

## Classification

Each unit in a pop is represented by:  $(\underline{x}^1, \dots, \underline{x}^p)$   $\xrightarrow{\text{Labels } \in \{1, \dots, g\}}$   
 vector of features  $\underline{x} = (x_1, \dots, x_p) \in X$  (e.g.  $R^p$ )

## Definition (classification)

Learning:  $\delta: X \rightarrow \{1, 2, \dots, g\}$

$\delta(\underline{x}) = i \iff$  the unit for which  $\underline{x}$  is observed belongs to group  $i$ .

## Supervised Learning (discriminant analysis)

Training data:

	$x_1$	...	$x_p$	$L$	Goals:
sample 1	$x_{11}$	...	$x_{1p}$	$l_1$	$\xrightarrow{\text{Learn}} \delta$
2	$x_{21}$	...	$x_{2p}$	$l_2$	
:	:		:		
n	$x_{n1}$	...	$x_{np}$	$l_n$	

## Unsupervised Learning (cluster analysis)

	$x_1$	...	$x_p$	$L$
1	$x_{11}$	...	$x_{1p}$	$\hat{l}_1$
2	$x_{21}$	...	$x_{2p}$	$\hat{l}_2$
:	:		:	
n	$x_{n1}$	...	$x_{np}$	$\hat{l}_n$

Goals:
1. Estimate $\hat{l}_1, \dots, \hat{l}_n$ & $\hat{\gamma}$
2. $1 \left[ \begin{array}{cccc} x_{11} & \dots & x_{1p} \\ x_{21} & \dots & x_{2p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{array} \right] \left[ \begin{array}{c} \hat{l}_1 \\ \vdots \\ \hat{l}_n \end{array} \right] \xrightarrow{\text{Learn}} \delta$

## Supervised learning

3 basic ingredients

1. Distribution of features:  $\underline{x} | L=i \sim f_i(\underline{x}) \quad i=1, \dots, g$  density of  $\underline{x}$  in group  $i$

2. Prior distribution:  $P[L=i] = p_i \quad i=1, \dots, g \quad p_i \geq 0, \quad \sum_i p_i = 1$  (context depend.)

3. Costs of misclassification:  $C(i|i) \quad i, j = 1, \dots, g$

Cost of attributing to group  $i$  a unit  $\underline{x}$ , belonging to group  $j$

- a.  $C(i|i) = 0 \quad i=1, \dots, g$
- b.  $C(i|j) > 0 \quad i \neq j$

CONTEXT DEPENDENT

## Optimal classifier

Obs:  $\delta: X \rightarrow \{1, \dots, g\}$  is equivalent to a partition  $\{R_1, \dots, R_g\}$  of  $X$ .

Partitions  $R_i \subseteq X \quad i=1, \dots, g$  s.t.:
 

- a.  $R_i \cap R_j = \emptyset \quad \text{if } i \neq j$
- b.  $\bigcup_{i=1}^g R_i = X$

$$R_i = \delta^{-1}(i) = \{\underline{x} \in X : \delta(\underline{x}) = i\}$$

Optimality criterion: minimize Expected Cost of Misclassification (ECM)

## Example

$g=2$  (discretization classifier)

$$\delta \longleftrightarrow \{R_1, R_2\}$$

$$\delta(\underline{x}) = 1 \iff \underline{x} \in R_1$$

$$\delta(\underline{x}) = 2 \quad \text{otherwise}$$

$$\delta: X \rightarrow \{1, 2\}$$

(indeed  $R_1 = R_1^\complement$ )

$$\text{ECM}(\delta) = \int_{R_2} C(2|1) f_1(\underline{x}) p_1 d\underline{x} + \int_{R_1} C(1|2) f_2(\underline{x}) p_2 d\underline{x} =$$

$$= \int_X C(2|1) f_1(\underline{x}) p_1 d\underline{x} - \int_{R_2} C(2|1) f_2(\underline{x}) p_2 d\underline{x} + \int_{R_2} C(1|2) f_2(\underline{x}) p_2 d\underline{x} =$$

$$= C(2|1) p_2 - \int_{R_2} (C(2|1) f_2(\underline{x}) p_2 - C(1|2) f_2(\underline{x}) p_2) d\underline{x} =$$

$$\text{Optimal } \delta: \quad R_1 = \{\underline{x} : C(2|1) f_2(\underline{x}) p_2 \leq C(1|2) f_1(\underline{x}) p_1\}$$

$$R_2 = \{\underline{x} \in X : C(2|1) f_2(\underline{x}) p_2 > C(1|2) f_1(\underline{x}) p_1\}$$

