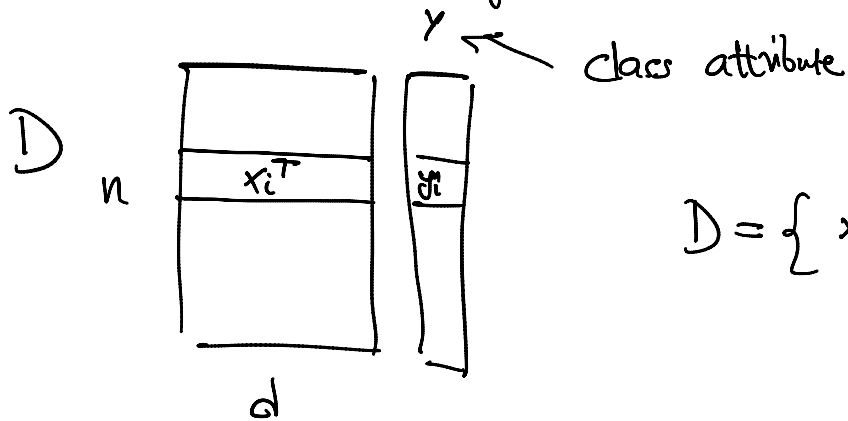


Lecture 6

Thursday, September 14, 2023 9:55 AM

Linear Discriminant Analysis



$$D = \{ x_i^T, y_i \}$$

↑
class

$$y_i \in \{c_1, c_2, \dots, c_k\}$$

y_i is binary

$$y_i \in \{1, -1\}$$

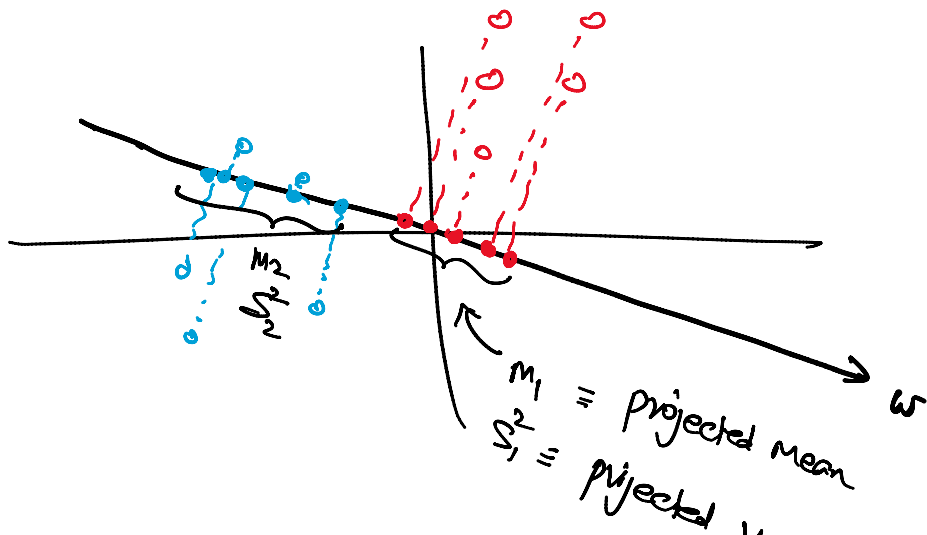
$$\begin{matrix} +1, -1 \\ 0 \quad 1 \end{matrix}$$

$$D_1 = \{ x_i^T : y_i = c_1 \}$$

$$D_2 = \{ x_i^T : y_i = c_2 \}$$

Q: find \vec{w} to separate the two classes

$\mathbb{R}^d \Rightarrow \mathbb{R}$



\rightarrow = projected variance (mean \rightarrow variance \rightarrow projected variance) $(s_1^2 + s_2^2)$

