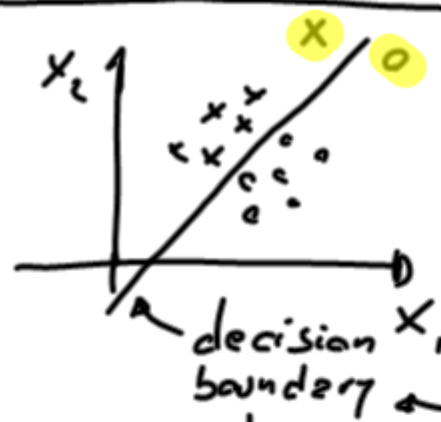
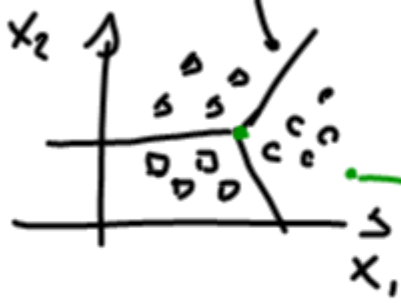


Linear models for classification



• we saw that we can divide the input space in regions, and assign a label to each of them

when linear linear methods for classification



$$\hat{f}_0(x) = \hat{f}_a(x) = \hat{f}_c(x)$$

$$\hat{f}_0(x) > \hat{f}_a(x) \wedge \hat{f}_c(x) > \hat{f}_0(x)$$

until now: classification based on linear regression

$$\forall k \quad \hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

decision boundary between two classes k and e

$$x: \hat{f}_k(x) = \hat{f}_e(x)$$

discriminant functions $\delta_k(x)$

$$\left. \begin{aligned} f_A(x) &= \Pr[G=A | X=x] \\ f_B(x) &= \Pr[G=B | X=x] \\ f_C(x) &= \Pr[G=C | X=x] \end{aligned} \right\}$$

posterior probabilities

it is a member of a class of methods, namely methods based on discriminant functions

If $\delta_k(x)$ or $\Pr[G=k | X=x]$ are linear in $x \rightarrow$ linear decision boundaries

Actually, we only need monotone transformation of $\delta_k(x)$ or $\Pr[G=k|X=x]$ to be linear

Examples

(ii) $\hat{\delta}_k(x) = \hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_{k1}x_1 + \hat{\beta}_{k2}x_2 + \hat{\beta}_{k3}\underline{x_1} + \hat{\beta}_{k4}\underline{x_2}$

the relationship is linear in the augmented ^{input} space, but the decision boundaries are quadratic in the original space

(iii) when there are two classes

$$\Pr[G=1|X=x] = \frac{\exp(\beta_0 + \beta^T x)}{1 + \exp(\beta_0 + \beta^T x)}$$

$$\Pr[G=2|X=x] = \frac{1}{1 + \exp(\beta_0 + \beta^T x)}$$

logit transformation

$$\log \frac{\Pr[G=1|X=x]}{\Pr[G=2|X=x]} = \beta_0 + \beta^T x$$

Linear regression of an Indicator Matrix
- codify each of the classes $1, \dots, K$ with an indicator variable

ex. $K=3 \Rightarrow$

	class 1	2	3	
\Rightarrow	$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$	$= Y$	$N \times K$	

1st obs. is of class 1
2nd and 3rd " " 2
4th " 5th " 3

→ linear regression: $\hat{Y} = X (X^T X)^{-1} X^T Y$

$N \times K$ $N \times p$ $p \times p$ $N \times p$ $p \times 1$ $N \times K$

new observation x_{new} $\hat{y} = \hat{f}(x_{\text{new}}) = x_0 \hat{\beta}$

$1 \times p$ $p \times K$

$$\left(\hat{f}_1(x_{\text{new}}) \quad \hat{f}_2(x_{\text{new}}) \quad \hat{f}_3(x_{\text{new}}) \right)$$

$$\hat{G}(x_{\text{new}}) = \underset{k}{\operatorname{argmax}} \hat{f}_k(x)$$

Why does it work

$$E[Y_k | X=x] = Pr[G=k | X=x]$$

Are $\hat{f}_k(x)$ reasonable estimates of $Pr[G=k | X=x]$?

