

Generalized Linear Models

* The exponential family

⇒ We say that a class of distributions is in the **exponential family** if it can be written in the form:

$$P(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

→ η is called the **natural parameter**

→ $T(y)$ is the **sufficient statistic** (mostly $T(y) = y$)

→ $a(\eta)$ is the **log partition function**

→ $b(y)$ is called **base measure**

⇒ A fixed choice of T , a and b defines a family of distributions that is parameterized by η .

↳

As we vary η , we then get different distributions within this family.

⇒ Bernoulli and Gaussian distributions are examples of exponential family distribution

① Bernoulli distribution

$$P(y; \phi) = \phi^y (1-\phi)^{1-y}$$

* Properties of Exponential family

① MLE is convex

$$\textcircled{2} \quad E[y; n] = \frac{\partial}{\partial \eta} a(\eta)$$

$$\textcircled{3} \quad \text{Var}[y; n] = \frac{\partial^2}{\partial \eta^2} a(\eta)$$

Date _____

Page _____



$$P(y; \phi) = \exp \left(\log \left(\frac{\phi}{1-\phi} \right) y + \log(1-\phi) \right)$$

- $b(y) = 1$

- $T(y) = y$

- $\eta = \log \left(\frac{\phi}{1-\phi} \right)$

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

② Gaussian distribution (fixed variance $\sigma^2 = 1$)

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} (y - \mu)^2 \right)$$

$$= \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} y^2 \right) \exp \left(\mu y - \frac{1}{2} \mu^2 \right)$$

- $b(y) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{y^2}{2} \right)$

- $T(y) = y$

- $\eta = \mu$

- $a(\eta) = \frac{\eta^2}{2}$

\Rightarrow There are many other distributions that are members of the exponential family:

- multinomial
- beta

- Poisson

- Dirichlet

- gamma

- ... etc ...

- exponential



Constructing GLMs

- Consider a classification or regression problem where we would like to predict the value of some random variable y as a function of x .
- To derive a GLM for this problem, we will make the following three assumptions about the conditional distribution of y given x and about our model:

① $y|x; \theta \sim \text{Exponential Family } (\eta)$

② Given x , our goal is to predict the expected value of $T(y)$ given x :

→ In most cases, $T(y) = y$

→ So this means we would like the prediction $h(x)$ output by our learned hypothesis h to satisfy,

$$h(x) = E[y|x; \theta]$$

③ The natural parameters η and the output input x are related linearly:

$$\eta = \theta^T x$$

⇒ These three assumptions/design choices will allow us to derive a very elegant class of learning algorithms, namely GLMs.

* Ordinary Least Squares

$$y|x; \theta \sim N(\mu, \sigma^2)$$

$$\eta = \mu = \theta^T x$$

$$h_{\theta}(x) = E[y|x; \theta]$$

$$= \mu$$

$$\boxed{h_{\theta}(x) = \theta^T x}$$

* Logistic Regression

$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

$$\Rightarrow \eta = \ln \left(\frac{\phi}{\phi-1} \right)$$

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

$$E[y|x; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$$

$$\boxed{h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}}$$

→ To introduce or little more terminology, ~~the~~
~~fitting~~

→ The function g given the distribution
 mean as a function of model parameter
 $g(\mu) = E[T(a)/N]$

is the **canonical response function**.

→ Its inverse g^{-1} is called the
canonical link function.

* Softmax Regression

→ Consider a classification problem in
 which the response variable y can take
 on only one of K values, so $y \in \{1, 2, \dots, K\}$

→ We will thus model it as distributed
 according to a multinomial distribution.

→ Let's derive a GLM for modelling this type of
 multinomial data.

→ To parameterize a multinomial over K
 possible outcomes, one could use K parameters
 ϕ_1, \dots, ϕ_K specifying the probability of
 each of the outcomes.

⇒ However, these parameters would be redundant, since knowing any $K-1$ of the ϕ_i s uniquely determines the last one.

$$\sum_{i=1}^K \phi_i = 1$$

⇒ So, we will instead parameterize the multinomial with only $K-1$ parameters

$$\phi_1, \dots, \phi_{K-1}$$

$$\phi_i = P(y=i; \phi)$$

$$P(y=k; \phi) = 1 - \sum_{i=1}^{K-1} \phi_i$$

⇒ To express the multinomial as an exponential family distribution, we will define $T(y) \in \mathbb{R}^{K-1}$ as follows

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \dots \quad T(K) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

⇒ Let $I\{\cdot\}$ be an indicator function that takes on value 1 if argument is true, and 0 otherwise.

$$\Rightarrow \text{So } (T(y))_i = I\{y=i\}$$

\Rightarrow Further, we have $E[(T(y))_i] = P(y=i) = \phi_i$

$$P(y; \phi) = \phi_1^{1\{y=1\}} \cdot \phi_2^{1\{y=2\}} \cdots \phi_k^{1\{y=k\}}$$

\downarrow \downarrow \downarrow

$$T(y)_1 \quad T(y)_2 \quad 1 - \sum_{i=1}^{k-1} 1\{y=i\}$$

\downarrow

$$T(y)_k$$

$$\Rightarrow \exp \left(T(y)_1 \log(\phi_1) + T(y)_2 \log(\phi_2) + \right.$$

$$\left. \left(1 - \sum_{i=1}^{k-1} T(y)_i \right) \log(\phi_k) \right)$$

$$\Rightarrow \exp \left(T(y)_1 \log \left(\frac{\phi_1}{\phi_k} \right) + \right.$$

$$\left. T(y)_2 \log \left(\frac{\phi_2}{\phi_k} \right) + \right.$$

$$\left. \left(1 - \sum_{i=1}^{k-1} T(y)_i \right) \log \left(\frac{\phi_{k-1}}{\phi_k} \right) + \right.$$

$$\left. \log(\phi_k) \right)$$

$$\Rightarrow P(y; \phi) = b(y) \exp(n^T T(y) - a(n))$$

where,

$$n = \begin{bmatrix} \log(\phi_1/\phi_K) \\ \log(\phi_2/\phi_K) \\ \vdots \\ \log(\phi_{K-1}/\phi_K) \end{bmatrix}$$

$$a(n) = -\log(\phi_K)$$

$$b(y) = 1$$

\Rightarrow The link function is given (for $i=1, \dots, K$) by

$$n_i = \log \frac{\phi_i}{\phi_K}$$

$$e^{n_i} = \frac{\phi_i}{\phi_K}$$

$$\phi_K e^{n_i} = \phi_i$$

$$\phi_K \sum_{i=1}^K e^{n_i} = \sum_{i=1}^K \phi_i = 1$$

$$\phi_K = \frac{1}{\sum_{i=1}^K e^{n_i}}$$

$$\Rightarrow \boxed{\phi_i = \frac{e^{n_i}}{\sum_{i=1}^K e^{n_i}}} \quad \{ \text{Response function} \}$$

⇒ Lastly, let's discuss parameter fitting

$$l = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \prod_{l=1}^K \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^K e^{\theta_j^T x^{(i)}}} \right)^{1\{y^{(i)}=l\}}$$

⇒ We can now obtain the maximum likelihood estimate of the parameters by maximizing $l(\theta)$ in terms of θ , using a method such as gradient ascent or Newton's method.