S_i^2 = Scatter
total squared deviation

σ_i^2 = Variance
avg. squared deviation

$$\max_{\vec{w}} J = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$x_i^T \in D_1$
 $x_i^T \in D_2$

$$a_i = \frac{\vec{w}^T x_i}{\sqrt{\vec{w}^T \vec{w}}} = \vec{w}^T x_i$$

$\vec{w}^T \vec{w} = 1$

projected point

$$m_1 = \frac{1}{n_1} \sum_{x_i \in D_1} a_i = \frac{1}{n} \sum_{x_i \in D_1} \vec{w}^T x_i = \vec{w}^T \left(\frac{1}{n} \sum_{x_i \in D_1} x_i \right)$$

$$m_1 = \vec{w}^T \vec{\mu}_1$$

projected mean \rightarrow mean of D_1

$$m_2 = \vec{w}^T \vec{\mu}_2$$

$$\vec{\mu}_1 \equiv \text{mean}(D_1) \in \mathbb{R}^d$$

$$\vec{\mu}_2 \equiv \text{mean}(D_2) \in \mathbb{R}^d$$

$$\vec{w} \in \mathbb{R}^d$$

$$(m_1 - m_2)^2 = (\vec{w}^T \vec{\mu}_1 - \vec{w}^T \vec{\mu}_2)^2$$

$$\begin{aligned}
&= (\vec{w}^T (\vec{\mu}_1 - \vec{\mu}_2))^2 \\
&= \vec{w}^T (\vec{\mu}_1 - \vec{\mu}_2) (\vec{\mu}_1 - \vec{\mu}_2)^T \vec{w} \\
&= \vec{w}^T \begin{pmatrix} B \end{pmatrix} \vec{w}
\end{aligned}$$

$B \equiv$ between class scatter

outer product of the difference vector

$(\vec{\mu}_1 - \vec{\mu}_2)$ with itself

$$\underbrace{S_1^2}_{\text{projected space}} = \text{total deviation} = \sum_{x_i \in D_1} (a_i - \mu_1)^2$$

$$= \sum_{x_i \in D_1} \left(\vec{w}^T \vec{x}_i - \vec{w}^T \vec{\mu}_1 \right)^2$$

$$= \sum_{x_i} \left(\vec{w}^T (\vec{x}_i - \vec{\mu}_1) \right)^2$$

$$= \vec{w}^T \left(\underbrace{\sum (\vec{x}_i - \vec{\mu}_1) (\vec{x}_i - \vec{\mu}_1)^T}_{S_1 = n_1 \cdot \Sigma_1} \right) \vec{w}$$

$$\begin{aligned}
S_2^2 &= \vec{w}^T (S_2) \vec{w} \\
&= \vec{w}^T (n_2 \Sigma_2) \vec{w}
\end{aligned}$$

$$S_1^2 + S_2^2 = \vec{w}^T \left(\underline{S_1} + \underline{S_2} \right) \vec{w}$$

$$= \vec{w}^T S \vec{w}$$

↑
pooled within-class scatter

$$\max_{\vec{w}} \quad J = \frac{\vec{w}^T B \vec{w}}{\vec{w}^T S \vec{w}} = \frac{f}{g} \quad \left(\frac{f'g - g'f}{g^2} \right)$$

$$\nabla_{\vec{w}} = \frac{\partial J}{\partial \vec{w}} = \frac{2B\vec{w}(\vec{w}^T S \vec{w}) - (2S\vec{w})(\vec{w}^T B \vec{w})}{(\vec{w}^T S \vec{w})^2} = 0$$

↑
gradient

(direction where J increases
w.r.t \vec{w})

$$\Rightarrow B\vec{w}(\vec{w}^T S \vec{w}) = S\vec{w}(\vec{w}^T B \vec{w})$$

divide by $\vec{w}^T S \vec{w}$ on both sides

$$\Rightarrow B\vec{w} = \left(\frac{\vec{w}^T B \vec{w}}{\vec{w}^T S \vec{w}} \right) S\vec{w}$$

$$J = \lambda$$

$$\Rightarrow B\vec{w} = J S\vec{w}$$

$$\Rightarrow B\vec{w} = \lambda S\vec{w}$$

generalized eigen equation

$$B\vec{w} = \lambda \vec{w}$$

$$B\vec{w} = \lambda [S] \vec{w}$$

If S^{-1} exists

$$\Rightarrow S^{-1} B \vec{w} = \lambda S^{-1} S \vec{w}$$

$$B\vec{w} = \lambda S \vec{w}$$

$$\Rightarrow \vec{S}^T B \vec{w} = \lambda \vec{S}^T S \vec{w}$$

$$\Rightarrow \boxed{(\vec{S}^T B) \vec{w} = \lambda \vec{w}}$$

Solution: the optimal direction is
the dominant eigenvector of $\vec{S}^T B$ matrix



