

Def. Let  $X$  be an  $m$ -dimensional random vector with sufficient statistics  $\phi(x) \in \mathbb{R}^d$ . An exponential family density with  $d$ -dimensional natural parameter vector  $\eta$  has the form:

$$p(x) = h(x) \exp\{\eta^T \phi(x) - A(\eta)\}$$

with base measure  $h(x)$  and log-partition function:

$$A(\eta) = \log \int \exp\{\eta^T \phi(x)\} h(x) dx$$

ASIDE

$$\sigma^{-2} = \exp(-\log \sigma)$$

$$= \exp(-\frac{1}{2} \log \sigma^2)$$

Ex. GAUSSIAN

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

where,

$$\eta = \begin{pmatrix} \sigma^{-2} \mu \\ -\frac{1}{2}\sigma^{-2} \end{pmatrix}, \quad \underbrace{h(x) = \frac{1}{\sqrt{2\pi}}}_{\text{constant in } x}, \quad A(\eta) = \frac{1}{2}\mu^2 \sigma^{-2} + \frac{1}{2} \log \sigma^2$$

MEAN PARAMETERS

- Alternate  $d$ -dimensional set of mean parameters given by:

$$\mu_i = \mathbb{E}[\phi_i(x)] \quad \text{for } i=1, \dots, d$$

-  $\eta(\mu)$  is an invertible mapping from mean to natural parameters  
 ↳ if minimal !!

## EX. GAUSSIAN PARAMETERS

$$M = \begin{pmatrix} \mathbb{E}[x] \\ \mathbb{E}[x^2] \end{pmatrix} = \begin{pmatrix} m \\ \sigma^2 + m^2 \end{pmatrix} \iff \eta(\mu) = \begin{pmatrix} \sigma^{-2} m \\ -\frac{1}{2} \sigma^{-2} \end{pmatrix}$$

- ASIDE: Gaussian mean params more conveniently written w/ variance:  
 $M' = (\mathbb{E}[x], \mathbb{E}[x^2] - \mathbb{E}[x]^2)^T = (m, \sigma^2)^T$

Def. An expfam  $p_\eta(x)$  is minimal if both of the following hold:  
a) Natural parameters are linearly independent,  
b) Suff. stats are linearly independent

Def. An expfam  $p_\eta(x)$  is overcomplete if it is not minimal

## EX. OVERCOMPLETE CATEGORICAL:

- Discrete RV  $X \in \{1, \dots, k\}$ , mean parameters  $\mu$  s.t.  
 $\{\mu \in \mathbb{R}^k : \mu_1 + \dots + \mu_k = 1, \mu_i \geq 0 \forall i = 1, \dots, k\}$

$$p(x) = \prod_{i=1}^k \mu_i^{\mathbb{I}(x=i)} = \exp\left(\sum_i \mathbb{I}(x=i) \log \mu_i\right)$$

But,

$$\mu_k = 1 - \sum_{i=1}^{k-1} \mu_i \quad \& \quad \mathbb{I}(x=k) = 1 - \sum_{i=1}^{k-1} \mathbb{I}(x=i)$$

Minimal Parameterization:

$$p(x) = \exp\left\{\sum_{i=1}^{k-1} \mathbb{I}(x=i) \log \mu_i + \left(1 - \sum_{i=1}^{k-1} \mathbb{I}(x=i)\right) \log\left(1 - \sum_{i=1}^{k-1} \mu_i\right)\right\}$$

## EXP FAM PROPERTIES

⇒ Set of natural parameters  $\eta \in H$  is convex:

Pf: Show for any  $\eta_1, \eta_2 \in H$  &  $\lambda \in [0, 1]$ :

$$\int \exp(\lambda \eta_1^T + (1-\lambda) \eta_2^T) \phi(x) h(x) dx < \infty$$

We have:

$$\int \exp\{\lambda \eta_1^T + (1-\lambda) \eta_2^T\} \phi(x) h(x) dx$$

$$= \int e^{\lambda \eta_1^T \phi(x)} e^{(1-\lambda) \eta_2^T \phi(x)} h(x) dx$$

$$\leq \left( \int e^{\eta_1^T \phi(x)} h(x) dx \right)^\lambda \left( \int e^{\eta_2^T \phi(x)} h(x) dx \right)^{(1-\lambda)} \quad (2)$$

ASIDE: See gradient properties of  $A(\eta)$  for alt. proof. (1)

(By Holder's inequality:  $\int |f(z)g(z)| dz \leq \left( \int |f(z)|^{\frac{1}{\lambda}} dz \right)^\lambda \left( \int |g(z)|^{\frac{1}{1-\lambda}} dz \right)^{1-\lambda}$ )  
 $< \infty$  (since  $\eta_1, \eta_2 \in H$ )

□

⇒ Log-partition  $A(\eta)$  is convex in  $\eta$ :

