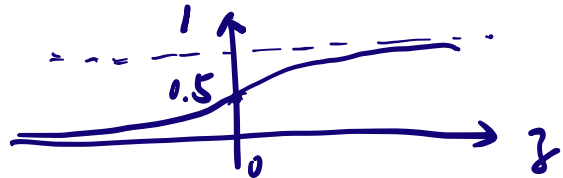Sep 26, 2022 Mon

$p(y|x)$ ~ Bernoulli $(\theta)$

$\theta = $ sigmoid $(w^T x)$

sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$p(y|x)$ ~ Ber $(\theta)$

$$= Ber(\sigma(z))$$

$$= Ber\left(\frac{1}{1 + \exp(-w^T x)}\right)$$

Training of Logistic Regression

Goal : estimate $w$, Given $\{x_i, y_i\}$

Testing of LR

Given a new $x^*$.

$$\theta = \sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \exp(-\hat{w}^T x^*)}$$

if $\theta > 0.5$. then $y^* = 1$

else $\quad\quad\quad\quad\quad\quad y^* = 0$

$$P(y^* = 1) = \theta_i = \frac{1}{1 + \exp(-w^T x_i^*)}$$

$$P(y^* = 0) = 1 - \theta_i$$

# Learning Parameters

## MLE approach.

Likelihood $\quad P(D|\theta) = \prod_{i=1}^{N} \theta_i^{y_i} (1-\theta_i)^{1-y_i}$

$NLL = -\log P(D|\theta)$

$\qquad = \sum_{i=1}^{N} -y_i \log \theta_i - (1-y_i) \log (1-\theta_i)$

Cross-entropy loss $\quad H(p,q) = -\sum_k p \log(q)$

log-loss

$MLE = \min NLL$

## Gradient Descent

$$W_{k+1} = W_k - \eta \frac{dNLL}{dw}$$

$\theta_i = \sigma(z) = \frac{1}{1+\exp(-w^Tx)}$

$\frac{d \log \theta_i}{w} = \frac{1}{\theta_i} \cdot \frac{d\theta_i}{dw} = \frac{1}{\theta_i} \frac{d\sigma}{dw}$

$\dfrac{dNLL}{dw}$

$\frac{d\sigma}{dw} = -1 (1+\exp(-w^Tx))^{-2} \exp(-w^Tx) \cdot (-X)$

$\qquad = \dfrac{\exp(-w^Tx)}{(1+\exp(-w^Tx))^2} \cdot X$

$= \sum_{i=1}^{N} -y_i \frac{1}{\theta_i} \theta_i (1-\theta_i) \cdot x_i$

$+ (1-y_i) \frac{1}{1-\theta_i} (1-\theta_i) \theta_i x_i$
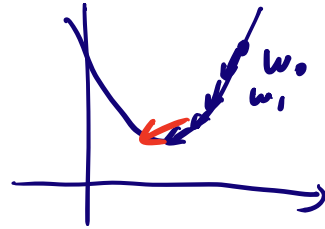
$\qquad = \sigma(1-\sigma) \cdot X$

$= \sum_{i=1}^{N} \left( y_i \theta_i x_i - y_i x_i + \theta_i x_i - y_i \theta_i x_i \right)$

$= \sum_{i=1}^{N} (\theta_i - y_i) x_i$

$$w_{k+1} = w_k - \eta \frac{dNLL}{dw} \qquad \eta \ \text{learning rate}$$

$$\frac{dNLL}{dw} = \sum_{i=1}^{N} (\theta_i - y_i) x_i$$



## Newton's method

$$w_{k+1} = w_k - H_k^{-1} \frac{dNLL}{dw_k}$$

H : Hessian matrix

Second order derivative of a function.

$$f(x) \qquad f'(x) = \frac{df(x)}{dx} \qquad f''(x) = \frac{d^2 f(x)}{dx^2}$$

$$f(x) = x^3 + 2x^2$$
$$f'(x) = 3x^2 + 4x$$
$$f''(x) = 6x + 4$$

$$f(w) \qquad\qquad w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}_{dx1}$$

$$\frac{df(w)}{dw} = \begin{bmatrix} \dfrac{df(w)}{dw_1} \\ \vdots \\ \dfrac{df(w)}{dw_d} \end{bmatrix}_{dx1}$$

$$H = \frac{d^2 f(w)}{dw^2} = \begin{bmatrix} \dfrac{d^2 f(w)}{dw_1^2} & \cdots & \dfrac{d^2 f(w)}{dw_1 dw_d} \\ \vdots & & \\ \dfrac{d^2 f(w)}{dw_d dw_1} & \cdots & \dfrac{d^2 f(w)}{dw_d^2} \end{bmatrix}_{dxd}$$

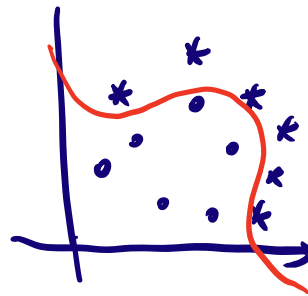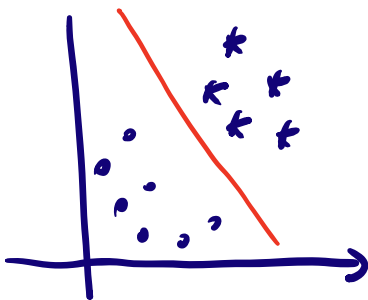$$f(w) = w_1 w_2 + 2w_1^2 + 3w_1^2 w_2 + 4$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_{2 \times 1}$$

$$\frac{df(w)}{dw} = \begin{bmatrix} w_2 + 4w_1 + 6w_1 w_2 \\ w_1 + 3w_1^2 \end{bmatrix}_{2 \times 1}$$

$$H = \begin{bmatrix} 4 + 6w_2 & 1 + 6w_1 \\ 1 + 6w_1 & 0 \end{bmatrix}_{2 \times 2}$$

For logistic Regression.

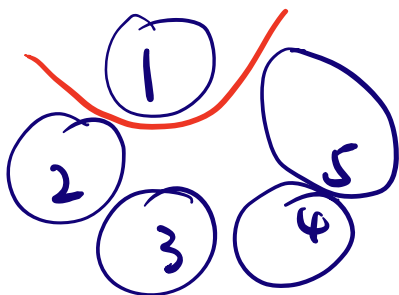$$H = -\sum_{i=1}^{N} \theta_i (1 - \theta_i) x_i x_i^T$$



$$\phi = [1, x, x^2 \cdots x^d]$$

$$\theta_i = \nabla(z) = \frac{1}{1 + \exp(-w^T \phi(x))}$$

$$L(w) = NLL(w) + \frac{1}{2} \lambda w^T w$$

$$\frac{dL(w)}{dw} = \frac{dNLL}{dw} + \lambda w$$

$$\hat{H} = H + \lambda I$$

**Multiple Classes**, multi class / multi label

$P(y | x) \sim \text{Multinoulli}(\theta)$    $\theta = [\theta_1 \cdots \theta_C]$

$$\theta_j = \frac{\exp(w^T_j x)}{\sum\limits_{k=1}^{C} \exp(w^T_k x)}$$    Softmax function

$C = 2$

$$\theta_1 = \frac{\exp(w_1^T x)}{\exp(w_1^T x) + \exp(w_2^T x)}$$

$$= \frac{1}{1 + \exp((w_2 - w_1)^T x)}$$

$$= \frac{1}{1 + \exp(-\hat{w}^T x)}$$    $\hat{w} = -(w_2 - w_1)$

Sigmoid function