

## Bayes classifier:

1. Prior probabilities

of class 0 and class 1

$$P(Y=0) \xrightarrow{P_0}, P(Y=1) \xrightarrow{P_1}$$

2. Likelihood / Class conditional

densities

$$f_0(x) \quad f_1(x)$$

$$P(X|Y=0), P(X|Y=1).$$

Class conditional densities.

Distribution of features

Given a class.

3. Loss function.

$$L(Y, h(x))$$



$$L(Y, h(x)) = 1 \text{ if } Y \neq h(x)$$

$$= 0 \text{ otherwise.}$$

(Zero - one Loss functions)

Perfect information case

Bayes Classifier

$$P(Y=1|x), P(Y=0|x)$$

$$q_{Y_1}(x)$$

$$q_{Y_0}(x)$$

$$q_{Y_1}(x) = \frac{P_1 f_1(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

$$q_{Y_0}(x) = \frac{P_0 f_0(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

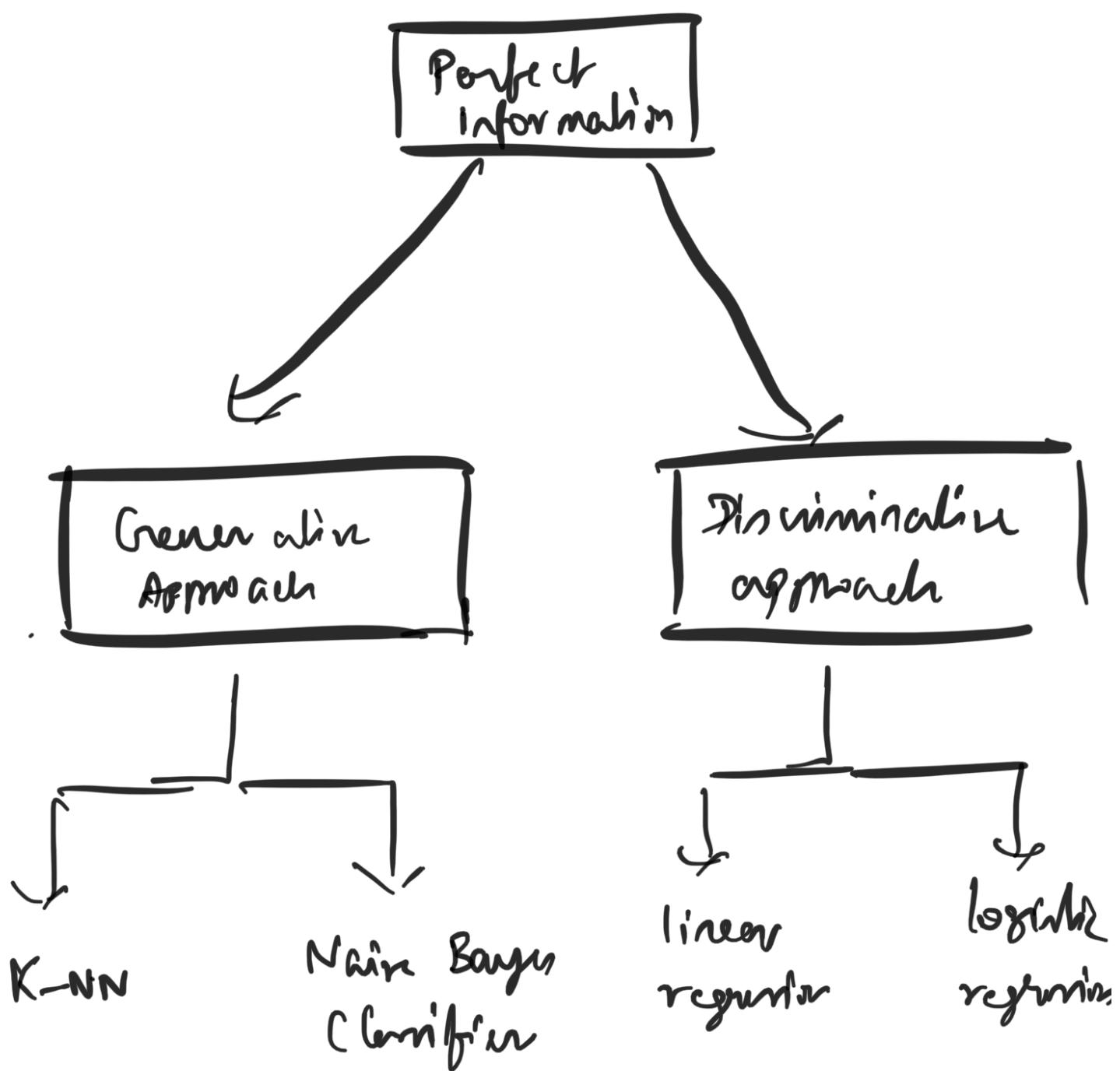
$$\text{if } q_{Y_1}(x) > q_{Y_0}(x)$$

$$h_B(x) = 1$$

else  $q_1(n) \leq q_0(n)$

$$h_B(n) = 0$$

---



## Generative Approach

1. Estimate prior probabilities from data
  2. Estimate class conditional densities from data.
