

Bayes Classifier

2 - class problems.

Ex: D - Default
 R - Return

Feature:

Ex: H - High risk loaners
 L - Low risk loaners.

\mathcal{X} - Feature space.

$x \in \mathcal{X}$ is called a feature.



Prior probabilities?

$P(C_0) \rightarrow$ probability of class 0

$P(C_1) \rightarrow$ probability of class 1

$Y = 0 \rightarrow$ class 0

$$Y = 1 \longrightarrow \text{class 1.}$$

$$P(Y = 0) , \quad P(Y = 1)$$

$$\downarrow$$

$$P_0$$

$$\downarrow$$

$$P_1$$

Likelihood or the class - conditional densities:

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

$$\downarrow$$

$$f_0(x)$$

$$\downarrow$$

$$f_1(x)$$

This gives distribution of the feature
conditioned on a class.

Ex: $P(H | D) , P(L | D) \rightarrow \text{class 0}$

$$P(H | R) , P(L | R) \rightarrow \text{class 1}$$

Posterior probabilities:

$$P(Y = 1 | X = x)$$

$$P(Y = 0 | X = x)$$

$$\downarrow$$

$$q_1(x)$$

$$\downarrow$$

$$q_0(x)$$

Ex: $P(D|H), P(R|H)$
 $P(D|L), P(R|L)$ } Can be computed using Bayes rule.

$$q_0(x) = \frac{P_0 f_0(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

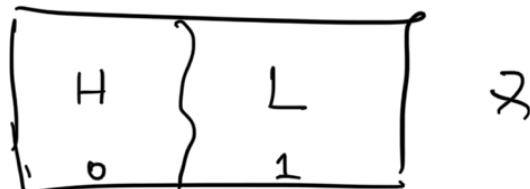
$$q_1(x) = \frac{P_1 f_1(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

Bayes classifier rule:

$$h_B(x) = \begin{cases} 1 & \text{if } q_1(x) > q_0(x) \\ 0 & \text{if } q_1(x) < q_0(x) \end{cases}$$



Ex:



Bayes Rule:

$$\text{Posterior} = \frac{\boxed{\text{likelihood}}}{\text{evidence}} \propto \text{Prior}$$

likelihood — How likely is the feature given a class.

Prior — what is the prior probabilities for each class

Posterior — what is the posterior probability of a class given a feature value

evidence — How likely is the feature.

$$\text{err}(h_B) = \frac{P(Y=1, h_B(x)=0) + P(Y=0, h_B(x)=1)}{1}$$

$$= P(Y=1) P(h_B(x)=0 | Y=1) + P(Y=0) P(h_B(x)=1 | Y=0)$$

$$= P(Y=1) P(x \in R_0 | Y=1)$$

$$+ P(Y=0) P(x \in R_1 | Y=0)$$

What is a classifier?

A classifier is a rule or a mapping from feature space $\mathcal{X} \longrightarrow \mathcal{Y} = \{0, 1\}$.



$$\begin{aligned} \text{err}(h) = & P(h(x) = 0, Y=1) \\ & + P(h(x) = 1, Y=0) \end{aligned}$$

h_B is the optimal classifier because it has the least error.

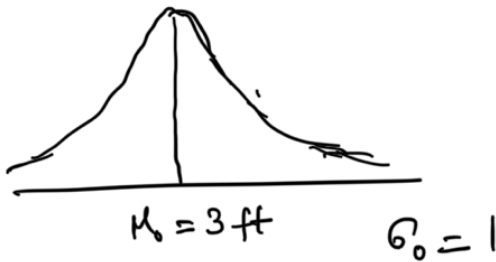
Height of the population is your feature

Adult }
child } Two classes



$$P(\text{Adult}) = 0.5$$

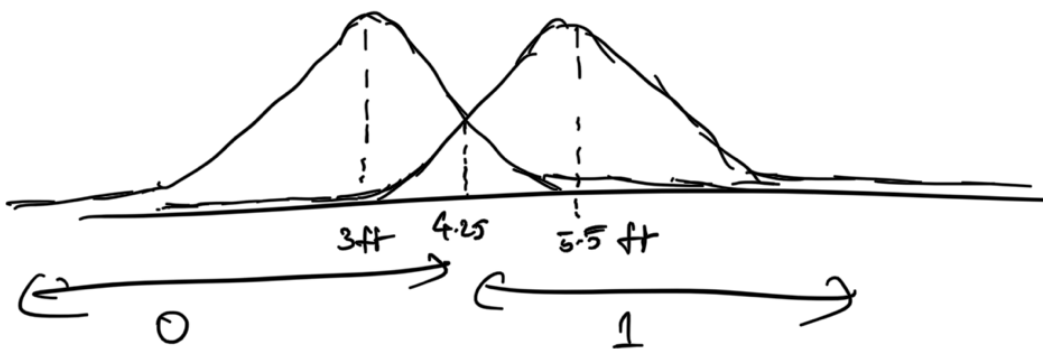
$$\mu_1 = 5.5 \text{ ft} \quad \sigma_1 = 1$$



$$P(\text{child}) = 0.5$$

$$f_0(x) = \frac{1}{\sigma_0 \sqrt{2\pi}} \exp\left(-\frac{(x - \mu_0)^2}{2\sigma_0^2}\right)$$

$$f_1(x) = \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{(x - \mu_1)^2}{2\sigma_1^2}\right)$$



$$P(Y=1|h) \rightarrow q_1(h)$$

$$P(Y=0|h) \rightarrow q_0(h)$$

$$\frac{q_1(x)}{q_0(x)} > 1 \Rightarrow \text{class Adult}$$

$$\frac{q_1(x)}{q_0(x)} < 1 \Rightarrow \text{class child}$$

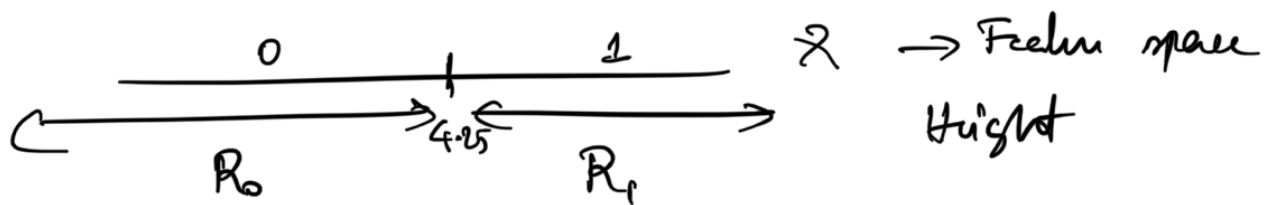
$$q_1(x) = \frac{P_1 f_1(x)}{P_1 f_1(x) + P_0 f_0(x)}$$

$$P_0 f_0(x) + P_1 f_1(x)$$

$$q_0(x) = \frac{P_0 f_0(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

$$\frac{q_1(x)}{q_0(x)} = \frac{P_1 f_1(x)}{P_0 f_0(x)} = \frac{f_1(x)}{f_0(x)} > 1$$

\Rightarrow Claim 2



$$R_0 = \{x : x < 4.25\}$$

$$R_1 = \{x : x \geq 4.25\}$$

$$h_B(x) = 0 \quad \text{if } x \in R_0$$

$$= 1 \quad \text{if } x \in R_1$$

When can max error weights matter?

Loss function:

$$L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

↓
True

↓
Prediction

$L(0, 0) \rightarrow$ Loss incurred when classifier says 0 class and the true class is 0

$L(0, 1) \rightarrow$ Loss incurred when classifier says 1 class and the true class is 0.

