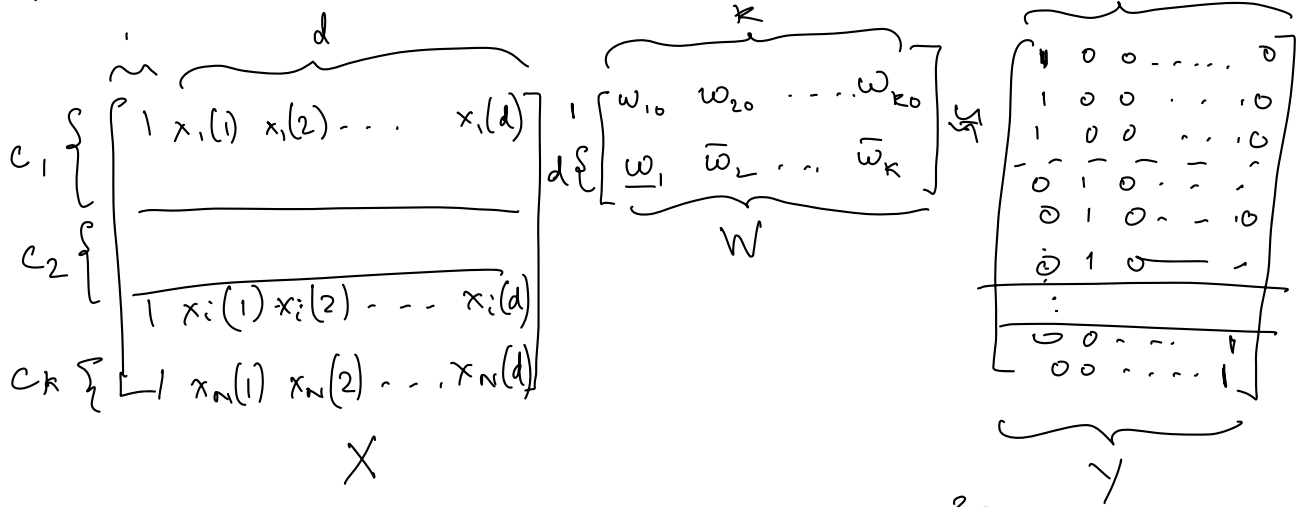# Classification : Regression Approaches

$(x_i, y_i)$ , $x_i \in \mathbb{R}^d$ , $y_i \in \mathbb{R}^k$ (k-class classification)

$$x = \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix} \Big\} d \quad , \quad y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \Big\} k \quad , \quad i = 1, 2, \ldots N$$

$\longrightarrow$ j-th position (encodes $x_i \in C_j$)

$\longrightarrow$ j-th column of the identity matrix

$$y(x) = w_0 + w_1 x(1) + w_2 x(2) + \ldots w_d x(d) \qquad - \text{ Linear fit}$$

Linear Discriminants : $\bar{w}_j^T x + w_{j_0}$ for j-th class



$$X W \backsim Y \qquad \Rightarrow \min_W \frac{1}{2} \| XW - Y \|_F^2$$

$$(X^T X) W^* = X^T Y \qquad\qquad W \in \mathbb{R}^{(d+1) \times k}$$

$$\boxed{W^* = (X^T X)^{-1} X^T Y}$$

$X W^*$ is prediction on training data

$$x \in \mathbb{R}^{d+1}$$
$$\underset{\shortparallel}{} \begin{bmatrix} 1 \\ x(1) \\ \vdots \\ x(d) \end{bmatrix}$$

$$x^T W^* = \begin{bmatrix} 0.01 & 0.02 \ldots & 0.98 \ldots \end{bmatrix} \qquad \text{if } x_i \in C_j$$

$\uparrow$ j-th position

k-dimensional

However, the above $W^*$ can lead to

$$x^T W^* \backsim \begin{bmatrix} -5.1 & -3.1 & 0 & 3.2 & 3.1 \ldots & 0.9 \end{bmatrix}$$

and clearly, this is not a good approximation to

an indicator vector $\rightarrow [0 \ldots 0 \ 1 \ldots 0]$ (1-hot vector)

Just Least Squares for Classification has obvious drawbacks

## Logistic Regression

We had modeled each class as a Gaussian with covariance $\underline{\underline{\Sigma}}$:

$$\log \frac{p(C_i|x)}{p(C_j|x)} = \log \frac{p(C_i)}{p(C_j)} - \frac{1}{2}(m_i + m_j)^T \Sigma^{-1}(m_i - m_j) \underbrace{+ x^T \Sigma^{-1}(m_i - m_j)}$$

$$\boxed{w_0 + w^T x}$$

### K-class problem

$$\log\left(\frac{p(C_1|x)}{p(C_k|x)}\right) = w_0 + \tilde{w}_1^T \tilde{x} = w_1^T x \quad -①$$

$$\log\left(\frac{p(C_2|x)}{p(C_k|x)}\right) = w_2^T x \quad -② \qquad\qquad , w_i \in R^{d+1}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad x \in R^{d+1}$$

$$\vdots \qquad\qquad \vdots \qquad\qquad \vdots$$

$$\log \frac{p(C_{k-1}|x)}{p(C_k|x)} = w_{k-1}^T x \quad -\boxed{k-1}$$

Write $p(C_i|x)$ as $p_i$

$$\log \frac{p_1}{p_k} = w_1^T x \quad\Rightarrow\quad \frac{p_1}{p_k} = e^{w_1^T x} \quad -①$$

$$\frac{p_2}{p_k} = e^{w_2^T x} \quad -②$$

$$\vdots$$

$$\frac{p_{k-1}}{p_k} = e^{w_{k-1}^T x} \quad -\boxed{k-1}$$

$① + ② + \ldots \boxed{k-1}$

$$\frac{p_1}{p_k} + \frac{p_2}{p_k} + \ldots + \frac{p_{k-1}}{p_k} = e^{w_1^T x} + e^{w_2^T x} + \ldots e^{w_{k-1}^T x}$$

