

today: exp. families
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) $h(x) d\mu(x)$ → reference measure

reference density base measure

→ counting measure (discrete R.V.) → \sum_x pmf
Lebesgue " " (cts R.V.) → \int_X pdf

II) $T: X \rightarrow \mathbb{R}^p$ called "sufficient statistics" vector
aka. feature vector

members of the family will have pmf/pdf

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

"canonical parameter"

log-normalization or cumulant generating functions

log partition fn.

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

• if \mathbb{X} is discrete, then $p(x; n)$ is a pmf

" " cts. ; " " " " a pdf

* want $1 = \int_X p(x; n) d\mu(x) = \int_X \exp(n^T T(x)) e^{-A(n)} h(x) d\mu(x)$

[discrete: $\sum_x p(x; n)$]

$$\Rightarrow A(n) \triangleq \log \left(\int_X \exp(n^T T(x)) h(x) d\mu(x) \right) z(n)$$

clash of notation:

$$\mathcal{Q} \neq \mathcal{Q}_X \subseteq X \\ \subseteq \mathbb{R}^p$$

domain $\mathcal{Q} \triangleq \{n \in \mathbb{R}^p \mid A(n) < \infty\}$

set of valid canonical parameters

note: $A(n)$ is convex in n $\Rightarrow \mathcal{Q}$ is convex

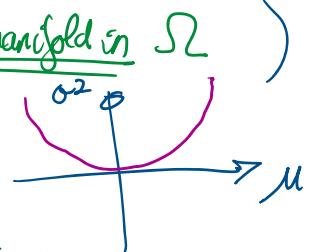
* more generally, consider a reparameterization of a subset of the flat exponential family

$$\dots \rightarrow m: \mathbb{R} \rightarrow \mathcal{Q}$$

⊗ more generally, consider a reparameterization of a subset of the flat exponential family by defining $\eta: \mathbb{H} \rightarrow \Omega$
 new set of parameters

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

(get a "curved exponential family" if $\eta(\mathbb{H})$ is curved manifold in Ω)
 ↳ 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 $\eta(x) = p(x)$

two examples of family not an exp family:

- $\text{Unif}(0, \theta)$
- mixture of Gaussians
 (latent variable model)

Example 1: (multinomial)

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

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

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

$$\begin{aligned} p(x; \pi) &= \prod_{j=1}^k \pi_j^{x_j} = \exp\left(\sum_{j=1}^k x_j \log \pi_j\right) \\ &\text{"think as "0"} \quad = \exp\left(\sum_{j=1}^k x_j \log \pi_j - 0\right) \end{aligned}$$

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

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

$$\delta_{\pi}(x) = \text{counting measure on } X \quad h(x) = \mathbb{1}\{x \in \Omega_X\} = \mathbb{1}\{x \in \Delta_k \cap X\}$$

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

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

$$\mathbb{H} \rightarrow \dim k-1$$

$$\eta(\mathbb{H}) \rightarrow \dots \dots \dots$$

$$\mathcal{L}_X \rightarrow \parallel K$$

we do not have a "minimal exp family"

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

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

affine linear dep between components of T

\Rightarrow multiple n 's give rise to same distribution
 \rightarrow "overparametrization"

$$\text{e.g. } (n + 1\alpha)^T T(x) = n T(x) + \alpha T(x)^T \underset{\perp}{\mathbf{1}} \quad \begin{matrix} \downarrow \\ \text{not a} \\ \text{minimal} \\ \text{exp family} \end{matrix}$$

* for a multinomial; a min. exp. family

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

$$Z(n) = \sum_{x \in \mathcal{L}_X} \exp(n T(x)) = \sum_{j=1}^{K-1} e^{n \eta_j} + 1$$

$$\boxed{p(x; n) = \exp \left(\sum_{j=1}^{K-1} \eta_j x_j - \log \left(1 + \sum_{j=1}^{K-1} e^{\eta_j} \right) \right)}$$

$$\text{recall: } D_n A(n) = \mathbb{E}_{p(x; n)} [T(x)] \quad (\text{valid for } n \in \text{int}(\mathbb{R}))$$

$$\text{for multinomial; } \frac{\partial A}{\partial \eta_j} = \frac{1}{Z(n)} e^{\eta_j} = p(\text{"}x=j\text{"}|n) \\ = \mathbb{E}_{p(x; n)} [T(x)] \quad \text{as required} //$$

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

gamma dist $\text{Ga}(\alpha, \beta)$ 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

example 2 : 1d Gaussian

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

$$\begin{aligned} p(x; (\mu, \sigma^2)) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{x^2}{2} \left[\frac{1}{\sigma^2}\right] + x \left[\frac{\mu}{\sigma^2}\right] - \left[\frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right]\right) \end{aligned}$$

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix} \quad n(\theta) = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$$

$$n_2 = \frac{1}{\sigma^2} = \text{precision} > 0$$

$$n_1 = n_2 \cdot \mu$$

$$\Omega = \{(n_1, n_2) : n_2 > 0, n_1 \in \mathbb{R}\}$$

$h(x) = 1$ (but some people use $h(x) = \frac{1}{J_{2\pi}}$ for Gaussian)

[we'll see later: Multivariate Gaussian $T(x) = \begin{bmatrix} x \\ -\frac{xx^T}{2} \end{bmatrix}$ "n₁" = $\mu = \sum_i \mu_i$, "n₂" = $\Sigma = \sum_i \Sigma_i$]

