# Quadratic and Linear Discriminant Analysis (QDA and LDA)

we discussed: reduce the size of NN training set by only keeping representative inst.
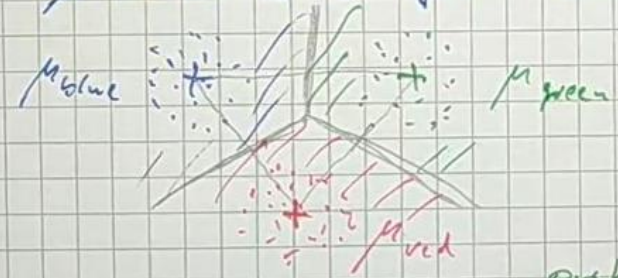
$\Rightarrow$ to the extreme: one representative per class

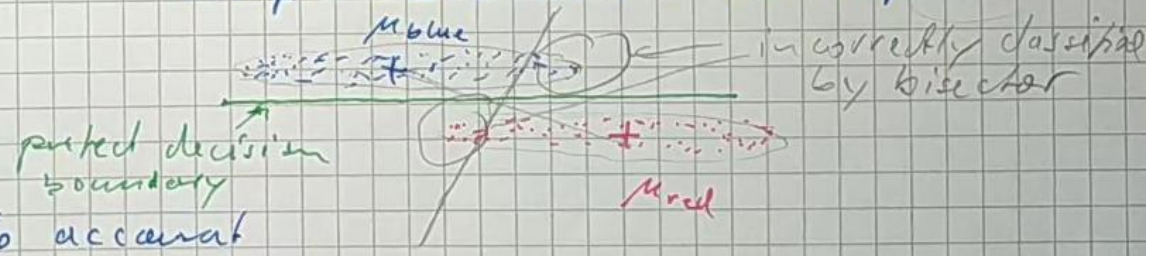- basic idea: use the class mean: $k = 1, ..., C$: $N_k$: # instances in class $k$

$$\mu_k = \frac{1}{N_k} \sum_{i:\, Y_i = k} X_i$$

classify query $X$: $\quad \hat{Y} = \arg\min_{k} d(X, \mu_k)$

this works, when data form clusters (one cluster per class) and
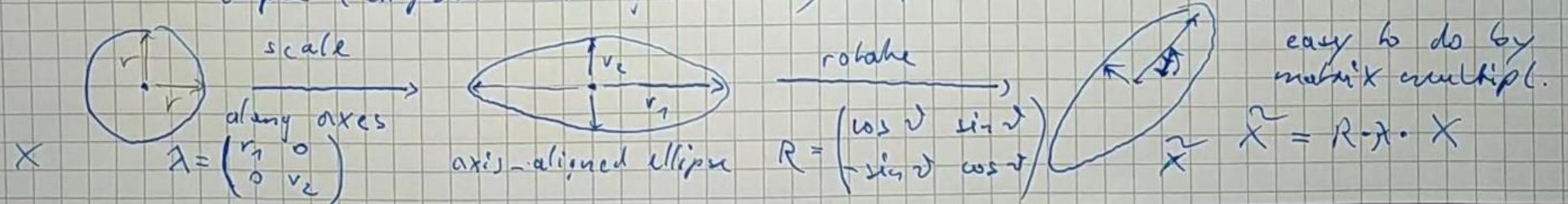
a mean is a good cluster representative



$\mu_{blue}$ $\mu_{green}$

$\mu_{red}$

perfect decision boundary

in practice, data are rarely so nice,

example where QDA would help:

$\mu_{blue}$

$\mu_{red}$

incorrectly classified by bisector
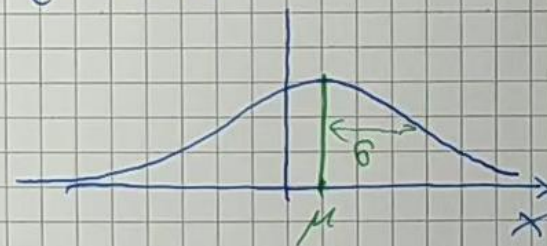
$\Rightarrow$ need to take cluster shape into account

nearest mean works for circular clusters, next complicated geometric shape

ellipse (ellipsoid in higher dims.). from circle to ellipse:



$X$

scale along axes $\Lambda = \begin{pmatrix} r_1 & 0 \\ 0 & r_2 \end{pmatrix}$

axis-aligned ellipse $\quad R = \begin{pmatrix} \cos \vartheta & \sin \vartheta \\ -\sin \vartheta & \cos \vartheta \end{pmatrix}$ rotate

easy to do by matrix multipl.