why? because $\max_{\vec{w}} J$

$\lambda_1 = J$ is the optimal objective value

Simpler solution (w/o eigenvalues)

$$B\vec{w} = \lambda S \vec{w}$$

Solve for \vec{w}

Sol 1: $\vec{S}^T B$ and
find dominant
eigenvector
 \vec{w}

Sol 2

$$B = (\vec{\mu}_1 - \vec{\mu}_2)(\vec{\mu}_1 - \vec{\mu}_2)^T$$

$$\underbrace{\vec{s}_1 \cdot \vec{s}_1^T}$$

Outer product

$$B \vec{v} = \alpha \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$

α \swarrow Some scalar

$$B \vec{v} = \lambda \vec{v}$$

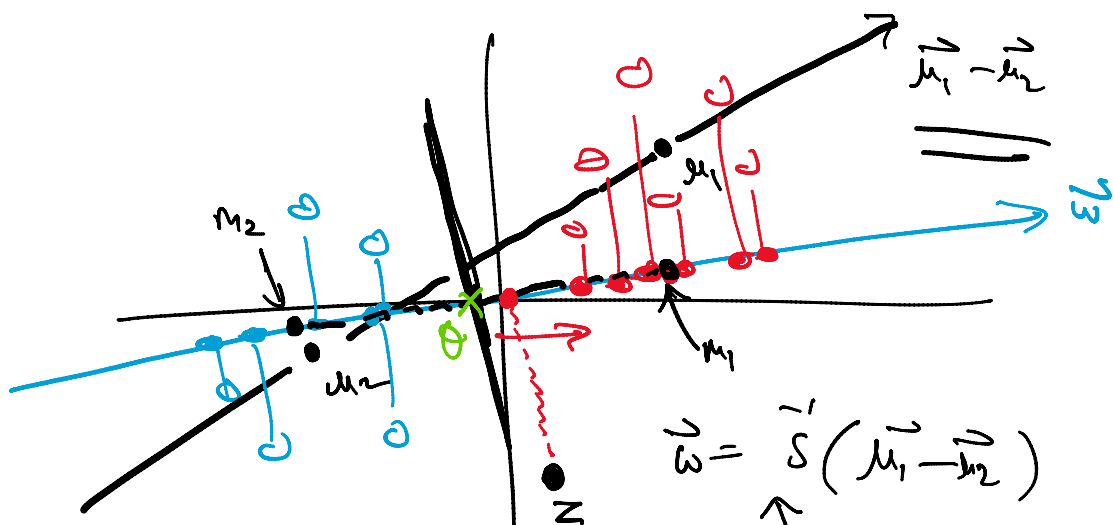
$$\alpha \vec{v} = \lambda \vec{v}$$

$$(\vec{u}_1 - \vec{u}_2) = \frac{\lambda}{\alpha} \vec{v}$$

$$\cancel{\frac{\alpha}{\lambda}} \vec{v} (\vec{u}_1 - \vec{u}_2) = \vec{v}$$

$$\vec{w} = \vec{v} (\vec{u}_1 - \vec{u}_2)$$

the normalize \vec{w} to be a unit vector!





$$\vec{w} = s(\vec{\mu}_1 - \vec{\mu}_2)$$

rotated/scaled difference vector

Solution: \vec{w} : Linear Discriminant Direction

Use to classify : what is the label for \vec{z} (test point)

classifier: $\theta = \frac{1}{2}(\mu_1 - \mu_2) = \frac{1}{2} \vec{w}^T (\vec{\mu}_1 - \vec{\mu}_2)$

$$M(\vec{z}) = \begin{cases} \text{red if } \vec{w}^T \vec{z} \geq \theta \\ \text{blue if } \vec{w}^T \vec{z} < \theta \end{cases}$$