$L(1, 0) \rightarrow$ "

$L(1, 1) \rightarrow$ "

Zero/one loss function

$$L(0, 0) = 0$$

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

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

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

$E_{x,y} [L(y, h(x))] \rightarrow$ Expected loss or Risk of a

clerical

$$R(h) = \mathbb{E}_{\underline{x, y}} [L(y, h(x))].$$

$$= L(0, 0) P(y=0, x \in R_0)$$

$$+ L(0, 1) P(y=0, x \in R_1)$$

$$+ L(1, 0) P(y=1, x \in R_0)$$

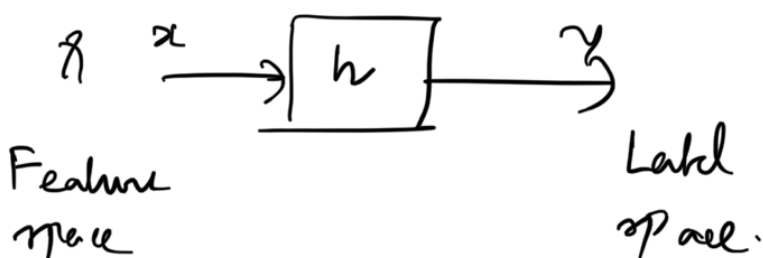
$$+ L(1, 1) P(y=1, x \in R_1)$$

Enron aircraft problem;

$$\underline{L}(0, 1) = 1$$

$$\underline{L}(1, 0) = 10$$

what is a classifier?



How do we measure the performance of a classifier?

$$R(h) = E_{x,y} [L(y, h(x))]$$

h_B^* = Bayes classifier is the one that minimizes the risk.

Perfect information

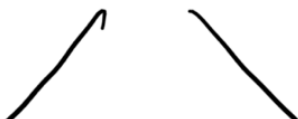
1. prior probabilities
2. class conditional densities
3. Loss function \rightarrow Zero/one loss function

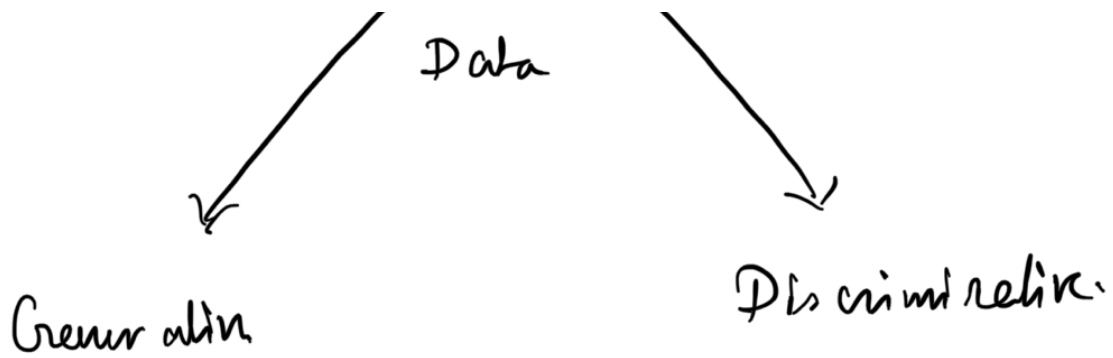


Bayes Classifier



optimal.





1. Naïve Bayes

2. K - Nearest Neighbors

1. Logistic regression

2. SVM (Support vector machines),

3. Linear discriminant analysis

Data: $(x_i, y_i)_{i=1}^n \longrightarrow \text{Dataset}$

Ex! $n = 6$

1. $x_1 = 7, y = 1$

2. $x_2 = 3, y = 0$

3. $x_3 = 3.5, y = 0$

4. $x_4 = 5, y = 0$

5. $x_5 = 6, y = 1$

6. $X_6 = 2, Y = 0$

What is generative approach for classifier design?

1. Estimate prior probabilities from data
2. Estimate class conditional densities from data

k - Nearest neighbour Classifier:

$x \rightarrow$ feature

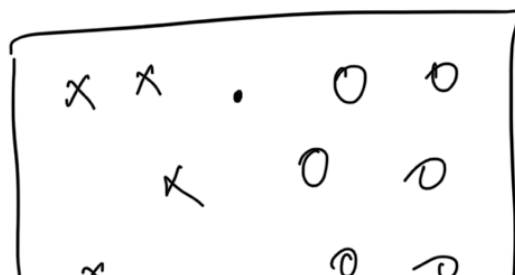
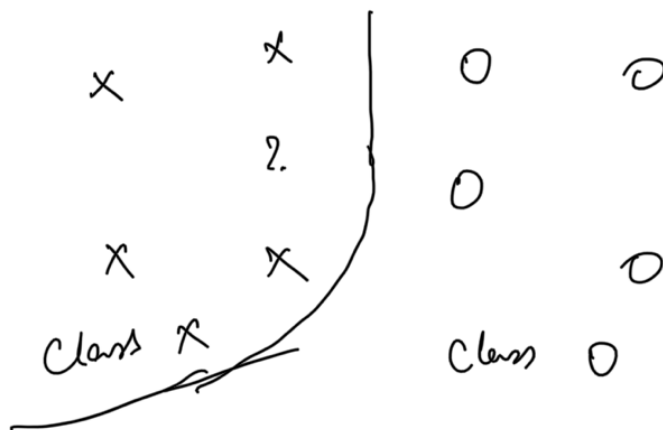
What is the label corresponding to this?

