

Linear Classification Methods

$\mathcal{Y} =$ discrete target space $\{1, \dots, K\}$ say.

$G(x) =$ predictor (i.e. $h_{\theta}(x)$)

For many methods the decision boundaries are linear in input variables - these are "linear methods".

If I fit linear model on K^{th} class,

$$\hat{f}_k(x) \text{ of } \hat{\theta}_{k0} + \vec{\theta}_k^T \vec{x}$$

$$\text{then } \{x \mid \hat{f}_k = \hat{f}_l\} = \{x \mid (\hat{\theta}_{k0} - \hat{\theta}_{l0}) + (\vec{\theta}_k - \vec{\theta}_l)^T x = 0\}$$

which is an affine space

i.e. the predictors are divided into regions separated by hyperplanes (a translated hyperplane).

One approach is to do linear regression on indicator matrix, namely,

$$\vec{Y} \in \mathbb{R}^K \quad \text{and} \quad Y = \begin{bmatrix} - & - & - \\ - & Y^{(1)T} & - \\ - & Y^{(2)T} & - \\ \vdots & \vdots & \vdots \\ - & Y^{(m)T} & - \end{bmatrix}$$

indicator response in design matrix of zeros & ones.

$$\text{Perform simple regression } \hat{\vec{Y}} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T Y$$

$\hat{\vec{Y}} = (p+1) \times K$
each column of Y is a response vector
and each column of $\hat{\vec{Y}}$ is a predicted response vector

$$\underline{X} = \begin{bmatrix} 1 & -x^{(1)T} \\ \vdots & \vdots \\ 1 & -x^{(k)T} \\ \vdots & \vdots \\ 1 & -x^{(n)T} \end{bmatrix}$$

$$\text{Set } \underline{\hat{B}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$

then Given a new \vec{x} , classify according to

(1) Compute $\hat{f}(\vec{x})^T = (1, x^T) \hat{B} \in \mathbb{R}^K$

(2) i.e. project \vec{x} using \hat{B}

(2) Output largest component

$$\hat{G}(\vec{x}) = \underset{k \in G}{\text{argmax}} \hat{f}_k(\vec{x})$$

Why? $E[Y_k | X=x] = P(G=k | X=x)$

and, we saw earlier, that linear regression matches this part of model framework.

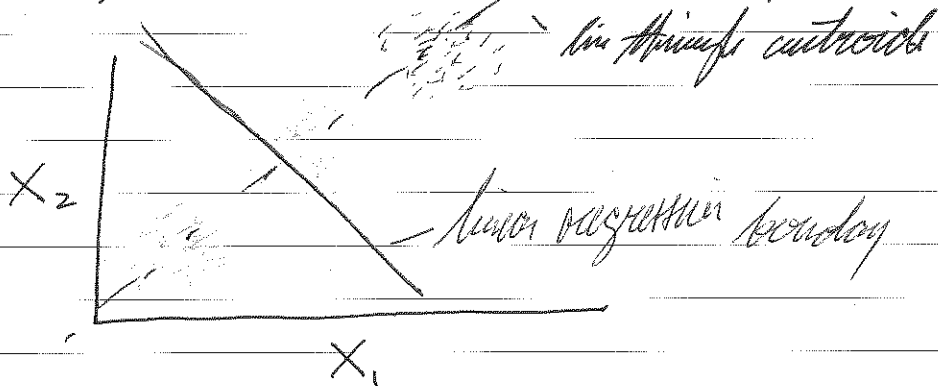
Of course $|\hat{f}_k(x)| > 1$ could happen, but that's not necessarily problematic.

Equivalently, can fit $\min_{\underline{B}} \sum_i \|y_i - ((1, x^{(i)T}) \underline{B})^T\|^2$

and then $\hat{G}(x) = \underset{k}{\text{argmin}} \| \hat{f}(x) - \vec{e}_k \|^2$

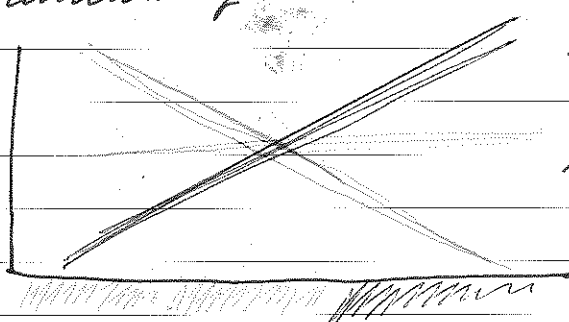
When $K = \# \text{ classes}$ is big a subtle problem can occur.

newly classes can be "invaded" by other classes.



linear regression failed to identify the middle class even though data is linearly separable.

Not in projected space, projecting onto line through the 3 centroids of data we have.



Linear Discriminant Analysis LDA

Suppose want to know $P(G_k | X)$ where $f_k(x)$

is class-conditional ($f_k(x) = P(X | G_k)$) & $\pi_k = \text{class prior}$, $\sum \pi_k = 1$

$$\text{Then } P(G_k | X = x) = \frac{f_k(x) \pi_k}{\sum_{j=1}^K f_j(x) \pi_j}$$

ie. knowing $f_k(x)$ is, basically, solving the problem.