  3. Using these estimated quantities, if we implement a Bayes classifier, then we call the approach Generative.
- 

## Discriminative approach.

---

Directly approaching posterior from the data rather than estimating conditional

Prior and class conditional  
densities

$$R(h) = \underset{x, y}{E} [L(y, h(x))]$$

→ Risk of a classifier.

$$h^* = \min_h \underset{x, y}{E} [L(y, h(x))]$$

A surrogate function to be optimized

$$h^* = \min_h$$

$$\frac{\sum_{i=1}^n L(y_i, h(x_i))}{n}$$

↑  
Empirical risk minimization.

Loss function

Data

...  $x_1$  ...  $x_n$

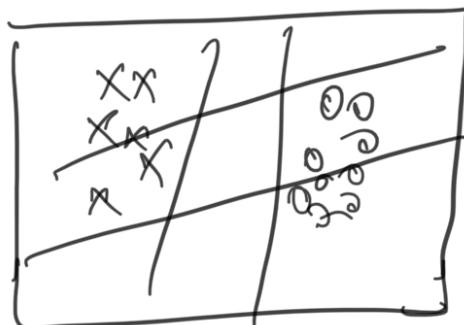
$h \rightarrow$  5ft thin adult  
thin kid

5ft Adult

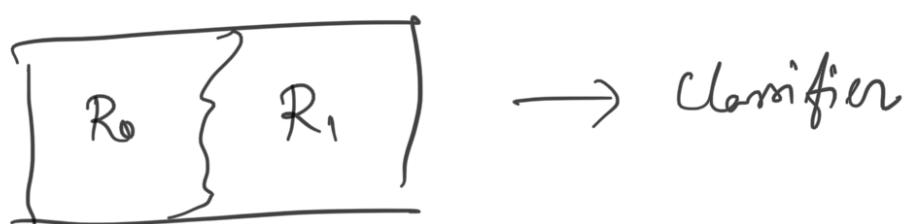
<u>5 ft</u>	man	...	
6 ft	Adult	6 ft	Adult ✓
<u>4.5 ft</u>	Kid	4.5 ft	Kid ✓
<u>4.5 ft</u>	Adult	4.5 ft	Kid ✗

$H = \{ w^T x : \text{Linear in the features} \}$

$$w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$



2 feature  
space.