1. Given a feature x find the k closest prototype to this example, or this feature
2. The class label will then correspond

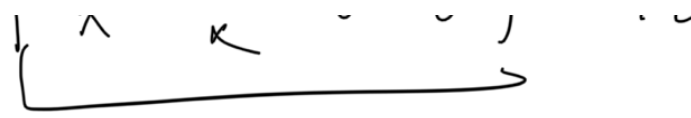
to majority of the labels in
the K prototypes.

Two important components:

1. distance metric.
2. How many prototypes to choose?



2



Voronoi regions
or voronoi diagram

1. Distance metric
→ Euclidean metric

$$x = (x_1, x_2, \dots, x_d)$$

$$y = (y_1, y_2, \dots, y_d)$$

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_d - y_d)^2}$$

2. $K \rightarrow$ odd number

$$K = 3, 5 \text{ etc.}$$

$K = \text{nn}$ Classifier error cannot be more
than $2 \times$ Bayes error.

Ex:

Height, weight

(7. ft, 10 kg)

$x = (x_1, x_2)$

... ..

$$y = (y_1, y_2)$$

$$d(x, y) = \sqrt{\frac{(x_1 - y_1)^2}{\sigma_1^2}} + \sqrt{\frac{(x_2 - y_2)^2}{\sigma_2^2}}$$

1.	(0.1, 100)	(-0.3, -3)
2.	(0.2, 101)	(-0.2, -2)
3.	(0.3, 104)	(-0.1, 1)
4.	(0.7, 99)	(0.3, -4)
5.	(0.6, 105)	(0.2, 2)

$$(0.15, 102)$$

$$\sigma_i^2 = \frac{\sum_{i=1}^n (x_i - \mu_i)^2}{n} \quad \mu_i = \frac{\sum_{i=1}^n x_i}{n}$$

$$h(x) \rightarrow \{0, 1\}$$

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

$$x \in \mathbb{R}$$

$$n : x \rightarrow \infty$$

$$+ \quad h(n) > 0$$

$$- \quad h(n) < 0$$

$$R(h) = E_{x,y} [L(y, \underline{h(x)})]$$

Example:

$$H \rightarrow D$$

$$1. \quad L \rightarrow D$$

$$2. \quad \begin{array}{l} H \rightarrow D \\ L \rightarrow R \end{array} \quad \left. \vphantom{\begin{array}{l} H \\ L \end{array}} \right\}$$

$$3. \quad \begin{array}{l} H \rightarrow R \\ L \rightarrow D \end{array} \quad \left. \vphantom{\begin{array}{l} H \\ L \end{array}} \right\}$$

$$4. \quad \begin{array}{l} H \rightarrow R \\ L \rightarrow R \end{array}$$

$$\boxed{\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} x}$$

$$\mathcal{H} = \{ \text{Hypothesis space} \}$$

$$\mathcal{H} = \{ \underline{w^t x}, w \in \mathbb{R}^d \}$$

$$w^t x = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_d x_d$$

Discriminative Models!

