

- today:
- exp. family
 - estimation PGM
 - sampling

Exponential family

a (flat/canonical) exponential family on X

is a parametric family of dist. defined by two quantities

I) $h(x) d\mu(x)$ → reference measure

reference density base measure

counting measure (discrete R.V.)
Lebesgue measure (cts. R.V.)

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

members of the family will have dist.

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

defining pieces ($+ \Omega_X$)

"canonical parameter" → log normalization or cumulant gen. fct.
or log partition fct.

if Ω_X is discrete, then $p(x; \mathcal{N})$ is a pmf
" " cts. " " " pdf

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

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

$\mathcal{Z}(\eta)$

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 family

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

consider $p(x; \theta) \triangleq p(x; \eta(\theta))$ for $\theta \in \Theta$

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

↳ eg could consider Gaussians where $N(\mu, \sigma^2)$



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

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

Example: (Multinomial)

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

$\Omega_X = \Delta_k \cap \mathcal{X}$ (one hot encoding)

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

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

think as \log

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

we have $\eta_j(\pi) = \log \pi_j \quad \Omega = \mathbb{R}^k$

$T(x) = x$

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

$h(x) = \mathbb{1}\{x \in \Omega_X\} = \mathbb{1}\{x \in \Delta_X \cap \mathcal{X}\}$

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

$\Theta \rightarrow$ dimension $k-1$
 $\eta(\Theta) \rightarrow$ " $k-1$
 $\Omega \rightarrow$ $k \quad k$

we do not have a "minimal exp family"

note: for any x st. $h(x) \neq 0$

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

affine linear dep. between of T

\Rightarrow multiple η 's give rise to same distribution
"overparameterization"

\hookrightarrow not a "minimal" exp family

(*) for multinomial, minimal exp family

$$T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix}$$

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

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

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

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

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

14h29

example 2: 1d Gaussian

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

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

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

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

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

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

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

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

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

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

Example 3: discrete UGM?

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

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

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

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

$$T(x) = \begin{pmatrix} \vdots \\ \mathbb{1}\{x_c = y_c\} \end{pmatrix}_{\substack{y_c \in X_c \\ c \in \mathcal{C}}}$$

$$X_c = \{(y_i)_{i \in c} : y_i \in X_i\}$$

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

[not a minimal representation]

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

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

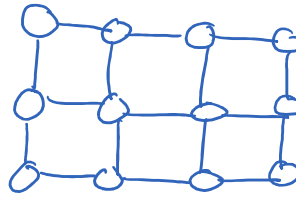
for example: suppose x is sentence
 x_i is a word

feature on x_i : $\{x_{i \in c}\}$ e.g. $\mathbb{1}\{x_i \text{ is a verb}\}$

$x_i \rightarrow$ is a noun

c) binary Ising model

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



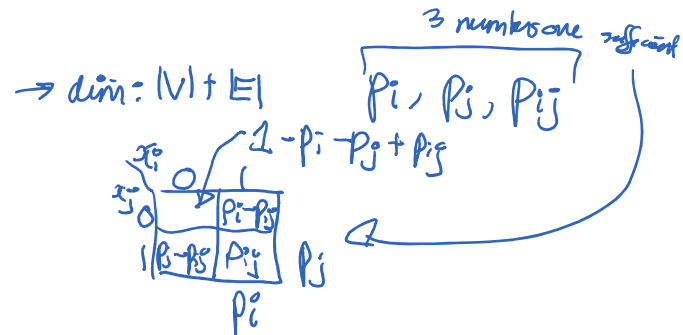
suppose use nodes & pairs (edges) as cliques

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

$\sum_{y_C} T_{C, y_C}(x) = 1$ for any $C \Rightarrow$ not a min. exp families

* a minimal representation:

$T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{\{i, j\} \in E} \end{pmatrix}$
 \downarrow
 $\mathbb{1}_{\{x_i=1, x_j=1\}}$



properties of A_θ

- $\mathbb{D}_n A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)] \triangleq \mu(\eta)$ "moment vector" (for $\eta \in \text{int}(\Omega)$)
- $\left(\frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j} \right)_{ij} = \mathbb{E}_{p(x; \eta)} [(T(x) - \mu(\eta))(T(x) - \mu(\eta))^T] = \text{cov}(T(x))$
 (proof as exercise)

Estimation of parameters in DGM

DGM: parametric family $\mathcal{P}_\Theta = \{ f_\theta(x) = \prod_i p(x_i | x_{\pi_i}, \theta_i) \}$

$\Theta = (\theta_1, \dots, \theta_{|V|})$

$\in \mathcal{H} = \mathcal{H}_1 \times \mathcal{H}_2 \times \dots \times \mathcal{H}_{|V|}$

ie. no tying of parameters

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

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

... $\sum_{i=1}^n \dots \sum_{j=1}^n \dots$

$$\log [] = \sum_{j=1}^n \underbrace{\left(\sum_{i=1}^n \log p(x_j^{(i)} | x_{-j}^{(i)}; \theta_j) \right)}_{f_j(\theta_j)}$$

example: for discrete RV. $\Rightarrow \theta_j^{MLE} = \frac{\text{proportion of observations } \#(x_j = k, x_{-j} = \text{something})}{\#(x_j = \text{something})}$

(fully observed DGM is relatively easy)

⊕ if have latent variable (i.e. unobserved variables)
 \Rightarrow use EM. (like HMM)

UGM:

example for exp family

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

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

gradient ascent on log likelihood

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

$$\nabla_{n_c} [] = \hat{\mu}_c - \mu_c(n)$$

$\hookrightarrow \mathbb{E}_{p(x;n)} [T_c(x_c)]$
 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 approx. inference
 sampling
 variational method