

# ① Review of exponential family

- $p(x|\eta) = h(x) \exp(\eta^T t(x) - a(\eta))$ 
  - $\eta \in \mathbb{R}^d$  "natural parameter"
  - $t(x) \in \mathbb{R}^d$  "sufficient statistics"
  - $h(x) \in \mathbb{R}$  "base measure"
  - $a(\eta) \in \mathbb{R}$  "log normalizer"
- $h(x) \exp(\eta^T t(x))$  is the "kernel"
- $a(\eta) = \log \int \underbrace{h(x) \exp(\eta^T t(x))}_{\text{kernel}} dx$

typically there is a "standard" parameterization

- $P(x|\theta) = h(x) \exp(g(\theta)^T t(x) - a(\theta))$
- $g(\theta) = \eta$  is the "link function"
- $g^{-1}(\eta) = \theta$  is the "inverse link"

there is also a "mean parameterization"

$$\begin{aligned}
 \nabla_{\eta} a(\eta) &= \nabla_{\eta} \log \int h(x) \exp(\eta^T t(x)) dx \\
 &= \frac{\nabla_{\eta} \int \dots dx}{\int \dots dx} = \frac{\int h(x) \nabla_{\eta} \exp(\eta^T t(x)) dx}{\int \dots dx} \\
 &= \int t(x) \frac{h(x) \exp(\eta^T t(x)) dx}{\int \dots dx} \\
 &\stackrel{\eta}{=} \mathbb{E}[t(x)] = \mu
 \end{aligned}$$

Example: Binomial (for fixed  $n$ )

$$X \sim \text{Binom}(N, \theta)$$

$$p(X | \theta, N) = \binom{N}{X} \theta^X (1-\theta)^{N-X}$$

$$= \binom{N}{X} \exp \left( X \log \theta + (N-X) \log (1-\theta) \right)$$

$$= \binom{N}{X} \exp \left( X \log \frac{\theta}{1-\theta} + \underbrace{N \log (1-\theta)}_{\stackrel{\downarrow}{g(\theta)}} \right)$$
$$\stackrel{\downarrow}{h(x)} \quad \stackrel{\downarrow}{t(x)} \quad \stackrel{\downarrow}{g(\theta)} \quad \stackrel{\triangle}{=} -a(\theta)$$

$$= h(x) \exp \left( t(x)^T g(\theta) - a(\theta) \right)$$

$$g(\theta) = \text{logit}(\theta) = \log \frac{\theta}{1-\theta} = \eta$$

$$g^{-1}(\eta) = \text{logistic}(\eta) = \frac{e^\eta}{1+e^\eta} = \theta$$

$$= h(x) \exp \left( t(x)^T \eta - \underbrace{N \log (1+e^\eta)}_{\equiv a(\eta)} \right)$$

confirm that  $E(X) = N\theta$ :

$$M = \mathbb{E}_\eta a(\eta) = \mathbb{E}_\eta N \log (1+e^\eta)$$

$$= N \frac{\mathbb{E}_\eta ((1+e^\eta))}{(1+e^\eta)} = \frac{N}{1+e^\eta} = N\theta$$

## ② Conjugacy

Recall beta-binomial conjugacy:

$$\theta \sim \text{Beta}(a, b)$$

$$x_i \stackrel{\text{ind.}}{\sim} \text{Binomial}(N_i, \theta), \quad i=1 \dots n$$

$$P(\theta | x_1, \dots, x_n, a, b) \propto \theta^{a-1} (1-\theta)^{b-1} \prod_{i=1}^n \theta^{x_i} (1-\theta)^{N_i - x_i}$$

$$\propto \theta^{\underbrace{a + \sum_{i=1}^n x_i}_{\triangleq a'}} (1-\theta)^{\underbrace{b + \sum_{i=1}^n (N_i - x_i)}_{\triangleq b'}}$$

This is an instance of a general "pattern w/ explicit dists."