Can be generalized to k classes $= \{C_1, C_2, \dots, C_k\}$

Solution: k eigenvector of $(S^{-1} B')$

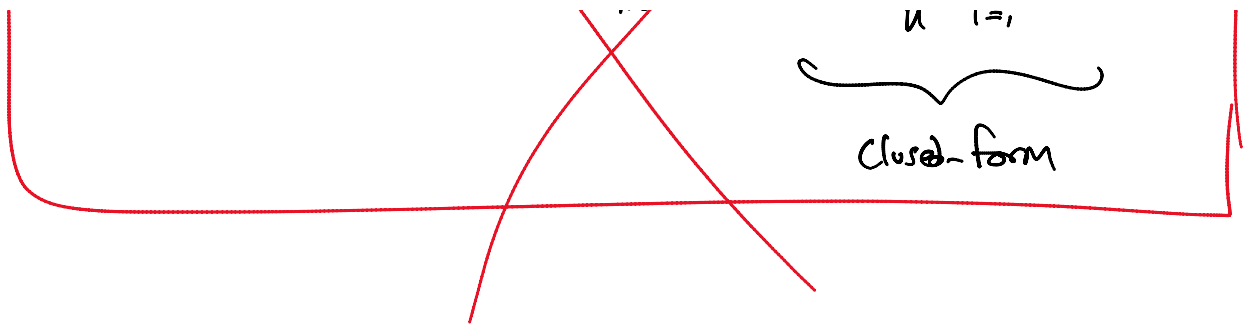
gradient descent algorithm

$$\min_{\vec{z}} J = \sum_{i=1}^n \|x_i - \vec{z}\|^2$$

find an estimate for the mean

$$\nabla_{\vec{z}} J = \frac{\partial J}{\partial \vec{z}} = \vec{0}$$

~~Solution~~ $= \vec{z} = \frac{1}{n} \sum_{i=1}^n x_i$



n is extremely large

∇_z

gradient vector

for all n points

all of D

∇_z

batch gradient descent

$B = n$

$\leftarrow I.$

batch size

$$J = \sum_{i=1}^B \frac{1}{2} \|x_i - z\|^2$$

$B = 1$

$I.$

stochastic gradient descent

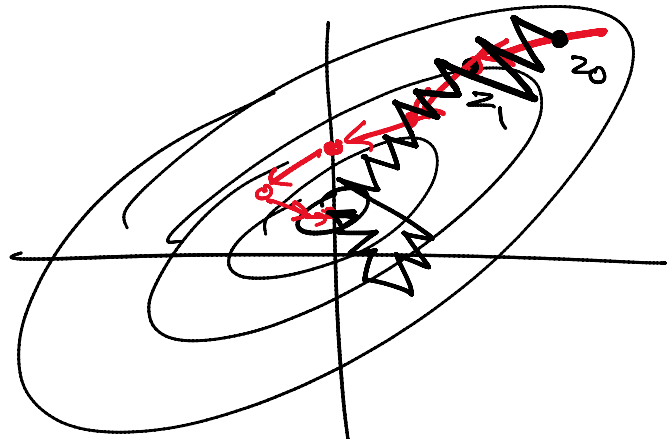
$$J = \frac{1}{2} \|\hat{x}_i - z\|^2$$

Some random point from D

$\nabla_z \leftarrow$ gradient by looking at a single random point

$$\vec{z}_t = \vec{z}_{t-1} - \underbrace{\eta \cdot \vec{\nabla}_z}_{\text{opposite to the gradient by a small step } \eta}$$

opposite to
the gradient
by a small
step η



(when η is very large)

Q1 how many iterations? t

Q2, how much does each iteration cost?

Batch Gradient descent:

$B = n$ (full dataset)

$$\vec{\nabla}_z = O(n)$$

total cost: $O(n) \cdot T$

T : total # of steps

$$T = 10$$

$$[10, O(n)]$$

$$n = 10^9$$

$$\text{total} = 10^{10}$$

SGD:
✓

$O(1)$ per update

T is arising to be large but $T = 10^4$

✓
T is going to be large but $T = 10^4$

total cor = 10^4

mini-batch SGD

$B = 128$

$\frac{\partial J}{\partial \vec{z}} = \nabla_{\vec{z}} = \vec{z} - \vec{x}_i \leftarrow B = 1$