## General case

$$\delta: X \rightarrow \{1, 2, \dots, g\}$$

$$\delta \longleftrightarrow \{R_1, \dots, R_g\} \text{ part of } X$$

$$\delta(\underline{x}) = \begin{cases} 1 & \underline{x} \in R_1 \\ 2 & \underline{x} \in R_2 \\ \vdots & \vdots \\ g & \underline{x} \in R_g \end{cases}$$



$$E_{\text{CCE}}(\delta) = \sum_{k \neq i} \int_{R_k} c(k|k) f_k(x) p_k dx + \sum_{k \neq i} \int_{R_k} c(k|k) f_k(x) p_k dx + \dots + \sum_{k \neq i} \int_{R_k} c(k|k) f_k(x) p_k dx = \\ = \int_{R_i} \sum_{k \neq i} c(k|i) f_k(x) p_k dx + \int_{R_i} \sum_{k \neq i} c(k|i) f_k(x) p_k dx + \dots + \int_{R_i} \sum_{k \neq i} c(k|i) f_k(x) p_k dx =$$

OPTIMAL  $\delta$ :  $R_{i,j} = \left\{ x \in X : \sum_{k \neq i} c(k|i) f_k(x) p_k \leq \sum_{k \neq j} c(k|j) f_k(x) p_k \quad \text{for } j \neq i \right\}$   
 $R_{i,j} = \left\{ x \in X : \sum_{k \neq i} c(k|i) f_k(x) p_k \leq \sum_{k \neq j} c(k|j) f_k(x) p_k \quad \text{for } j \neq i \right\}$

$$\delta(x) = i \iff x \in R_{i,i} = \left\{ x \in X : \sum_{k \neq i} c(k|i) f_k(x) p_k \leq \sum_{k \neq j} c(k|j) f_k(x) p_k \quad \text{for } j \neq i \right\}$$

Obs:  $\delta(x) = t \in \{1, 2, \dots, g\} \iff \frac{\sum_{k \neq t} c(t|k) f_k(x) p_k}{\sum_{k \neq t} f_k(x) p_k} \leq \frac{\sum_{k \neq t} c(j|k) f_k(x) p_k}{\sum_{k \neq t} f_k(x) p_k}$

$$\frac{f_k(x) p_k}{\sum_{k \neq t} f_k(x) p_k} = \frac{P[x=k | L=k] \cdot P[L=k]}{\sum_{k \neq t} P[x=k | L=i] P[L=i]} = \frac{P[x=k, L=k]}{P[x=k]} = \underbrace{P[L=k | x=x]}_{\text{posterior}} \quad (\text{Bayes rule})$$

$$\delta(x) = t \iff \sum_{k \neq t} c(t|k) P[L=k | x=x] \leq \sum_{k \neq t} c(j|k) P[L=k | x=x] \quad \text{optimal class.}$$

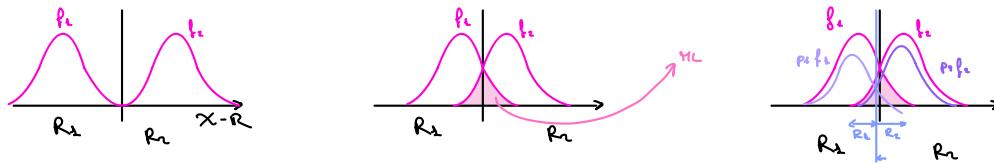
### Special cases

1.  $c(i|i) = \text{const} = d > 0 \quad \text{for } i \neq j$   
 $c(i|i) = 0 \quad i = 1, \dots, g$

Optimal classifier:  $\delta(x) = t \iff \sum_{k \neq t} d \cdot P[L=k | x=x] \leq \sum_{k \neq t} d \cdot P[L=k | x=x]$        $i \neq t$   
 $d - P[L=t | x=x] \leq 1 - P[L=j | x=x]$   
 $P[L=t | x=x] \geq P[L=j | x=x] \quad \text{BAYES CLASSIFIER}$

2.  $c(i|i) = d > 0 \quad i \neq j$   
 $c(i|i) = 0 \quad i = 1, \dots, g$   
 $\delta(x) = t \iff \frac{f_t(x) p_t}{\sum f_k(x) p_k} \geq \frac{f_j(x) p_j}{\sum f_k(x) p_k} \quad i \neq t \quad \text{ML CLASSIFIER}$

### Example



Obs.: Bayes cl is more flexible than what you might think!