---

Generally:

$$\eta \sim F(\cdot | \lambda)$$

$$x_i \stackrel{\text{ind.}}{\sim} G(\cdot | \eta) \quad i = 1 \dots n$$

conjugate if

$$P(\eta | x_{1:n}, \lambda) \propto f(\eta | \lambda) \prod_{i=1}^n g(x_i | \eta)$$

$$\propto F(\eta | \lambda_n)$$

(same family as prior)

All exp fams  $G(\cdot | \eta)$  have a conjugate prior.  
(likelihood)

$$P(x_i | \eta) = h_e(x_i) \exp \left( \eta^\top t_e(x_i) - a_e(\eta) \right)$$

$$P(\eta | \lambda) = h_c(\eta) \exp \left( \lambda_1^\top \eta + \lambda_2 (-a_e(\eta)) - a_c(\lambda) \right)$$

- $\lambda = [\lambda_1, \lambda_2] \in \mathbb{R}^{\dim(\eta) + 1}$  "natural parameter"

- $t_c(\eta) = [\eta, -a_e(\eta)]$  "suff stats"

Example: Binomial's conjugate prior

$$\eta = g(\theta) = \text{logit } \theta, \quad \dim(\eta) = 1$$

$$a_\ell(\eta) = N \log(1 + e^\eta)$$

Set  $N=1$  (for now):

$$p(\eta | \lambda) \propto h_c(\eta) \exp\left(\lambda_1 \eta - \lambda_2 \log(1 + e^\eta)\right)$$

notice that  $\eta = g_e(\theta)$  is the likelihood's natural parameterization at  $\theta$ . When deriving the prior, we will often want to think about  $\theta = g_e^{-1}(\eta)$ .

change of variables to  $\theta \in (0,1)$ :

$$p_\theta(\theta | \lambda) = \left| \frac{\partial}{\partial \theta} g_e(\theta) \right| p_\eta(g_e(\theta) | \lambda)$$

$$\frac{\partial}{\partial \theta} \left| \frac{\partial}{\partial \theta} \log \frac{\theta}{1-\theta} \right| \mathbb{1}_{(\theta \in (0,1))} h_c(g_e(\theta)) \exp\left(\lambda_1 \log \frac{\theta}{1-\theta} + \lambda_2 \log(1-\theta)\right)$$

$$\frac{\partial}{\partial \theta} \frac{1}{\theta(1-\theta)} \mathbb{1}_{(\theta \in (0,1))} h_c(g_e(\theta)) \exp\left(\lambda_1 \log \theta + [-\lambda_1 + \lambda_2] \log(1-\theta)\right)$$

We recognize the kernel of a Beta:

Let  $h_c(g_e(\theta)) = 1, \alpha = \lambda_1, \beta = -\lambda_1 + \lambda_2$ :

$\alpha \theta^{\alpha-1} (1-\theta)^{\beta-1} \rightarrow$  conjugate prior: Beta.

③ Conjugate posterior under indep. sampling

$$\eta \sim F(\cdot | \lambda)$$

$$x_i \stackrel{\text{ind.}}{\sim} G(\cdot | \eta) \quad i=1 \dots n$$

$$P(\eta | x_1, \dots, x_n, \lambda) \propto \pi(\eta | \lambda) \prod_{i=1}^n g(x_i | \eta)$$

$$\propto \eta^{h_c(\eta)} \exp(\lambda_1^\top \eta + \lambda_2(-a_e(\eta)) - a_c(\lambda)) \\ \left[ \frac{1}{\pi} \prod_{i=1}^n h_e(x_i) \exp(\eta^\top t_e(x_i) - a_e(\eta)) \right]$$

$$\propto \eta^{h_c(\eta)} \exp(\lambda_1^\top \eta + \lambda_2(-a_e(\eta))) \\ \prod_{i=1}^n \exp(\eta^\top t_e(x_i) - a_e(\eta))$$

$$\propto \eta^{h_c(\eta)} \exp\left(\eta^\top \left[\lambda_1 + \sum_{i=1}^n t_e(x_i)\right] - a_e(\eta)[\lambda_2 + n]\right)$$

Same family as  $F(\cdot | \lambda)$  prior but with

updated natural parameters:

$$\lambda_h = \begin{bmatrix} \lambda_{h1} \\ \lambda_{h2} \end{bmatrix} = \begin{bmatrix} \lambda_1 + \sum_{i=1}^n t_e(x_i) \\ \lambda_2 + n \end{bmatrix}$$

↳ posterior updating = adding / incrementing!

