

Recall: $P(y=1|x;\theta) = h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

$$\ell(\theta) = \sum y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) \cdot x_j^{(i)}$$

Newton method:

$$\frac{f(\theta^{(0)}) - 0}{\theta^{(0)} - 0} = f'(\theta_0)$$

$$\Rightarrow \theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

to maximize $\ell(\theta)$, want $f'(\theta) = 0$

When θ is a vector:

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_{\theta} \ell(\theta)$$

H is hessian matrix

$$H_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j}$$

Short cutting: need to inverse hessian.

Generalized Linear Model:

$$P(y|x;\theta)$$

$y \in \mathbb{R}$: Gaussian \rightarrow least square.

$y \in \{0,1\}$: Bernoulli \rightarrow Logistic function

$$\text{Bernoulli}(\phi) : P(y=1;\phi) = \phi$$

exponential function

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

η : natural parameter

T : sufficient statistic.

(a,b,T) can determine a special instance of exp family.

$$\text{Ber}(\phi) \quad P(y=1;\phi) = \phi$$

$$\begin{aligned} P(y;\phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp(y \log \phi + (1-y) \log (1-\phi)) \\ &= \exp\left(\underbrace{\log \frac{\phi}{1-\phi}}_{\eta} y + \underbrace{\log (1-\phi)}_{-a(\eta)}\right) \end{aligned}$$

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

$$\begin{cases} a(\eta) = -\log(1-\phi) = \log(1+e^{\eta}) \\ T(y) = y \\ b(y) = 1 \end{cases}$$

Gaussian: $N(\mu, \sigma^2)$

Set $\sigma^2=1$ since it will take no effect on θ .

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

$$\Rightarrow: \underbrace{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right)}_{b(y)} \exp(y\mu - \frac{1}{2}\mu^2)$$

$$\mu = \eta, \quad T(y) = y, \quad a(\eta) = \frac{1}{2}\mu^2$$

Gamma, Poisson, exponential and β , Dirichlet, Wishart distribution is exponential distribution family.

GLM:

Assume:

$$1): y|x;\theta \sim \text{ExpFamily}(\eta)$$

2): Given x : goal is to output $E[T(y)|x]$

$$\text{want } h(x) = E[T(y)|x]$$

$$3): \eta = \theta^T x \quad // \quad \eta_i = \theta_i^T x$$

Bernouli:

$$y|x;\theta \sim \text{ExpFamily}(\mu)$$

For fixed x, θ , algorithm output

$$h_\theta(x) = E[y|x;\theta] = P(y=1|x;\theta)$$

$$= \phi$$

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

$$= \frac{1}{1+e^{-\theta^T x}}$$

$$g(\eta) = E[y;\eta] = \frac{1}{1+e^{-\eta}}$$

canonical response function

g^{-1} canonical link function.

example:

Multinomial: $y \in \{1, \dots, k\}$

Parameter: $\phi_1, \phi_2, \dots, \phi_k$

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

$$\phi_k = 1 - (\phi_1 + \dots + \phi_{k-1})$$

parameter $(\phi_1, \phi_2, \dots, \phi_{k-1})$

$$T(y) \neq y!$$

$$T(y) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{k-1} \quad T(y) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}_{k-1} \in \mathbb{R}^{k-1}$$

$$T(y) = [0, 0, \dots, 0]^T$$

$$1\{\text{true}\} = 1 \quad / \quad 1\{\text{false}\} = 0$$

$$T(y)_i = 1\{y=i\}$$

the ith element of $T(y)$

$$\begin{aligned} P(y) &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}} \\ &= \phi_1^{T(y)_1} \phi_2^{T(y)_2} \dots \phi_k^{1-\sum_{j=1}^{k-1} T(y)_j} \end{aligned}$$

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

$$\text{where } \eta = \begin{bmatrix} \log(\phi_1/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix}_{k-1} \quad a(\eta) = \log(\phi_k)$$

$$b(y) = 1$$

$$\phi_i = \frac{e^{\eta_i}}{1 + \sum_{i=1}^{k-1} e^{\eta_i}} \quad (i=1, \dots, k-1)$$

$$= \frac{e^{\theta_i^T x}}{1 + \sum_{j=1}^{k-1} e^{\theta_j^T x}}$$

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

$$= E \left[\begin{bmatrix} 1\{y=1\} \\ \vdots \\ 1\{y=k-1\} \end{bmatrix} \middle| x; \theta \right] = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{e^{\theta_1^T x}}{(1 + \sum_{j=1}^{k-1} e^{\theta_j^T x})} \\ \vdots \\ \frac{e^{\theta_{k-1}^T x}}{(1 + \sum_{j=1}^{k-1} e^{\theta_j^T x})} \end{bmatrix}$$

Soft max regression: ($2 \rightarrow k$)

train set: $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)}) \quad y \in \{1, \dots, k\}$

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \phi_1^{1\{y^{(i)}=1\}} \dots \phi_k^{1\{y^{(i)}=k\}}$$

$$\downarrow$$

$$\frac{e^{\theta_i^T x}}{1 + \sum_{j=1}^{k-1} e^{\theta_j^T x}}$$