

Linear Discrimination

classification $\Rightarrow \mathcal{X} = \{(x_i, y_i)\}_{i=1}^N, y_i \in \mathcal{Y}$

$\left. \begin{matrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_K(x) \end{matrix} \right\}$ score functions

$$\hat{y}_{N+1} = \arg \max_{c=1}^K g_c(x_{N+1})$$

$$g_c(x) = p(x | y=c) \Pr(y=c)$$

univariate
($x_i \in \mathbb{R}$)

$$N(x_i; \mu_c, \sigma_c^2)$$

$$\hat{\mu}_c, \hat{\sigma}_c^2$$

multivariate
($x_i \in \mathbb{R}^D$)

$$N(x_i; \mu_c, \Sigma_c)$$

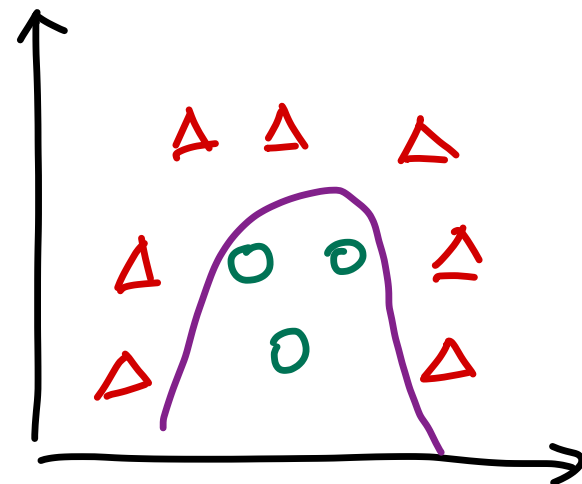
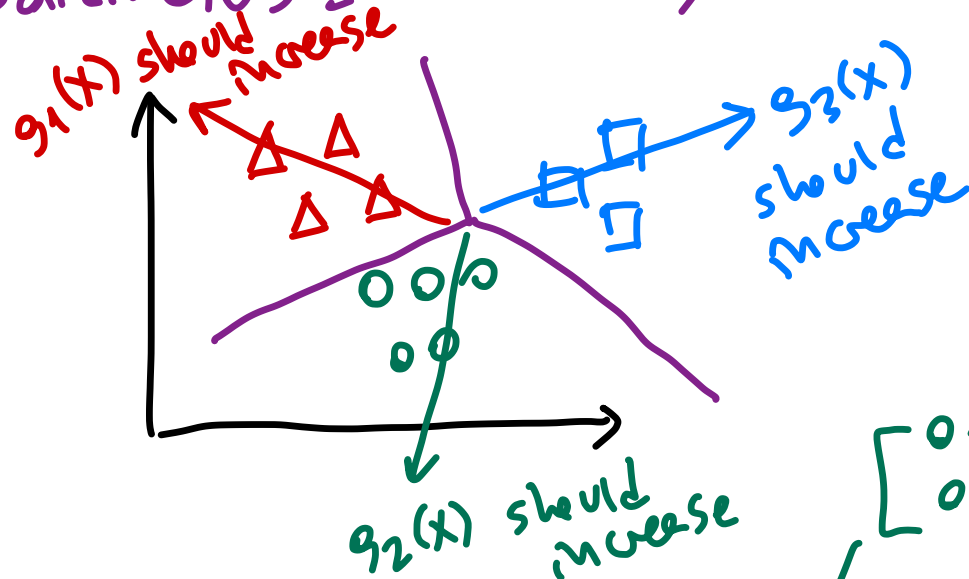
$$\hat{\mu}_c, \hat{\Sigma}_c$$