Example :

$$\theta \sim \text{Beta}(\alpha_1, \beta)$$

$$x_i \stackrel{\text{ind.}}{\sim} \text{Binom}(N_i, \theta) \quad i = 1, \dots, n$$

$$f(\theta | \cdot) \propto \theta^{\alpha + \sum_i x_i - 1} (1 - \theta)^{\beta + \sum_i (N_i - x_i) - 1}$$

$$\alpha_n = \alpha + \sum_{i=1}^n x_i \quad , \quad \beta_n = \beta + \sum_{i=1}^n (N_i - x_i)$$

From above:

$$\lambda_{n1} = \alpha_n \quad , \quad -\lambda_{n1} + \lambda_{n2} = \beta_n$$

So:  $\leftarrow (x_i \equiv t_\ell(x_i))$

$$\lambda_{n1} = \lambda_1 + \sum_{i=1}^n x_i$$

$$\lambda_{n2} = \beta_n + \alpha_n = \beta + \sum_{i=1}^n (N_i - x_i) + \alpha + \sum_{i=1}^n x_i$$

$$= \alpha + \beta + \sum_{i=1}^n N_i$$

$$= \lambda_1 + (\lambda_2 - \lambda_1) + \sum_{i=1}^n N_i$$

$$= \lambda_2 + \sum_{i=1}^n N_i$$

④ Likelihood principle

Note the "sample size" above is  $\sum_{i=1}^n N_i$ .

This is because  $x_i \sim \text{Binom}(n_i, \theta)$  is the same as...

$$b_{ij} \stackrel{iid}{\sim} \text{Bin}(1, \theta), \quad j = 1 \dots N_i$$

... from the perspective of D.

$$\prod_{i=1}^n \text{Binom}(x_i; N_i, \theta) \quad \text{and} \quad \prod_{i=1}^n \prod_{j=1}^{N_i} \text{Binom}(b_{ij}; 1, \theta)$$

We can also understand this as satisfying the likelihood principle: all information in the data about  $\theta$  should be contained in the likelihood function:

If two likelihood fractions are proportional in  $\theta$ , then they contain the same information.

( $\hookrightarrow$ ) posterior (Bayesian) inference satisfies the likelihood principle generically.

Another example:  $X_i \stackrel{\text{ind.}}{\sim} NB(r_i, \theta)$   $i = 1 \dots n$

$$\frac{n}{\prod_{i=1}^n} \text{NB}(x_i; r_i, \theta) \propto_\theta \theta^{\sum_{i=1}^n x_i} (1-\theta)^{\sum_{i=1}^n r_i}$$

(negative binomial)

↳ this "experiment" contains the same info about  $\theta$  as observing  $\sum_{i=1}^n x_i$  "successes" and  $\sum_{i=1}^n r_i$  "failures" in  $\sum_{i=1}^n (x_i + r_i)$  indep. Bernoulli trials

## ⑤ Posterior Predictive in conjugate models

$$\begin{aligned}
 & P(X_{n+1} | X_1, \dots, X_n) \\
 &= \int P(X_{n+1} | \eta) p(\eta | X_{1:n}) d\eta \\
 &= \int g(x_{n+1} | \eta) f(\eta | \lambda_n) d\eta \\
 &= \int h_e(x_{n+1}) \exp\left(t(x_{n+1})^\top \eta - a_e(\eta)\right) \\
 &\quad h_c(\eta) \exp\left(\eta^\top \lambda_{n1} + \lambda_{n2}(-a_e(\eta)) - a_c(\lambda_n)\right) d\eta \\
 &= h_e(x_{n+1}) \exp(-a_c(\lambda_n)) \underbrace{\int \exp\left(\eta^\top [t(x_{n+1}) + \lambda_{n1}] + (-a_e(\eta))[\lambda_{n2} + 1]\right) d\eta}_{=} \\
 &\quad = \exp(a_c([\lambda_{n1} + t(x_{n+1}), \lambda_{n2} + 1])) \\
 &= h_e(x_{n+1}) \frac{\exp(a_c([\lambda_{n1} + t(x_{n+1}), \lambda_{n2} + 1]))}{\exp(a_c([\lambda_{n1}, \lambda_{n2}]))}
 \end{aligned}$$

ratio of normalizers!

