

## Lecture 9: Discriminative Learning - Linear classifier

Generative model: 1) learn CCD from data  $p(x|y)$

2) BDR to get classifier

$$P_{\text{Lg}}(x) = \frac{p(x|y)p(y)}{e(x)}$$

Note: data is used only in step (1) to get the CCDs  
classifier is secondary.

Density estimator is an ill-posed problem

Gaussian, GMM, Gamma.

Vapnik advice: "when solving a given problem, avoid solving a more general problem as an intermediate step"

Discriminative solution: solve the decision boundary or  $p(y|x)$  directly

"use the data to learn to discriminative classes,  
rather than generatively data."

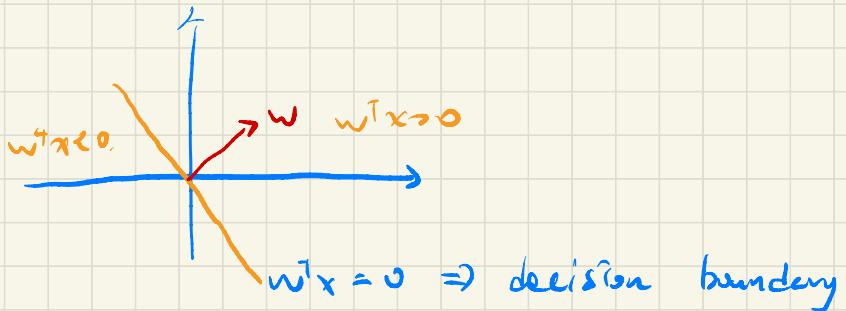
### Linear classifier

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

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

Linear function:  $f(x) = w^T x$

$w$  separates the space into 2 half-spaces.



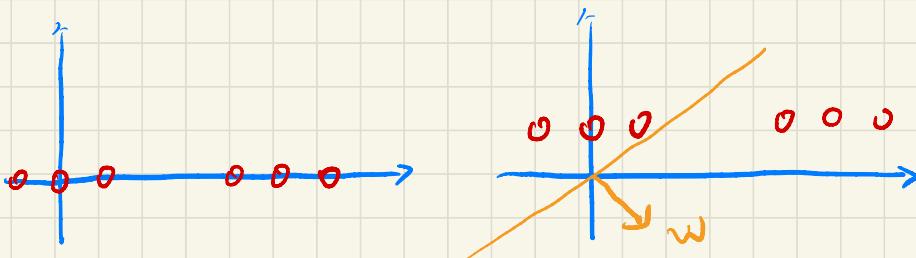
\*:  $w$  points into the positive half. (?)

Decision Rule:

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

Note: bias can be included as another dimension

$$\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad \tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix} \quad f(\tilde{x}) = \tilde{w}^T \tilde{x} = w^T x + b$$



Training set:  $D = \{(x_i, y_i)\}_{i=1}^n$

$$D = \{x, y\}, x = [x_1, \dots, x_n]$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Note:

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

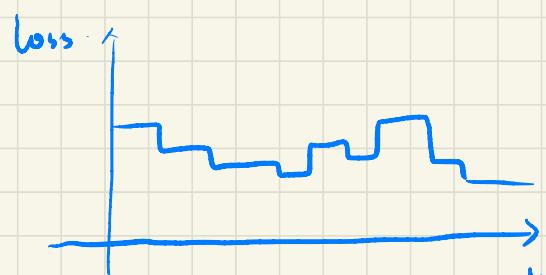
$y_i w^T x_i \leq 0 \Rightarrow$  misclassified  $x_i$

Ideal Case: 0-1 loss function:

Optimize the number of misclassified samples.

$$w^* = \arg \min_w \sum_{i=1}^n \begin{cases} 0, & y_i w^T x_i > 0 \\ 1, & y_i w^T x_i \leq 0 \end{cases}$$

0-1 loss function.



★: the gradient is 0  
or undefined. So it's  
difficult to optimize.

