

today: exp families
estimation in 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 d.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

$$p(x; \eta) d\mu(x) = \exp(\underbrace{\eta^T T(x)}_{\text{"canonical parameter"}} - \underbrace{A(\eta)}_{\text{log-normalization or cumulant generating function or log partition fct.}}) \underbrace{h(x) d\mu(x)}_{\text{defining pieces } (+ \Omega_X)}$$

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

" " cts. " " " pdf

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

discrete: $\left[\sum_x p(x, \eta) \right]$

$$\Rightarrow A(\eta) \triangleq \log \left(\underbrace{\int_X \exp(\eta^T T(x)) h(x) d\mu(x)}_{Z(\eta)} \right)$$

domain $\Omega \triangleq \{ \eta \in \mathbb{R}^p \mid A(\eta) < \infty \}$
set of valid canonical parameters

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

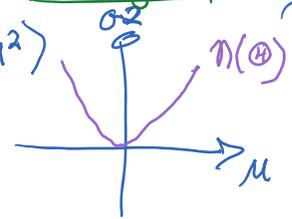
⊕ more generally, consider a reparameterization of a subset of the flat family

by defining $\eta: \Theta \rightarrow \Omega$
↑ new set of parameters

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

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

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



* note: any single dist. $p(x)$ can be put in an exponential family by using $\eta(x) = p(x)$

* two examples of family not an exp. family: • unif $(0, \theta)$
• mixture of Gaussians (latent variable model)

Example 1: (multinoulli)

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

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

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

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

think as "dots"

$$= \exp\left(\sum_{j=1}^k x_j \log \pi_j - 0\right)$$

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

$$\Omega_X \subseteq \mathbb{R}^k$$

$$T(x) = x$$

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

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

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

$$\Theta \rightarrow \text{dimension } k-1$$

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

$$\Omega_X \rightarrow \text{" " } k$$

we do not have a "minimal exp family"

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

$$\text{here } \underbrace{\sum_{j=1}^k T_j(x)}_j = \sum_{j=1}^k x_j = 1$$

offer linear dep between the components of T

⇒ multiple η 's give rise to same dist.
→ 'overparameterization'

↳ not a "minimal" exp. family

⊛ for a multinomial, min. exp family

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

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

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

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

for multinomial, $\frac{\partial A}{\partial \eta_j} = \frac{1}{Z(\eta)} e^{\eta_j} = p("x=j")(\eta)$
 $= \mathbb{E}_{p(x; \eta)} [T_j(x)]$ as required //

moment matching can be different from MLE in exp. family

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

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

16/07

example 2: 1d Gaussian

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

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

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

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

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

$\eta_1 = \eta_2 \cdot \mu$

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

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

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

[we'll see later : multivariate Gaussian]

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

Canonical parameters

$$\begin{cases} \eta = \mu \\ \Lambda = \Sigma^{-1} \end{cases}$$

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

example 3: discrete UGM

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

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

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

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

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

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

$$\eta(\theta) = \begin{pmatrix} \vdots \\ \log \psi_c(y_c) \\ \vdots \end{pmatrix} \leftarrow \begin{matrix} y_c \in \mathcal{X}_c \\ c \in \mathcal{C} \end{matrix}$$

[not a minimal representation?]

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

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

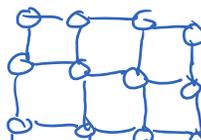
for example: suppose x is a sentence
 x_i is a word

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

~ much smaller set of parameters

c) binary Ising model

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



suppose use nodes & pairs (edges) as cliques



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

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

\Rightarrow not a min. exp family

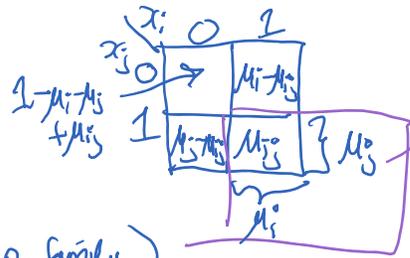
* a minimal representation

$$\bar{\theta} T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{\{i, j\} \in E} \end{pmatrix} \begin{matrix} \eta_i \\ \eta_{ij} \end{matrix}$$

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

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

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



\rightarrow numbers are sufficient to parameterize

properties of A : (for generic joint exp family)

$\cdot \nabla_n A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)] \triangleq \mu(\eta)$ "moment vector" (for $\eta \in \text{int}(\mathcal{R})$)

$\cdot \left(\frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j} \right)_{i, j} \approx \mathbb{E}_{p(x; \eta)} [(T(x) - \mu(\eta)) (T(x) - \mu(\eta))^T] = \text{cov}(T(x))$
(proof as exercise)

"cumulant generating fct."

Estimation of parameters DGM/UGM

DGM (fully observed)

parametric family $\mathcal{P}_{\mathcal{H}} = \{ p_{\theta}(x) = \prod_i p(x_i | x_{\mathcal{H}_i}, \theta_i) : \theta = (\theta_1, \dots, \theta_{|V|}) \}$
 $\mathcal{H}_{\mathcal{H}} = \mathcal{H}_1 \times \mathcal{H}_2 \times \dots \times \mathcal{H}_{|V|}$

independent parameterization i.e. no tying of parameters

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

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

$$\log(\quad) = \sum_{j=1}^{|V|} \underbrace{\left(\sum_{i=1}^n \log p(x_j^{(i)} | x_{\mathcal{H}_j}^{(i)}; \theta_j) \right)}_{\text{in } \mathcal{H}_j}$$

example: for discrete R.V. $\Rightarrow \theta_j^{MLE} = \frac{\sum_j f_j(x_j)}{\sum_j f_j(x_j)}$ = proportion of observations
 (multinomial conditions)

$$\hat{p}(x_j = k | x_{\text{rest}} = \text{stuff}) = \frac{\#(x_j = k, x_{\text{rest}} = \text{stuff})}{\#(x_{\text{rest}} = \text{stuff})}$$

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

⊗ if have latent variables (ie. unobserved variables)

\Rightarrow use E.M. (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 as $\sum_c f_c(\eta_c)$

gradient ascent on log-likelihood

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

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

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

to compute this, need inference

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

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

here need approximate inference \leftarrow sampling variations [mean field] [Gibbs sampling]