$\tilde{X} = R \cdot \Lambda \cdot X$

\# of degrees of freedom:
- circle: center ($D$ dimensions) radius ($1$ dimension)
- ellipse: —''— scaling + rotation
  $$\stackrel{1}{=} \text{symmetric matrix} \quad \frac{D}{2}(D+1)$$

$$O(D): (D+1) \quad \text{vs.} \quad O(D^2): D + \frac{D}{2}(D+1) \quad \text{d.o.f.}$$

- How to fit an ellipse to a data cluster
  - consider each class separately $\exists$ fit $C$ independent ellipses
    ( drop the class label for convenience )
  - central limit theorem of probability theory: the superposition of ~~many~~
    infinitely many random events converges to a Gaussian distribution
  - reminder: 1-D Gaussian
    $$N(\underset{\underset{\text{mean}}{\uparrow}}{\mu}, \underset{\underset{\text{variance}}{\uparrow}}{\sigma^2}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$



  - generalization to arbitrary dimensions:
    $\mu$: $D$-dimensional vector
    $\sigma^2 \rightarrow \Sigma$: $D \times D$ symmetric positive definite matrix
    " co-variance matrix "

compute square root: $\Sigma = \Sigma^{1/2} \Sigma^{1/2}$

$\Sigma^{1/2} = R \cdot \lambda$

$\sigma^2 = \sigma \cdot \sigma$

$\underset{\uparrow}{R}$ rotation, $\underset{\uparrow}{\lambda}$ axis-aligned scaling

$$N(\mu, \Sigma) = \frac{1}{\sqrt{\underset{\underset{\text{determinant}}{\uparrow}}{\det(2\pi\Sigma)}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

## Derivation of QDA (Quadratic Discriminant Analysis)

- intuition: data features are arranged in feature space such that we have an elliptic cluster for each class ⇒ fit a model that captures these clusters

- mathematics: arrangement of the features for a given class is expressed by the likelihood function: $p(X \mid Y = k)$ and prior $\hat{p}(Y = k) = \frac{N_k}{N}$

  ⇒ learn a model for these likelihoods,

  specifically: elliptic clusters ≝ multivariate Gaussian

  ⇒ generative classifier ⤳ one cluster per class

$$p(Y = k \mid X) = \frac{p(X \mid Y = k) \cdot p(Y = k)}{p(X)} \leftarrow \text{how often occurs each class?}$$

  ↑ generative classifier          ↖ normalization

- **maximum likelihood principle:**

  - we can treat each class in isolation ⇒ simplified $TS = \{\, \{X_i\}_{i=1}^{N_k}$ for fixed class $k$

  - assumption: $TS = \{X_i\}_{i=1}^{N_k}$ is a *typical* representation of the cluster / class $k$

    turn assumption around: search for model, such that observed data are "as typical as possible".

    mathematically: $p(TS \mid model) \rightarrow$ maximize, fit the model so "maximum likelihood fit" ≝ model where the data have max. prob.

  - simplify $p(TS \mid model)$ by the *i.i.d.* assumption ≝ assume that all training instances are independently drawn from the same model ≝ factorize prob.

$$\boxed{p(TS \mid model) = \prod_i p(X_i \mid model)} \Rightarrow \text{maximize via model fit}$$

- specifically when $p(x_i | \text{model})$ is a multivariate Gaussian
$\underset{\mu, \Sigma}{}$ mean and covariance matrix

$$p(x_i | \mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi \Sigma)}} \cdot \exp\left(-\frac{1}{2} (x_i - \mu) \Sigma^{-1} (x_i - \mu)^T\right)$$

( $x_i$ is a row vector, therefore $\mu$ also row vector $\underbrace{(x_i - \mu) \cdot (x_i - \mu)^T}_{\text{scalar product}}$

