

Exponential families

1. definition
2. examples - Poisson, Bernoulli
3. properties of exp. family
4. conjugate priors

Exponential family:

$$p(x|\eta) = \pi(x) \exp\{\eta \cdot t(x) - a(\eta)\}$$

sufficient statistic $x \rightarrow \mathbb{R}^d$
log normalizer
natural parameter $\eta \in \mathbb{R}^d$
base measure carrier distribution

log normalizer

$$a(\eta) = \log \int \pi(x) \exp\{\eta \cdot t(x)\} dx$$

$$p(x|\eta) = \frac{\pi(x) \exp\{\eta \cdot t(x)\}}{\int \pi(x') \exp\{\eta \cdot t(x')\} dx'}$$

integrates to 1

Bernoulli

$$p(x|\theta) = \theta^x (1-\theta)^{1-x} \quad x \in \{0,1\}, \theta \in [0,1] \quad \theta: \text{prob of "heads"}$$

$$= \exp\{x \log \theta + (1-x) \log (1-\theta)\}$$

Brown (1986)

$$= \exp\{x \log(\theta/(1-\theta)) + \log(1-\theta)\}$$

Efron (2023)

Wainwright + Jordan (2008)

$$t(x) = x$$

$$\pi(x) = 1$$

$$\eta = \log(\theta/(1-\theta)), \theta = 1/(1+e^{-\eta})$$

$$a(\eta) = -\log(1-\theta) = \log(1+e^{\eta})$$

exponential family
form of the
Bernoulli

Poisson

$$p(x|\lambda) = \frac{1}{x!} \lambda^x \exp\{-\lambda\} \quad \begin{array}{l} x \in \{0, 1, 2, \dots\} \\ \lambda > 0 \end{array}$$

$$= \frac{1}{x!} \exp\{x \log \lambda - \lambda\}$$

$$\eta = \log \lambda, \quad \lambda = \exp\{\eta\}$$

$$t(x) = x$$

$$\pi(x) = 1/x!$$

$$a(\eta) = \lambda = \exp\{\eta\}$$

Moments of an exp. family

$$\mathbb{E}[t(X)] = \nabla_{\eta} a(\eta)$$

$$\nabla_{\eta} a(\eta) = \nabla_{\eta} \left\{ \log \int \exp\{\eta \cdot t(x)\} \pi(x) dx \right\}$$

$$= \frac{\nabla_{\eta} \int \exp\{\eta \cdot t(x)\} \pi(x) dx}{\int \exp\{\eta \cdot t(x)\} \pi(x) dx}$$

$$= \int t(x) \frac{\exp\{\eta \cdot t(x)\} \pi(x)}{\int \exp\{\eta \cdot t(x')\} \pi(x') dx'} dx$$

← $p(x|\eta)$

$$= \int t(x) p(x|\eta) dx = \mathbb{E}[t(X)]$$

$t(x): \mathcal{X} \rightarrow d\text{-vector.}$
 $\eta: \mathbb{R}^d$

$$\frac{\partial^2 a(\eta)}{\partial \eta_i \partial \eta_j} = \mathbb{E}[t_i(X) t_j(X)] - \mathbb{E}[t_i(X)] \mathbb{E}[t_j(X)] = \text{Cov}(t_i(X), t_j(X))$$

Example: Poisson.

$$\mathbb{E}[X] = \lambda$$

$$\frac{da}{d\eta} = \exp\{\eta\} = \lambda$$

$$\frac{d^2 a}{d\eta^2} = \exp\{\eta\} = \lambda$$

mean parameterization

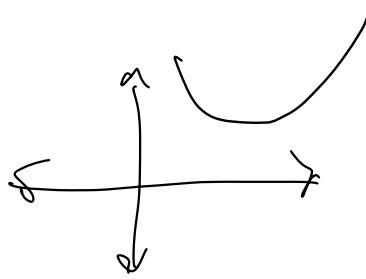
1-1 relationship b/w η and $\mathbb{E}_{\eta}[t(X)]$ ← mean parameterization.

why?

1d expfam.