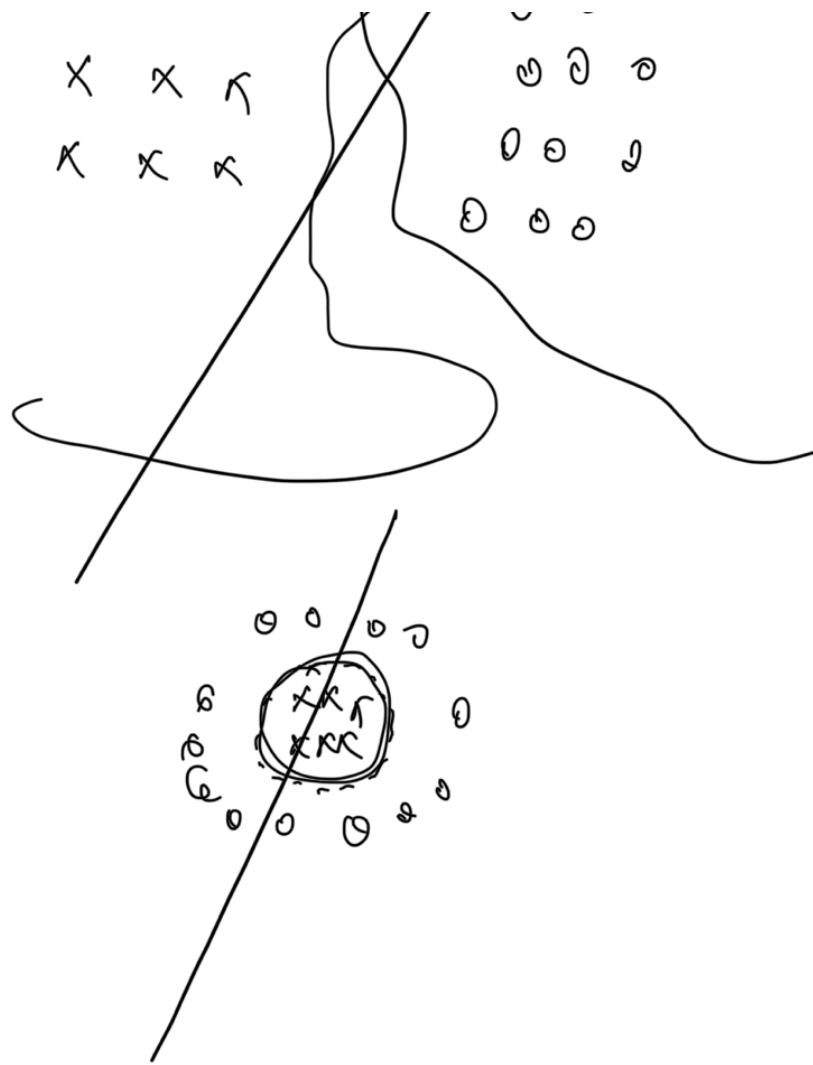
$$1. \quad h : \mathcal{X} \rightarrow \mathbb{R}$$

2. we are going to restrict the space of hypothesis from all possible mapping from $\mathcal{X} \rightarrow \mathbb{R}$ to a restricted class

$$\{ w^t x \} \quad w \in \mathbb{R}^d$$

$$x \times x \times \bigcup \bigcup$$

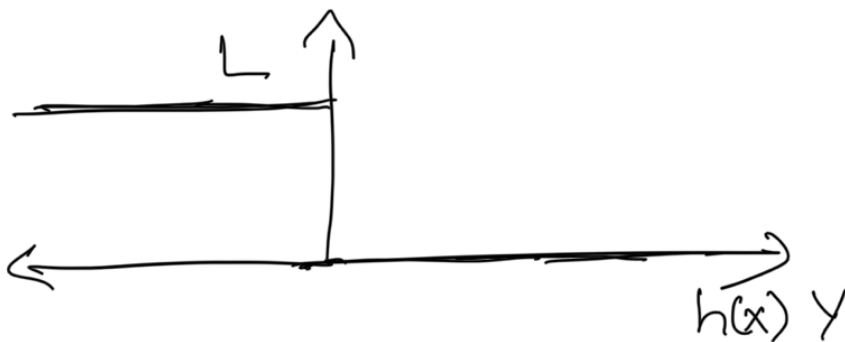
- - ?



$$\mathcal{X} = \mathbb{R}^2$$

$$\mathbb{R} \times \mathbb{R}$$

$$R(h) = \mathbb{E}_{x, y} (L(h(x), y))$$



→ Squared
Loss
function

