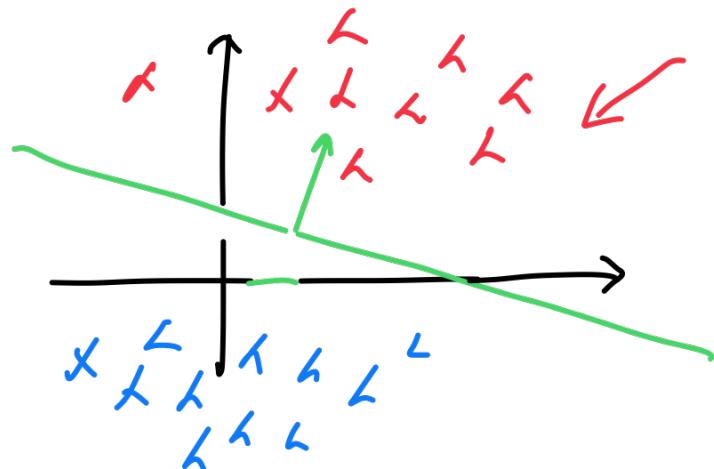


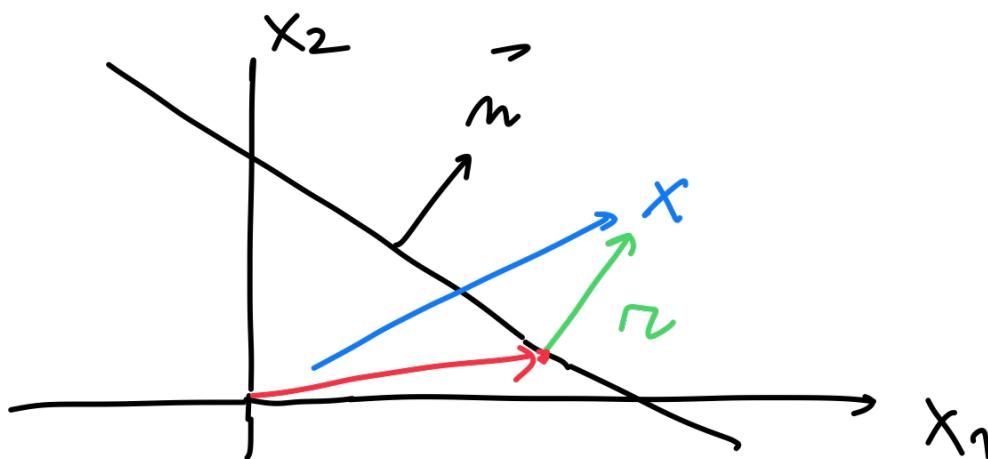
Classification  $\{t^{(i)}\}_{i=1}^N \quad t^{(i)} \in \{1, 2, 3, \dots, k\}$



Plane equation

normal to  $\pi$   
 $= (\beta_0, \beta_1, \dots, \beta_D)$

$$\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D = 0$$



if  $x$  above plane

$$\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D > 0$$

$x_1$  if  $x$  below plane

$$\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D < 0$$

$$x = x_{\perp} + \frac{\vec{m}}{\|\vec{m}\|} \gamma$$

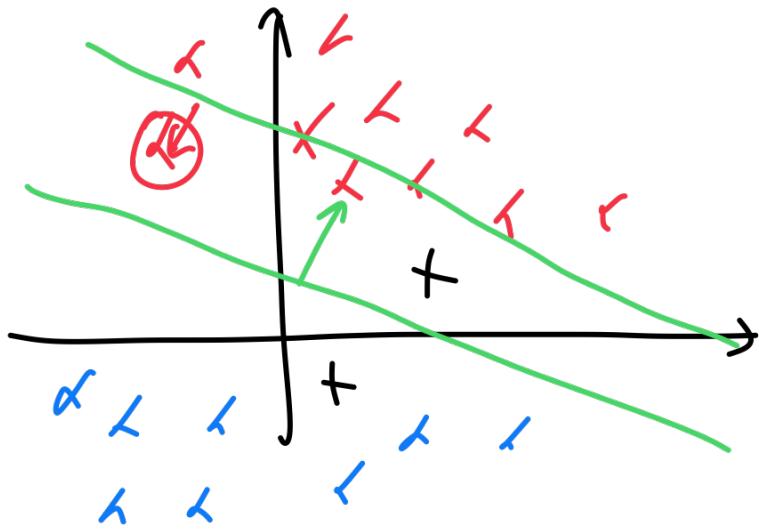
$$\beta_0 + \vec{m}^T \vec{x} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_D x_D$$

$$\beta_0 + \vec{m}^T x_{\perp} + \frac{\vec{m}^T \vec{m}}{\|\vec{m}\|} \gamma = \frac{\|\vec{m}\|^2}{\|\vec{m}\|} \gamma > 0$$

$$\beta_0 + \beta_1 x_{\perp,1} + \beta_2 x_{\perp,2} + \dots + \beta_D x_{\perp,D} = 0$$

1st approach: learn classifier through least squares

$$\min_{\beta} \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$



$$(1 - \left(\frac{7}{10}\right))^2$$

Possible extensions to Multi-class framework:

→ One vs rest →  $K-1$  linear classifiers

→ One vs one →  $\frac{K(K-1)}{2}$  classifiers

Multiple class discriminant      Assume K classes

if  $x^{(i)} \in C_k \Rightarrow$  store  $t^{(i)}$  as  $t^{(i)} = [0 \dots 0, \underset{k}{1}, 0 \dots 0]$

$$T = \begin{bmatrix} | & | & | \\ t^{(1)} & t^{(2)} & \dots & t^{(N)} \\ | & | & & | \end{bmatrix}^T \quad \tilde{x} = \begin{bmatrix} | & (\tilde{x}^{(1)})^T & | \\ | & (\tilde{x}^{(2)})^T & | \\ | & \vdots & | \\ | & (\tilde{x}^N)^T & | \end{bmatrix}$$

*k-th entry*

$$\beta = \begin{bmatrix} | & | & | \\ \beta^{(1)} & \beta^{(2)} & \dots & \beta^{(K)} \\ | & | & & | \end{bmatrix}$$

$$\underset{\beta}{\operatorname{argmin}} \quad \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K (t_k^{(i)} - (\tilde{x}^{(i)})^T \beta_k)^2$$

$$\underset{B}{\operatorname{argmin}} \frac{1}{NK} \operatorname{Tr} \left( \underbrace{(\mathbf{T} - \mathbf{X}\mathbf{B})^T}_{K} \underbrace{(\mathbf{T} - \mathbf{X}\mathbf{B})}_{N} \right) \quad (*)$$

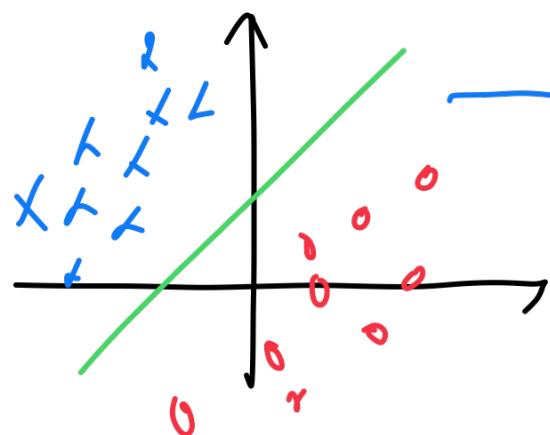
$$\mathbf{T} - \mathbf{X}\mathbf{B} = \underbrace{\begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{t}^{(2)} \\ \vdots & & \vdots \\ \mathbf{0} & \mathbf{t}^{(\omega)} & \mathbf{0} \end{bmatrix}}_K - \underbrace{\begin{bmatrix} \tilde{\mathbf{x}}^{(2)} \\ \vdots \\ \tilde{\mathbf{x}}^{(\omega)} \end{bmatrix}}_N \underbrace{\begin{bmatrix} \mathbf{B}^{(2)} \\ \vdots \\ \mathbf{B}^{(\omega)} \end{bmatrix}}_D$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}^T \begin{bmatrix} \frac{a_{11}}{a_{21}} & \frac{a_{12}}{a_{21}} & \frac{a_{13}}{a_{21}} \\ \frac{a_{12}}{a_{31}} & \frac{a_{22}}{a_{31}} & \frac{a_{32}}{a_{31}} \\ \frac{a_{13}}{a_{31}} & \frac{a_{23}}{a_{31}} & \frac{a_{33}}{a_{31}} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{bmatrix}$$

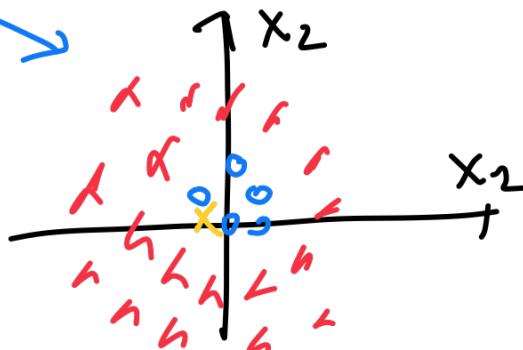
$$= \begin{bmatrix} \frac{a_{11}^2 + a_{21}^2 + a_{31}^2}{a_{12}^2 + a_{22}^2 + a_{32}^2} \\ \frac{a_{12}^2 + a_{22}^2 + a_{32}^2}{a_{13}^2 + a_{23}^2 + a_{33}^2} \end{bmatrix}$$

To find  $B$ , compute derivative of (\*) w.r.t  $B$  and set it to 0.

