

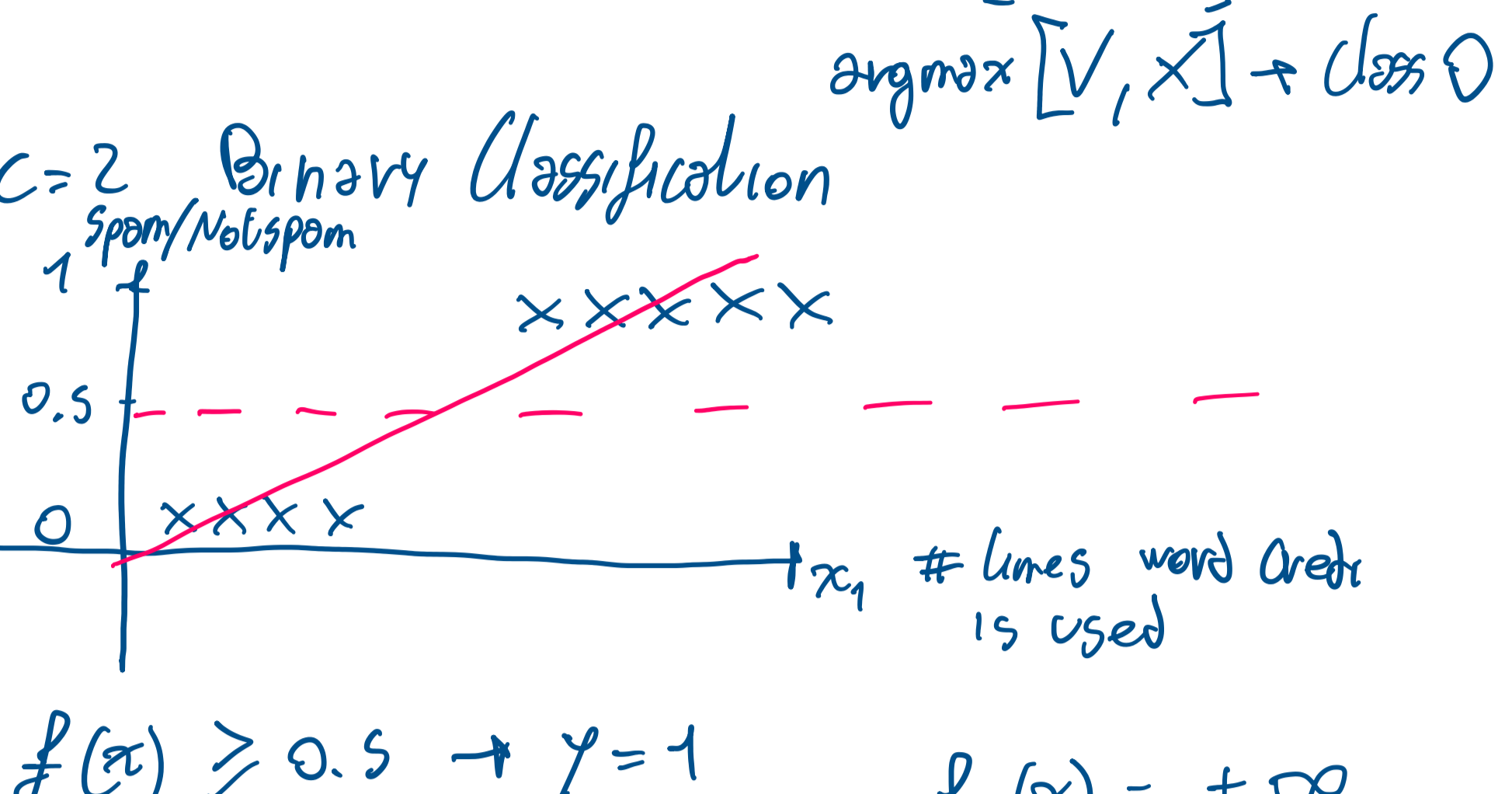
# Classification (Binary)

$y \in \mathbb{R}$   
 $y = \{0, 1\}$  Positive outcome  
 Negative Outcome

$y = \{0, 1, 2\}$  C = Symbol to denote the cardinality of classes  
 Not Spom  
 Unsure  
 Spom