Example 3 : discrete UGM

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

with $\mathcal{N}_c(x_c) \neq \emptyset \forall c, x_c$

$$\begin{aligned} p(x) &= \frac{1}{Z} \prod_{c \in G} \mathcal{N}_c(x_c) = \exp\left(\sum_c \log \mathcal{N}_c(x_c) - \log Z\right) \\ &= \exp\left(\sum_{c \in G} \sum_{y_c \in \mathcal{X}_c} \underbrace{\mathbb{I}_{\{y_c = x_c\}}}_{T_{c,y_c}(x)} \underbrace{\log \mathcal{N}_c(y_c)}_{n_{c,y_c}} - \log Z\right) \end{aligned}$$

$$T(x) = \left(\begin{array}{c} \vdots \\ \mathbb{I}_{\{x_c = y_c\}} \\ \vdots \end{array} \right)_{y_c \in \mathcal{X}_c, c \in G}$$

$$n(\theta) = \left(\begin{array}{c} \vdots \\ \log \mathcal{N}_c(y_c) \\ \vdots \end{array} \right)_{y_c \in \mathcal{X}_c, c \in G}$$

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

[not a minimal representation]

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

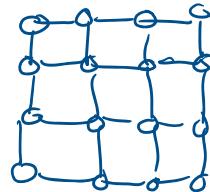
notes: a) $\text{Mult}(\pi)$ is a special case where have complete graph (1 big clique)
 b) feature perspective: instead of using all possible indicators $\{y_c = x_c\}$
 you could use a subset or a function of them
 See a task

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

feature on $x_i \in x_{i+1}$ e.g. $\{x_i \text{ is a verb}, x_{i+1} \text{ is a noun}\}$
 → much smaller set of parameters
 "parameter sharing"

c) binary Ising model

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



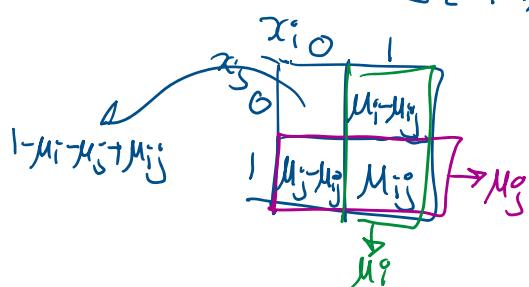
suppose we have nodes & pairs (edges)
 as cliques

→ dimension of $T(x)$ $2|V| + 4|E|$ → "overparameterized"
 $\sum_{y_c} T_{c,y_c}(x) = 1$ for any c exp. family
 → not a. min. exp. fam.

* a minimal representation

$$\text{is } T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{i, j \in E} \end{pmatrix} \begin{matrix} \rightarrow m_i \\ \downarrow \\ \begin{cases} 1 & x_i=1, x_j=1 \\ 0 & \text{else} \end{cases} \end{matrix} \rightarrow \dim |V| + |E|$$

$$\mathbb{E} T(x) = \begin{pmatrix} (M_i)_{i \in V} \\ (M_{ij})_{i, j \in E} \end{pmatrix} + p(x_i=1, x_j=1) M_{ij}$$



properties of A (for general flat exponential family)

properties of A (for general flat exponential family)

$$\cdot \nabla_n A(n) = \mathbb{E} p(x; n) [T(x)] \stackrel{\triangle}{=} \mu(n) \quad \text{"moment vector"} \quad (\text{for } n \in \text{int}(S))$$

$$\cdot \left(\frac{\partial^2 A}{\partial n_i \partial n_j} \right)_{(i,j) \in I \times I} = \mathbb{E} p(x; n) [(T(x) - \mu(n))(T(x) - \mu(n))^T] = \text{cov}(T(x))$$

(proof as exercise)

"cumulant generating fct."

Estimation of parameters DGM / UGM

DGM:
(fully observed)

$$\text{parametric family } P_{\Theta} = \left\{ p_{\theta}(x) = T_q(p(x_i; \theta_i); \theta) : \begin{array}{l} (\theta_1, \dots, \theta_M) \in \\ \Theta = \Theta_1 \times \Theta_2 \times \dots \times \Theta_M \end{array} \right\}$$

independent parameterization
i.e. no tying of parameters

\Rightarrow MLE decouple in $|V|$ independent MLE problems

$$\left\{ x^{(i)} \right\}_{i=1}^n \quad p(\text{data} | \theta) = \prod_{i=1}^n p(x^{(i)} | \theta) = \prod_{i=1}^n \prod_{j=1}^{|V|} p(x_j^{(i)} | x_{\pi_j}^{(i)}; \theta_j)$$

$$\log(p(\text{data} | \theta)) = \sum_{j=1}^{|V|} \left(\sum_{i=1}^n \underbrace{\log p(x_j^{(i)} | x_{\pi_j}^{(i)}; \theta_j)}_{f_j(\theta_j)} \right)$$

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

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

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

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

* if have latent variables (ie unobserved variables)

\Rightarrow use E.M. (Dek 14.14)

UGM:

example for exp. family

$$p(x; n) = \exp \left(\sum c_i \langle n_c, T_c(x_c) \rangle - A(n) \right)$$

$$p(x; n) = \exp\left(\sum_i \langle n_c, T_c(x_c) \rangle - A(n)\right)$$

\rightarrow unlike in a DGM, $\log p(x; n)$ does not separate $\sum_i f_c(n_c)$
gradient ascent on $\log p(x; n)$

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)} | n) = \sum_i n_c^T \left[\underbrace{\frac{1}{n} \sum_{i=1}^n T_c(x_c^{(i)})}_{\mu_c} \right] - n \frac{A(n)}{n}$$

$$\nabla_{n_c} \left[\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)} | n) \right] = \mu_c - \mu_c(n)$$

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

↑
to compute this, need inference

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

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

here need approximate inference

sampling [Gibbs Sampling]
variational
[mean field]