Least square classification (Label regression)

- Ignore the fact that  $y$  is discrete and just apply

## LS regression

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2 = \arg \min_w \|y - X^T w\|^2$$

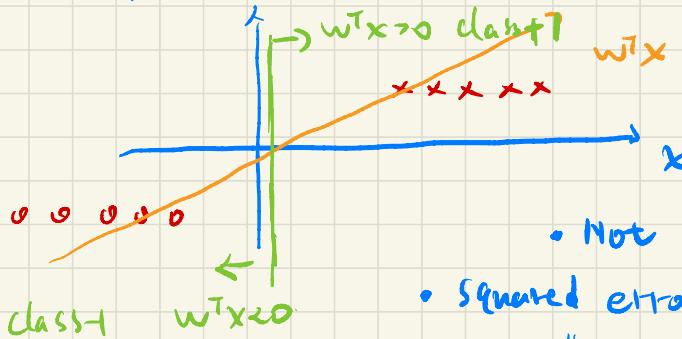
$$\|y - X^T w\|^2 = (y - X^T w)^T (y - X^T w)$$

$$= y^T y - y^T X^T w - (X^T w)^T y + w^T X^T X^T w$$

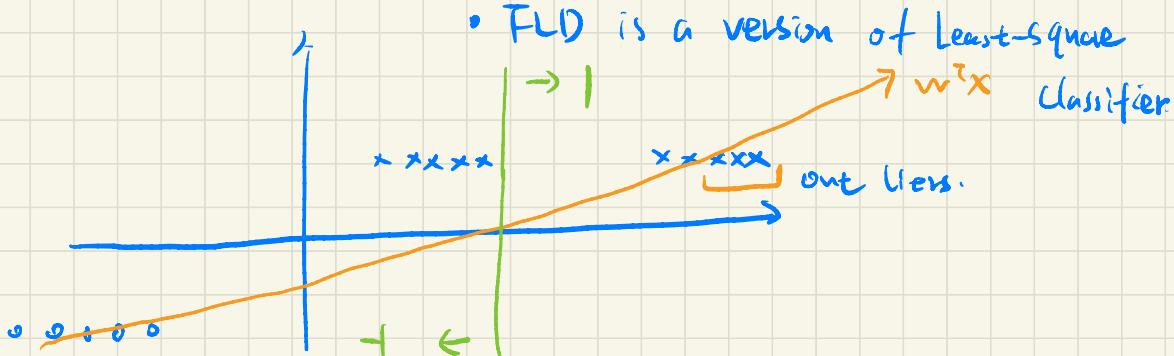
$$= y^T y - 2 y^T X^T w + w^T X^T X^T w$$

$$\frac{\partial}{\partial w} = -2 X y + 2 X^T w = 0 \Rightarrow w = (X^T X)^{-1} X^T y$$

Example:



- Not robust to outliers
- Squared error penalizes that are "two corrected"



# Perceptron

Perceptron criteria - only look at misclassified points.

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

higher loss for  $x_i$  that are  
badly misclassified i.e.  
 $y_i w^T x_i << 0$

M: misclassified points

$$E(w) = 0 \text{ when data correctly classified}$$

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 60's ...
- Apply "stochastic gradient descent" (SGD)
  - use one data point at a time :

$$w^{t+1} = w^t + \eta y_i x_i, \text{ for some } i \in M$$

$\eta$ : learning rate.

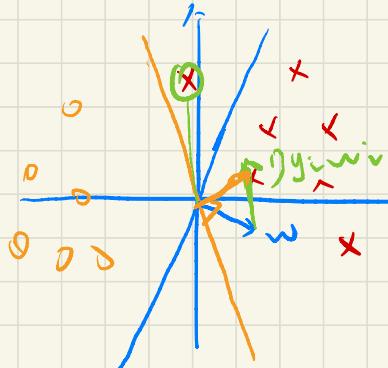
Rotates  $w$  to point towards

1) the positive class.