Assume  $c(i|i) = c_i \quad i \neq j$

$c(i|i) = 0 \quad i = 1, \dots, g$

Classifier:  $\delta(x) = t \iff \sum_{k \neq t} c(t|k) f_k(x) p_k \leq \sum_{k \neq t} c(s|k) f_k(x) p_k$   
 $\sum_{k \neq t} c_k f_k(x) p_k \leq \sum_{k \neq t} c_s f_k(x) p_k \implies \pi_k = \frac{c_k p_k}{\sum c_s p_s} \geq 0, \sum \pi_k = 1$   
 $\sum_{k \neq t} f_k(x) \pi_k \leq \sum_{k \neq t} f_k(x) \pi_k$

### Exercise

Work out the Bayes cl. when:

1.  $c(i|i) = c_i \quad i \neq j$

2.  $c(i|i) = c_i \cdot h_i \quad i \neq j$

( $c^{\neq i}, h^{\neq i}$ )

### Different costs

$c(t|k) = c_k > 0 \quad t, k = 1, \dots, g$

$c(t|t) = 0$

Opt. Class  $\implies$  Bayes cl with priors  $\pi_j = \frac{c_j p_j}{\sum c_k p_k} \quad j = 1, \dots, g$



$$\text{Ex: } c(i, s) = c^{\log p_i} \quad i, s = 1, \dots, g$$

Bayes CL when features are Gaussian

$$x | L=i \sim N_p(\mu_i, \Sigma_i)$$

$$\text{Bayes CL: } \delta(x) = t \iff P[L=t | x=x] \geq P[L=s | x=x] \quad s=1, \dots, g$$

$$f_t(x) p_t > f_s(x) p_s$$

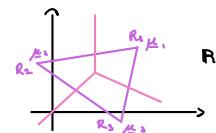
$$\frac{p_t}{\sqrt{(2\pi)^p |\Sigma_t|}} \exp\left[-\frac{1}{2}(x - \mu_t)' \Sigma_t^{-1} (x - \mu_t)\right] \geq \frac{p_s}{\sqrt{(2\pi)^p |\Sigma_s|}} \exp\left[-\frac{1}{2}(x - \mu_s)' \Sigma_s^{-1} (x - \mu_s)\right]$$

$$\log p_t - \frac{1}{2} \log |\Sigma_t| - \frac{1}{2} (x - \mu_t)' \Sigma_t^{-1} (x - \mu_t) \geq \log p_s - \frac{1}{2} \log |\Sigma_s| - \frac{1}{2} (x - \mu_s)' \Sigma_s^{-1} (x - \mu_s) \quad s=1, \dots, g$$

$$\text{Obs: If } p_1 = p_2 = \dots = p_g = \frac{1}{g}$$

$$\text{& } \Sigma_1 = \Sigma_2 = \dots = \Sigma_g$$

$$\text{Bayes CL} \Rightarrow x \in R_t \quad \delta(x) = t \iff d_{\Sigma_t}(x, \mu_t) \leq d_{\Sigma_s}(x, \mu_s)$$

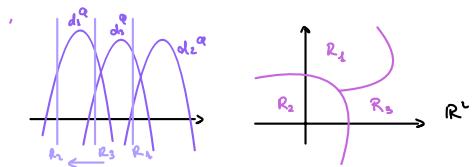


Definition (Quadratic discriminant function)

$$\text{for } t=1, \dots, g \quad d_t^Q(x) = \log p_t - \frac{1}{2} \log |\Sigma_t| - \frac{1}{2} (x - \mu_t)' \Sigma_t^{-1} (x - \mu_t)$$

$$\text{Bayes CL: } \delta(x) = t \iff x \in R_t = \{x \in R^p : d_t^Q(x) \geq d_s^Q(x) \quad s=1, \dots, g\}$$

QDA (quadratic discriminant analysis)



$$\text{What if: } \Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$$

$$\delta(x) = t, \quad t=1, \dots, g \iff \log p_t - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_t)' \Sigma^{-1} (x - \mu_t) \geq \log p_s - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_s)' \Sigma^{-1} (x - \mu_s)$$

$$\log p_t - \frac{1}{2} \mu_t' \Sigma^{-1} x + x' \Sigma^{-1} \mu_t - \frac{1}{2} \mu_t' \Sigma^{-1} \mu_t \geq \log p_s - \frac{1}{2} \mu_s' \Sigma^{-1} x + x' \Sigma^{-1} \mu_s - \frac{1}{2} \mu_s' \Sigma^{-1} \mu_s$$

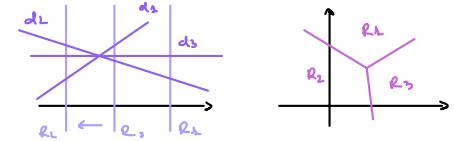
$$x' \Sigma^{-1} \mu_t + \log p_t - \frac{1}{2} \mu_t' \Sigma^{-1} \mu_t \geq x' \Sigma^{-1} \mu_s + \log p_s - \frac{1}{2} \mu_s' \Sigma^{-1} \mu_s \quad s=1, \dots, g$$

Definition (Linear discriminant function)

$$t=1, \dots, g \quad d_t(x) = x' \Sigma^{-1} \mu_t + \log p_t - \frac{1}{2} \mu_t' \Sigma^{-1} \mu_t$$

$$\text{Bayes CL: } \delta(x) = t \iff x \in R_t = \{x \in R^p : d_t(x) \geq d_s(x), \quad s=1, \dots, g\}$$

LDA (linear discriminant analysis)



Use the training data to estimate  $f_i$   $i=1, \dots, g$ , i.e. the distrib. of  $x | L=i$   $i=1, \dots, g$

$$\text{Data: } \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

EXAMPLE

LDA:  $\mu_1, \dots, \mu_g, \Sigma$  to be estimated from data

QDA:  $\mu_1, \dots, \mu_g, \Sigma_1, \dots, \Sigma_g$  "

Estimate:  $\mu_i$  with  $\bar{x}_i = \frac{1}{n_i} \sum_{j:e_j=i} x_j$   $n_i = \# \{j : e_j=i\} \quad i=1, \dots, g$

In QDA:  $\Sigma_i$  with  $\hat{\Sigma}_i = \frac{1}{n_i - 1} \sum_{j:e_j=i} (x_j - \bar{x}_i)(x_j - \bar{x}_i)'$

In LDA:  $\Sigma$  with  $\hat{\Sigma}_{pool} = \frac{1}{n-g} \sum_{i=1}^g (n_i - 1) \hat{\Sigma}_i$

To estimate  $\Sigma_1, \dots, \Sigma_g$  you need  $n_1, \dots, n_g$  to be large wrt.  $p$ .

Class that do not require large sample sizes (parametric  $\Sigma_1, \dots, \Sigma_g$ ):

• Naive Bayes

$$\text{Assume } \Sigma_i = \begin{bmatrix} \sigma_{ii}^{(x)} & 0 \\ 0 & \sigma_{ii}^{(e)} \end{bmatrix} \Rightarrow d_t^Q = \log p_t - \frac{1}{2} \sum_{i=1}^p \log \sigma_{ii}^{(x)} - \frac{1}{2} \sum_{i=1}^p \frac{(x_i - \bar{x}_{ti})^2}{\sqrt{\sigma_{ii}^{(x)}}}$$

$\bar{x}_{ti}$   $i$ -th comp. of  $\bar{x}_t$

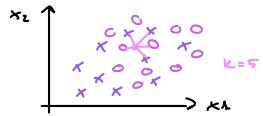
$$\sigma_{ii}^{(x)} = \frac{1}{n_i - 1} \sum_{j:e_j=i} (x_{ji} - \bar{x}_{ti})^2$$



## KNN Classifier

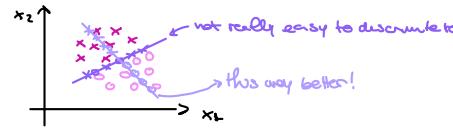
Fix  $k \geq 1$

$N_k(x) = \{x_i : x_i \text{ in the training set}\}$  ("closest" to  $x \rightarrow$  according to a specified dist)  
 $\delta(x) = t \iff \text{the majority of } x_i \in N_k(x) \text{ have class } t$



## Fisher's argument for LDA

- Robustness of LDA to Gaussianity ass
- Dimensionality reduction



Assumptions:  $i = 1, \dots, g$   $x | L=i \sim \mu_i, \Sigma$   
 $\Sigma$  same on in every group  $i$

$$\text{let } B = \frac{1}{g-1} \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T \quad \bar{\mu} = \frac{1}{g} \sum_{i=1}^g \mu_i$$

$$\Sigma \in \mathbb{R}^p \quad x | L=i \sim \mu_i, \Sigma$$

$$\mathbb{E}[x | L=i] = \mu_i$$

$$\text{Var}(x | L=i) = \Sigma$$

$$\begin{aligned} \text{Problem: } \arg \max_{B \in \mathbb{R}^{p \times p}} \frac{\alpha^T B \alpha}{\alpha^T \Sigma \alpha} &= \arg \max_{B \in \mathbb{R}^p} \frac{1}{g-1} \frac{\sum_{i=1}^g (\alpha^T \mu_i - \alpha^T \bar{\mu})^2}{\alpha^T \Sigma \alpha} \\ \frac{\alpha^T B \alpha}{\alpha^T \Sigma \alpha} &= \frac{\alpha^T \Sigma^{-1/2} B \Sigma^{-1/2} \alpha}{\alpha^T \alpha} \quad \alpha = \Sigma^{-1/2} \alpha \quad B = \Sigma^{-1/2} \alpha \Sigma^{-1/2} \\ \arg \max_{B \in \mathbb{R}^p} \frac{\alpha^T \Sigma^{-1/2} B \Sigma^{-1/2} \alpha}{\alpha^T \alpha} &= \alpha^T \alpha \quad \Sigma^{-1/2} B \Sigma^{-1/2} = \sum_{i=1}^g \lambda_i e_i e_i^T \\ \arg \max_{B \in \mathbb{R}^p} \frac{\alpha^T B \alpha}{\alpha^T \Sigma \alpha} &= \Sigma^{-1/2} \alpha^T \alpha \\ \alpha &= \Sigma^{-1/2} \alpha \\ \Sigma \alpha &= \Sigma^{-1/2} \alpha \Sigma \quad S = \min(g-1, p) \\ x \rightarrow \begin{bmatrix} \alpha^T x \\ \Sigma^{-1/2} x \\ \vdots \\ \alpha^T \Sigma^{-1/2} x \end{bmatrix} & \text{Fisher discriminant scores} \\ \text{Cov}(\alpha^T x, \alpha^T x) &= \alpha^T \Sigma \alpha = \alpha^T \Sigma^{-1/2} \Sigma^{-1/2} \alpha = \alpha^T \alpha = \left\langle \begin{array}{c} 0 \\ \vdots \\ 0 \end{array} \right\rangle \end{aligned}$$

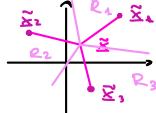
$$A = \begin{bmatrix} \alpha \\ \vdots \\ \alpha \end{bmatrix} \quad \text{Cov}(Ax) = I$$

## How to build a class

$\mu_i$  est  $\bar{\mu}_i$

$\Sigma$  est  $\hat{\Sigma}_{\text{pooled}}$

$$\bar{\mu}_i \rightarrow \begin{bmatrix} \alpha^T \bar{\mu}_i \\ \Sigma^{-1/2} \bar{\mu}_i \\ \vdots \\ \alpha^T \Sigma^{-1/2} \bar{\mu}_i \end{bmatrix} \quad k \leq S$$



$R^k$

$$\delta(x) = ? \quad x \rightarrow \begin{bmatrix} \alpha^T x \\ \Sigma^{-1/2} x \\ \vdots \\ \alpha^T \Sigma^{-1/2} x \end{bmatrix}$$

find the eigenvectors of:

$$\hat{\Sigma}_{\text{pooled}}^{-1/2} \hat{\Sigma}^{-1/2} \rightarrow \hat{B} = \frac{1}{g-1} \sum_{i=1}^g (\bar{\mu}_i - \bar{\mu})(\bar{\mu}_i - \bar{\mu})^T$$

$$\delta(x) = t \iff \sum_{i=1}^g (\bar{\mu}_i - \hat{x}_{i,i})^2 \leq \sum_{j=1}^g (\bar{\mu}_j - \hat{x}_{j,j})^2$$

## Evaluating a classifier

$$\begin{aligned} \mathcal{X} &= \begin{bmatrix} x_1 & l_1 \\ \vdots & \vdots \\ x_n & l_n \end{bmatrix} \quad l_i \in \{1, \dots, g\} \\ \text{train set} & \quad x_i \in \mathbb{R}^p \\ & \quad i = 1, \dots, n \end{aligned}$$

$$\delta: \mathbb{R}^p \rightarrow \{1, 2, \dots, g\} \quad \text{fitted to } \mathcal{X}$$

## 1. Actual error rate (AER)

$\delta \leftarrow \{R_1, \dots, R_g\}$  partition of  $\mathbb{R}^p$

$$\text{AER}(\delta) = \sum_{k=1}^g \int_{R_k} f_k(x) p_k dx + \sum_{k=2}^g \int_{R_{k-1}} f_k(x) p_{k-1} dx + \dots + \sum_{k=g}^g \int_{R_1} f_k(x) p_g dx$$

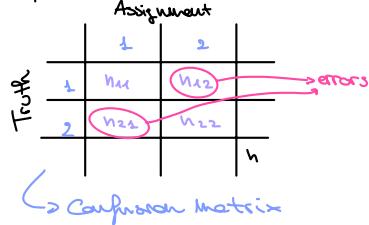
How to estimate AER using data?

we can let  $f=2$

$$\text{AER}(\delta) = \int_{R_2} f_2(x) p_2 dx + \int_{R_1} f_2(x) p_1 dx$$



Naive idea: Apply  $\delta$  to  $X$  and count the mistakes for  $i=1, \dots, n$   $\delta(x_i)$  and then check if  $\delta(x_i) = l_i$  or not.



Apparent error rate:

$$\hat{AER}(\delta) = APER(\delta) = \frac{n_{12} + n_{21}}{n}$$

TOO OPTIMISTIC, you tried to overfit the data to make it smaller, BAD!

### EXAMPLE

KNN with  $K=1$



That means that we cannot rely on the  $\overline{err}$  cause it's not the same.

$$\overline{err} = \frac{n_{12} + n_{21}}{n} = \frac{n_{12}}{n} + \frac{n_{21}}{n} = \frac{n_1}{n} \cdot \frac{n_{12}}{n_1} + \frac{n_2}{n} \cdot \frac{n_{21}}{n_2} =$$

$$= \hat{p}_1 \int_{R_2} f_1(x) dx + \hat{p}_2 \int_{R_1} f_2(x) dx$$

Obs: Dichotomous class

$$\delta(x) = \begin{cases} 1 & \text{positive} \\ 0 & \text{negative} \end{cases}$$

Confusion matrix:

		Assignment	
		1	0
True	1	True pos n <sub>11</sub>	False neg n <sub>10</sub>
	0	False pos n <sub>01</sub>	True neg n <sub>00</sub>
		n <sub>1.</sub>	n <sub>0.</sub>
			n

$$APER = \frac{n_{10} + n_{01}}{n}$$

$$\text{Precision: } \frac{n_{11}}{n_{1.}}$$

$$\text{Recall: } \frac{n_{11}}{n_{10}}$$

$$\text{Specificity: } \frac{n_{00}}{n_{0.}}$$

$$\text{Precision} \uparrow \quad \text{False Pos} \downarrow$$

$$\text{Recall} \uparrow \quad \text{False Neg} \downarrow$$

For a more realistic estimate of AER (and the other indices computed out of the confusion matrix) USE TEST DATA

$$\tilde{\mathbf{x}} = \begin{bmatrix} \tilde{x}_1 & \tilde{l}_1 \\ \vdots & \vdots \\ \tilde{x}_m & \tilde{l}_m \end{bmatrix} \quad \text{Apply } \delta_x \text{ (learned with *)} \quad \hat{AER}(\delta_x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} [\delta_x(\tilde{x}_i) \neq \tilde{l}_i] \quad \mathbb{1} [\delta_x(\tilde{x}_i) \neq \tilde{l}_i] = \begin{cases} 1 & \text{if } \delta_x(\tilde{x}_i) \neq \tilde{l}_i \\ 0 & \text{if } \delta_x(\tilde{x}_i) = \tilde{l}_i \end{cases}$$

### 2. CROSS VALIDATION

Leave one out

$$\mathbf{X} = \begin{bmatrix} x'_1 & l'_1 \\ \vdots & \vdots \\ x'_n & l'_n \end{bmatrix}$$

train data

For  $i = 1, \dots, n$

step 1: take unit  $i$  out

$$\mathbf{X}_{-i} = \begin{bmatrix} x'_1 & l'_1 \\ \vdots & \vdots \\ \cancel{x'_i} & \cancel{l'_i} \\ \vdots & \vdots \\ x'_n & l'_n \end{bmatrix}$$

Step 2: train  $\delta$  on  $\mathbf{X}_{-i} \Rightarrow \delta_{-i}: \mathbb{R}^d \rightarrow \{1, \dots, k\}$

Step 3: Evaluate  $\delta_{-i}$  on unit  $i$   $\delta_i = \mathbb{1} [\delta_{-i}(x'_i) \neq l'_i]$

Return:  $AER(\delta_x) = \frac{1}{n} \sum_{i=1}^n \delta_i \leftarrow \text{low bias - high variability}$

### 3. K-fold cross-validation

Fix  $K \geq 2$ , split the train set in  $K$ -subsections (usually 5/10, but it's up to us)

$$\mathbf{X} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad K=4 \quad \text{By permuting the rows of } \mathbf{X} \text{ you are randomly assigning units to the } K \text{ subsets}$$

Obs: try to balance the  $K$  subsets



Step 1: train  $\delta$  on  $K$ -parts  $\Rightarrow \delta_{\text{parts}}: \mathbb{R}^p \rightarrow \{z_1, \dots, z_J\}$   
 Step 2: Apply  $\delta_{\text{parts}}$  to part  $j \Rightarrow \text{Err}_j = \frac{1}{n_j} \sum_{i \in \text{parts } j} \varepsilon_i \quad \varepsilon_i = \begin{cases} 1 & [\delta_{\text{parts}}(x_i) \neq c_i] \\ 0 & \text{otherwise} \end{cases} \quad n_j = \#\{i: i \in \text{part } j\}$   
 Return:  $\hat{\text{AER}}(\delta_x) = \frac{1}{n} \sum_{j=1}^K n_j \text{Err}_j$

Obs: Can be used to estimate: Precision, Recall, ...

Obs: Leave-one-out  $\Leftrightarrow K=n$

Given  $K$ : repeat  $K$ -fold cross-val  $B$  times,  $B$  large. ( $\text{repeat} = \text{permute and then repeat, it changes!}$ )

$$\begin{aligned}
 \hat{\text{AER}}_1(\delta_x), \dots, \hat{\text{AER}}_B(\delta_x) &\Rightarrow \mathbb{E}[\hat{\text{AER}}(\delta_x)] = \hat{\text{AER}}_m(\delta_x) = \frac{1}{B} \sum_{i=1}^B \hat{\text{AER}}_i(\delta_x) \\
 &\Rightarrow \text{Var}[\hat{\text{AER}}(\delta_x)] = \frac{1}{B-1} \sum_{i=1}^B [\hat{\text{AER}}_i(\delta_x) - \hat{\text{AER}}_m(\delta_x)]^2 \\
 &\Rightarrow \text{CI}_{1-\alpha}(\mathbb{E}[\hat{\text{AER}}(\delta_x)]) = [\hat{\text{AER}}_m(\delta_x) \pm \sqrt{\text{Var}[\hat{\text{AER}}(\delta_x)]} z_{1-\frac{\alpha}{2}}]
 \end{aligned}$$

Obs: Leave-one-out  $\hat{\text{AER}}(\delta_x) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$

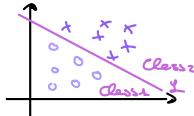
$$\frac{\text{Var}(\varepsilon_i)}{n} \approx \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i\right) \ll \text{Var}(\varepsilon_i)$$

$\varepsilon_{-1}, \varepsilon_{-2}, \dots, \varepsilon_{-n}$  strongly dependent!!  $\Rightarrow \text{Var}\left(\frac{1}{n} \sum \varepsilon_i\right) \approx \text{Var}(\varepsilon_i)$

## SUPPORT VECTOR MACHINES

Dichotomous Problem

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 \\ \vdots & \vdots \\ x_n & x_m \end{bmatrix} \quad x_i \in \mathbb{R}^p \quad i = 1, \dots, n$$



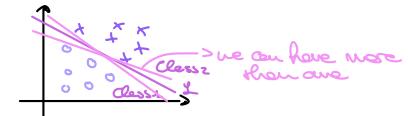
Q: Given  $A, B \subseteq \mathbb{R}^p$  when can they be separated by an hyperplane?

Assumption: Let  $\text{CH}(A)$  and  $\text{CH}(B)$  be the convex hull generated by  $A$  and  $B$

1.  $\text{CH}(A) \neq \emptyset$  and  $\text{CH}(B) \neq \emptyset$
2.  $\text{CH}(A) \cap \text{CH}(B) = \emptyset$   $\Rightarrow$   $\exists$  separating hyperplane
3. Either  $\text{CH}(A)$  or  $\text{CH}(B)$  is open

Obs: Condition 3 can be substituted by:

- 3'.  $\text{CH}(A)$  and  $\text{CH}(B)$  are closed and at least one of them is compact.  
(satisfied by the training data)

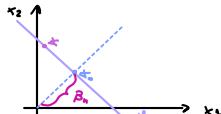


## Hyperplanes

$L$  hyperplane in  $\mathbb{R}^p$

affine subspace of dim  $p-1$

Let  $\beta \in \mathbb{R}^p$ ,  $\|\beta\|=1$  s.t.  $\beta \perp L$



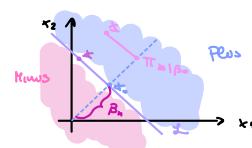
$\Delta_0 = \text{span}(\beta) \cap L$

$\beta_0 = \|\Delta_0\|$

$$\mathbb{R}^p \ni x \in L \iff \pi_{\Delta_0} x = x_0 \implies \frac{\beta^\top \beta}{\|\beta\|^2} x = x_0 \implies \boxed{\beta^\top x = \beta_0}$$

$$x \in L \iff \beta^\top x = \beta_0$$

$$\beta^\top x \geq \beta_0$$



Obs:  $\beta^\top x - \beta_0$  distance with sign of  $x$  from  $L$ .

Going back to classification

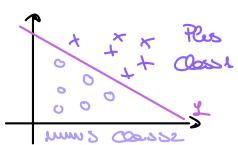
Assume  $\mathbf{x}$  are s.t.  $\exists L$  separating Class 1 and Class 2.