## Example : Beta binomial

- $a_c(\lambda_1, \lambda_2) = \mathcal{B}(\lambda_1, \lambda_2 - \lambda_1)$  beta function
- $h_e(x_{n+1}) = \binom{N}{x_{n+1}}$  binomial coefficient

$$\begin{aligned}
 & P(X_{n+1} | N_{n+1}, X_{1:n}, N_{1:n}) \\
 &= \binom{N_{n+1}}{X_{n+1}} \frac{\mathcal{B}(\lambda_{n1} + X_{n+1}, \lambda_{n2} - \lambda_{n1} + 1)}{\mathcal{B}(\lambda_{n1}, \lambda_{n2} - \lambda_{n1})} \\
 &\equiv \binom{N_{n+1}}{X_{n+1}} \frac{\mathcal{B}(\alpha + \sum_{i=1}^n x_i + X_{n+1}, \beta + \sum_{i=1}^n (N_i - x_i) + N_{n+1} - X_{n+1})}{\mathcal{B}(\alpha + \sum_{i=1}^n x_i, \beta + \sum_{i=1}^n (N_i - x_i))} \\
 &= \text{BetaBinom}(X_{n+1}; N_{n+1}, \alpha_n, \beta_n)
 \end{aligned}$$

⑥ "Prior predictive"

$$n=0: P(X_{n+1} | X_{1:n}) \equiv P(X_1) \\
 = h_e(x_1) \frac{\exp(a_c([\lambda_1 + t(x_1), \lambda_2 + 1]))}{\exp(a_c([\lambda_1, \lambda_2]))}$$

(this is the "marginal likelihood", or  
 "model evidence", or  
 "normalizing constant")

⑦ Alternative parameterization at conjugate prior

Let  $\gamma_1 = t_0 \kappa_0$ ,  $\gamma_2 = \kappa_0$ , so that the prior:

$$F(\eta | t_0, \kappa_0) \propto_{\eta} h_c(\eta) \exp(t_0 \kappa_0^T \eta - \kappa_0 a(\eta))$$

Then the posterior is:

$$P(\eta | t_0, \kappa_0, X_{1:n}) = f(\eta | t_n, \kappa_n)$$

where

$$\kappa_n = \kappa_0 + n$$

$$t_n = \frac{\kappa_0 t_0 + \sum_{i=1}^n t(x_i)}{\kappa_0 + n} = \left( \frac{\kappa_0}{\kappa_0 + n} \right) t_0 + \left( \frac{n}{\kappa_0 + n} \right) \frac{1}{n} \sum_{i=1}^n t(x_i)$$

Recall that the MLE is:

$$\underset{\eta}{\operatorname{argmax}} \sum_{i=1}^n \log G(x_i | \eta)$$

$$= \underset{\eta}{\operatorname{argmax}} \eta^T \left( \sum_{i=1}^n t(x_i) \right) - n a_e(\eta)$$

$\underbrace{\quad}_{\triangleq LL(\eta; x)}$

$$\nabla_{\eta} LL(\eta; x) = \sum_{i=1}^n t(x_i) - n \underbrace{\nabla_{\eta} a_e(\eta)}_{= \mathbb{E}_{\eta}[t(x)]} = 0$$

$\hookrightarrow \hat{\eta}^{\text{MLE}}: \frac{1}{n} \sum_{i=1}^n t(x_i) = \mathbb{E}_{\hat{\eta}^{\text{MLE}}} [t(x)]$  "moment matching"