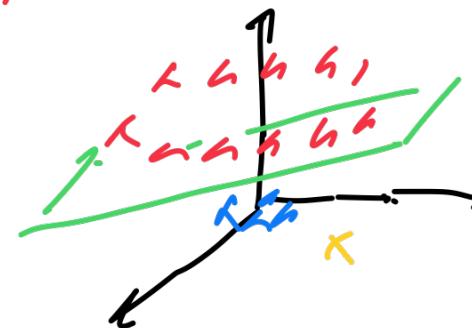
$\Rightarrow$



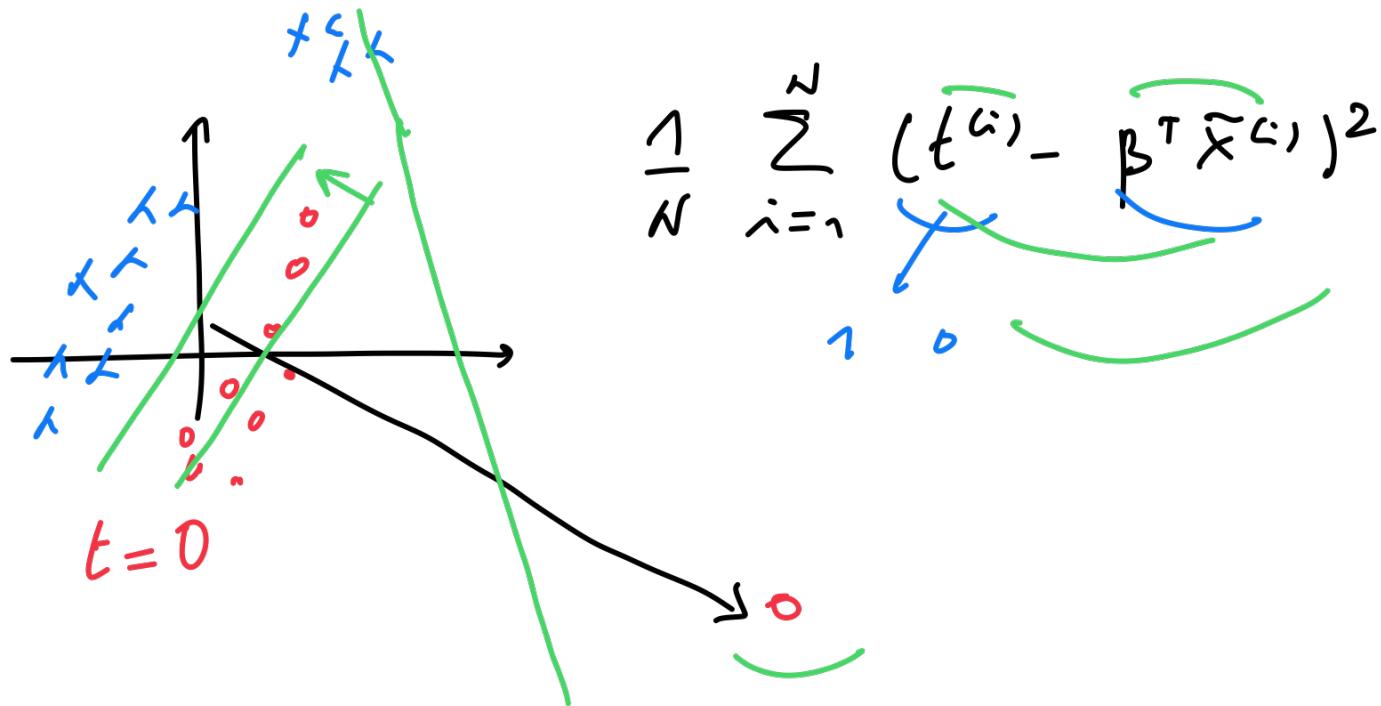
$$B = (X^T X)^{-1} X^T T$$



$$x_1, x_2, x_1^2 + x_2^2$$

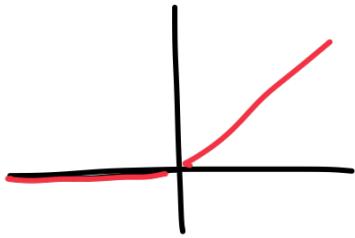


$$x_1, x_2, x_1^2 + x_2^2$$

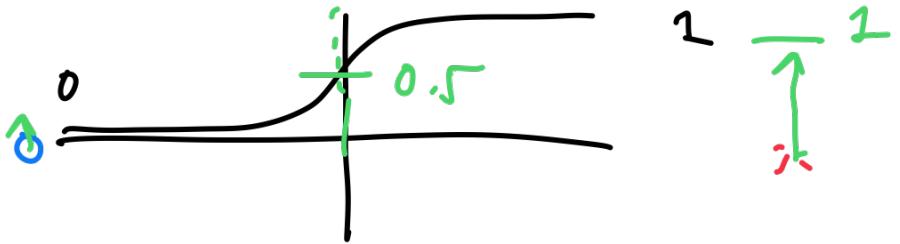


2 Motivations: presence of outliers  
 interpretability of classification  
 (quantitative criterion)

Solution: take an activation function



Sigmoid function



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

## Logistic regression

$$h_{\beta}(x) = \sigma(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_D x_D)$$

$$= \sigma(\beta^T \tilde{x}) \in [0, 1]$$

$$P(t^{(i)} = 1 | x^{(i)}) = \sigma(\beta^T \tilde{x}^{(i)}) \quad \{$$

$$P(t^{(i)} = 0 | x^{(i)}) = 1 - \sigma(\beta^T \tilde{x}^{(i)}) \quad \}$$