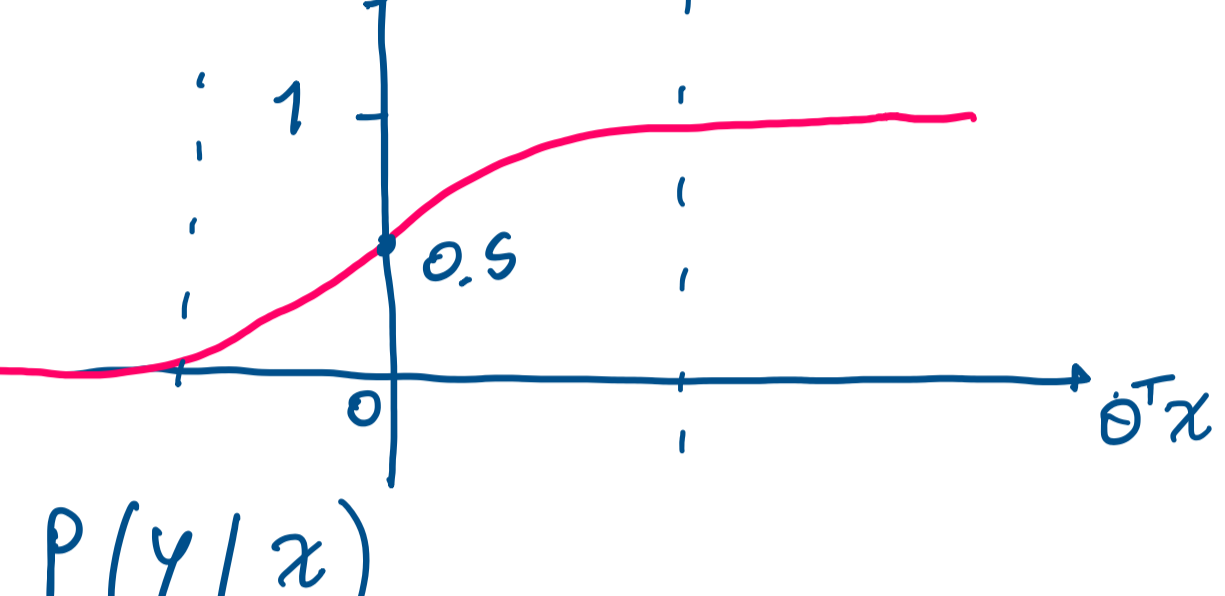
$\hat{y} = f_{\theta}(x) = h_{\theta}(x)$   
 $\hat{y}$  belongs to probability simplex  $\Delta_c$   
 $\gamma_i \geq 0$   $\sum_i \gamma_i = 1$   
 Single element  $\pi_i = \begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix}$

$f_{\theta}(x) = \hat{y} \in \Delta_c$  we can interpret  $\hat{y}$  as a categorical distribution  
 $C := \arg\max_i \hat{y}$   $[0, 1]$   
 $[0.7, 0.3]$   
 $\arg\max [V, x] \rightarrow \text{Class } 0$