## ⑧ Shannon entropy and information

"information content" of event  $X = x$

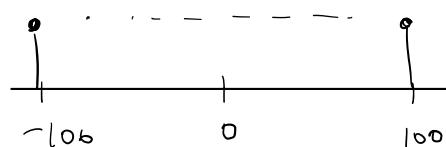
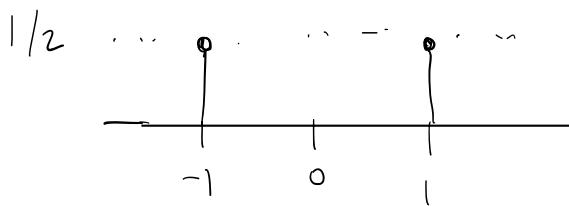
$$h(x) = \log \left[ \frac{1}{p(x)} \right]$$

"entropy" of random variable  $X$

$$H(X) = \mathbb{E}_X \left[ \log \left[ \frac{1}{p(x)} \right] \right]$$

measure of uncertainty

equal entropy



⑨ Expfam dists are uniquely maximum entropy

• consider an expfam dist for  $x \in S$ :  
 $\leftarrow \text{support}$

$$P_\eta(x) = 1(x \in S) h(x) \exp(\eta^T t(x) - a(\eta))$$

• define the set of distributions

$$\mathcal{S}_\alpha \stackrel{\text{def}}{=} \left\{ P : \mathbb{E}_{\substack{x \sim p(x)}} [t(x)] = \alpha, \quad x \in S \right\}$$

all distributions, with support  $S$ , that satisfy  
the moment constraint  $\mathbb{E}[t(x)] = \alpha$ .

• THEOREM: if  $P_\eta(x) \in \mathcal{S}_\alpha$ , then:

$$H(P_\eta) > H(p') \quad \forall p' \in \mathcal{S}_\alpha$$

" $P_\eta$  is the unique maximum entropy dist  
in the set  $\mathcal{S}_\alpha$ "

• Proof: ① show  $H(P_\eta) \geq H(p') \quad \forall p' \in \mathcal{S}_\alpha$   
② show uniqueness

Show  $H(p_n) \geq H(p')$  &  $p' \in S$

- for simplicity, discrete  $X \in S$ :

$$p(X=k) = p(k), \quad \sum_{k \in S} p(k) = 1$$

- introduce Lagrangian, with Lagrange multipliers

$$\eta = (\eta_1, \dots, \eta_d), \eta_0$$

$$\Lambda(p, \eta, \eta_0) = -H(p)$$

$$+ \eta^T (\alpha - \mathbb{E}_{X \sim p} [t(X)]) \quad \begin{matrix} \leftarrow \text{moment} \\ \text{constraints} \end{matrix}$$

$$+ \eta_0 \left( \sum_{k \in S} p(k) - 1 \right)$$

$\nwarrow$  normalization constraint

- argmin  $\Lambda(p, \eta, \eta_0)$  will maximize  $H(p)$

$p$  subject to constraints

- Derivative wrt a single component of  $p(j)$ ,  $j \in S$ :

$$\frac{\partial}{\partial p(j)} \Lambda(p, \eta, \eta_0) = \frac{\partial}{\partial p(j)} \sum_{k \in S} p(k) \log p(k) + \log p(j) + 1$$
$$+ \frac{\partial}{\partial p(j)} \eta^T (\alpha - \sum_k p(k) t(k)) = -\eta^T t(j)$$
$$+ \frac{\partial}{\partial p(j)} \eta_0 \left( \sum_k p(k) - 1 \right) = \eta_0$$

- set to zero

$$\log p(j) + 1 - \eta^T t(j) + \eta_0 = 0$$

$$p(j) = \exp(\eta^T t(j) - \eta_0 - 1)$$

- $p(j) \equiv p_{\eta_0, \eta}(j)$  is parameterized by  $\eta_0, \eta$

- We can drop  $\eta_0$  by setting it to a value that satisfies the constraint it enforces

$$\eta_0 \triangleq a(\eta) - 1$$

- $p(j) = p_\eta(j) = \exp(\eta^T t(j) - a(\eta))$

$$\hookrightarrow \sum_{k \in S} p_\eta(k) = 1 \quad (\text{by defn of } a(\eta))$$

- Now if  $\exists \eta$  s.t.  $\mathbb{E}_{x \sim p_\eta}[t(x)] = \alpha$ , then:

$$p_\eta \in \mathcal{P}_\alpha \quad \text{and} \quad H(p_\eta) \geq H(p'), \quad p' \in \mathcal{S}_\alpha$$

$\hookrightarrow p_\eta$  is a max ent dist in  $\mathcal{S}_\alpha$ .