$\frac{N_c}{N} \rightarrow$ # of data points in class c
 \rightarrow total # of data points

$$g_c(x | w_c, w_{c0}) = w_c^T \cdot x + w_{c0}$$

$$= [w_{c1} \ w_{c2} \ \dots \ w_{cD}] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} + w_{c0} = \sum_{d=1}^D w_{cd} x_d + w_{c0}$$

of parameters = $K \cdot (D+1)$



$$g_c(x | W_c, w_c, w_{c0}) = x^T \cdot W_c \cdot x + \underbrace{w_c^T \cdot x}_{\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}} + \underbrace{w_{c0}}_{\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}}$$

of parameters = $K \left(\frac{D \cdot (D+1)}{2} + D+1 \right)$

$$a^2 + 2ab + b^2$$

$$[a \ b] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

Binary Classification ($K=2$)

$$\begin{array}{ll} 83 & g_1(x) \\ 79 & g_2(x) \end{array} \left. \vphantom{\begin{array}{l} 83 \\ 79 \end{array}} \right\} \begin{array}{ll} \text{if } g_1(x) > g_2(x) & \Rightarrow \hat{y} = 1 \\ \text{if } g_2(x) > g_1(x) & \Rightarrow \hat{y} = 2 \end{array}$$

+4
-3

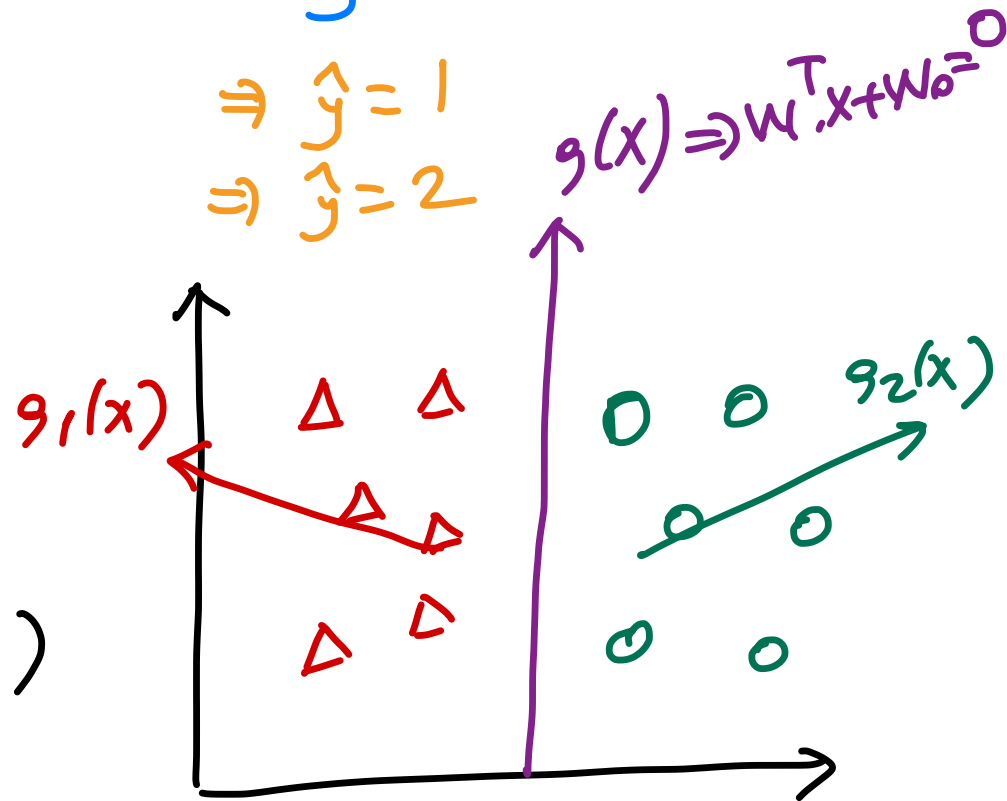
$$\begin{array}{ll} \text{if } g(x) > 0 & \Rightarrow \hat{y} = 1 \\ \text{if } g(x) < 0 & \Rightarrow \hat{y} = 2 \end{array}$$

$$\begin{array}{ll} \text{if } g(x) > 0 & \Rightarrow \hat{y} = 1 \\ \text{if } g(x) < 0 & \Rightarrow \hat{y} = 2 \end{array}$$

$$g_1(x) = w_1^T \cdot x + w_{10}$$

$$g_2(x) = w_2^T \cdot x + w_{20}$$

$$g_1(x) - g_2(x) = \underbrace{(w_1 - w_2)}_w^T \cdot x + \underbrace{(w_{10} - w_{20})}_{w_0}$$



Multiclass Classification ($K > 2$)