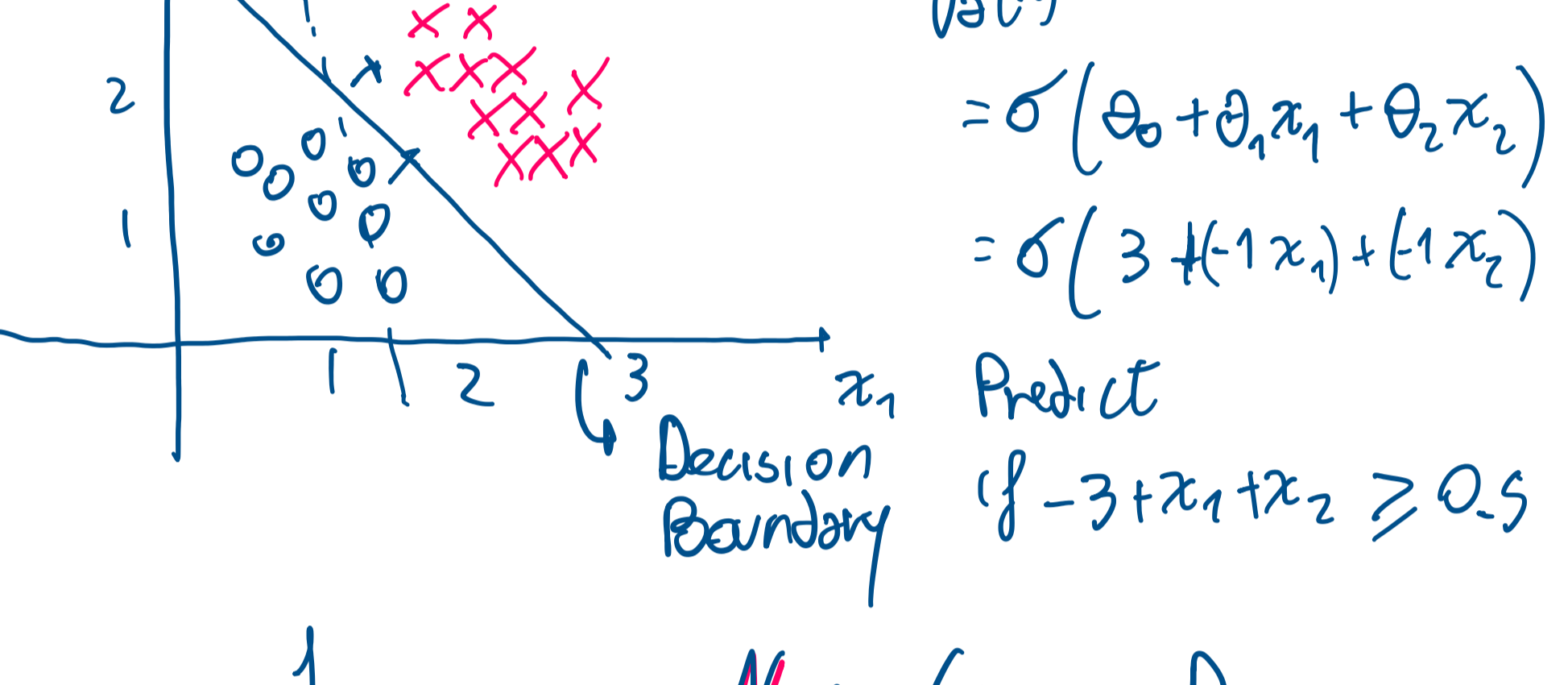


$f(x) \geq 0.5 \rightarrow y=1$   
 $< 0.5 \rightarrow y=0$   $f_{\theta}(x) = \pm \infty$

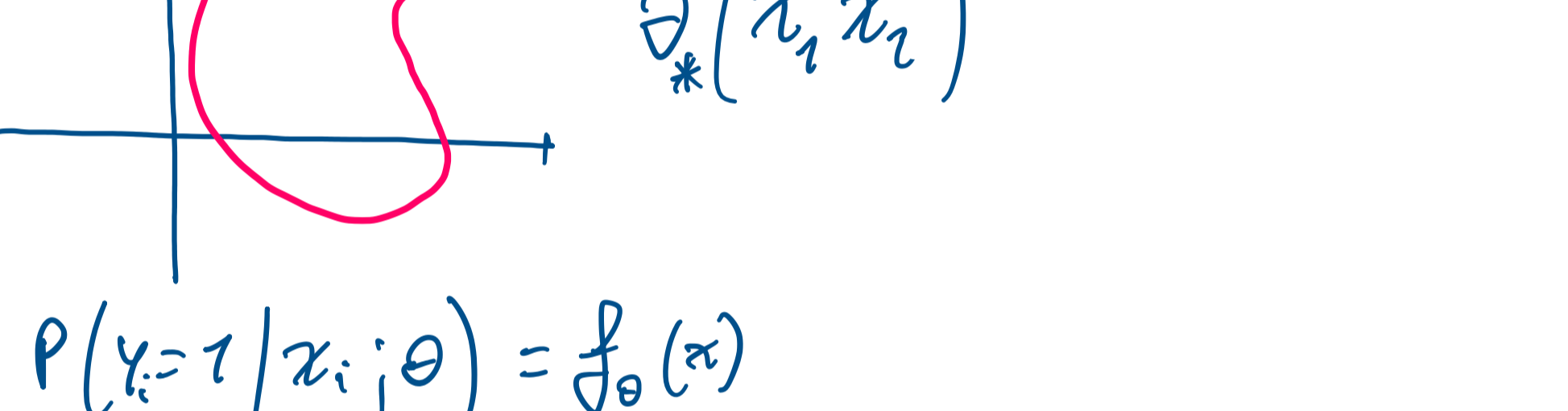
let apply a function  $g(f_{\theta}(x)) \rightarrow \Delta_c$   
 $g(\cdot) = \text{Sigmoid function} = \sigma(\theta^T x_i) = \frac{1}{1 + e^{-\theta^T x}}$



