# Lec 08: Logistic Regression

$$P\left(y_i \mid x_i, W\right)$$

$$\begin{bmatrix} P(y_i = 1 \mid x_i, W) \\ P(y_i = 0 \mid x_i, W) \end{bmatrix}$$

$$P\left(y_i = 1 \mid x_i, W\right) = \frac{e^{W_1^T x_i}}{e^{W_1^T x_i} + e^{W_2^T x_i}}$$

$$f\left(W, x_i\right) = \sigma\left(W^T x_i\right) = \frac{1}{1 + e^{-W^T x_i}}, \text{ where } W = W_1 - W_2$$

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

$$\begin{bmatrix} e \; e^{W_1^T x_i} \\ e \; e^{W_2^T x_i} \end{bmatrix}$$

$$e = \frac{1}{e^{W_1^T x_i} + e^{W_2^T x_i}}$$

## Binary LR:
$$W_{LR}^* \in \mathbb{R}^d$$

$$W_{LR}^* = \arg\max \prod_{i=1}^{n} P(y_i | x_i, w)$$

$$= \arg\max \sum_{i=1}^{n} \log P(y_i | x_i, w)$$

$$= \arg\min \left\{ -\sum_{i=1}^{n} \log P(y_i | x_i, w) \right\}$$

$$NLL_i(w) = -\log P(y_i | x_i, w) \quad : \text{negative log likelihood}$$

$$\hookrightarrow \log P(y_i = 1 | x_i, w), y_i = 1 \quad \text{or} \quad \log(1 - P(y_i = 1 | x_i, w)^{y_i}_{y_i = 0}$$

$$-\log P(y_i | x_i, w)$$

$$= -y_i \log \left[ \sigma(w^T x_i) \right]$$

$$- (1 - y_i) \log \left[ 1 - \sigma(w^T x_i) \right]$$

$$\boxed{\text{Cross-entropy loss}}$$

$$H(p, q) = -\sum_{x \in X} p(x) \log q(x)$$

$$NLL_i(w) = y_i \log\left(1 + e^{-w^T x_i}\right) - (1 - y_i) \log\left(\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}}\right)$$

$$= y_i \log\left(1 + e^{-w^T x_i}\right) + (1 - y_i) w^T x_i + (1 - y_i) \log\left(1 + e^{-w^T x_i}\right)$$

$$= \log\left(1 + e^{-w^T x_i}\right) + (1 - y_i) w^T x_i$$

$$\nabla_w NLL_i(w) = -\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}} \cdot x_i + (1 - y_i) \cdot x_i \quad \in \mathbb{R}^d$$

$$= -\left(y_i - \sigma\left(w^T x_i\right)\right) x_i$$

Apply GD :

$$\boxed{w_{t+1} \leftarrow w_t - \eta \sum_{i=1}^{n} \nabla_w NLL_i(w)}$$

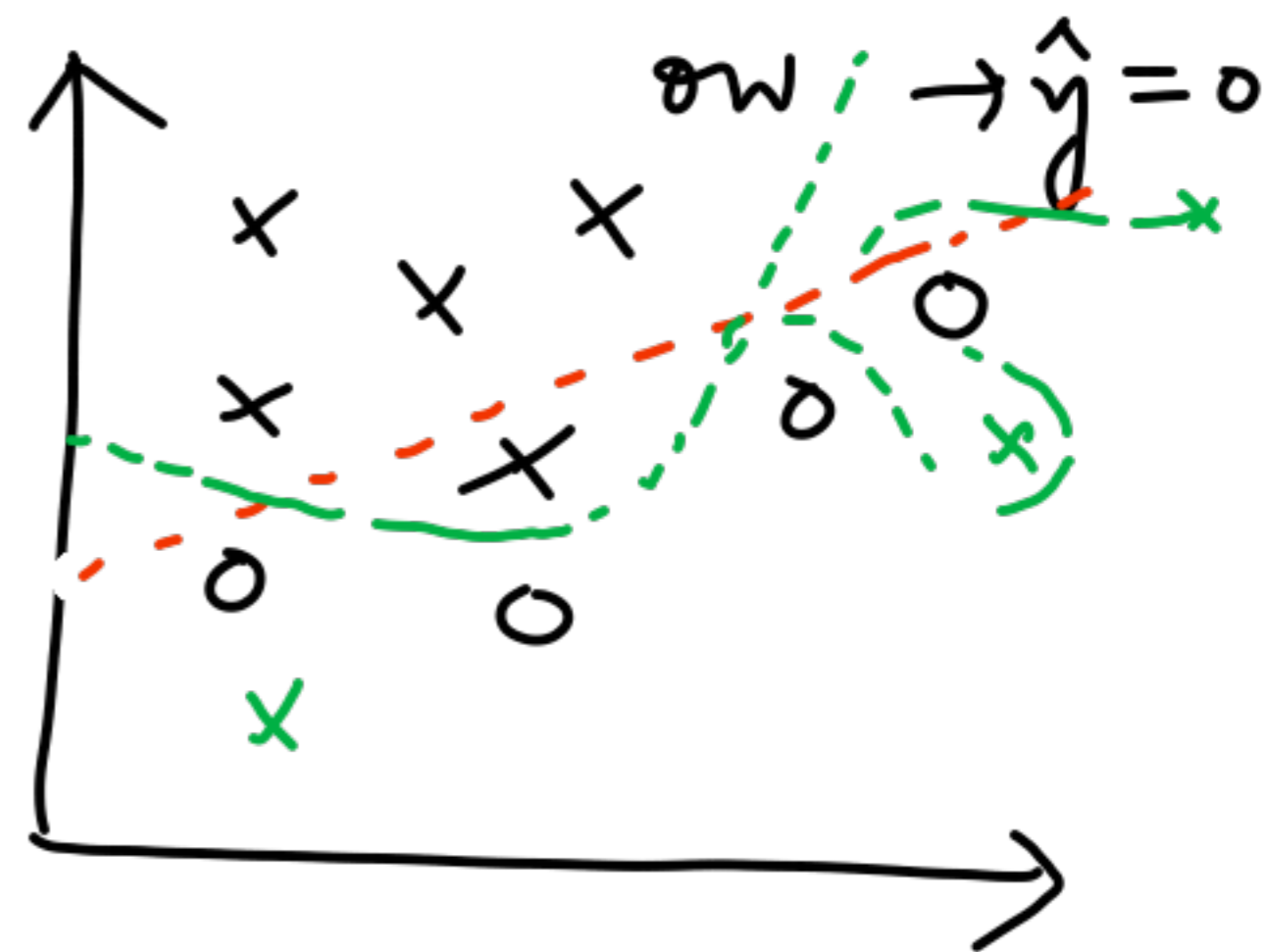$$\frac{P(y=1 \mid \hat{x}, w)}{P(y=0 \mid \hat{x}, w)} > 1 \rightarrow \hat{y}=1$$

$$\text{ow} \rightarrow \hat{y}=0 \Bigg\} \Rightarrow w^T \hat{x} > 0 \rightarrow \hat{y}=1$$