73 tennis players

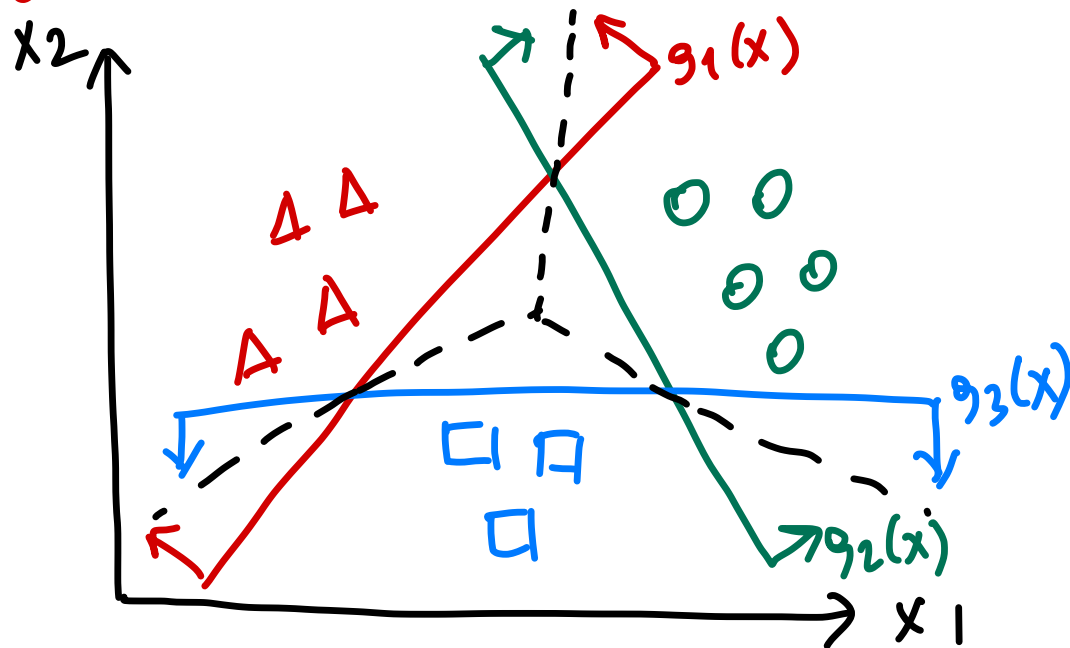
$g_1(x)$
 $g_2(x)$
 \vdots
 $g_K(x)$

assume $K=3$

if $g_1(x) > g_2(x) \wedge g_1(x) > g_3(x) \Rightarrow \hat{y} = 1$
 if $g_2(x) > g_1(x) \wedge g_2(x) > g_3(x) \Rightarrow \hat{y} = 2$
 if $g_3(x) > g_1(x) \wedge g_3(x) > g_2(x) \Rightarrow \hat{y} = 3$

