrt{2\pi}})$$

[we'll see later: multivariate Gaussian  $T(x) = \begin{bmatrix} x \\ -\frac{xx^T}{2} \end{bmatrix}$  "for Gaussian"  $\eta_1 = \mu = \Sigma^{-1} \mu$   $\eta_2 = \Sigma^{-1}$ ]

Example 3: discrete UGM

let  $p \in \mathcal{P}(\mathcal{X})$   $\mathcal{X}$  is undirected

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

$$p(z) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(x_c) = \exp\left(\sum_{c \in \mathcal{C}} \log \psi_c(x_c) - \log z\right)$$

$$= \exp\left(\sum_{c \in \mathcal{C}} \sum_{y_c \in \mathcal{X}_c} \underbrace{\mathbb{1}\{y_c = x_c\}}_{T_{c, y_c}(x)} \log \psi_c(y_c) - \log z\right)$$

$$T(x) = \begin{pmatrix} \vdots \\ \mathbb{1}\{x_c = y_c\} \end{pmatrix} \leftarrow \begin{matrix} y_c \in \mathcal{X}_c \\ c \in \mathcal{C} \end{matrix}$$

$$\mathcal{X}_c = \{ (y_i)_{i \in \mathcal{C}} : y_i \in \mathcal{X}_i \}$$

$$\eta(\theta) = \begin{pmatrix} \vdots \\ \log \psi_c(y_c) \end{pmatrix} \leftarrow "$$

[not a minimal representation]

notes: a) mult(x) is a special case where have complete graph (1 big clique)

b) feature perspective: instead of using all possible indicators  $\mathbb{1}\{y_c = x_c\}$  you could use a subset or a function of them for a task

for example: suppose  $x$  is a sentence  
 $x_i$  is a word

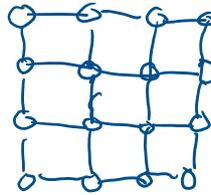
feature on  $x_i$  &  $x_{i+1}$  e.g.  $\mathbb{1} \{ \begin{matrix} x_i \text{ is a verb} \\ x_{i+1} \text{ is a noun} \end{matrix} \}$

→ much smaller set of parameters

"parameter sharing"

c) binary Ising model

$x_i \in \{-1, 1\}$   $|V| \leq 2$



suppose use nodes & pairs (edges)  
as cliques

⇒ dimension of  $T(x)$   $2|V| + 4|E|$

"overparameterized exp. family"

$\sum_{y \in C} T_{C,y}(x) = 1$  for any  $C$

⇒ not a min. exp. fam

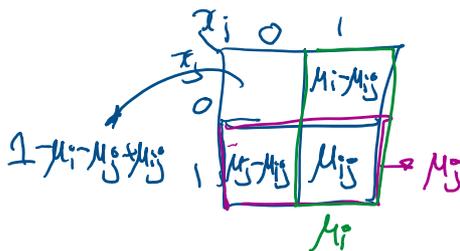
\* a minimal representation

$$B T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{\{i,j\} \in E} \end{pmatrix} \begin{matrix} \int n_i \\ \int n_{ij} \end{matrix} \rightarrow \dim |V| + |E|$$

$$E T(x) = \begin{pmatrix} (\mu_i)_{i \in V} \\ (\mu_{ij})_{\{i,j\} \in E} \end{pmatrix}$$
  

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

$$p(x_i = 1, x_j = 1)$$



properties of  $A$  (for generic flat exponential family)

$\bullet \nabla_n A(n) = \mathbb{E}_{p(x;n)} [T(x)] \triangleq \mu(n)$  "moment vector" (for  $n \in \text{int}(\Omega)$ )

$\bullet \left( \frac{\partial^2 A(n)}{\partial n_i \partial n_j} \right)_{(i,j) \in (1:p)^2} = \mathbb{E}_{p(x;n)} [ (T(x) - \mu(n)) (T(x) - \mu(n))^T ] = \text{cov}(T(x))$   
(proof as exercise)

"cumulant generating fct."

Estimation of parameters DBM/UGM

independent + . . .

# Estimation of parameters DGM/UGM

DGM  
(fully observed)

parametric family  $\mathcal{P}_{\Theta} = \{p_{\theta}(x) = \prod_i p(x_i | x_{\pi_i}, \theta_i)\} \equiv$   
 $(\theta_1, \dots, \theta_M) \in \Theta$  } se. no. of parameters  
 independent parametrization  
 $\Theta = \Theta_1 \times \Theta_2 \times \dots \times \Theta_M$

$\Rightarrow$  MLE decouples in  $M$  independent MLE problems:

$$p(\{x^{(i)}\}_{i=1}^n | \theta) = \prod_{i=1}^n p(x^{(i)} | \theta) = \prod_{i=1}^n \prod_{j=1}^M p(x_j^{(i)} | x_{\pi_j}^{(i)}, \theta_j)$$

$$\log(\quad) = \sum_{j=1}^M \underbrace{\left( \sum_{i=1}^n \log p(x_j^{(i)} | x_{\pi_j}^{(i)}, \theta_j) \right)}_{f_j(\theta_j)}$$

example: for discrete R.V.  $\Rightarrow \hat{\theta}_j^{MLE} =$  proportion of observations  
 (multiplicity conditions)  $\hat{p}(x_j = k | x_{\pi_j} = \text{stuff})$   
 $= \frac{\#(x_j = k, x_{\pi_j} = \text{stuff})}{\#(x_{\pi_j} = \text{stuff})}$

(fully observed DGM is relatively easy)  
often in closed form?

$\oplus$  if have latent variable (i.e. unobserved variables)

$\rightarrow$  use EM. (like HMM)

## UGM

example for exp. family

$$p(x; \eta) = \exp\left(\sum_c \langle \eta_c, T_c(x_c) \rangle - A(\eta)\right)$$

$\rightarrow$  unlike in a DGM,  $\log p(x; \eta)$  does not separate  $\sum_c f_c(\eta_c)$

gradient ascent on log-likelihood

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)} | \eta) = \sum_c \eta_c^T \left[ \frac{1}{n} \sum_{i=1}^n T_c(x^{(i)}) \right] - \frac{A(\eta)}{n}$$

$$\nabla_{\eta_c} [\quad] = \hat{\mu} - \mu_c(\eta)$$

$\downarrow$   
 $\mathbb{E}_{p(x; \eta)} [T_c(x_c)]$

to compute this, need inference

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

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

here need approximate inference

← sampling [Gibbs sampling]  
variation [mean field]