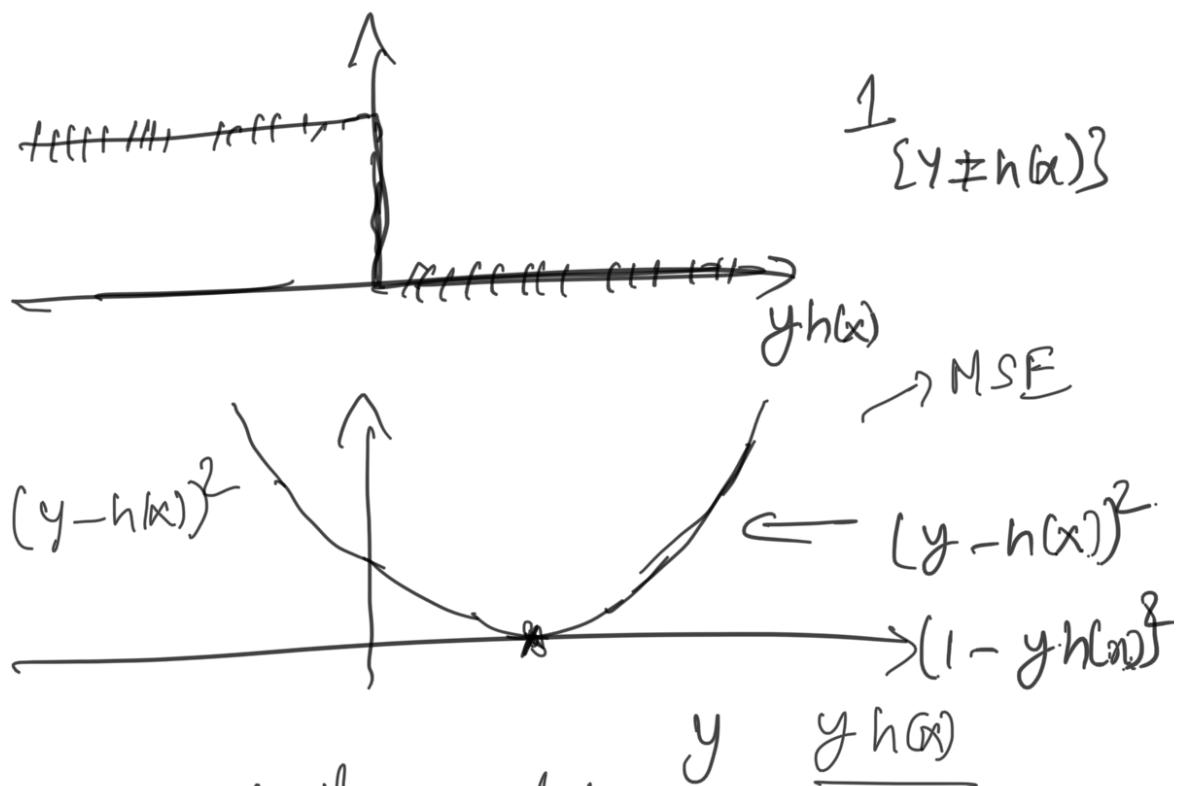


1.  $E \rightarrow$  arrange over the data given to  $w$ .

2. Instead of considering all possible classifiers, we will consider a

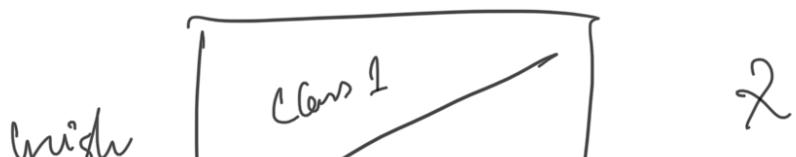
restricted hypothesis space

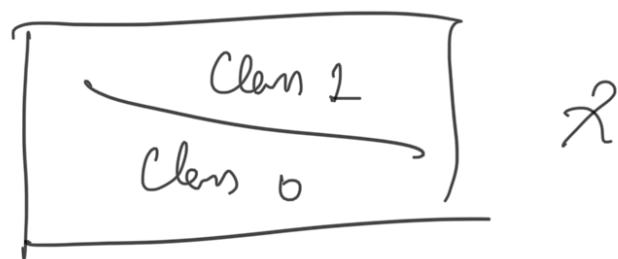
### 3. Nature of Loss function



$h(x) \rightarrow \text{Classifier Prediction}$

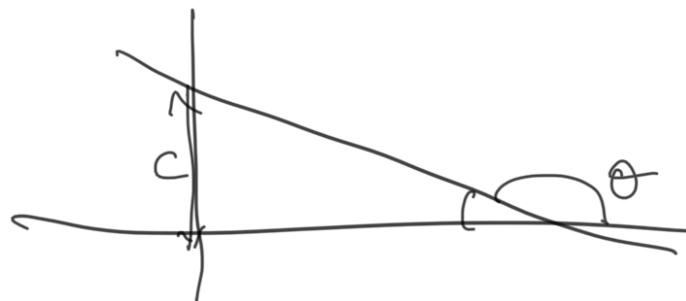
$y \rightarrow \text{True Label}$





$$\{ w_1 x_1 + w_2 x_2 + \dots + w_d x_d \}$$

$$\{ m x + c \}$$



$$h : \chi \rightarrow \mathbb{R}$$

$$h(x) \rightarrow \text{Score}$$

$$h(x) > 0 \quad \text{Class } 1$$

$$h(x) < 0 \quad \text{Class } -1$$

$$h(x) > \frac{1}{2} \quad \text{Class } 1$$

$$h(x) < \frac{1}{2} \quad \text{Class } 0$$

Objektiv:

$$(x_i; y_i)_{i=1}^n \quad y_i = +1 \\ = -1$$

$$\min_w \sum_{i=1}^n (\underline{w^T x_i} - y_i)^2$$

n

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \end{pmatrix}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

$$(x_n^T / y_n)$$

and also

$$Xw = \begin{pmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{pmatrix}$$

$$\frac{1}{n} \|Xw - Y\|^2 \Leftrightarrow \frac{\sum_{i=1}^n (w^T x_i - y_i)^2}{n}$$

$$f(w) = \frac{1}{n} \|Xw - Y\|^2$$

↓

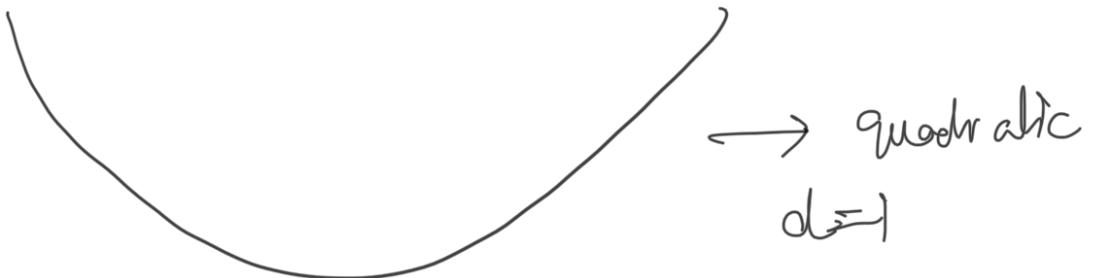
$$= \frac{1}{n} \underbrace{(Xw - Y)^T (Xw - Y)}_{\frac{1}{n} (Xw - Y)^T (Xw - Y)}$$

$$= \frac{1}{n} (w^T X^T - y^T) (Xw - Y)$$

$$= \frac{1}{n} (w^T X^T w - \underbrace{y^T X w}_{+ y^T y} - \underbrace{w^T X^T y}_{+ y^T y})$$

$$= \frac{1}{n} \left( w^T \underbrace{(x^T x)}_{d \times d} w - 2(y^T x)w + \underbrace{y^T y}_{\text{constant}} \right)$$

Concave minima.