Suppose each class conditional is multivariate Gaussian

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}$$

Assume a common Covariance Matrix $\Sigma_k^{-1} = \Sigma^{-1} \forall k$

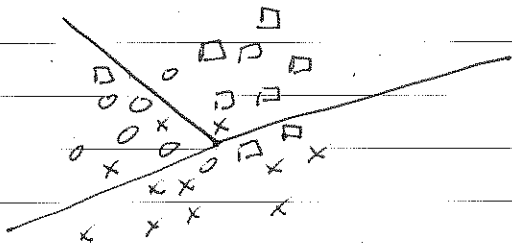
then can look at the log-odds function

(in general $\frac{P}{1-P}$ is "odds" so $\log \frac{P}{1-P}$ is "log-odds")

$$\log \frac{P(G=k|X=x)}{P(G=l|X=x)} = \log \frac{\pi_k}{\pi_l} + \log \frac{f_k(x)}{f_l(x)}$$

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + x^T \Sigma^{-1} (\mu_k - \mu_l)$$

But then we get the decision boundary in linear
i.e. hyperplane decision boundaries.



- Decision boundaries are NOT
the \perp bisectors of
the centroids.

Leads to linear discriminant function

$$S_k(x) = x^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k + \log \pi_k$$

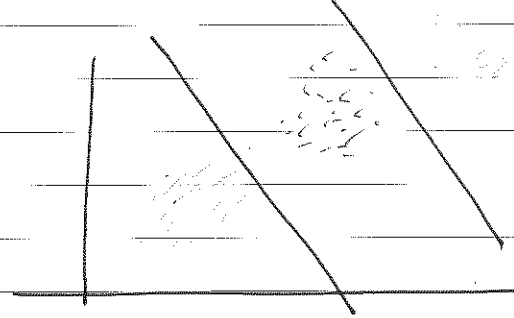
So that $G(x) = \underset{k}{\operatorname{argmax}} S_k(x)$

Of course, need to estimate the parameters

- $\hat{\pi}_k = \frac{M_k}{M} = \% \text{ of observations in } k^{\text{th}} \text{ class.}$

- $\hat{\mu}_k = \frac{1}{M_k} \sum_{g_i=k} x_i$

- $\hat{\Sigma}_k = \frac{1}{M - K} \sum_{k=1}^K \sum_{g_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$



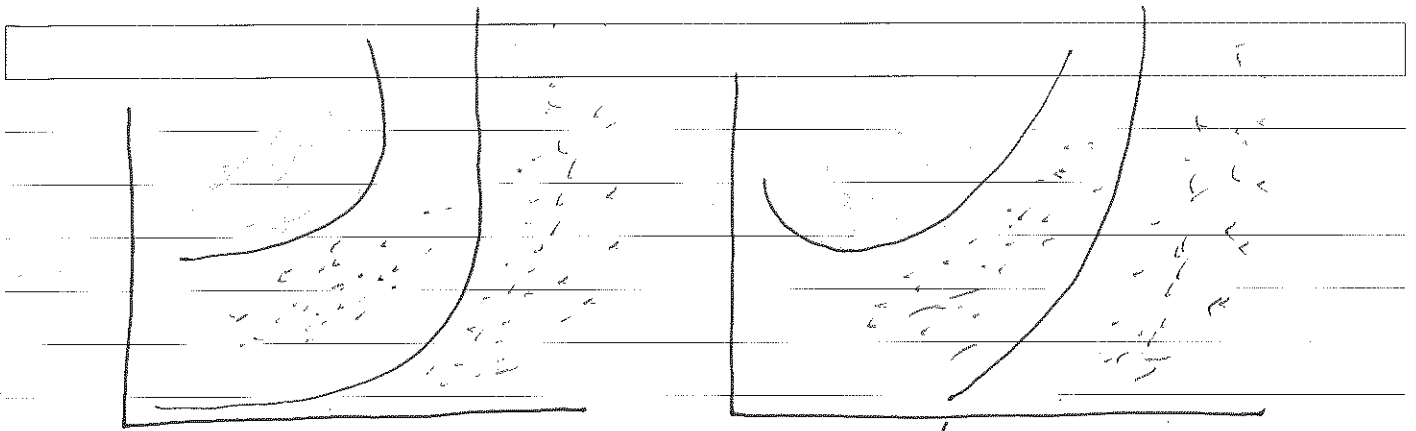
- LDA boundaries avoid masking

If not all Σ_k 's are assumed same, then we don't get the cancellation from before & instead we get

$$S_k(x) = \log \pi_k - \frac{1}{2} \log \det \Sigma_k - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)$$

so that $\{x | S_k(x) = S_l(x)\}$ is quadratic

is QDA quadratic discriminant analysis



QDA decision boundaries

LDA boundaries
in 5-dimensional
($x_1, x_2, x_1^2, x_2^2, x_1 x_2$)
space.

Very similar decision boundaries, but QDA
is advised in general.

Why use LDA/QDA?

Data not Gaussian & it still works well.
That's bias-variance. Accept Gaussian
bias to reduce variance in more complicated
models.

Can Regularize LDA/QDA via

$$\hat{\Sigma}_k(\alpha, \gamma) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}_k(\gamma)$$

$$\hat{\Sigma}_k(\gamma) = \gamma \hat{\Sigma}_k + (1-\gamma) \sigma^2 \mathbf{I} \quad \gamma \in [0, 1]$$

This gives a nice tradeoff between LDA & QDA