Show uniqueness:  $H(p_\eta) > H(p')$ ,  $p' \in \mathcal{S}_\alpha$

- consider any  $p \in \mathcal{S}_\alpha$

$$\bullet H(p) = - \sum_k p(k) \log p(k)$$

$$= - \sum_k p(k) \log \left[ p(k) \frac{p_\eta(k)}{p_\eta(k)} \right]$$

$$= - \sum_k p(k) \log \frac{p(k)}{p_\eta(k)} - \underbrace{\sum_k p(k) \log p_\eta(k)}_{= \eta^T t(k) - a(\eta)}$$

$$= -KL(p \parallel p_\eta) - \sum_k p(k) [\eta^T t(x) - a(\eta)]$$

$$= -KL(p \parallel p_\eta) - \underbrace{[\eta^T \mathbb{E}_{x \sim p}[t(x)] - a(\eta)]}_{\text{since } p \in \mathcal{J}_\epsilon, \mathbb{E}_p[t(x)] = \mathbb{E}_{p_\eta}[t(x)]}$$

$$= -KL(p \parallel p_\eta) - \underbrace{\sum_k p_\eta(k) \log p_\eta(k)}_{= H(p_\eta)}$$

so:

$$H(p) = H(p_\eta) - \underbrace{KL(p \parallel p_\eta)}_{\geq 0 \text{ IFF } p = p_\eta}$$

## ⑩ Back to Shannon Information and Entropy

"Game of Submarine" (MacKay Chap 4.1)

- submarine is somewhere in 64-cell search grid

- $P(Z=k) = \frac{1}{64}$  (no prior info)

- $P(\text{hit on first guess}) = P(X_1=1) = \frac{1}{64}$

- info content of event  $X_1=1$ ?

$$h(X_1=1) = \log_2 64 = 6 \text{ bits}$$

- info content of event  $X_1=0$ ?

$$h(X_1=0) = \log_2 \frac{64}{63} \approx 0.0227 \text{ bits}$$

- prob. guessing correctly in second guess:

$$P(X_2=1 | X_1=0) = \frac{1}{63} \leftarrow \text{one cell eliminated by } X_1$$

$$\hookrightarrow h(X_2=0 | X_1=0) = \log_2 \frac{63}{62} \approx 0.023 \text{ bits}$$

- total info content in  $X_1=0, X_2=0$ :

$$h(X_2=0, X_1=0) = \log_2 \frac{1}{P(X_1=0, X_2=0)}$$

$$= \log_2 \frac{64}{63} + \log_2 \frac{63}{62}$$

- total info content in 32 misses:  $X_1=0 \dots X_{32}=0$

$$h(X_1=0 \dots X_{32}=0) = \log_2 \frac{64}{63} + \log_2 \frac{63}{62} \dots + \log_2 \frac{33}{32}$$

$$= \log_2 \frac{64}{32} = \log_2 2 = 1 \text{ bit}$$

- this makes sense since:

$$p(z \in \{1, \dots, 32\}) = \frac{32}{64} = \frac{1}{2}$$

- similarly:

$$\begin{aligned} h(X_{(1)}=0, \dots, X_{(64)}=0, X_{(65)}=1) &= \\ = \log_2 \frac{64}{16} + \log_2 16 &= \underline{6 \text{ bits}} \end{aligned}$$

No matter when you find the sub, the total info gained is 6 bits.

$$\hookrightarrow H(p) = \sum_{k=1}^{64} p(z=k) \log_2 \frac{1}{p(z=k)} = \sum \frac{1}{64} \log_2 64 = 6.$$