$$\hat{y} = \arg \max_{c=1}^K g_c(x)$$

ONE-VERSUS-ALL (OVA) APPROACH



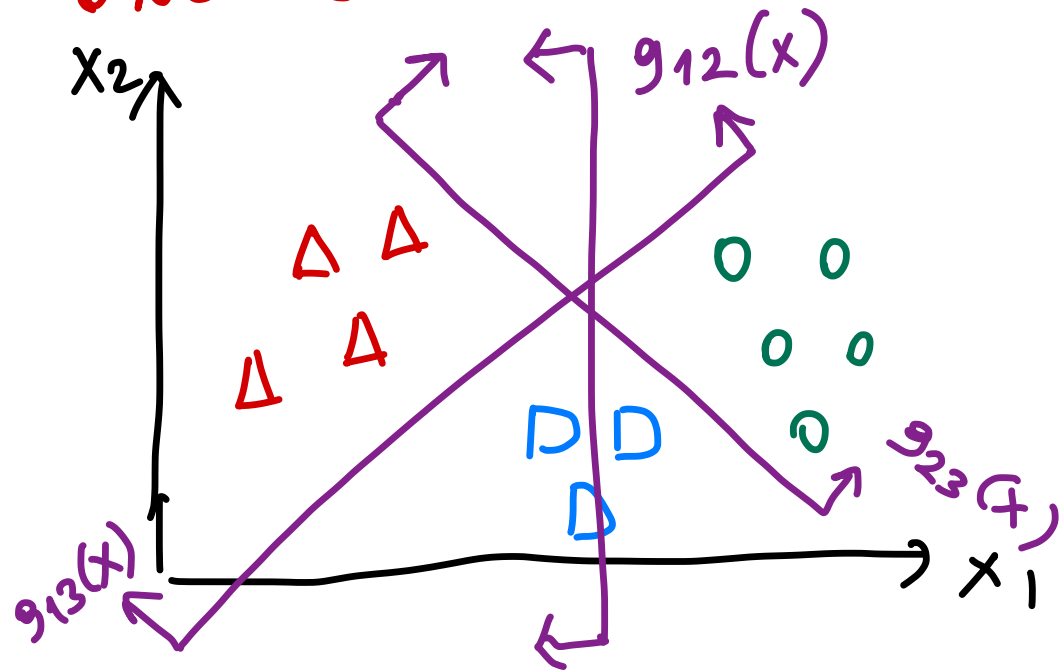
+	-
$\Delta \Delta \Delta$	$\bigcirc \bigcirc \bigcirc \bigcirc \square \square \square$
$\bigcirc \bigcirc \bigcirc \bigcirc$	$\Delta \Delta \Delta \Delta \square \square \square$
$\square \square \square$	$\Delta \Delta \Delta \Delta \bigcirc \bigcirc \bigcirc \bigcirc$

of parameters = $K(D+1)$

data set size = N

of score functions = K

ONE-VERSUS-OTHER (OVO) APPROACH



$\underbrace{\quad\quad\quad}_+$	$\underbrace{\quad\quad\quad}_-$
$\Delta\Delta\Delta\Delta$	$\circ\circ\circ\circ\circ$
$\Delta\Delta\Delta\Delta$	$\Diamond\Diamond\Diamond$
$\circ\circ\circ\circ\circ$	$\Diamond\Diamond\Diamond$

$$\# \text{ of parameters} = \frac{K(K-1)}{2} \cdot (D+1)$$

$$\text{data set size} \approx \frac{2N}{K}$$

$$\# \text{ of score functions} = \frac{K(K-1)}{2}$$

x	1	2	3
$g_{12}(x)$	1	0	\emptyset
$g_{13}(x)$	1	\emptyset	0
$g_{23}(x)$	\emptyset	0	1
$\# \text{ of wms}$	2	0	1

K=2

$$\Pr(y=1|x) = \delta$$

$$\Pr(y=2|x) = 1-\delta$$

odds ratio \leftarrow if $\delta > 0.5 \Rightarrow \hat{y} = 1$
log odds ratio \leftarrow if $\frac{\delta}{1-\delta} > 1 \Rightarrow \hat{y} = 1$
if $\log\left[\frac{\delta}{1-\delta}\right] > 0 \Rightarrow \hat{y} = 1$

$$\log\left[\frac{\Pr(y=1|x)}{\Pr(y=2|x)}\right] = \log\left[\frac{\frac{p(x|y=1) \Pr(y=1)}{\cancel{p(x)}}}{\frac{p(x|y=2) \Pr(y=2)}{\cancel{p(x)}}}\right]$$