$P(y|x)$   
 $P(y|x; \theta)$  a core of  $f_{\theta}(x)$   
 $P(y=1|x; \theta) = f_{\theta}(x)$   
 $P(y=0|x; \theta) = 1 - f_{\theta}(x)$



$f_{\theta}(x) = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$   
 $= \sigma(3 + (-1)x_1 + (1)x_2)$   
 Predict if  $-3 + x_1 + x_2 \geq 0.5$   
**Non Linear Decision Boundary.**  
 $f_{\theta}(x) = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$   
 $-1 + x_1^2 + x_2^2 \geq 0 \rightarrow y=1$   
 $x_1^2 + x_2^2 = 1 \quad < 0 \rightarrow y=0$



**Can decision boundary?**  
 $\theta^*(x_1, x_2^2)$   
 $P(y_i=1|x_i; \theta) = f_{\theta}(x)$   
 $P(y_i=0|x_i; \theta) = 1 - f_{\theta}(x)$   
 $P(y_i|x_i; \theta) = \begin{cases} f_{\theta}(x_i) & \text{if } y_i=1 \\ 1 - f_{\theta}(x_i) & \text{if } y_i=0 \end{cases} = f_{\theta}(x_i)^{y_i} (1 - f_{\theta}(x_i))^{1-y_i}$   
 $P(Y|X; \theta) = \prod f_{\theta}(x_i)^{y_i} (1 - f_{\theta}(x_i))^{1-y_i}$   
 vector Matrix Likelihood

Log Likelihood  $\log L(\theta) = \frac{1}{n} \sum y_i \log f_{\theta}(x_i) + (1-y_i) \log (1 - f_{\theta}(x_i))$   
 In practice we minimize the Negative Log Likelihood

Max Likelihood = Min NLL  
 From a numerical point of view ensures that as the probability for the right class increases the loss decreases

$y_i=1 \quad \log(\sigma(\theta^T x_i))$   
 $\rightarrow$  close to 1

$NLL(\theta) = -\frac{1}{n} \sum y_i \log f_{\theta}(x_i) + (1-y_i) \log (1 - f_{\theta}(x_i))$

$y_i \log(\sigma(\theta^T x_i))$   
 $\downarrow$  Chain Rule  $\frac{\partial \log(\sigma(z))}{\partial z} = \frac{1}{\sigma(z)} \cdot \sigma'(z)$

$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$   
 $\frac{\partial \theta^T x_i}{\partial \theta} = x_i$   
 $\frac{\partial y_i \log(\sigma(\theta^T x_i))}{\partial \theta} = y_i \cdot \frac{1}{\sigma(\theta^T x_i)} (\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i))) x_i$   
 $= y_i (1 - \sigma(\theta^T x_i)) x_i$

Chain Rule  $\frac{\partial \log(1 - \sigma(z))}{\partial z} = \frac{-1}{(1 - \sigma(z))} \cdot \sigma'(z)$   
 $= (1 - y_i) \frac{-1}{(1 - \sigma(\theta^T x_i))} \sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i)) x_i$   
 $= -(1 - y_i) \sigma(\theta^T x_i) x_i$

$\frac{\partial NLL(\theta)}{\partial \theta} = -\frac{1}{n} \sum (y_i - \sigma(\theta^T x_i)) x_i$

$\theta + \alpha \frac{1}{n} \sum (y_i - \sigma(\theta^T x_i)) x_i$   
 Difference between Real Label and the predicted one  
 if error is large  $\rightarrow$  the update will be large  
 if " is small  $\rightarrow$  " " " " small

$(0 - 0.3) \cdot x_i$   
 $(1 - 0.7) \cdot x_i$   
 $(0 - 0.7) x_i$

**Replaced by Newton's Method**

$x$  velocity  $x''$  acceleration  
 $x'$ : how quickly (change) in which direction  $x''$ : rate of change of first derivative

Use this information to adjust the step size

$\theta \leftarrow \theta - H^{-1} \nabla NLL(\theta)$   
 $H = \sigma(\theta^T x_i) (1 - \sigma(\theta^T x_i)) x_i x_i^T$   
 $\rightarrow H^{-1}$  is the inverse  $\rightarrow$  Costly and not always guaranteed

Algorithm BFGS: it is a Quasi-Newton Method  
 $\rightarrow$  Compute approximation of the Hessian.