$J = \frac{1}{2} \|\vec{x}_i - \vec{z}\|^2$

$$\begin{cases} J = \frac{1}{2} \sum_{i=1}^B \|\vec{x}_i - \vec{z}\|^2 \leftarrow B \\ \nabla_{\vec{z}} = \frac{\partial J}{\partial \vec{z}} = \sum_{i=1}^B (\vec{z} - \vec{x}_i) = B \cdot \vec{z} - \sum_{i=1}^B \vec{x}_i \end{cases}$$

$\vec{z}_{t+1} = \vec{z}_t - \eta \cdot \nabla_{\vec{z}}$

$\|\vec{z}_{t+1} - \vec{z}_t\| \leq \epsilon$

stop!

large n

10^{-2}



$\epsilon = 10^{-2}$

10^{-2}

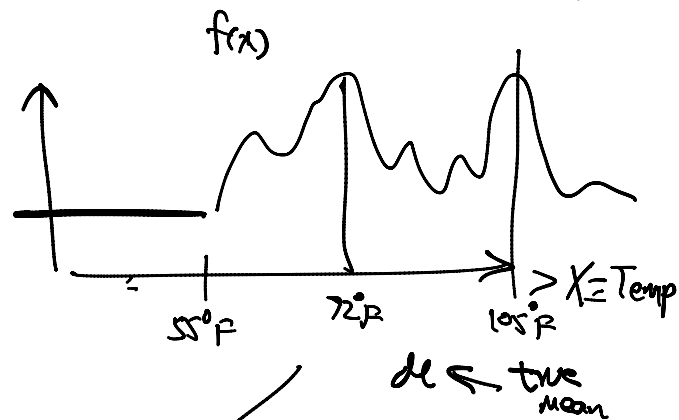


very small n
✓



Probability Distributions

Normal distribution



Unknown prob density function

Sample 1

$$\begin{cases} x_1 = 71 \\ x_2 = 65.5 \\ x_3 = 60.2 \end{cases}$$

Sample 2

⋮

Sample k

$$\Rightarrow \hat{\mu}_1$$

$$\Rightarrow \hat{\mu}_2$$

$$\Rightarrow \hat{\mu}_k$$

$$\mu \in (lb, ub)$$

99% interval

Confidence

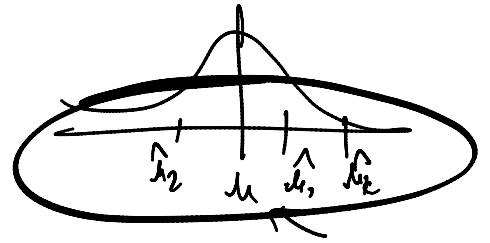
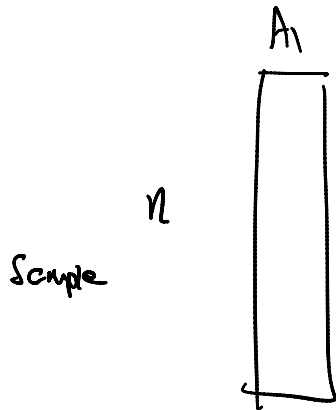
true mean

True
Mean
(Unknown)

Confidence
Interval

Sample $k \Rightarrow \hat{x}_{ik}$

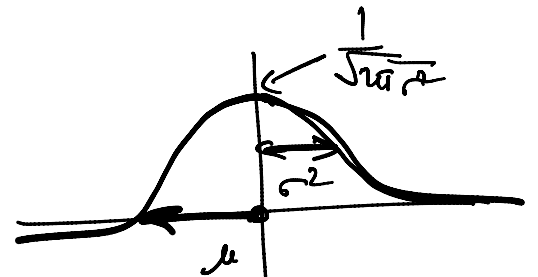
Univariate (1d)



$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

\uparrow mean \uparrow variance
 parameters

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



d-dims (Multivariate)

$$f(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(\sqrt{2\pi})^d (\det(\Sigma))^{1/2}} e^{-\frac{(\vec{x}-\vec{\mu})^T \Sigma^{-1} (\vec{x}-\vec{\mu})}{2}}$$

$\mathbb{R}^d \nearrow$ $\mathbb{R}^d \uparrow$ $\mathbb{R}^{d \times d} \uparrow$

$\Sigma^{-1} \equiv$ Connection to PCA