$\rightarrow \mathbb{R}$

$$\underline{f(w)} = \underline{w^T A w} + b^T w + c$$

$$\nabla f(w) = 0$$

$$A = \underbrace{x^T x}_{d \times d}, \quad c = y^T y$$

$$b = -2y^T x$$

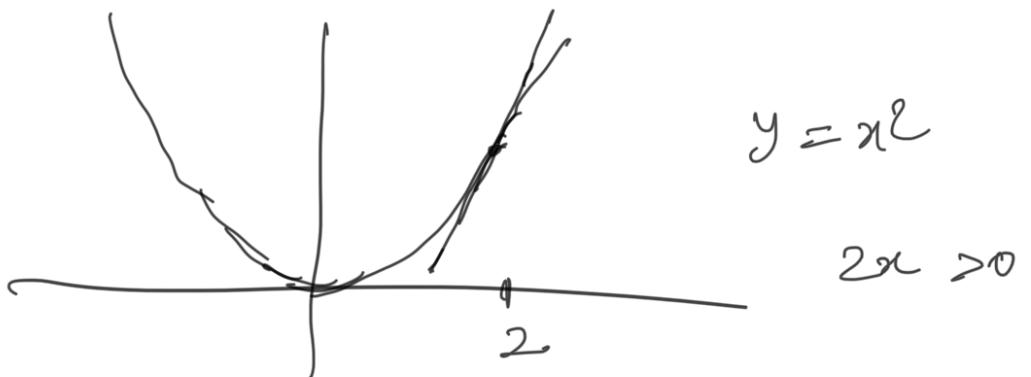
$$2Aw + b = 0$$

$$w = \frac{1}{2} A^{-1} b$$

$$\boxed{w^* = \underbrace{(x^T x)^{-1}}_{d \times d} y^T x}$$

We found the minima of the loss function

## Gradient descent Algorithm:

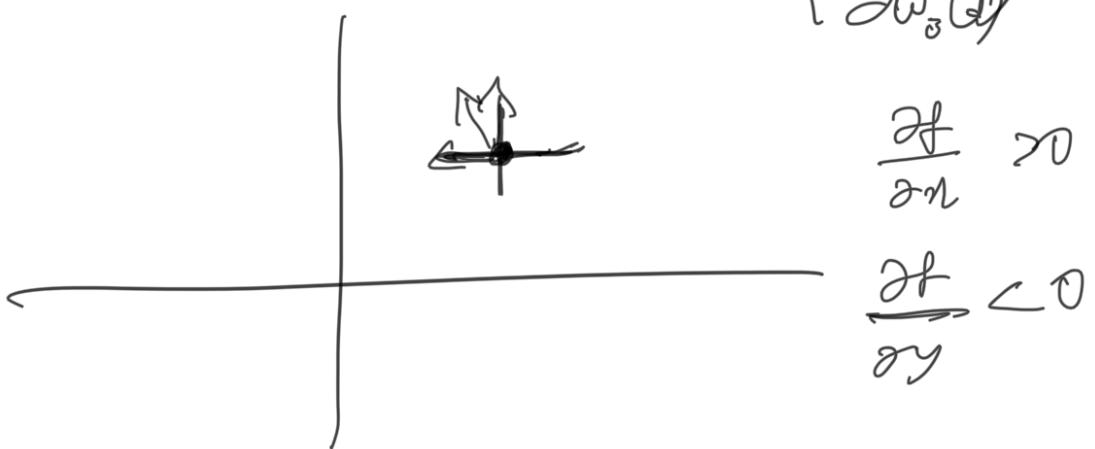


$$w_0 \rightarrow w_1$$

$$w_1 = w_0 - \alpha \left( \frac{dl}{dw} \right)$$

slope

$$\begin{pmatrix} w_1(1) \\ w_1(2) \\ \vdots \\ w_1(d) \end{pmatrix} = \begin{pmatrix} w_0(1) \\ w_0(2) \\ \vdots \\ w_0(d) \end{pmatrix} - \alpha \begin{pmatrix} \frac{\partial f}{\partial w_0(1)} \\ \frac{\partial f}{\partial w_0(2)} \\ \vdots \\ \frac{\partial f}{\partial w_0(d)} \end{pmatrix}$$



$$f(w) = \frac{\|x^T w - y\|^2}{2} + \frac{w^T A w + b^T w}{2} + c$$

$$\nabla f(w) = 2Aw + b$$

$$w_{t+1} = w_t - \alpha (2Aw_t + b)$$

$$w_t \rightarrow w^*$$

Stochastic gradient descent:

$$f(w) = \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n L_i$$

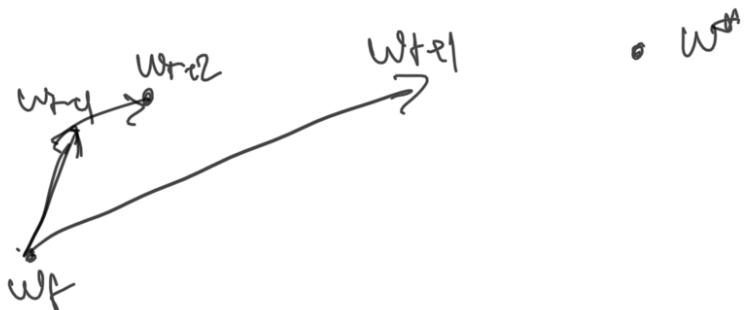
$$L_i = (w^T x_i - y_i)^2$$

$$w_{t+1} = w_t - \alpha_t \left[ (w_t^T x_i - y_i) x_i \right]$$

$$w_t \rightarrow w^* \quad \text{error}_i \text{ (scalar)}$$

$$w_{t+1} = w_t - \alpha_f \left( \frac{w_t^T x_i}{\text{Predictor}} - y_i \right) x_i$$

Predictor      Target




---

Linear regression based classification

---

1. Loss function  $\rightarrow$  MSE

2.  $f(w) = \sum_{i=1}^n (w^T x_i - y_i)^2$

Algorithm for finding best  $w^*$

---

1.  $\nabla f(w^*) = 0$

$$w^* = \underbrace{(X^T X)^{-1}}_{n \times n} \underbrace{X^T Y}_{n \times 1}$$

$$d\alpha d \beta \cup (\alpha \beta)$$

2. Incremental way:

$$w_{t+1} = w_t - \alpha_t (A w_t + b).$$

3. Stochastic gradient descent

$$w_{t+1} = w_t - \alpha (w_t^T x_i - y_i) x_i$$

4. Mini batch gradient descent: