

# Lecture 8 Discriminative Learning - Linear Classification

## Generative model

⇒ 1) learn CCD from training set ( $p(x|y)$ )

2) use BDT & BDR to get  $p(y|x)$   
 $g(x) = \arg\max_y p(y|x)$

- The data is used to learn the CCD, the decision rule is secondary.
- Density estimation is an ill-posed problem.  
• which density to use? Gaussian, Laplacian, GMM, KDE

Vapnik's Advice: "When solving a given problem, try to avoid solving a more general problem as an intermediate step."

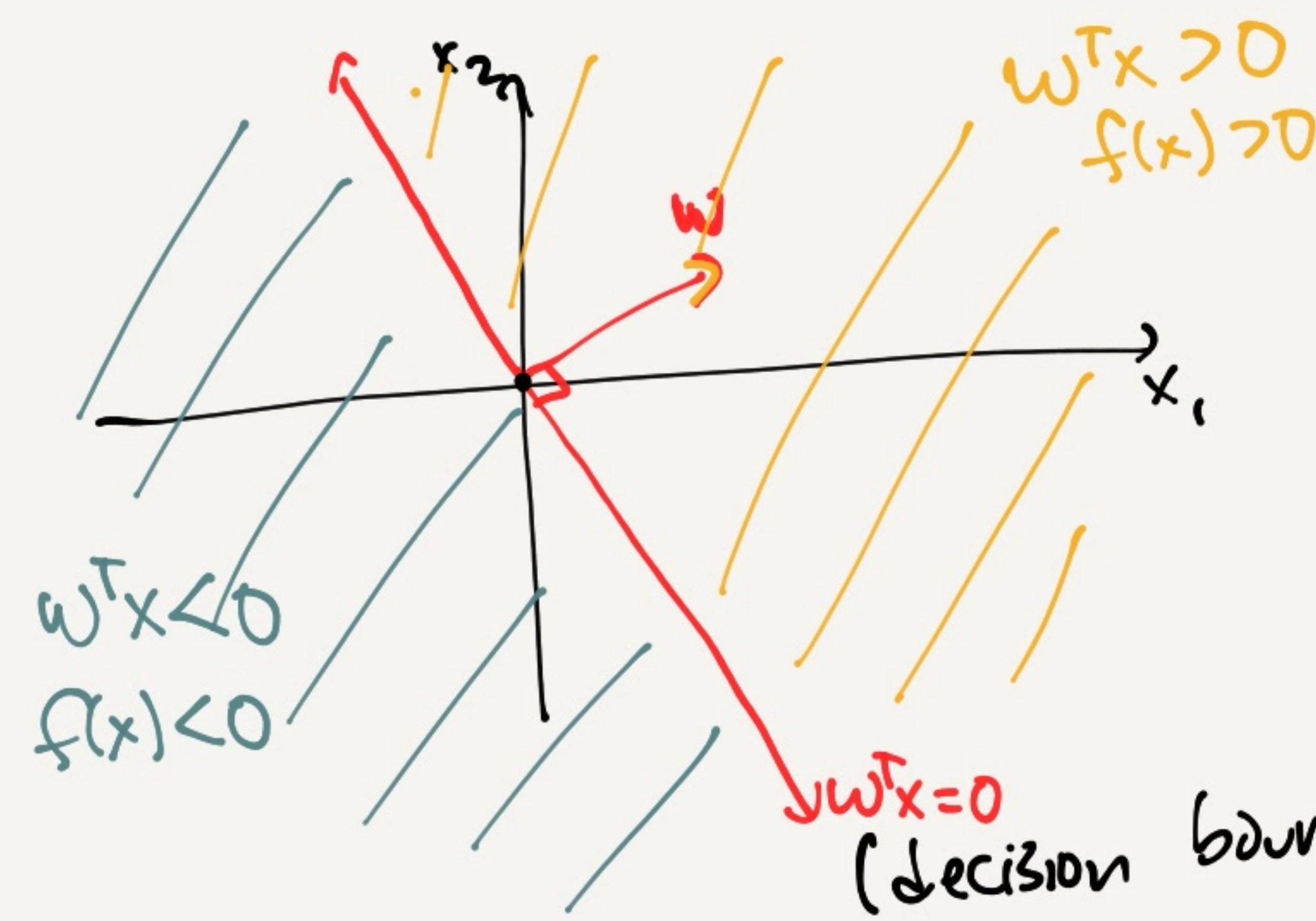
Discriminative Soln: use the data to directly estimate the decision rule.  
•  $g(x)$  or  $p(y|x)$

## Linear classifier

input:  $x \in \mathbb{R}^d$

output:  $y \in \{+1, -1\}$  (binary class)

linear function:  $f(x) = w^T x$ ,  $w \in \mathbb{R}^d$  parameters of classifier.



$w$  separates the space into 2 half-spaces.  $w$  points into the "positive" space.

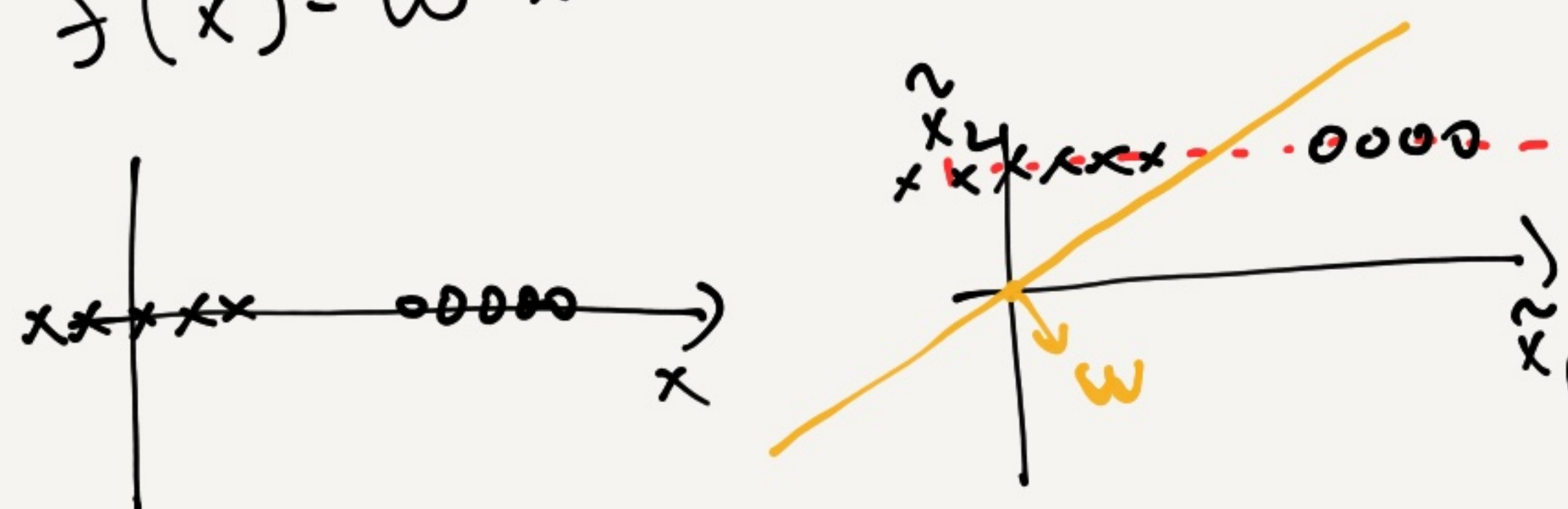
Decision Rule:

$$y^* = \text{sign}(w^T x) = \begin{cases} +1, & w^T x \geq 0 \\ -1, & w^T x < 0 \end{cases}$$

Note: bias term can be included into input.

$$\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad \tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$f(\tilde{x}) = \tilde{w}^T \tilde{x} = w^T x + b$$





Training Set

$$D = \{(x_i, y_i)\}_{i=1}^N \\ = \{X, y\}$$

$$X = \begin{bmatrix} x_1 & \dots & x_N \end{bmatrix} \\ y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Note:

Given a  $w$ :  $y_i w^T x_i > 0 \Rightarrow x_i$  correctly classified

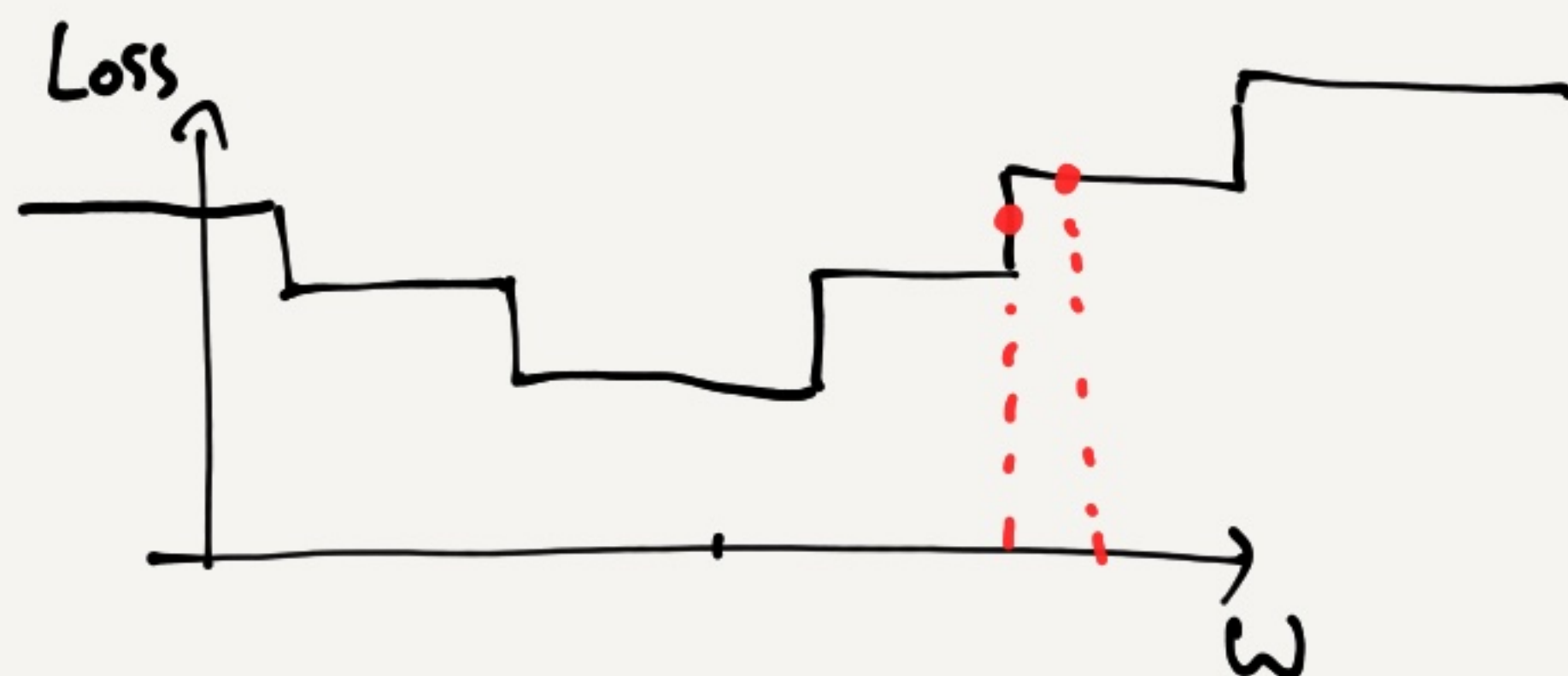
$$\begin{matrix} +1 & > 0 \\ -1 & < 0 \end{matrix}$$

$y_i w^T x_i < 0 \Rightarrow x_i$  is misclassified.

Ideal Case: 0-1 loss function

Optimize the # of misclassifications.

$$w^* = \arg\min_w \sum_{i=1}^N \begin{cases} 0, & y_i w^T x_i \geq 0 \text{ (correct)} \\ 1, & y_i w^T x_i < 0 \text{ (incorrect)} \end{cases}$$



Difficult to optimize!

Gradient is either 0, or undefined.