$$\text{ow} \rightarrow \hat{y}=0$$

Basis functions: $\Phi(x) \rightarrow$
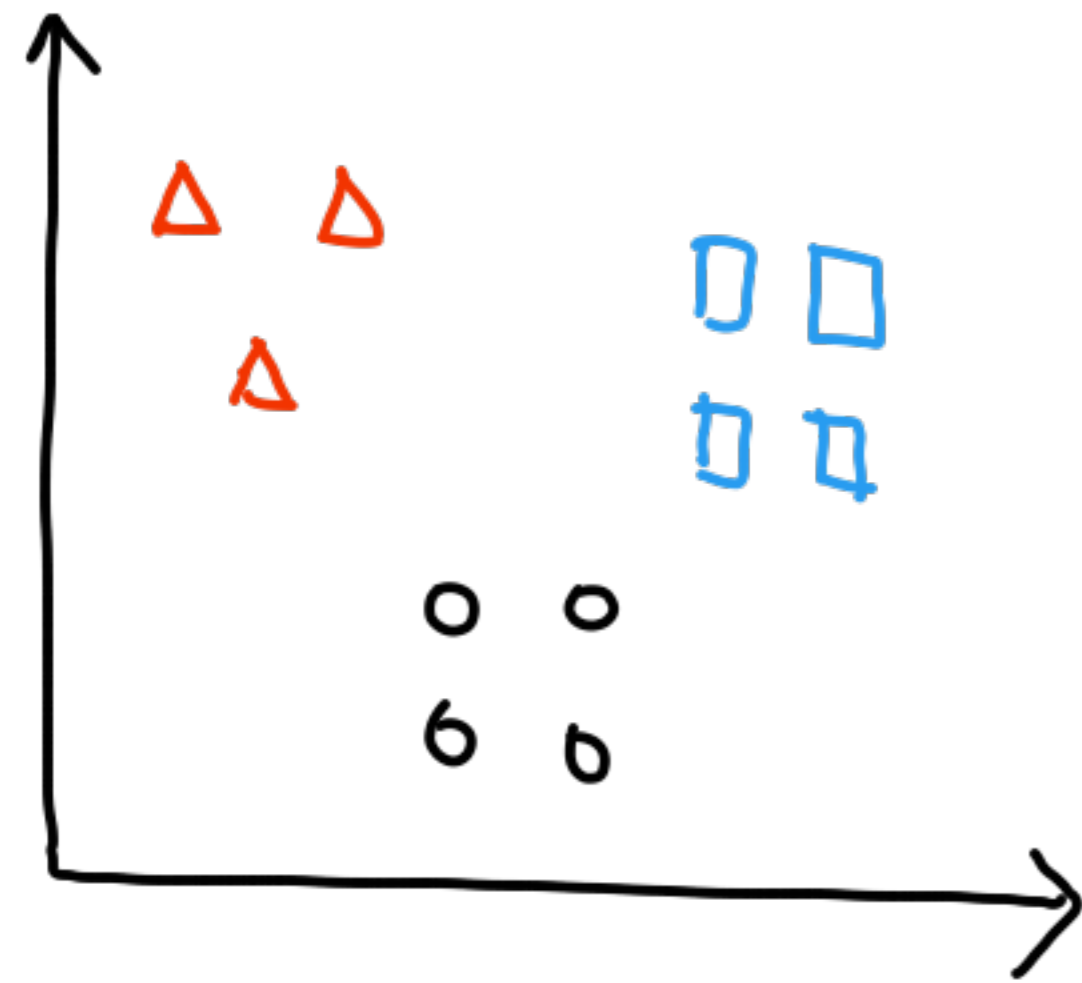
$$\sigma(w^T \phi)$$