$L \Rightarrow$  defines  $R_1$  and  $R_2$  half-spaces

$\Rightarrow$  classifier

let  $y_i = +1$  if unit  $i$  belongs to class 1.

$y_i = -1$  if unit  $i$  belongs to class 2.



Note:  $y_i \cdot (\beta^T x_i - \beta_0) \geq 0$

distance between  $x_i$  and  $L$

$$H_L = \min \{ y_i (\beta^T x_i - \beta_0) : i=1, \dots, n \}$$

↳ Problem: find  $L$  (i.e.  $\beta, \beta_0$ ) s.t.  $H_L$  is max

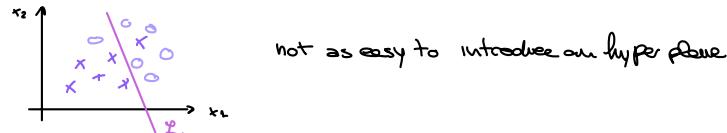
$$\begin{aligned} \max_{\beta, \beta_0} H & \text{ s.t. } \|\beta\| = 1 \\ & y_i (\beta^T x_i - \beta_0) \geq H \text{ for } i=1, \dots, n \end{aligned}$$

(OPTIMIZATION PROBLEM)

Solution:  $\hat{\beta} = \sum_{i=1}^n \lambda_i y_i x_i \quad \hat{\beta}_0 \in \mathbb{R}$   
 $f(x) = \hat{\beta}^T x - \hat{\beta}_0 = \sum_{i=1}^n \lambda_i y_i x^T x - \hat{\beta}_0$

classifier:  $C(x) = \operatorname{sgn}(f(x))$

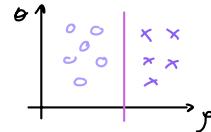
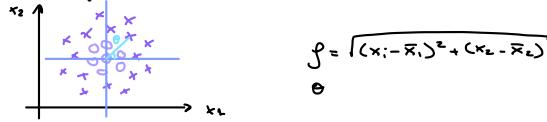
What if:



Soft solution:  $\max_{\beta, \beta_0} H \quad \text{s.t. } \|\beta\| = 1$   
 $y_i (\beta^T x_i - \beta_0) \geq (1 - \varepsilon_i) y_i$   
 $\varepsilon_i \geq 0, \sum_{i=1}^n \varepsilon_i \leq C \rightarrow \text{Budget cost}$

- Extra parameters:
- $C$  budget
  - $C=0 \Rightarrow$  hard problem
  - $C$  very large  $\Rightarrow$  any  $L$  would do

what if?



Take  $h_1: \mathbb{R}^d \rightarrow \mathbb{R}, \dots, h_m: \mathbb{R}^d \rightarrow \mathbb{R}$

and let  $h(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_m(\mathbf{x}))^T \in \mathbb{R}^m$

Now training set:

$$\tilde{\mathbf{X}} = \begin{bmatrix} h_1(\mathbf{x}_1) & y_1 \\ \vdots & \vdots \\ h_1(\mathbf{x}_n) & y_n \end{bmatrix} \Rightarrow \text{find soft or hard separation hyperplane for } \tilde{\mathbf{X}}$$

Solution:  $\hat{\beta} = \sum_{i=1}^n \lambda_i y_i h(\mathbf{x}_i)$

$$x \in \mathbb{R}^d \quad f(x) = \hat{\beta}^T h(x) - \hat{\beta}_0 = \sum_{i=1}^n \lambda_i y_i h(\mathbf{x}_i)^T h(x) - \hat{\beta}_0$$

$$C(x) = \operatorname{sign}(f(x))$$

Idea:  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  s.t.  $K(\mathbf{x}_i, \mathbf{x}_j) = \langle h(\mathbf{x}_i), h(\mathbf{x}_j) \rangle = h(\mathbf{x}_i)^T h(\mathbf{x}_j) \quad \forall \mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$

$$\Rightarrow f(x) = \sum_{i=1}^n \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}) \quad \text{RKHS}$$

### Examples of kernels

- $K(\mathbf{x}, \mathbf{w}) = [\mathbf{1} + \mathbf{x}^T \mathbf{w}]^\alpha$  d-degree polynomials
- $K(\mathbf{x}, \mathbf{w}) = \exp[-\gamma \|\mathbf{x} - \mathbf{w}\|^2]$  radial basis
- $K(\mathbf{x}, \mathbf{w}) = \tanh[\kappa_1 \mathbf{x}^T \mathbf{w} + \kappa_0]$  neural networks