2)  $w$  gets longer as we iterate

$\Rightarrow$  each sample has diminishing effect

example:

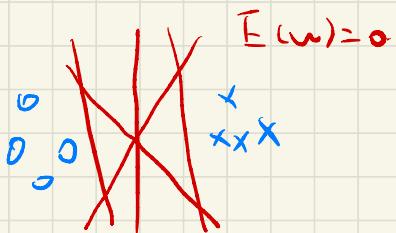


- Rosenblatt proved that perceptron algorithm converges in  $(\frac{R}{\gamma})^2$  iterations if the data is linearly separable

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

$\gamma$  = "margin" =  $\|\hat{w}\|^2 - |y_i \hat{w}^T x_i|$   
how separate is the data

- will not converge if the data is not linearly separable
- many possible solutions with  $\text{loss} = 0$ , depends on initialization



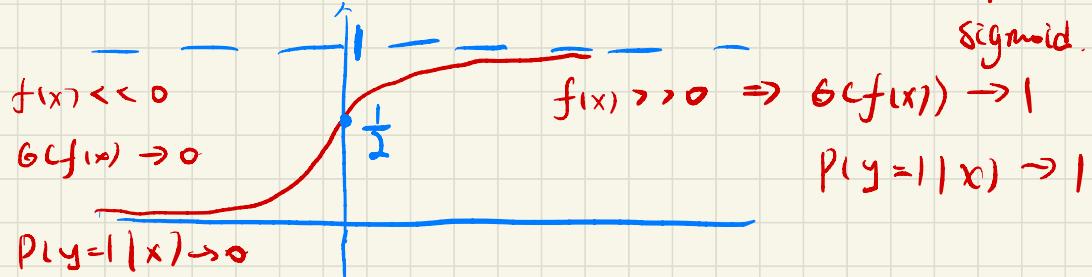
- Note: learning rate does not matter.

## Logistic Regression (probabilistic approach)

- Consider two class problem:  $y \in \{0, 1\}$
- Ps 6-7: fit CDF of Gaussians  $\Rightarrow$  Posterior  $P(y|x)$  is a sigmoid function.

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

↑  
sigmoid.



- When CCD Gaussians have same covariance  $\Rightarrow$   
 $f(x)$  is a linear function.
- with BDR,  $f(x)$  is determined from the  
CCD parameters.  
 $\Rightarrow$  now, learn the parameter of  $f(x) = w^T x$   
directly

Set up:

$$f(x) = w^T x \quad P(y=1|x) = \frac{1}{1+e^{-w^T x}} = g(w^T x) = \tau$$

Decision rule:

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

Look at the number of parameters

LR:  $w \in \mathbb{R}^d \rightarrow d$  parameters

BDR:  $w \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \rightarrow O(\frac{d^2}{2})$  parameters (?)

much better for high dimensional data (need less data compared to BDR)

Learning:  $D = \{(x_i, y_i)\}_{i=1}^n$

$$\text{let } \pi_i = g(w^T x_i) = P(y_i=1 | x_i)$$

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

log-like likelihood of data:

$$\begin{aligned} L(w) &= \sum_i \log p(y_i | x_i, w) \\ &= \sum_i y_i \log \pi_i + (1-y_i) \log (1-\pi_i) \end{aligned}$$

MLE:

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

Alternatively:

$$w^* = \underset{w}{\operatorname{arg\,min}} -L(w)$$

$$= \underset{w}{\operatorname{arg\,min}} \sum_i -y_i \log z_i w - (1-y_i) \log (1-z_i)$$

cross-entropy loss

- Minimize  $L(w)$
- Apply Newton-Raphson method.

$$w^{\text{new}} = w^{\text{old}} - [\nabla^2 L(w)]^{-1} \nabla L(w)$$

: Hessian      Gradient  
 { iteration

$$w^{\text{new}} = (X R^T X)^{-1} X R z \leftarrow \begin{array}{l} \text{weighted least} \\ \text{squares with } R, z, X \end{array}$$

where: iterative reweighted least squares

(IRWLs) (IRLS)

$R = \text{diag}(\pi_1(1-\pi_1), \dots, \pi_n(1-\pi_n)) \leftarrow \begin{array}{l} \text{weighted depend} \\ \text{on current } w \end{array}$

$z = \underbrace{X^T w^{\text{old}}}_{\text{current prediction}} - R^T \underbrace{(\pi - y)}_{\text{error}} \leftarrow \begin{array}{l} \text{target depends on} \\ \text{current } w \end{array}$

## Comparison of loss (error) functions

All of these methods are of the form:

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n L(f(x_i), y_i)$$

loss function  
empirical risk

"empirical risk minimization" → all about training error

Let  $z_i = y_i w^\top x_i \Rightarrow \begin{cases} z_i > 0 \rightarrow \text{classify correctly} \\ z_i < 0 \rightarrow \text{misclassify} \end{cases}$

Loss function:

Ideal 0-1 loss:  $L(z) = \begin{cases} 0, & z > 0 \\ 1, & z \leq 0 \end{cases}$

LS classifier:  $L(z) = (z - 1)^2$

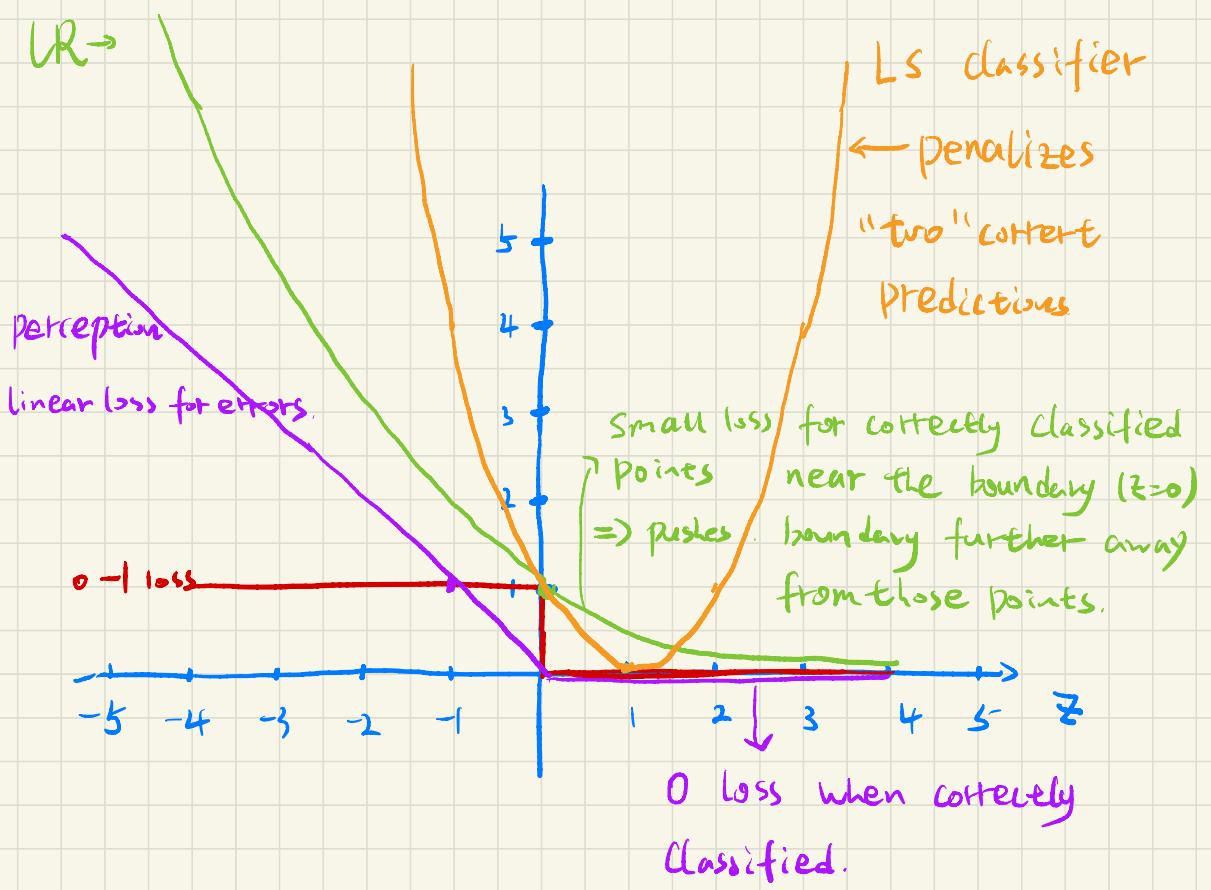
Perception:

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

logistic regression:

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

$$z=0, L(z)=1$$



LSC and LR loss are convex approx to the ideal 0-1 loss. (7)