

$$P(y_i|x_i, \omega) = \frac{e^{\omega_i^T x_i}}{e^{\omega_i^T x_i} + e^{\omega_2^T x_i}}$$

$$\begin{aligned} \rightarrow f(\omega, x_i) &= \frac{1}{1 + e^{-\omega^T x_i}} \quad \text{where } \omega = \omega_1 - \omega_2 \\ &= \sigma(\omega^T x_i) \quad (\text{Sigmoid fn: } \sigma(z) = \frac{1}{1 + e^{-z}}) \end{aligned}$$

$$\begin{cases} P(y_i=1|x_i, \omega) \\ P(y_i=0|x_i, \omega) \end{cases} \quad \begin{array}{l} \text{for class 1, } \omega_i^T x_i \\ \text{for class 2, } \omega_2^T x_i \end{array}$$

$$\begin{bmatrix} \tau e^{\omega_i^T x_i} \\ \tau e^{\omega_2^T x_i} \end{bmatrix}$$

$$\text{where } \tau = \frac{1}{e^{\omega_i^T x_i} + e^{\omega_2^T x_i}}$$

## Binary Logistic Regression

We've  $\omega_{LR}^* \in \mathbb{R}^d$ , such that  $\omega_{LR}^* = \arg \max_{\omega} \prod_{i=1}^n P(y_i|x_i, \omega)$

$$\begin{aligned} &= \arg \max_{\omega} \sum_{i=1}^n \log P(y_i|x_i, \omega) \\ &= \arg \min_{\omega} \left( - \sum_{i=1}^n \log P(y_i|x_i, \omega) \right) \end{aligned}$$

We've Negative log likelihood,

$$\begin{aligned} NLL_i &= -\log P(y_i|x_i, \omega) \\ &= -\log P(y_i=1|x_i, \omega), y_i=1 \quad \text{or} \quad -\log (1 - P(y_i=1|x_i, \omega)), y_i=0 \\ &= -y_i \log (\sigma(\omega^T x_i)) - (1-y_i) \log (1 - \sigma(\omega^T x_i)) \end{aligned}$$

(Cross Entropy Loss)

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

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

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

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

We are going to apply gradient descent on  $NLL_i$  to minimize it.

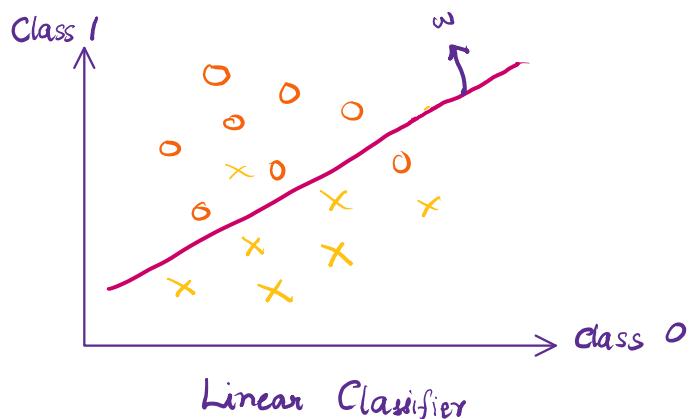
$$\rightarrow \dots \rightarrow -\frac{\partial}{\partial \omega^T x_i} \dots \rightarrow \dots$$

We are going to apply gradient descent on  $NLL_i(\omega)$  to minimize it -

$$\begin{aligned}\therefore \nabla_{\omega} NLL_i(\omega) &= \frac{-e^{-\omega^T x_i}}{1+e^{-\omega^T x_i}} x_i + (1-y_i) x_i \\ &= \frac{1}{1+e^{\omega^T x_i}} x_i - x_i y_i \in \mathbb{R}^d \\ &= -(y_i - \sigma(\omega^T x_i)) x_i\end{aligned}$$

Applying gradient descent:  $\omega_{t+1} \leftarrow \omega_t - \eta \sum_{i=1}^n \nabla_{\omega} NLL_i(\omega)$

Recall that  $\frac{P(y=1|x, \omega)}{P(y=0|x, \omega)} > 0 \rightarrow \hat{y}=1$  or,  $\hat{y}=0 \Rightarrow \omega^T x > 0 \rightarrow \hat{y}=1$  or,  $\hat{y}=0$



To accommodate a better classification, we can use a basis function, (give place to outliers)