$t(x)$ scalar, $\eta \in \mathbb{R}$

$$\frac{d^2 a}{d\eta^2} = \text{Var}(X)$$



Variance > 0

$a(\eta)$ convex

1-1 relationship b/w $\frac{da}{d\eta}$ and η 
" \nwarrow 1-1
 $\mathbb{E}[t(X)]$

notation $\mu \triangleq \mathbb{E}[t(X)]$

$\eta_{\mu}(\mu)$ maps μ to η

$\mu_{\eta}(\eta)$ maps η to μ

Maximum likelihood estimation

$$X_i \sim \text{expfam}(\eta)$$

$$\hat{\eta}_{\text{MLE}} = \arg \max_{\eta} \sum_{i=1}^n \log p(x_i | \eta)$$

- 1- MLE for expfam
- 2- conjugate priors for EF
- 3- GLMs
- 4- hierarchical GLMs

anonclass.com

Code: 5596
1 782

$$\begin{aligned}
 \ell(\eta) &= \sum_{i=1}^n \log p(x_i | \eta) \\
 &= \sum_{i=1}^n (\log \pi(x_i) + \eta \cdot t(x_i) - a(\eta)) \\
 &= \sum_{i=1}^n \log \pi(x_i) + \eta \left(\sum_{i=1}^n t(x_i) \right) - n a(\eta)
 \end{aligned}$$

$$\nabla_{\eta} \ell(\eta) = \left(\sum_i t(x_i) \right) - n \nabla_{\eta} a(\eta)$$

$$\mu_{\eta}(\eta) \triangleq \nabla_{\eta} a(\eta) = \mathbb{E}[t(x) | \eta]$$

$$\left(\sum_i t(x_i) \right) - n \mu = 0$$

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_i t(x_i)$$

$$\hat{\eta}_{MLE} = \eta_{\mu}(\hat{\mu}_{MLE})$$

Conjugacy

expfam EDM:

$$\eta \sim f(\eta; \lambda)$$

$$x_i | \eta \sim \text{expfam}(\eta)$$

observe $x_{1:n}$, goal: $p(\eta | x_{1:n}; \lambda)$

$$p(\eta | \mathbf{x}; \lambda) \propto p(\eta; \lambda) \prod_{i=1}^n p(x_i | \eta)$$

Suppose $p(\eta | \mathbf{x}; \lambda) = f(\eta; \hat{\lambda})$ where $\hat{\lambda}$ is a function of λ, \mathbf{x}

$$p(x | \eta) = \pi_{\eta}(x) \exp\{\eta \cdot t_{\eta}(x) - a_{\eta}(\eta)\} \quad \text{likelihood model}$$

$$p(\eta; \lambda) = \pi_{\eta}(\eta) \exp\{\lambda_1 \cdot \eta + \lambda_2 (-a_{\eta}(\eta)) - a_{\eta}(\lambda_1, \lambda_2)\}$$

$$t_{\eta}(\eta) = (\underbrace{\eta}_{\substack{\uparrow \\ \text{d-vector}}}, \underbrace{-a_{\eta}(\eta)}_{\substack{\uparrow \\ \text{scalar}}}) \quad \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$$

λ_1 : d-vector
 λ_2 : scalar
 $\lambda = (\lambda_1, \lambda_2)$

$$p(\eta | \mathbf{x}; \lambda) \propto \pi_{\eta}(\eta) \exp\{\lambda_1 \cdot \eta + \lambda_2 (-a_{\eta}(\eta))\} \exp\{\eta \cdot (\sum_{i=1}^n t(x_i)) - n a_{\eta}(\eta)\}$$

$$= \pi_{\eta}(\eta) \exp\{(\lambda_1 + \sum_{i=1}^n t(x_i)) \cdot \eta + (\lambda_2 + n) (-a_{\eta}(\eta))\} \quad \leftarrow \text{unnormalized}$$

$$t(\eta) = (\eta, -a_{\eta}(\eta)) \quad a_{\eta}(\hat{\lambda}_1, \hat{\lambda}_2)$$

$$\hat{\lambda}_1 = \lambda_1 + \sum_{i=1}^n t(x_i)$$

$$\hat{\lambda}_2 = \lambda_2 + n$$

$$\mu(\eta)$$

$$\mathbb{E}_{\lambda}[\mu(\eta)] = \lambda_1 / \lambda_2 \quad \text{prior expectation of the mean parameter}$$

$$\mathbb{E}[\mu(\eta) | \mathbf{x}; \lambda] = \frac{\hat{\lambda}_1}{\hat{\lambda}_2} = \frac{\lambda_1 + \sum_{i=1}^n t(x_i)}{\lambda_2 + n}$$