Pf: Exponentiating both sides of (1) & (2) yields:

$$e^{A(\lambda \eta_1 + (1-\lambda) \eta_2)} \leq e^{\lambda A(\eta_1) + (1-\lambda) A(\eta_2)}$$

taking the log yields:

$$A(\lambda \eta_1 + (1-\lambda) \eta_2) \leq \lambda A(\eta_1) + (1-\lambda) A(\eta_2)$$

□

ASIDE: exp and log are strictly increasing monotonic fun's

⇒ Closed under multiplication:

Pf: Let  $X_i \stackrel{iid}{\sim} P_{\eta_i}(x)$  for  $i=1, \dots, N$  and  $P_{\eta}(\cdot)$  is exp fam w/ suff. stats  $\phi(x)$ . The joint density is then,

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p_{\eta}(x_i) = \prod_i h(x_i) \exp\{\eta^T \phi(x_i) - A(\eta)\}$$

$$= \left(\prod_i h(x_i)\right) \exp\left\{\eta^T \left(\sum_i \phi(x_i)\right) - nA(\eta)\right\}$$

↑ Add suff. stats. □

Pf: Let  $p_{\eta_i}(x)$  be expfam w/ parameters  $\eta_1, \dots, \eta_n$  then the product is a density in the expfam:

$$p(x) \propto \prod_{i=1}^n p_{\eta_i}(x) = h(x)^n \exp\left\{\phi(x)^T \left(\sum_i \eta_i\right) - \sum_i A(\eta_i)\right\}$$

w/ Natural parameters  $\eta = \sum_i \eta_i$  □

⇒ Differentiating log-partition w.r.t natural parameters yields non-central moments:

$$\underbrace{\frac{\partial A}{\partial \eta_i} = \mathbb{E}[\phi_i(x)] = \mu_i}_{\text{General inverse mapping between moment-natural params.}}, \quad \frac{\partial^2 A}{\partial \eta_i \partial \eta_j} = \mathbb{E}[\phi_i(x) \phi_j(x)]$$

Pf: We just show the first derivative result,

$$\frac{\partial A}{\partial \eta_i} = \frac{\partial}{\partial \eta_i} \log\left(\int h(x) \exp\{\langle \eta, \phi(x) \rangle\} dx\right) = e^{-A(\eta)} \frac{\partial}{\partial \eta} \exp\{\langle \eta, \phi(x) \rangle\}$$

$$= \int \exp\{\langle \eta, \phi(x) \rangle - A(\eta)\} \phi_i(x) dx = \mathbb{E}[\phi_i(x)] = \mu_i$$

→ Note: 2<sup>nd</sup> derivative result also proves convexity of  $A(\eta)$

$$\text{Since } \frac{\partial^2 A}{\partial \eta_i^2} - \left(\frac{\partial A}{\partial \eta_i}\right)^2 = \text{VAR}(\phi_i(x)) \geq 0$$

(4)

## CONJUGACY:

Ex. Beta-Bernoulli:

- Consider  $\theta \sim \text{Beta}(\alpha, \beta)$ ,  $X|\theta \sim \text{Bernoulli}(X|\theta)$  for  $X \in \{0, 1\}$

$$p(\theta)p(X|\theta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^X (1-\theta)^{(1-X)}$$

$$= \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \theta^{\alpha-1+X} (1-\theta)^{(\beta-1+1-X)}$$

$$\propto \text{Beta}(\theta | \alpha+X, \beta+1-X)$$

$\Rightarrow$  All members of the exponential family have conjugate priors:

PF: Let  $\{x_i\}_{i=1}^N \stackrel{iid}{\sim} p(\cdot|\eta)$  be from an exponential family w/ PDF:

$$p(x_i|\eta) = h(x_i) \exp\{\eta^T \phi(x_i) - A(\eta)\}$$

Let  $\eta \sim q(\eta|\lambda)$  have prior density,

$$q(\eta|\lambda) = g(\eta) \exp\{\lambda_1^T \eta - \lambda_2 A(\eta) - B(\lambda)\}$$

which is conjugate since,

$$p(\eta|x_1, \dots, x_N) \propto q(\eta|\lambda) \prod_{i=1}^N p(x_i|\eta)$$

$$= g(\eta) \exp\{\lambda_1^T \eta - \lambda_2 A(\eta) - B(\lambda)\} \cdot$$

$$\prod_{i=1}^N h(x_i) \exp\{\eta^T \phi(x_i) - A(\eta)\}$$

$$\propto g(\eta) \exp\{\eta^T (\lambda_1 + \sum_{i=1}^N \phi(x_i)) + (\lambda_2 + N)(-A(\eta))\}$$

$$\propto q(\eta | (\lambda_1 + \sum_{i=1}^N \phi(x_i), \lambda_2 + N))$$

ASIDE:

Previously showed  
this is explain  
due to closure

# INFORMATION THEORY:

(1) Kullback-Liebler divergence:

$$KL(p||q) = H_p(q(x)) - H_p(x) = \mathbb{E} \left[ \log \frac{p(x)}{q(x)} \right]$$

