

09-Exponential Family and GLIMs

Exponential Family

Def (exponential family): A distribution belonging to the exponential family has the following form

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{T}(\mathbf{x}) - A(\boldsymbol{\eta}) \right\}$$

where $A(\boldsymbol{\eta}) = \log \int h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{T}(\mathbf{x}) \right\} d\mathbf{x}$

is a log-normalizing factor. In other words,

$$p(\mathbf{x}|\boldsymbol{\eta}) = \frac{1}{Z} h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{T}(\mathbf{x}) \right\}$$

where $Z = \int h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{T}(\mathbf{x}) \right\} d\mathbf{x}$ [assuming convergent]

Here, $\mathbf{T}(\mathbf{x})$ are the sufficient statistics.

$\boldsymbol{\eta}$ are the parameters, and are called the natural parameters.

Basically, exponential family is log-linear in the sufficient statistics of \mathbf{x} , with some residual effect $h(\mathbf{x})$.

A variety of common distributions are exponential family

- Gaussian

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

Thus $\mathbf{T}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \mathbf{x}^2 \end{bmatrix}$

$$\boldsymbol{\eta} = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix} \quad \text{Let } \boldsymbol{\eta} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix}. \quad \text{Then} \quad \begin{aligned} \mu &= -\frac{\eta_1}{2\eta_2} \\ \sigma^2 &= -\frac{1}{2\eta_2} \end{aligned}$$

$$A(\boldsymbol{\eta}) = \frac{\mu^2}{2\sigma^2} + \log \sigma = -\frac{\eta_1^2}{4\eta_2} + \log \sqrt{-\frac{1}{2\eta_2}}$$

$$h(\mathbf{x}) = \frac{1}{\sqrt{2\pi}}$$

- Bernoulli

$$p(\mathbf{x}; \pi) = \pi^x (1-\pi)^{1-x}$$

$$\begin{aligned}
 &= \exp \left\{ x \log \pi + (1-x) \log (1-\pi) \right\} \\
 &= \exp \left\{ [\log \pi - \log (1-\pi)] \cdot x + \log (1-\pi) \right\}
 \end{aligned}$$

Thus, $T(x) = x$

$$\eta = \log \frac{\pi}{1-\pi} \Rightarrow \pi = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}}$$

$$A(\eta) = -\log(1-\pi) = \log(1+e^\eta)$$

$$h(x) = 1$$

For the Gaussian distribution, $\mu = \mathbb{E}[x]$, $\sigma^2 = \mathbb{E}[x^2]$

For the Bernoulli distribution, $\pi = \mathbb{E}[x]$

They are known as moment parameters, and have some correspondence with natural parameters.

$$\psi: \mu \mapsto \eta \quad \psi^{-1}: \eta \mapsto \mu$$

A few useful properties of exponential family

- $A(\eta)$ is convex in natural parameters η .
- $\frac{\partial A(\eta)}{\partial \eta} = \mathbb{E}[T(x)] \quad \frac{\partial^2 A(\eta)}{\partial \eta^2} = \text{Var}[T(x)]$
- Maximum likelihood estimation \Leftrightarrow Moment matching

$$\begin{aligned}
 \text{log-likelihood } l(\eta) &= \log \prod_{m=1}^M [h(x^{(m)}) \exp \{ \eta^T T(x^{(m)}) - A(\eta) \}] \\
 &= \sum_{m=1}^M [\log h(x^{(m)})] + \eta^T \sum_{m=1}^M T(x^{(m)}) - M A(\eta)
 \end{aligned}$$

Closed-form solution

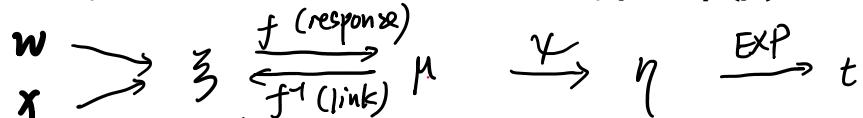
$$\begin{aligned}
 \frac{\partial l(\eta)}{\partial \eta} &= \sum_{m=1}^M T(x^{(m)}) - M \frac{\partial A(\eta)}{\partial \eta} \stackrel{\text{set}}{=} 0 \\
 \Rightarrow \frac{\partial A(\eta)}{\partial \eta} &= \frac{1}{M} \sum_{m=1}^M T(x^{(m)}) \\
 \Rightarrow \hat{\mu}_{MLE} &= \frac{1}{M} \sum_{m=1}^M T(x^{(m)})
 \end{aligned}$$

Generalized Linear Models

In **generalized linear models** (GLM, GLIM), we assume

- x enters the model by linear (or affine) transformation, yielding $\xi = w^T x$
- A response function f transforms ξ to the mean parameters μ
- Natural parameters are thus obtained by $\eta = \eta(\mu)$

• Natural parameters are thus obtained by $\mu = \psi(\eta)$



Usually, we set $f = \psi^{-1}$, in which case f is known as the natural response. The goal of prediction is to estimate the mean of t .

- For Gaussian distribution,

$$f(x) = \mathbb{E}[t|x] = \mu = \eta = w^T x$$

- For Bernoulli distribution,

$$f(x) = \mathbb{E}[t|x] = \mu = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-w^T x}}$$

MLE for GLIM with natural response

Log-likelihood

$$\begin{aligned} l(w, t) &= \sum_{m=1}^M \log h(t^{(m)}) + \sum_{m=1}^M ((w^T x^{(m)}) t^{(m)} - A(\eta)) \\ \frac{\partial l}{\partial w} &= \sum_{m=1}^M \left[x^{(m)} y^{(m)} - \frac{\partial A(\eta^{(m)})}{\partial \eta^{(m)}} \cdot \frac{\partial \eta^{(m)}}{\partial w} \right] \\ &= \sum_{m=1}^M (y^{(m)} - \mu^{(m)}) x^{(m)} \\ &= X^T(t - y) \end{aligned}$$

Summary

- We showed that exponential family is a useful family of distributions with nice properties.
- With further assumptions of generalized linear models and the natural response function, we may derive the formulas for linear regression, logistic regression, and softmax regression (HW). And they share the same formula of the derivative in MLE.
- It should be emphasized that GLIM and natural response that lead to the sigmoid or softmax is still an assumption. Other assumptions may also be acceptable and used.