$$\Rightarrow \quad \frac{p_1 + p_2 + \cdots + p_{k-1}}{p_k} = \sum_{i=1}^{k-1} e^{w_i^T x}$$

$$\frac{1 - p_k}{p_k} = \sum_{i=1}^{k-1} e^{w_i^T x}$$

$$1 - p_k = p_k \sum_{j=1}^{k-1} e^{w_j^T x}$$

$$1 = p_k \left( 1 + \sum_{j=1}^{k-1} e^{w_j^T x} \right)$$

$$\Rightarrow \quad \boxed{ p_k = \frac{1}{1 + \sum_{j=1}^{k-1} e^{w_j^T x}} }$$

$$\frac{p_i}{p_k} = e^{w_i^T x}$$

$$\Rightarrow \quad \boxed{ p_i = \frac{e^{w_i^T x}}{1 + \sum_{j=1}^{k-1} e^{w_j^T x}} } \quad , \quad i = 1, 2, \cdots k-1$$

$w_1, w_2, w_3 \cdots, w_{k-1}$ are all parameters.

How do we find them?

Parameters in logistic regression $\{ w_i \in \mathbb{R}^{d+1} \}_{i=1}^{k}$ are usually fit by maximum likelihood.

Consider the 2-class problem

$$p(C_1 | x) = p \qquad\qquad p(C_1 | x) + p(C_2 | x) = 1$$
$$p(C_2 | x) = 1 - p \qquad\qquad p = f(w)$$

$(x_i, y_i)$ is training data, $i = 1, 2, \cdots N$

Let $y_i = 1$ when $x_i \in C_1$
$\quad\quad y_i = 0$ when $x_i \in C_2$

$$\boxed{ \text{Data Likelihood} = \prod_{i=1}^{N} p^{y_i} (1-p)^{1-y_i} }$$

$(x_i, y_i)$
If $x_i \in C_1, y_i = 1$
$p^{y_i} (1-p)^{1-y_i}$

$= p'(1-p)^{1-1} = p \cdot (1-p)^0$
$\qquad\qquad = p$

Log-likelihood $\ell(\omega)$

$$\ell(\omega) = \sum_{i=1}^{N} \log\left(p^{y_i}(1-p)^{1-y_i}\right)$$

$$= \sum_{i=1}^{N} \log p^{y_i} + \log(1-p)^{1-y_i}$$

$$= \sum_{i=1}^{N} \left[ y_i \log p + (1-y_i)\log(1-p) \right]$$

Maximum log-likelihood

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad p = p(\omega)$

$$\max_{\omega} \ell(\omega) = \boxed{\max_{\omega}} \; \boxed{\left[ \sum_{i=1}^{N} \{ y_i \log p + (1-y_i)\log(1-p) \} \right]}$$

$$p = p(\omega) = \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} \qquad , \quad 1-p = 1 - \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} = \frac{1}{1+e^{\omega^T x_i}}$$

$$\log p = \omega^T x_i - \log(1+e^{\omega^T x_i}) \qquad \log(1-p) = -\log(1+e^{\omega^T x_i})$$

$$\nabla_\omega \ell(\omega) = 0$$
$$\underbrace{\qquad\qquad}_{\ell}$$

will involve $\nabla_\omega \log p$ & $\nabla_\omega \log(1-p)$

$$\boxed{\nabla_\omega \log p} = x_i - \frac{1}{1+e^{\omega^T x_i}} e^{\omega^T x_i} \cdot x_i = \boxed{x_i \left( 1 - \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} \right)}$$

$$\boxed{\nabla_\omega \log(1-p)} = \boxed{-\frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} x_i}$$

$$\nabla_\omega \ell(\omega) = \sum_{i=1}^{N} \left[ y_i \nabla_\omega \log p + (1-y_i) \nabla_\omega \log(1-p) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i x_i \left( 1 - \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} \right) + (1-y_i)x_i\left(-\frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}}\right) \right]$$

$$= \sum_{i=1}^{N} x_i \left[ y_i - y_i \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} - \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} + y_i \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} \right]$$

$$\Rightarrow \quad \nabla_\omega \ell(\omega) = \sum_{i=1}^{N} x_i \left[ y_i - \frac{e^{\omega^T x_i}}{1+e^{\omega^T x_i}} \right] \longrightarrow p(C_1 | x)$$

To find $w$ that maximizes $l(w)$

$$\nabla_w l(w) = 0$$

$$\nabla_w l(w) = \sum_{i=1}^{N} x_i \left[ y_i - \frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \right] = 0 \quad , x_i \in \mathbb{R}^{d+1}$$
$$w \in \mathbb{R}^{d+1}$$

$(d+1)$ parameters

$(d+1)$ equations since $x_i \in \mathbb{R}^{d+1}$

nonlinear equations

There is no closed form solution for above

How do we solve for $w$?

Need to use optimization methods to solve for these parameters

$\hookrightarrow$ Gradient Descent / Ascent

Newton's Method

$$w_{j+1} \leftarrow w_j - \eta \nabla_w l(w) \quad - \text{Gradient Descent}$$

$$w_{j+1} \leftarrow w_j - \eta \left( \nabla^2 l(w) \right)^{-1} \nabla l(w) \quad - \text{Newton's Method}$$

Drawback of these methods (Gradient Descent, Newton) is that each step requires $O(N)$ computation

Not feasible when $N$ is very large

Stochastic Gradient Descent (SGD)

Regularization: $\lambda \|w\|_2^2$ or $\lambda \|w\|_1$ should be used