Least squares classification (Label regression)

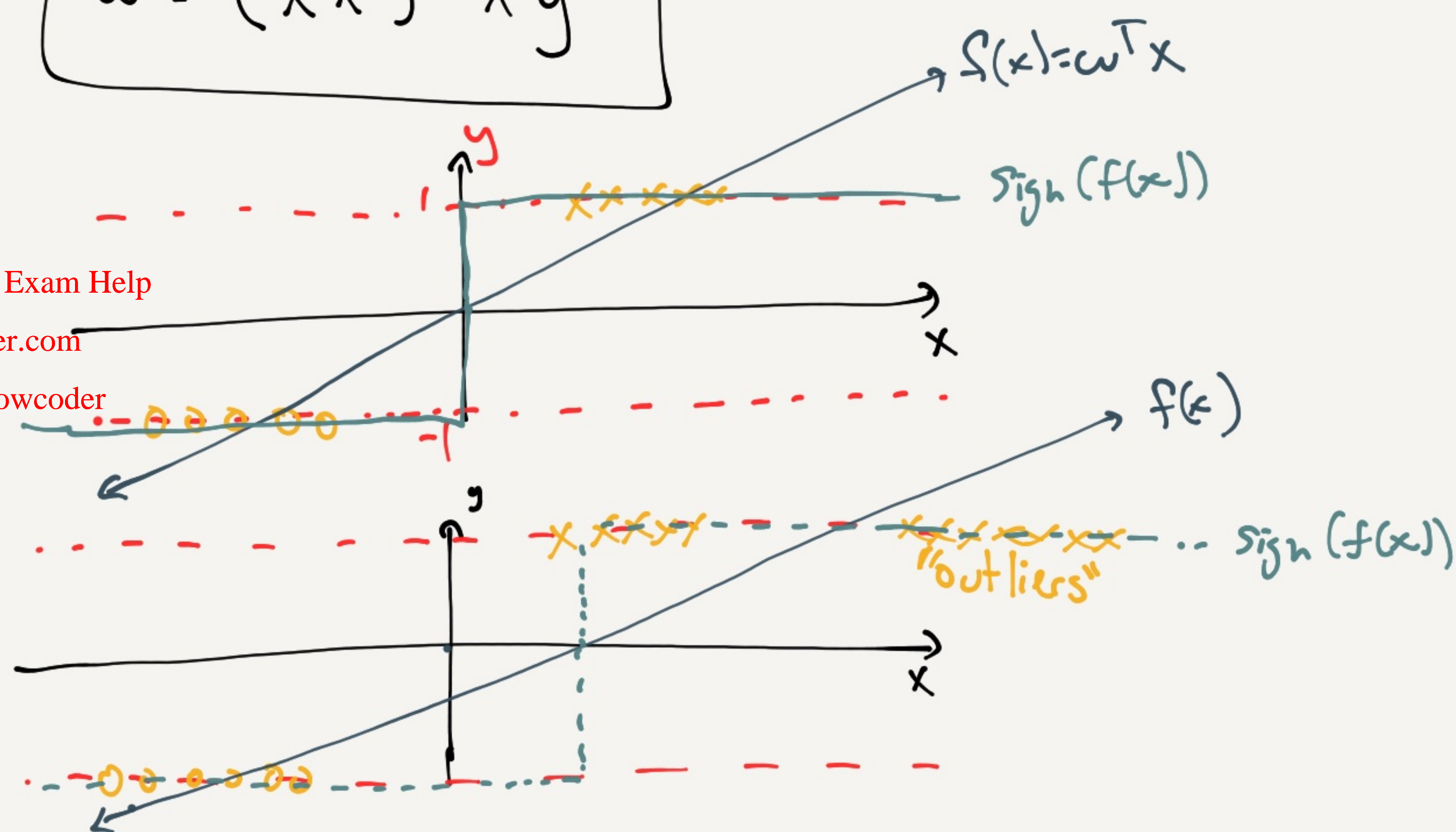
- Ignore the fact that  $y$  is discrete, & just apply LS regression.

• Given  $D = \{X, y\}$

$$\hat{w} = \arg\min_w \sum_{i=1}^N \|w^T x_i - y_i\|^2$$

$$= \arg\min_w \|X^T w - y\|^2$$

$$\hat{w} = (X X^T)^{-1} X y$$



Note: not robust to "outliers"

- the squared penalizes predictions that are "too correct"

- FLD is a version of LSC (PS 7.7)



## Perceptron (Rosenblatt 1962)

Perceptron criteria - only look at misclassified points.

$$E(w) = \sum_{i \in M} -y_i w^T x_i, \quad M = \{\text{misclassified points}\} \\ = \{i \mid y_i w^T x_i < 0\}$$

$E(w) = 0$  when all points correctly classified.