- posterior variance goes down as $1/n$

- posterior predictive distribution.

$$p(x_{\text{new}} | \mathbf{x}; \lambda) = \pi_{\eta}(x_{\text{new}}) \exp\{a_{\eta}(\hat{\lambda}_1 + t(x_{\text{new}}), \hat{\lambda}_2 + 1) - a_{\eta}(\hat{\lambda}_1, \hat{\lambda}_2)\}$$

- Beta-Bernoulli

Generalized linear models (GLMs)

$$\mathbb{E}[Y | X=x_i, \beta] = f(\beta \cdot x_i)$$

↑ response ↑ features ↑ coefficients ("weights")

in linear regression $f(\eta) = \eta$

logistic regression $f(\eta) = \sigma(\eta)$

GLM.

anon class - 4455
4501

given features x_i

$$w_i = \beta \cdot x_i$$

$$\mu_i = f(w_i)$$

$$\eta_i = \eta_\mu(\mu_i)$$

$$y_i \sim \text{expfam}(\eta_i)$$

f : link function or response function

$$\mu_i \triangleq \mathbb{E}[Y | X=x_i, \beta]$$

$$\begin{array}{c} \beta \\ x_i \end{array} \rightarrow w_i \xrightarrow{f = \eta_\mu^{-1}} \mu_i \xrightarrow{\eta_\mu} \eta_i \rightarrow y_i$$

$$\begin{array}{c} \beta \\ x_i \end{array} \rightarrow \eta_i \rightarrow y_i$$

choices: ① link function

② response distribution

put a prior on coefficients $\beta \rightarrow$ Bayesian GLM.

$$\beta \sim N_p(0, \lambda^2)$$

$$y_i | x_i, \beta \sim \text{GLM}(x_i, \beta) \quad i=1 \dots n$$

special link: canonical link $f = \eta_\mu^{-1}$

$$\eta_i = \beta \cdot x_i$$

$$y_i | \eta_i \sim \text{expfam}(\eta_i)$$

MAP inference

$$\beta \sim N_p(0, \lambda^2)$$

$$y_i | x_i, \beta \sim \text{GLM}(x_i, \beta) \quad i=1 \dots n$$

$$\mathcal{L}(\beta) = -\frac{1}{2\lambda^2} \sum_{j=1}^p \beta_j^2 + \sum_{i=1}^n \underbrace{(\eta_i \cdot y_i - a(\eta_i))}_{\log p(y_i | x_i, \beta)}$$

$$l_i \triangleq \eta_i \cdot y_i - a(\eta_i)$$

$$\eta_i = \beta \cdot x_i$$

$$\nabla_{\beta} l_i = \nabla_{\eta} l_i \nabla_{\beta} \eta_i$$

$$= (y_i - \nabla_{\eta} a(\eta_i)) \nabla_{\beta} \eta_i$$

$$= (y_i - \mathbb{E}[Y | X=x_i, \beta]) x_i$$

$$\nabla_{\beta} \mathcal{L} = -\frac{1}{\lambda^2} \beta + \sum_{i=1}^n (y_i - \mathbb{E}[Y | X=x_i, \beta]) x_i$$

signed residual

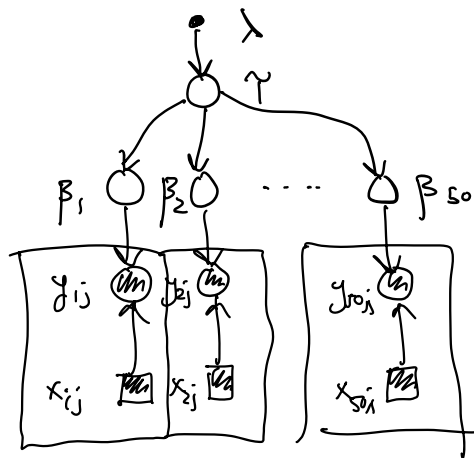
general link function: $\eta_i = \eta_{\mu}(f(\beta \cdot x_i))$

Hierarchical GLMs

data from each state $\mathcal{D}_i = \{(x_{ij}, y_{ij})\}_{j=1}^{n_i}$

How to analyze?

- separate GLMs for each state
- pool all the data, one GLM
- hierarchical model



hierarchical GLM.

$$\gamma \sim p(\gamma | \lambda)$$

For each $i \in \{1, \dots, 50\}$

$$\beta_i \sim p(\beta | \gamma)$$

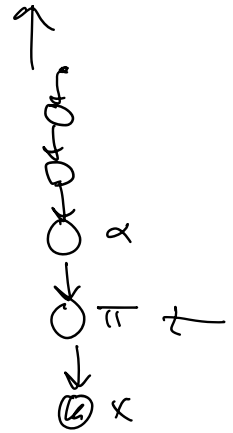
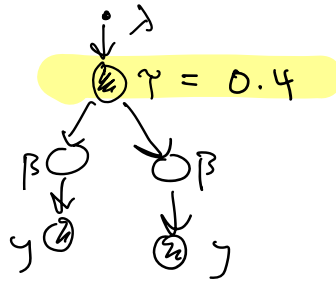
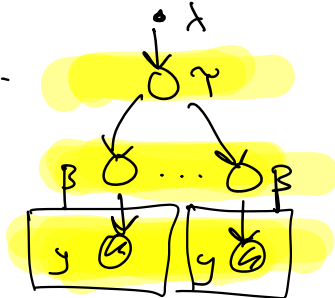
For each $j \in \{1 \dots n_i\}$

$$y_{ij} | x_{ij}, \beta_i \sim \text{GLM}(\beta_i)$$

$$\log p(-) = \log p(\gamma; \lambda) + \sum_{i=1} (\log p(\beta_i | \gamma) + \sum_{j=1} \log p(y_{ij} | \beta_i, x_{ij}))$$

$$p(\beta_i | \mathbf{x}, \mathbf{y}) \propto \mathbb{E}_{\gamma} [p(\beta_i | \gamma) | \mathbf{x}_{-i}, \mathbf{y}_{-i}] p(y_i | \mathbf{x}_i, \beta_i)$$

Empirical Bayes.
Hierarchical Bayes.
↳ Gelman & Hill.



Gaussian process : f.w. $GP(\mu(\cdot), K(\cdot, \cdot))$ — Bayesian nonparametrics.

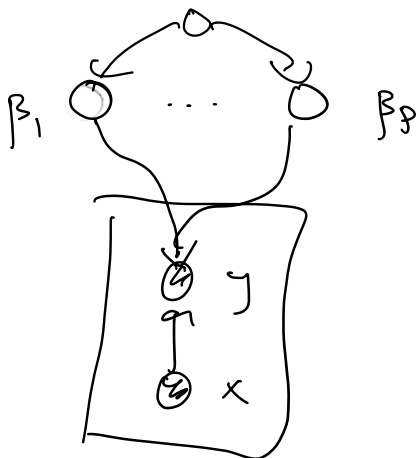
$$\mathbb{E}_{\mathcal{D}} [d(\beta_{1:50}^*, \hat{\beta}_{1:50})] = R(\beta^*)$$

Risk.

$$\mathbb{E}_{\gamma^*} [\mathbb{E}_{\mathcal{D}} [d(\beta_{1:50}^*, \hat{\beta}_{1:50})]] = R_{\text{Bayes}}$$

$$R_{\text{Bayes}}^{\text{MLE}} > R_{\text{Bayes}}^{\text{Shrinkage}}$$

Robust, Efron, James-Stein...



— ARD
— Horseshoe priors.