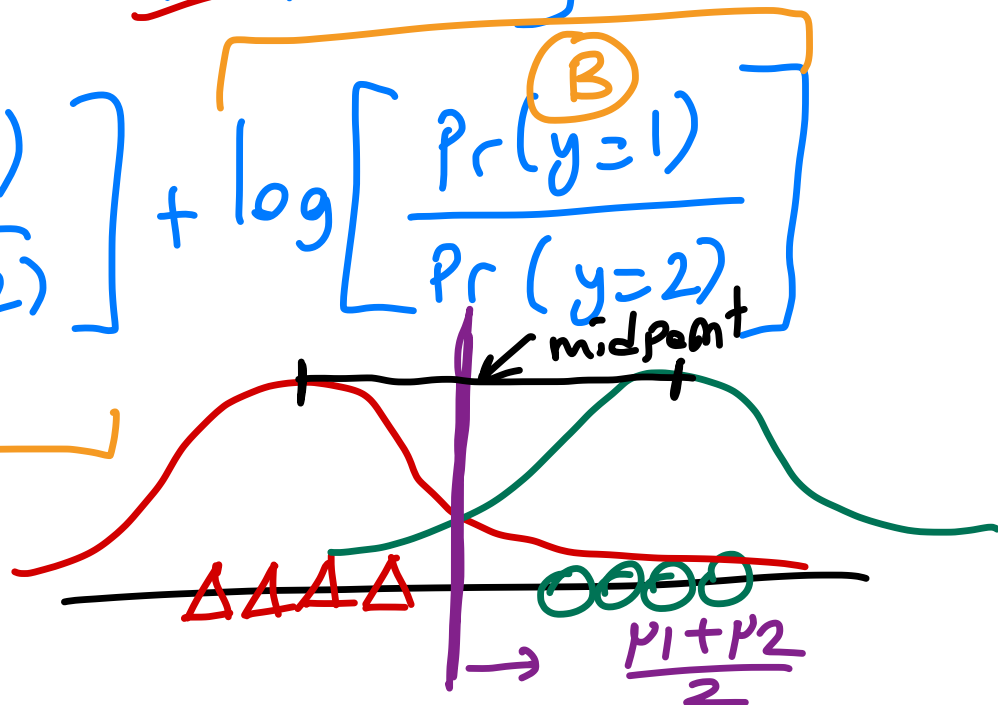
$$\log\left(\frac{a \cdot c}{b \cdot d}\right) = \log\left(\frac{a}{b}\right) + \log\left(\frac{c}{d}\right)$$

$$= \log\left[\frac{p(x|y=1)}{p(x|y=2)}\right] + \log\left[\frac{\Pr(y=1)}{\Pr(y=2)}\right]$$

(A) (B)

assuming

$$\Sigma = \Sigma_1 = \Sigma_2, \Pr(y=1) = \Pr(y=2)$$



Shared covariance assumption $\Sigma_1 = \Sigma_2 = \Sigma$

$$\frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \cdot \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] = N(x; \mu, \Sigma)$$

$$\log \left[\frac{\frac{1}{\sqrt{(2\pi)^D |\Sigma_1|}} \cdot \exp \left[-\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right]}{\frac{1}{\sqrt{(2\pi)^D |\Sigma_2|}} \cdot \exp \left[-\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right]} \right] + \log \left[\frac{\Pr(y=1)}{\Pr(y=2)} \right]$$

$\frac{\exp(a)}{\exp(b)} = \exp(a-b)$

$$= -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) + \log \left[\frac{\Pr(y=1)}{\Pr(y=2)} \right]$$

$$= \underbrace{-\frac{1}{2} x^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} x}_{\text{linear}} - \underbrace{\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1}_{\text{constant}} + \underbrace{\frac{1}{2} x^T \Sigma^{-1} x - \mu_2^T \Sigma^{-1} x + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2}_{\text{linear}} + \log \left[\frac{\Pr(y=1)}{\Pr(y=2)} \right]_{\text{constant}}$$

$$= (\mu_1 - \mu_2)^T \Sigma^{-1} x + \left[-\frac{1}{2} (\mu_1 + \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) + \log \left[\frac{\Pr(y=1)}{\Pr(y=2)} \right] \right]$$

$$= \underbrace{W^T \cdot X}_{D \times 1} + \underbrace{W_0}_{D \times 1}$$

$$\underbrace{W}_{D \times 1} = \underbrace{\hat{\Sigma}^{-1}}_{D \times D} \cdot \underbrace{(\hat{\mu}_1 - \hat{\mu}_2)}_{D \times 1}$$

$\hat{\Sigma}$ = shared covariance estimate
(calculated on the whole dataset)

$$\underbrace{W_0}_{1 \times 1} = -\frac{1}{2} \underbrace{(\hat{\mu}_1 + \hat{\mu}_2)^T}_{1 \times D} \underbrace{\hat{\Sigma}^{-1}}_{D \times D} \cdot \underbrace{(\hat{\mu}_1 - \hat{\mu}_2)}_{D \times 1} + \underbrace{\log \left[\frac{\hat{p}_r(y=1)}{\hat{p}_r(y=2)} \right]}_{1 \times 1}$$