## Logistic Regression:

Training Set $= \{(x^{(i)}, y^{(i)})\}_{i=1}^{n}$

$\in \mathbb{R}^{d+1}$ ↰     ⮑ $\in \{0,1\}$

**Goal:** Find $\theta \in \mathbb{R}^{d+1}$ such that

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} l_{nll}(\theta, x_{i}^{(i)}, y^{(i)}) \text{ is minimized}$$

$$l_{nll}(\theta, x, y) = -[y \log h_{\theta}(x) + (1-y) \log(1 - h_{\theta}(x))]$$

Maximum Likelihood Estimation

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

Intuitively $P[y=1 | x; \theta] = h_{\theta}(x)$.

Gradient descent updates

$$\theta_{j} = \theta_{j} - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_{\theta}(x^{(i)}) - y^{(i)}) x_{j}^{(i)}]$$

The "same" as linear regression

---

## Exponential Family of Distributions

Distribution ↰   ⟶ parameter
over                             $(\eta^{T} T(y) - a(\eta))$

$$p(y; \eta) = b(y) e^{(\eta^{T} T(y) - a(\eta))}$$

↓ p.d.f /
p.m.f

$\eta$ — natural parameter
$T$ — sufficient statistic
$a$ — log partition.

# Generalized Linear Models

(i) $P(y \mid x; \theta) \sim$ Exponential Family $(\eta)$

(ii) Goal of the learning algorithm is to
$$E[T(y) \mid x; \theta]$$
$$= y$$

(iii) $\eta = \theta^T x$ when $\eta \in \mathbb{R}$

$\eta_i = \theta^{(i)T} x$ when $\eta \in \mathbb{R}^k$.

We get $\theta$ by maximized likelihood
maximizing log likelihood
minimizing negative log likelihood

## Logistic Regression

$\rightarrow$ depends on $x$ & $\theta$.

$$P[y=1 \mid x; \theta] \sim \text{Bernouli}(\phi)$$

~~P.m.f~~ P.m.f of Bernouli $(\phi)$

$$p(y; \phi) = \phi^y (1-\phi)^{(1-y)} \qquad ] \rightarrow \exp(\partial) = e^\partial$$

$$= \exp[y \ln \phi + (1-y) \ln(1-\phi)]$$

$$= \exp[y \ln \frac{\phi}{1-\phi} + \ln(1-\phi)]$$

$b(y) = 1$, $T(y) = y$, $\eta = \ln \frac{\phi}{1-\phi}$ $\qquad a(\eta) = \ln(1-\phi)$

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

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

$$\phi = \text{sigmoid}(\eta)$$

$$a(\eta) = \ln(1-\phi) = \ln\left[1 - \frac{1}{1+e^\eta}\right]$$

Learning algorithm of Logistic Regression found

$$\rightarrow h_\theta(x) = p[y=1 \mid x; \theta]$$

$$\stackrel{?}{=} E[y \mid x; \theta]$$

$$= \phi$$

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

$z \sim$ Bernouli $(\phi)$

$$E(z) = \sum_{z=i} i \cdot p(z=i)$$

$$= 1 \cdot p(z=1) + 0 \cdot p[z=0]$$

$$= p(z=1) = \phi$$

## Softmax Regression

~~Output~~ Output $\in \{1, 2, \ldots k\}$.

$$p[y=i \mid x; \Theta] \sim \text{Multinomial Dist} (\phi_1, \phi_2 \ldots \phi_k)$$

$\hookrightarrow$ prob that output is 2.

$$= \phi_i$$

$$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$$

$$p[y=i \mid x; \Theta] \sim \text{Multinomial Dist} (\phi_1 \ldots \phi_{k-1})$$

One-shot encoding.

$$T(y) = \begin{bmatrix} 1\{y=1\} \\ 1\{y=2\} \\ \vdots \\ 1\{y=i\} \\ \vdots \\ 1\{y=k-1\} \end{bmatrix} \in \{0,1\}^{k-1}$$

Example: $k = 4$.

$$T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$T(4) = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$p[y \mid x; \Theta] = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \cdots \phi_{k-1}^{1\{y=k-1\}} \phi_k^{1\{y=k\}} \qquad \left[ \phi_k = 1 - \sum_{i=1}^{k-1} \phi_i \right] \qquad \textcircled{4}$$

$$= \phi_1^{T(y)_1} \phi_2^{T(y)_2} \cdots \phi_{k-1}^{T(y)_{k-1}} \phi_k^{1 - \sum_{i=1}^{k-1} T(y)_i}$$

$$= \exp\left[ T(y)_1 \ln \phi_1 + T(y)_2 \ln \phi_2 + \cdots T(y)_{k-1} \ln \phi_{k-1} + \left(1 - \sum T(y)_i\right) \ln \phi_k \right]$$

$$= \exp\left\{ T(y_1) \ln \frac{\phi_1}{\phi_k} + T(y_2) \ln \frac{\phi_2}{\phi_k} + \cdots T(y)_{k-1} \ln \frac{\phi_{k+1}}{\phi_k} + \ln \phi_k \right\}$$

$$\underbrace{\qquad\qquad} = \eta^T T(y)$$

$$b(y) = 1, \qquad T(y)$$

$$\eta = \begin{bmatrix} \ln \frac{\phi_1}{\phi_k} \\ \vdots \\ \ln \frac{\phi_{k-1}}{\phi_k} \end{bmatrix} \qquad a(\eta) = -\ln(\phi_k)$$

$$\eta_i = \ln \frac{\phi_i}{\phi_k} \quad \Rightarrow \quad \phi_i = \phi_k e^{\eta_i} \qquad \left[ \eta_k = 0 \right.$$

$$1 = \sum \phi_i = \phi_k \sum_{i=1}^{k} e^{\eta_i}, \qquad \eta_k = 0.$$

$$\phi_k = \frac{1}{1 + \sum_{i=1}^{k-1} e^{\eta_i}}$$

$$\phi_i = \frac{e^{\eta_i}}{1 + \sum_{i=1}^{k-1} e^{\eta_i}}$$

$$T(y) = \begin{bmatrix} 1(y=1) \\ 1(y=2) \\ \vdots \\ 1(y=k-1) \end{bmatrix}$$

$$h_\theta(x) = E\left[ T(y) \mid x; \theta \right]$$

$$= \begin{bmatrix} E[1(y=1)] \\ \vdots \\ E[1(y=k-1)] \end{bmatrix} = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$\eta_i = \theta^{(i)T} x.$$

$$= \begin{bmatrix} \dfrac{e^{\eta_1}}{1 + \sum e^{\eta_i}} \\ \vdots \\ \dfrac{e^{\eta_{k-1}}}{1 + \sum e^{\eta_i}} \end{bmatrix} = \begin{bmatrix} \dfrac{e^{\theta^{(1)T} x}}{1 + \sum e^{\theta^{(i)T} x}} \\ \vdots \\ \dfrac{e^{\theta^{(k-1)T} x}}{1 + \sum e^{\theta^{(i)T} x}} \end{bmatrix}$$