Yes and no

→ if intercept is included $\sum_{k=1}^K \hat{f}_k(x) = 1$

→ $\hat{f}_k(x)$ can be < 0 or > 1 → problems happen when the new observation is outside the training hull, due to the rigidity of linear regression

many times it works despite this issue

A bigger problem is the so called masking effect

- only if $K \geq 3$

As we saw in Fig 4.3, when $K=3$, it is sufficient to have a quadratic rule.

More generally, if we have K classes, we need a curve of degree $K-1$

LINEAR DISCRIMINANT ANALYSIS

- From the decision theory (Section 2.4) we know that for optimal classification we need to know the class posterior $\Pr(G=k|X=x)$

Suppose:

- $f_k(x)$ is the density of x conditional to class $G=k$
 $\Pr[X=x|G=k]$
- $\pi_k(x)$ is the prior probability to be in class k , $\Pr[G=k]$

Then

$$\Pr[G=k|X=x] = \frac{\Pr[X=x|G=k] \Pr[G=k]}{\Pr[X=x]} = \frac{f_k(x) \pi_k(x)}{\sum_{e=1}^K f_e(x) \pi_e(x)}$$

$= \frac{\sum \Pr[X=x|G=k] \Pr[G=k]}{\sum \Pr[X=x|G=k] \Pr[G=k]}$

We can choose $f_k(x)$ and $\pi_k(x)$ as we prefer...

- when $f_k(x)$ is from a multivariate Gaussian distribution, then Linear Discriminant Analysis (LDA)
Quadratic Discriminant Analysis (QDA)

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

In particular, for LDA we suppose $\Sigma_k = \Sigma \quad \forall k$

We can then compare two classes by the log-ratio

$$\log \frac{Pr[G=k|X=x]}{Pr[G=e|X=x]} = \log \frac{\frac{1}{(\hat{n}_k)^K |\Sigma|^{K/2}} \exp\left\{-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right\} \hat{n}_k / D}{\frac{1}{(\hat{n}_e)^K |\Sigma|^{K/2}} \exp\left\{-\frac{1}{2} (x - \mu_e)^T \Sigma^{-1} (x - \mu_e)\right\} \hat{n}_e / D}$$

$$= \log \frac{\hat{n}_k}{\hat{n}_e} - \frac{1}{2} \left(\cancel{x^T \Sigma^{-1} x} - \cancel{2 x^T \Sigma^{-1} \mu_k} + \mu_k^T \Sigma^{-1} \mu_k - \cancel{x^T \Sigma^{-1} x} + \cancel{2 x^T \Sigma^{-1} \mu_e} - \mu_e^T \Sigma^{-1} \mu_e \right)$$

Scalar
= $\mu_k^T \Sigma^{-1} \mu_k$

$$= \log \frac{\hat{n}_k}{\hat{n}_e} - \frac{1}{2} (\mu_k + \mu_e)^T \Sigma^{-1} (\mu_k - \mu_e) + x^T \Sigma^{-1} (\mu_k - \mu_e)$$

Note

- the decision boundaries are NOT the perpendicular bisectors of the segments joining the centroids (happens only if $\Sigma = \sigma^2 I$)
- the linear discriminant function $\delta_k(x)$ is

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \hat{n}_k$$

$$\rightarrow \hat{G}(x) = \arg \max_k \delta_k(x)$$

Since we do not know the values of the parameters \hat{n}_k, μ_k, Σ , we need to estimate them

$$\bullet \hat{n}_k = N_k / N = \frac{\text{\# observations in class } k}{\text{total \# of observations}}$$

$$\bullet \hat{\mu}_k = \sum_{g_i=k} x_i / N_k$$

$$\bullet \hat{\Sigma} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \mu_k) (x_i - \mu_k)^T / (N - K)$$

With two classes, there is a simple correspondence between LDA and classification by linear regression (\rightarrow Ex. 4.2)

With $K > 2$, there are substantial differences

LDA does NOT suffer from the masking effect.