$\tilde{=}$ simplify: take the negative logarithm of likelihood

likelihood $\to$ max $\hat{=}$ log likelihood $\to$ max $\hat{=}$ $-$log likelihood $\to$ min

$$-\log \prod_{i=1}^{N} p(x_i | \mu, \Sigma) = \sum_{i=1}^{N} -\log p(x_i | \mu, \Sigma)$$

$$= \sum_{i=1}^{N} -\left[ \log \cdot \frac{1}{\sqrt{\det(2\pi \Sigma}} + \left(-\frac{1}{2}(x_i - \mu) \Sigma^{-1}(x_i - \mu)^T\right)\right]$$

$$\text{Loss}(TS | \mu, \Sigma) = \sum_{i=1}^{N} \log \sqrt{\det(2\pi \Sigma)} + \frac{1}{2} (x_i - \mu) \Sigma^{-1} (x_i - \mu)^T \to \text{min}$$

- to find the minimum, take the derivative of the loss w.r.t. parameters and
set it to zero $\Rightarrow$ solve for the parameters $\Rightarrow$ ~~most~~ estimated params $\hat{\mu}, \hat{\Sigma}$

$$\frac{\partial \text{Loss}}{\partial \mu} = \cancel{\sum \cancel{\log \frac{}{}}} - \sum_{i=1}^{N} \Sigma^{-1}(x_i - \mu)^T \overset{!}{=} 0 \quad | -\Sigma \cdot \left[\text{matrix calculus } \frac{\partial v A v^T}{\partial v} = 2 A v^T\right]$$

$$\underset{\text{left multiply}}{}$$

$$\sum_{i=1}^{N} \underbrace{\Sigma \cdot \Sigma^{-1}}_{=\, \mathbb{I}}(x_i - \mu)^T = -\Sigma \cdot 0 = 0 \qquad \overset{=N}{\underline{\phantom{x}}}$$

$$\Rightarrow \boxed{\mu = \frac{1}{N} \sum_{i=1}^{N} x_i} \qquad \sum_{i=1}^{N} x_i = \sum_{i=1}^{N} \mu = \mu \underset{i=1}{\overset{N}{\sum}} \cdot 1 = \mu \cdot N$$

$\frac{\partial \text{Loss}}{\partial \Sigma}$ is a bit hard, because we actually have $\Sigma^{-1}$

$\Rightarrow$ introduce abbreviation $\quad K = \Sigma^{-1}$ precision matrix $\quad S = \sum_{i=1}^{N} z_i^T z_i$

$$z_i = x_i - \mu$$

$$= \sum_{i=1}^{N} (x_i - \mu)^T (x_i - \mu)$$

scatter matrix

$\Rightarrow \frac{\partial \text{Loss}}{\partial K} = \frac{\partial}{\partial K} \left[ \sum_{i=1}^{N} \left[ \log \sqrt{\det(2\pi K^{-1})'} + \frac{1}{2} z_i K z_i^T \right] \right]$

$$= \frac{1}{\det(2\pi K)}$$

$$-\frac{1}{2} \log \det (2\pi K) = -\frac{1}{2} \log (2\pi)^D \cdot \det(K) = -\frac{1}{2} \log (2\pi)^D - \frac{1}{2} \log \det (K)$$

$$= \frac{\partial}{\partial K} \left[ -\frac{1}{2} \log (2\pi)^D - \frac{1}{2} \log \det (K) + \frac{1}{2} z_i \right]$$

$$= \frac{\partial}{\partial K} \sum_{i=1}^{N} \left( -\frac{1}{2} \log (2\pi)^D - \frac{1}{2} \log \det (K) + \frac{1}{2} z_i K z_i^T \right)$$

$\left[ \text{matrix calculus:} \quad \frac{\partial v A v^T}{\partial A} = v^T v \qquad \frac{\partial \log \det A}{\partial A} = (A^T)^{-1} \right]$

$\frac{\partial \text{Loss}}{\partial K} = \sum_{i=1}^{N} \left( -\frac{1}{2} (K^T)^{-1} + \frac{1}{2} z_i^T z_i \right) \overset{!}{=} 0 \qquad K^T = K \quad (K \text{ and } \Sigma \text{ are symmetric}$