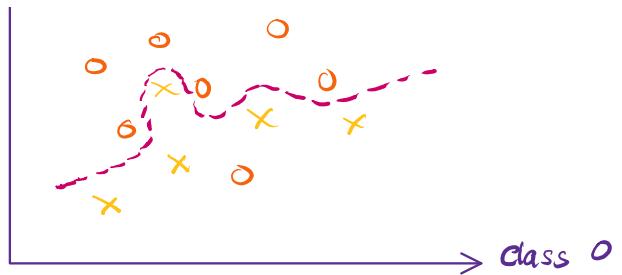
Basis function:  $\Phi(x)$

Logistic function:  $\sigma(\omega^T \Phi)$

Now we check,  $\omega^T \Phi > 0$  instead of  $\omega^T x > 0$

This may be linear in higher dimensions, but non-linear in lower





Classification using basis

We can accomodate for overfitting by using a regularizer:

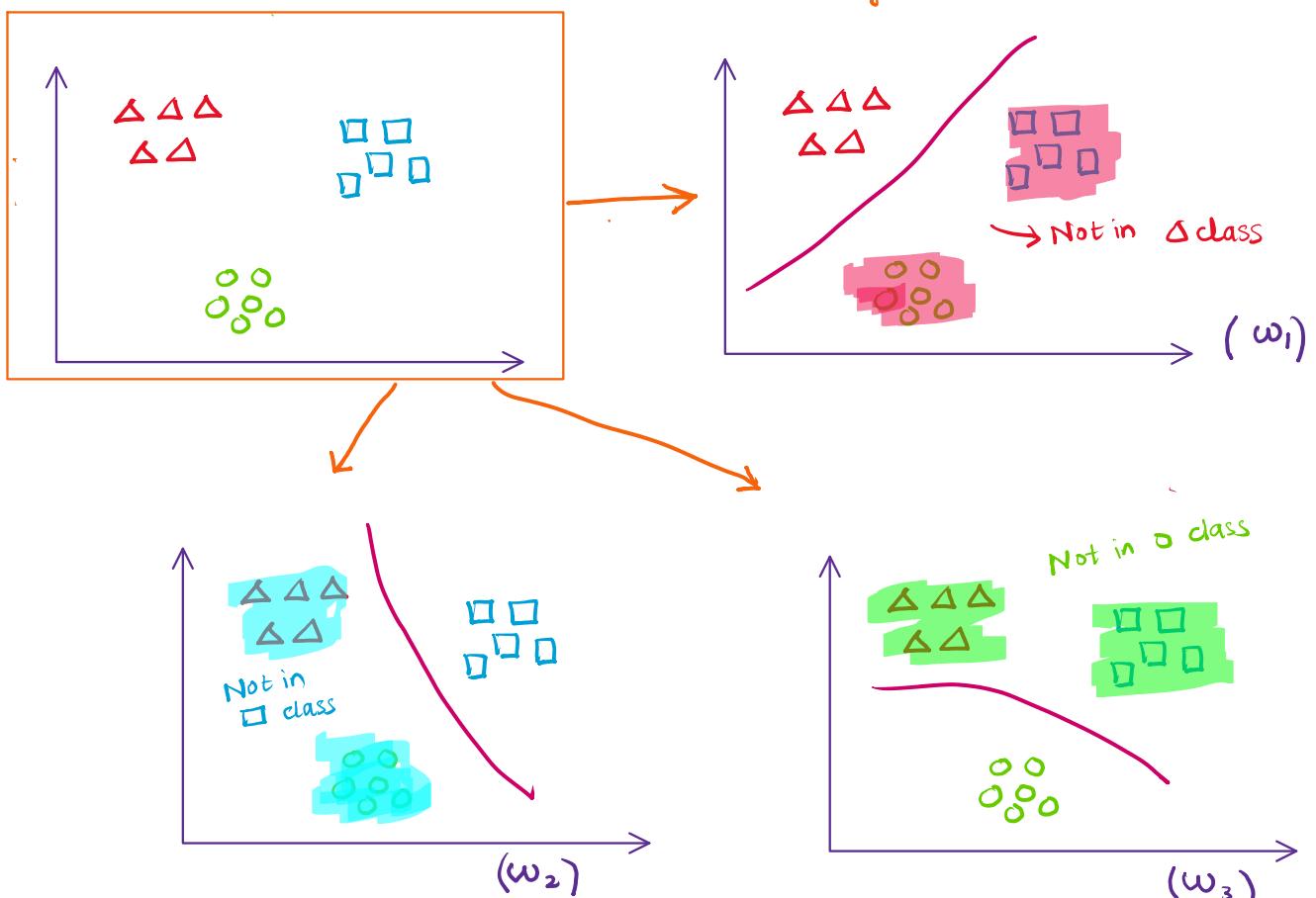
$$NLL_{reg}(w) = NLL_{LR}(w) + \lambda \|w\|^2$$

(Minimize using gradient descent)

## Multi-class Classification

### ① One-vs-rest Classifier

We can use any binary classifier



Let  $\hat{x}$  be test point

$$\sigma(w_1^\top \hat{x}) \rightarrow \text{class 1}$$

$$\sigma(w_2^\top \hat{x}) \rightarrow \text{class 2}$$

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

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

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

$\therefore$  Final prediction,  $\hat{y} = \operatorname{argmax}_k \sigma(w_k^T \hat{x})$

## Softmax Regression:

$$w_i \rightarrow P(y=i|x, w) \quad (\text{Probability that } y \text{ belong to } i^{\text{th}} \text{ class})$$

Like logistic regression we've

$$P(y=i|x, w) = \tau e^{w_i^T x} : \text{where } \tau = \frac{1}{\sum_{j=1}^k e^{w_j^T x}} \quad (\text{Assume k classes})$$

$$= \frac{e^{w_i^T x}}{\sum_{j=1}^k e^{w_j^T x}}$$

Softmax fn

$$f(x, w) = \begin{bmatrix} P(y=1|x, w) \\ P(y=2|x, w) \\ \vdots \\ P(y=k|x, w) \end{bmatrix} = \begin{bmatrix} \tau e^{w_1^T x} \\ \tau e^{w_2^T x} \\ \vdots \\ \tau e^{w_k^T x} \end{bmatrix} \text{ where, } w = \begin{bmatrix} -w_1^T \\ -w_2^T \\ \vdots \\ -w_k^T \end{bmatrix}_{k \times d}$$

Here,  $NLL(w) = -\sum_{i=1}^n P(y_i|x, W) = -\sum_{i=1}^n \sum_{k=1}^K \mathbb{I}\{y_i=k\} \log \frac{e^{w_k^T x}}{\sum_{j=1}^k e^{w_j^T x}}$  Indicator fn.

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

Again, we will perform gradient descent:

$$-\nabla_{w_k} NLL_i(w) = \begin{cases} x - \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}, & y_i = k \\ -\frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}, & y_i \neq k \end{cases}$$

$$\therefore \nabla_{w_k} \text{NLL}_i(\omega) = - \left[ \mathbb{I}\{y_i=k\} - f_k(x, \omega) \right] x$$

One-Hot Representation:

Instead of indicator function,

we can represent  $y_i$  as a vector with  $k$  entries such that if  $y_i \in \text{class } \alpha$ , everywhere in the vector it is zero except  $\alpha^{\text{th}}$  position

$$y_i = [0 \ 0 \ \dots \underset{\alpha}{1} \ \dots \ 0]$$

Gradient Descent Step:

$$\begin{bmatrix} w_1^\top \\ w_2^\top \\ \vdots \\ w_k^\top \end{bmatrix}_{t+1} \leftarrow \begin{bmatrix} w_1^\top \\ w_2^\top \\ \vdots \\ w_k^\top \end{bmatrix}_t - \eta \begin{bmatrix} -\nabla_{w_1} \text{NLL}(\omega) \\ -\nabla_{w_2} \text{NLL}(\omega) \\ \vdots \\ -\nabla_{w_k} \text{NLL}(\omega) \end{bmatrix}$$

where,  
 $\nabla_{w_k} \text{NLL}(\omega) = \sum_{i=1}^n \nabla_{w_k} \text{NLL}_i(\omega)$

NB vs LR

$$\text{Naive Bayes} : \arg \max_y P(y=y) \prod_{i=1}^d \underbrace{P(x_i | y=y)}_{\substack{\text{inter class} \\ \text{distribution}}} \quad \text{given by data distribution}$$

Models like NB are called generative model which uses data distribution to infer class distribution

$$\text{Logistic Regression} : \arg \max_y P(y|x, \omega) \quad \text{given by } \sigma(w^\top x)$$

Models like LR are called discriminative models.

Gaussian Naive Bayes: Special case of LR

## Gaussian Naive Bayes: Special case of LR

$$\underset{y_k}{\operatorname{argmax}} \ P(y=y_k) \prod_{i=1}^d P(x_i | Y=y_k)$$

$$\text{where, } P(x_i | y_k) \sim N(\mu_{ik}, \sigma_{ik}^2)$$

with following assumptions :

1.  $x_i, x_j$  are conditionally independent
2.  $P(y=1) = \pi, P(y=0) = 1-\pi$
3.  $P(x_i | y=0) \sim N(\mu_{i0}, \sigma_i^2)$   
 $P(x_i | y=1) \sim N(\mu_{i1}, \sigma_i^2)$

$$\begin{aligned} \therefore \text{We've } P(y_i=1 | x) &= \frac{P(x|y_i=1) P(y_i=1)}{P(x|y_i=1) P(y_i=1) + P(x|y_i=0) P(y_i=0)} \\ &= \frac{1}{1 + \exp\{w_0 + \sum w_j x_j\}} \quad \left( \begin{array}{l} \text{(What are } w_0, w_j \text{'s?)} \\ (\text{HW}) \end{array} \right) \end{aligned}$$