$\uparrow$                        $\uparrow$   
 CROSS ENTROPY      ENTROPY

- Closest expfam  $q(x)$  approx. of  $p(x)$  in KL sense:

$$q^*(x) = \operatorname{argmin}_q KL(p||q)$$

$$= \operatorname{argmin} H_p(q(x)) + \text{const.}$$

ASIDE: Channel coding interpretation:  $KL(p||q)$  is loss incurred if I code  $x \sim p$  as  $x \sim q$ .

$$H_p(q(x)) = \mathbb{E}_p[-\log q(x)]$$

$$= \mathbb{E}_p[-\log h(x) - \eta^T \phi(x) + A(\eta)]$$

- Take derivative wrt  $\eta_i$ , set to zero, solve:

$$\frac{\partial H}{\partial \eta_i} = -\mathbb{E}_p[\phi_i(x)] + \mathbb{E}_{q_\eta}[\phi_i(x)] = 0 \Rightarrow$$

ASIDE: Recall  $\frac{\partial A}{\partial \eta_i} = \mathbb{E}_q[\phi_i(x)]$

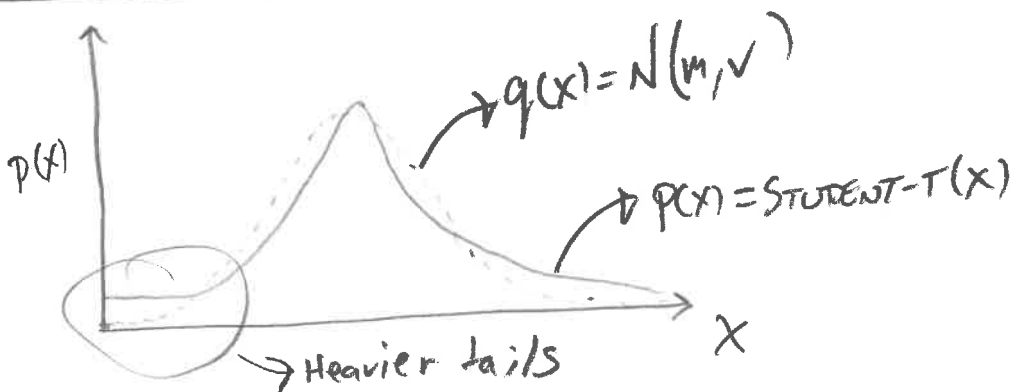
$$\mathbb{E}_{q_\eta}[\phi_i(x)] = \mathbb{E}_p[\phi_i(x)]$$

- This is the well-known "moment-matching" property

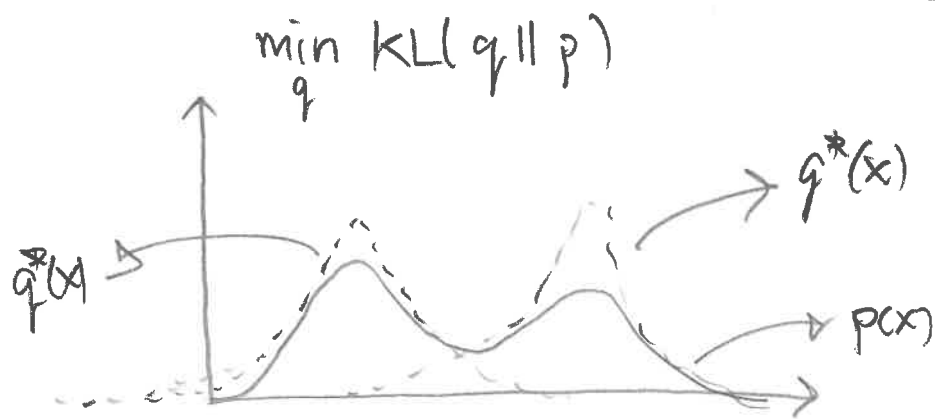
- Note: The sol'n is unique since:

$$\nabla_\eta^2 H = \nabla_\eta^2 A(\eta) = \mathbb{E}_q[\phi(x)\phi(x)^T] \succeq 0 \text{ PSD}$$

## EX. GAUSSIAN PROJECTION



② Reverse KL projection is non-convex in general:



- $KL(q||p)$  called "Exclusive KL" because there is high penalty when  $q(x)$  high in regions  $p(x)$  low
- $\min KL(q||p)$  tends to capture modes of  $p(x)$
- We will see that  $KL(q||p)$  important in variational inference

③ Maximum entropy principle

Thm. Given RV  $X \sim p$  and statistics  $\{\phi_i(x)\}_{i=1}^d$  w/ known expectations  $E_p[\phi_i(x)] = \mu_i$ , the dist.  $q(x)$  of maximum entropy  $q(x) = \arg \max H_q(x)$ , subject to moment constraints, is in the expfam w/ suff. stats.  $\{\phi_i(x)\}_{i=1}^d$ .

Pf Overview: Can be proven by forming Lagrangian w/ Lagrange multipliers  $\vec{\eta}$  for moment constraints.

Problem is convex and sol'n is  $q_{\vec{\eta}}(x)$

↳ Lagrange multipliers are natural params