$\sum_{i=1}^{N} -K^{-1} + \underbrace{\sum_{i=1}^{N} z_i^T z_i}_{S} = 0 \quad \Rightarrow \quad S = N \cdot K^{-1}$

$$\boxed{\hat{\Sigma} = \frac{1}{N} S}$$

$\frac{1}{N} S$ : empirical covariance matrix

$\Rightarrow$ training: for each class $k$, compute $\hat{\mu}_k = \frac{1}{N_k} \sum_{i: Y_i = k} x_i$, $\hat{\Sigma} = \frac{1}{N} \sum_k \sum_{i: Y_i = k} (x_i - \mu_k)^T$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{i: Y_i = k} (x_i - \mu_k)^T \cdot (x_i - \mu_k)$$

· prediction: for a new instance $x$, compute the likelihood of $x$ for each class,
weight with the class prior, return the most probable class

$$\hat{Y} = \arg\max_k p(Y = k \mid x) = \arg\max_k p(x \mid Y = k)\, p(Y = k)$$

$$= \arg\min_k -\log p(x \mid Y = k) - \log p(Y = k)$$

$$= \arg\min_k \underbrace{\tfrac{1}{2} \log \det(2\pi \hat{\Sigma}_k)}_{\text{independent of } X} + \tfrac{1}{2}(x - \hat{\mu}_k) \hat{\Sigma}_k^{-1}(x - \hat{\mu}_k)^T - \underbrace{\log p(Y = k)}_{\substack{\text{independent}\\ \text{of } X}}$$

$$b_k = \tfrac{1}{2} \log \det(2\pi \hat{\Sigma}_k) - \log p(Y = k)$$

adjust for cluster shape

QAA: $\boxed{\hat{Y} = \arg\min_k \tfrac{1}{2}(x - \hat{\mu}_k) \hat{\Sigma}_k^{-1}(x - \hat{\mu}_k)^T + b_k}$

squared Mahalanobis distance

$\Big[$ Scatter matrix is an outer product:

$A = v^T v$

$\llcorner$ row vector

$A_{ij} = v_i \cdot v_j$ $\Big]$

## LDA (linear discriminant analysis)

- simplification of QDA when all clusters have the same shape

$$\Sigma_k = \Sigma_{k'} := \Sigma_w \qquad \text{"within-class covariance"}$$

$\Rightarrow$ we can get rid of the quadratic terms in QDA, only linear terms remain

~~replace eq.~~ rewrite eq. with precision matrix

$$\arg\min_k \frac{1}{2}(x - \hat{\mu}_k)\hat{K}_w(x - \hat{\mu}_k)^T + b_k$$

$$\underbrace{x\hat{K}_w x^T}_{\text{independent of } k} - 2\hat{\mu}_k\hat{K}_w x^T + \underbrace{\hat{\mu}_k\hat{K}_w\hat{\mu}_k^T}_{\text{independent of } x}$$

independent of $k$
$\Rightarrow$ has no effect on the $\arg\min_k$
$\Rightarrow$ drop

independent of $x$
$\Rightarrow$ absorb into $b_k' = b_k + \frac{1}{2}\hat{\mu}_k\hat{K}_w\hat{\mu}_k^T$

$$\arg\min_k -\underbrace{\hat{\mu}_k\hat{K}_w}_{w_k} x^T + b_k' = \boxed{\arg\max_k w_k x^T - b_k' = \hat{y}} \quad \text{LDA}$$

- interpretation: case $C = 2$ $\Rightarrow$ $w_1$ for class 1, $w_0$ for class 0
$b_1'$ $b_0'$

$$\arg\max_{k=0,1}\left(w_1 x^T - b_1', \ w_0 x^T - b_0'\right) \iff (w_1 - w_0)x^T - (b_1' - b_0') = \begin{cases} \geq 0 \Rightarrow \hat{y} = 1 \\ < 0 \Rightarrow \hat{y} = 0 \end{cases}$$
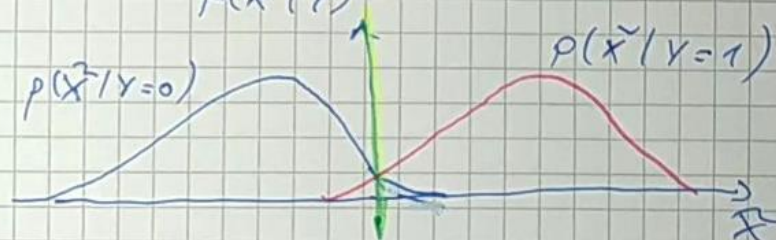
simplified rule: $\boxed{\hat{y} = \text{sign}(w x^T + b)}$ $\qquad w = w_1 - w_0 \quad b = b_0' - b_1'$
LDA, 2 classes