$$\underline{w^T \phi > 0}$$



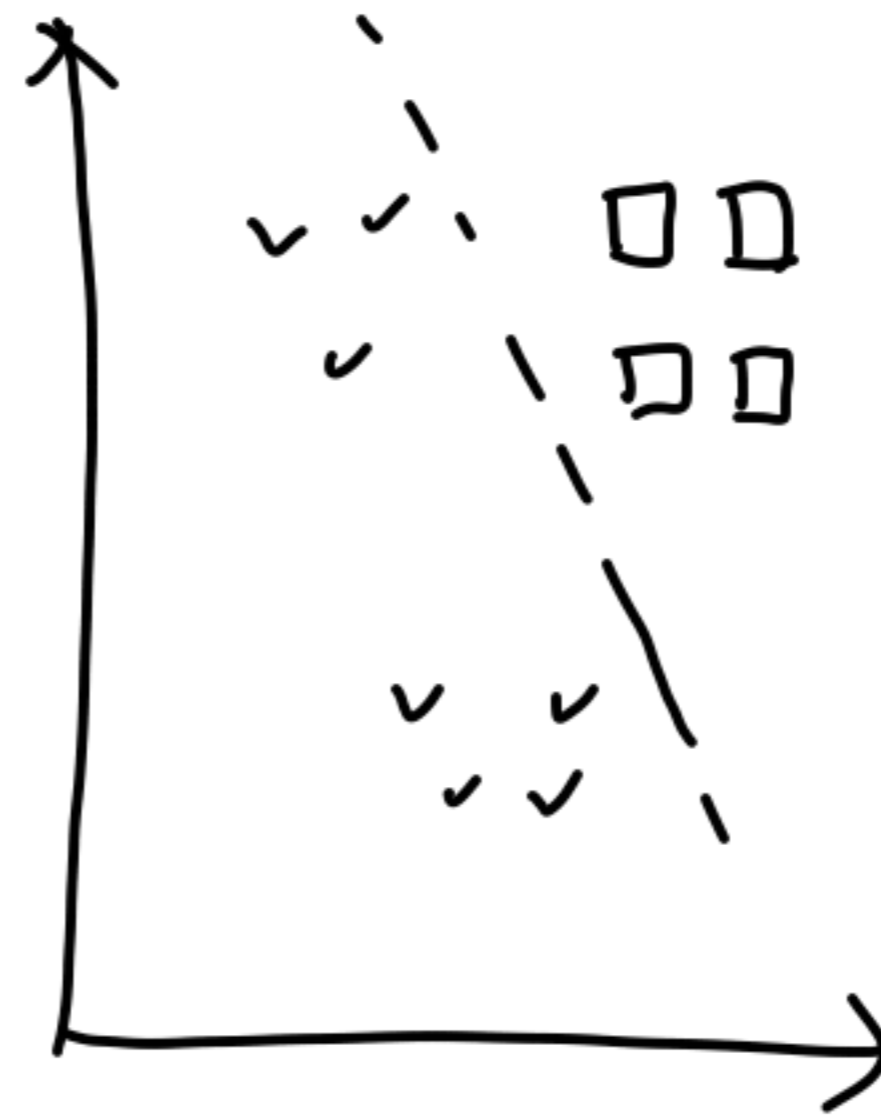$$NLL_{Reg}(w) = NLL_{LR}(w) + \lambda \|w\|^2 \quad : \text{minimize using}$$

$$GD.$$

# Multi-class classification:

## ① One-vs-rest classifier



Any binary classifier can be used.

test point $\hat{x}$

$$\sigma(w_1^T \hat{x}) \to \text{class 1}$$

$$\sigma(w_2^T \hat{x}) \to \text{class 2}$$

$$\sigma(w_3^T \hat{x}) \to \text{class 3}$$

Final prediction

$$\hat{y} = \text{argmax}_k \sigma(w_k^T \hat{x})$$

# Softmax Regression

$$1 \quad \varrho \, e^{w_1^T x} \longrightarrow P(y = 1 \mid x, w_1)$$

$$2 \quad \varrho \, e^{w_2^T x} \longrightarrow P(y = 2 \mid x, w_2)$$

$$3 \quad \varrho \, e^{w_3^T x}$$

$$\vdots$$

$$k \quad \varrho \, e^{w_k^T x}$$

$$\varrho = \frac{1}{\sum\limits_{j=1}^{k} e^{w_j^T x}}$$

$$\boxed{P(y = j \mid x, w) = \frac{e^{w_j^T x}}{\sum\limits_{i=1}^{k} e^{w_i^T x}}}$$

$$P(y = k \mid x, w_k)$$

$$f(x, W) = \begin{bmatrix} P(y=1 \mid x, W) \\ P(y=2 \mid x, w) \\ \vdots \\ \vdots \\ P(y=K \mid x, W) \end{bmatrix}$$

$$= \begin{bmatrix} \varrho \, e^{w_1^T x} \\ \vdots \\ \vdots \\ \varrho \, e^{w_k^T x} \end{bmatrix}$$

$$= \text{Softmax}(x, w)$$

$$W = \begin{bmatrix} - & w_1^T & - \\ - & w_2^T & - \\ & \vdots & \\ - & w_k^T & - \end{bmatrix}_{K \times d}$$

$$NLL(w) = -\sum_{i=1}^{n} \log P(y_i \mid x_i, W)$$

$$= -\sum_{i=1}^{n} \sum_{k=1}^{K} \mathbb{I}\{y_i = k\} \log \frac{e^{w_k^T x}}{\sum_{j=1}^{K} e^{w_j^T x}}$$

$\uparrow$

indicator $\{ \qquad \}$

$$NLL_i(w) = -\sum_{k=1}^{K} \mathbb{I}\{y_i = k\} \left[ w_k^T x - \log \sum_{j=1}^{K} e^{w_j^T x} \right]$$

$$-\nabla_{w_k} NLL_i(w) = \begin{cases} x - \dfrac{e^{w_k^T x}}{\sum e^{w_j^T x}} \cdot x \; ; & y_i = k \\[4mm] - \dfrac{e^{w_k^T x}}{\sum e^{w_j^T x}} \cdot x \; ; & y_i \neq k \end{cases}$$

$$\nabla_{w_k} NLL_i(w) = -\left[\mathbb{I}\{y_i = k\} - f_k(x, w)\right] \cdot x$$

$$y_i \in \{1, 2, 3, \ldots, K\}$$

$$\nabla_{w_k} NLL(w)$$
$$= \sum_{i=1}^{n} \nabla_{w_k} NLL_i(w)$$

$$y_i^{(k)}$$

$$y_i = \begin{bmatrix} 0 & 0 & \cdots & 1 & 0 & 0 \end{bmatrix} : \text{one-hot representation}$$
$$\uparrow$$
$$y_i^{(k)}$$

GD Step

$$\begin{bmatrix} - & w_1^T & - \\ & \vdots & \\ - & w_k^T & - \end{bmatrix}_{t+1} \leftarrow \begin{bmatrix} w_1^T \\ \vdots \\ w_k^T \end{bmatrix}_t - \eta \begin{bmatrix} \cdots & \nabla_{w_1} NLL(w) & \cdots \\ & \vdots & \\ \cdots & \nabla_{w_k} NLL(w) & \cdots \end{bmatrix} \to w^*$$

## NB vs LR

$$\begin{cases} \underset{y}{\text{argmax}} \quad \underbrace{P(Y=y)}_{\text{infer class distri-bution}} \prod_{i=1}^{d} \underbrace{P(x_i \mid Y=y)}_{\text{infer data distribution}} \end{cases} : \text{Naive Bayes}$$

Generative models

Logistic regression $\longrightarrow$ $\underset{y}{\text{argmax}} \quad P(y \mid x, w) \longrightarrow$ Discriminative models.

# Gaussian Naive Bayes : Special case of LR

$$\underset{y_k}{\text{argmax}} \quad P(Y=y_k) \prod_{i=1}^{d} P(x_i \mid Y=y_k)$$

continuous

HW: find $w_0, w_j\text{'s}$

$$P(x_j \mid y_k) \sim N\left(\mu_{jk}, \sigma_{jk}^2\right)$$

① $x_i, x_j$ are conditionally indept.

② $P(y=1) = \pi$ , $P(y=0) = 1-\pi$

③ $P(x_i \mid y=0) \sim N(\mu_{i0}, \sigma_i^2)$
   $P(x_i \mid y=1) \sim N(\mu_{i1}, \sigma_i^2)$

$$P(y_i=1 \mid x)$$

$$= \frac{P(x \mid y_i=1) P(y_i=1)}{P(x \mid y_i=1) P(y_i=1) + P(x \mid y_i=0) P(y_i=0)}$$

$$\vdots$$

$$\implies \frac{1}{1 + \exp\{w_0 + \sum w_j x_j\}}$$