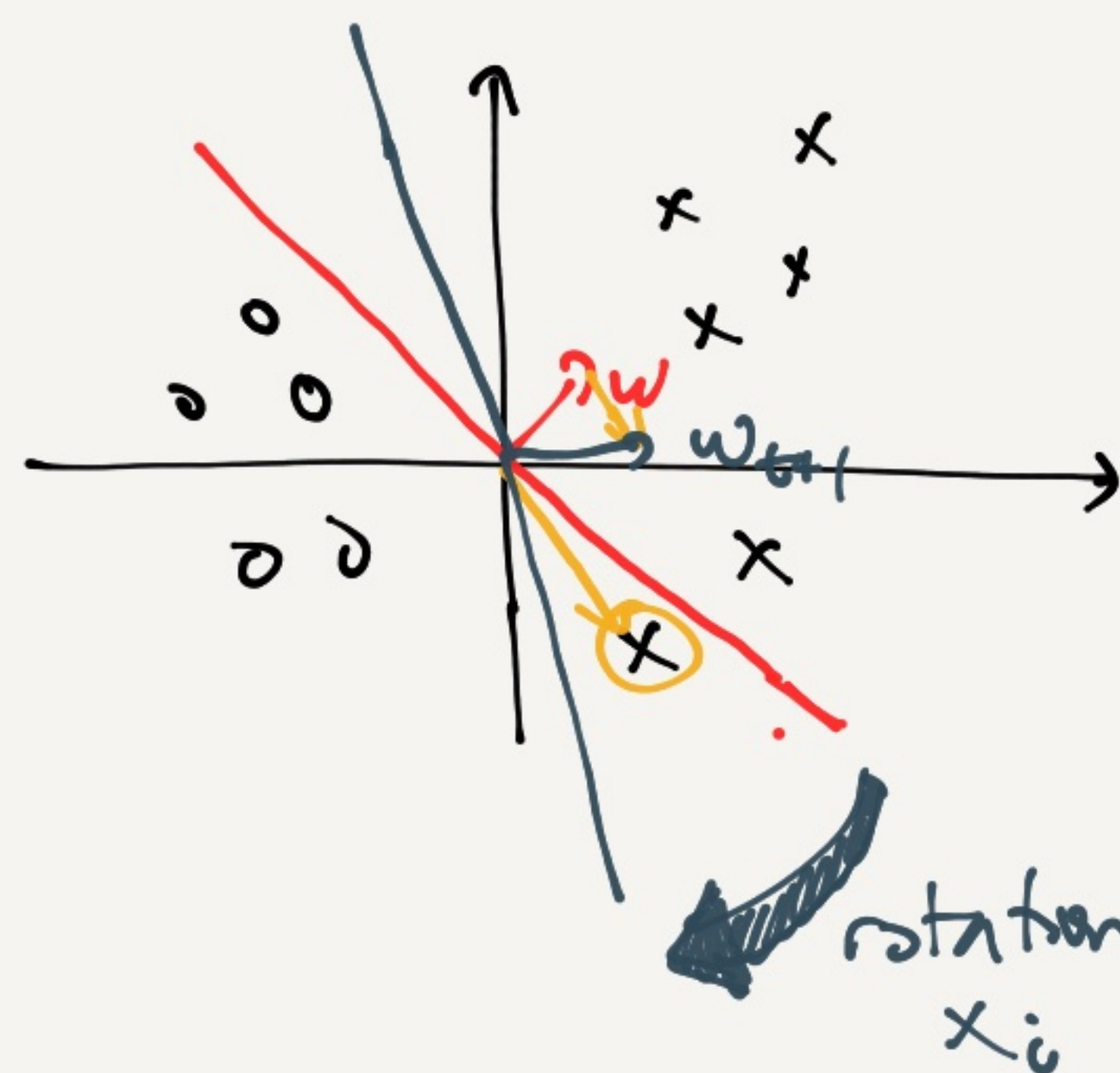
## Perceptron Algorithm

$$w^* = \underset{w}{\operatorname{argmin}} E(w) = \underset{w}{\operatorname{argmin}} \sum_{i \in M} -y_i w^T x_i$$

- computers were slow in 60s...
- apply Stochastic gradient descent (SGD)
  - use one datapoint at a time

$$\Rightarrow w_{t+1} = w_t + \eta y_i x_i, \quad \text{for some } i \in M$$

rotate  $w$  to point towards the misclassified point ( $y_i = 1$ ), & vice versa

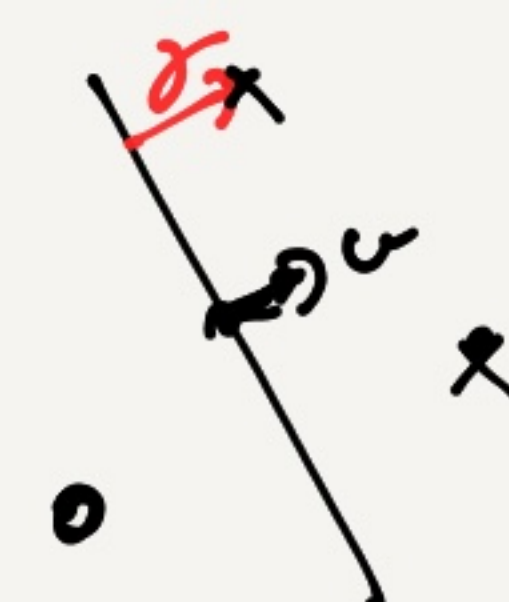


- How to set  $\eta$  the learning rate?

- Rosenblatt proved that SGD converges in  $(\frac{R}{\gamma})^2$  iterations, if the data is linearly separable.

$$R = \max_i \|x_i\|$$

$\gamma$  = "margin"  $\Rightarrow \|\hat{w}\|^2 = 1, y_i w^T x_i \geq \gamma \quad \forall i$   
measurement of how "separable" the data is.



- it will not converge if the data is not linearly separable.



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- many possible solutions, depending on the initialization.



## Logistic Regression (probabilistic approach)

- Consider 2-class binary problem:  $y \in \{0, 1\}$

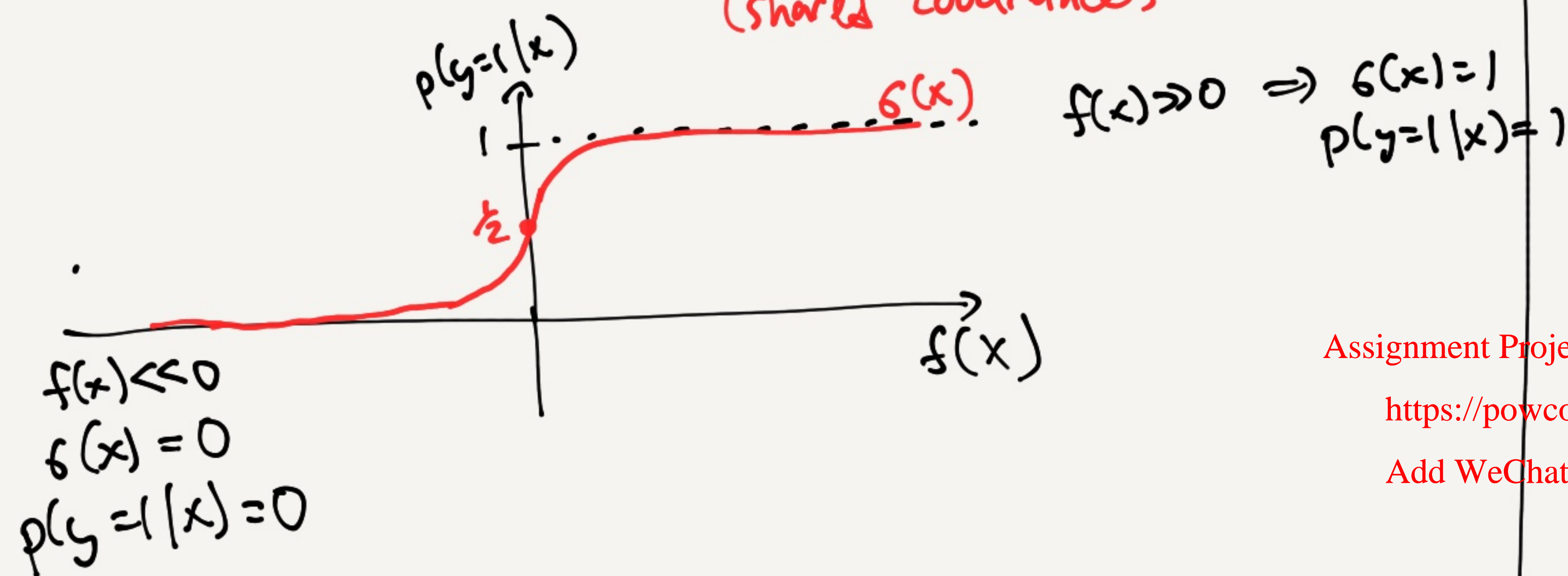
- From PS 6-7, if the CCDs are Gaussian

$$p(x|y) = N(x|\mu_y, \Sigma_y)$$

then the posterior distribution  $p(y|x)$  is a sigmoid func.

$$p(y=1|x) = \frac{1}{1 + e^{-f(x)}} = \sigma(x)$$

$f(x)$  is linear function if  $\Sigma_y = \Sigma$   
(shared covariance)



• with BDR,  $f(x)$  is determined by the learned Gaussian parameters.

• Now, learn  $f(x)$  directly.

Assume  $f(x)$  is linear

$$f(x) = w^T x, \quad p(y=1|x) = \sigma(w^T x) = \pi$$

Decision Rule:

$$\hat{y} = \begin{cases} 1, & p(y=1|x) > \frac{1}{2} \\ 0, & \text{otherwise} \end{cases} = \begin{cases} 1, & w^T x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

## Learning (2-class problem)

$$\text{Dataset } D = \{(x_i, y_i)\}_{i=1}^N$$

$$\text{let } \pi_i = \sigma(w^T x_i) \quad (\text{probability of class 1 given } x_i)$$

$$p(y_i|x_i, w) = \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

BDR = Note we model  $P(x|y)$

conditional prob. given  $x_i$  &  $w$

log-likelihood of data

$$\ell(w) = \sum_{i=1}^N \log p(y_i|x_i, w)$$

$$= \sum_i [y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i)]$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

MLE

$$w^* = \underset{w}{\operatorname{argmax}} \ell(w)$$

Error (loss) function

$$w^* = \underset{w}{\operatorname{argmin}} -\ell(w)$$

$$= \underset{w}{\operatorname{argmin}} \sum_i \underbrace{-y_i \log \pi_i - (1 - y_i) \log (1 - \pi_i)}_{\text{cross-entropy loss}}$$



• maximize  $l(w)$

Apply Newton-Raphson method: (more in tutorial)

Iterate

$$w^{(new)} = w^{(old)} - \underbrace{[\nabla^2 l(w)]^{-1}}_{\text{Hessian}} \underbrace{[\nabla l(w)]}_{\text{gradient}}$$

$$w^{(new)} = (XRX^T)^{-1} X R z$$

← weighted least-squares using  $R, z, X$

where

$$\begin{cases} R = \text{diag}(\pi_1(1-\pi_1), \dots, \pi_N(1-\pi_N)) \\ z = X^T w^{(old)} - R^{-1}(\pi - y) \end{cases}$$

← depend current  $w$  & predicted  $\pi_i$ . higher weight to  $\pi_i = \frac{1}{2}$

← target depends on  $w$

←  $f^{old}(x)$  weight on each point error b/w prediction & class label

"iterative reweighted least squares"

IRWLS / IRLS

## Comparison of loss/error functions

All these methods optimize a similar form:

$$w^* = \underset{w}{\text{argmin}} E(w) = \underset{w}{\text{argmin}} \sum_i \underbrace{L(f(x_i), y_i)}_{\text{loss function}}$$

"empirical risk"

"empirical risk minimization" - all about training error.

Let  $z = y w^T x$

Ideal 0-1:  $L = \begin{cases} 0, & z > 0 \\ 1, & z < 0 \end{cases}$

LSC:  $L = (z-1)^2$

Perceptron:  $L = \begin{cases} 0, & z > 0 \\ -z, & z < 0 \end{cases} = \max(0, -z)$

Log. Regr:  $L = \log(1 + e^{-z}) \cdot \frac{1}{\log 2}$