- in practice, this also works when clusters are not elliptic $\Rightarrow$ choose $b$ by cross-validation
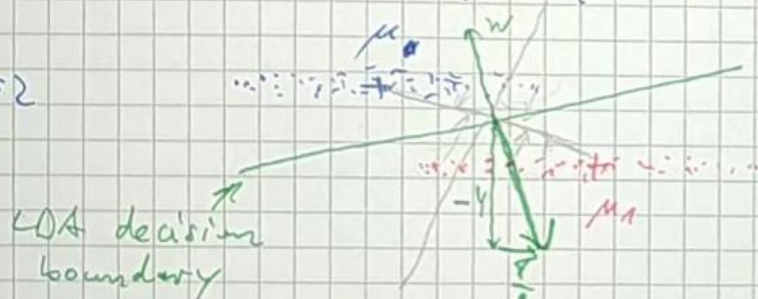
an expression of the form $w \cdot X^T + b$ has an intuitive effect:

- projects the data $X^T$ onto the vector $w$ $\Rightarrow$ 1-dimensional feature $w \cdot X^T$

- $+b$ shifts the origin of the new feature such that the decision boundary is at the origin
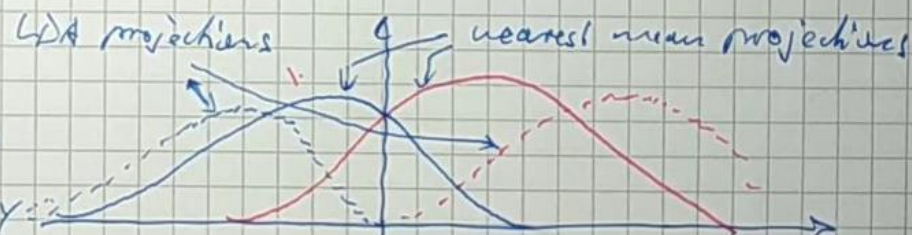
$$\tilde{X} = \underset{\mathbb{R}^D}{w} \cdot \underset{\mathbb{R}^P}{X^T} + b \in \mathbb{R}$$



$D=2$

LDA decision boundary

$$\boxed{w = (\hat{\mu}_1 - \hat{\mu}_0) \cdot \hat{\Sigma}_w^{-1}}$$

effect: rotates the normal
of the decision boundary
(relative to the line through the
means) such that the cluster shape is considered $\Rightarrow$ less overlap

nearest mean classifier $w = \mu_1 - \mu_0$
$\Rightarrow$ data are projected on the line through the means

LDA projections      nearest mean projections

numerical example:

$\mu_1^T = \begin{pmatrix} 3 \\ 0 \end{pmatrix}$  $\mu_0^T = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  $\Sigma = \begin{pmatrix} (3/2)^2 & 0 \\ 0 & (\frac{1}{2})^2 \end{pmatrix}$  $w = (2 \;\; -1) \begin{pmatrix} \frac{4}{9} & 0 \\ 0 & 4 \end{pmatrix} = (\frac{8}{9} \;\; -4)$

$(\mu_1 - \mu_0)^T = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$

$\Sigma^{-1} = \begin{pmatrix} \frac{4}{9} & 0 \\ 0 & 4 \end{pmatrix}$

LDA training alg: • compute mean of each class: $\hat{\mu}_u = \frac{1}{N_u} \sum_{i : y_i = u} x_i$

• compute within-class covariance:

$$\Sigma_w = \frac{1}{N} \sum_{i=1}^{n} (x_i - \mu_{y_i})^T (x_i - \mu_{y_i})$$

always subtract the mean of the appropriate class

• in 2-class case: $\hat{w} = (\hat{\mu_1} - \hat{\mu_0}) \cdot \Sigma_w^{-1}$

determine offset $\hat{b}$ such that the training error is minimized for $\tilde{x} = w x^T + b$