

14. Classification

- qualitative variable, ~~same~~ also referred to as ~~quantitative~~ categorical variable.
- classification: approach for predicting qualitative responses.
- classifier: classification technique
 - ↳ 1. Logistic regression
 - 2. Linear discriminant analysis
 - 3. Quadratic discriminant analysis
 - 4. Naive Bayes
 - 5. K-nearest neighbors

4.2 Why not linear regression?

- the categorical data is not ordered \nparallel linear regression is applied using coded dummy variables, different set of predictions will be produced for each encoding.
- The above is especially the case when there are more than two categories which do not have any natural ordering.

Reasons :

- (a) A regression method cannot accommodate a qualitative response with more than two classes.
- (b) a regression method will not provide meaningful estimates of $P(y|x)$ - even for two classes.
ex: $P(y_1|x_0) > 1$ or $P(y_1|x_0) < 0$ for some x_0 .

4.3. Logistic regression

Let, $y = \begin{cases} 1 & ; \text{category 1} \\ 0 & ; \text{category 2} \end{cases}$

$P(y=1|x) = p(x)$: response, x : predictor,

If consider: $p(x) = \beta_0 + \beta_1 x$, \rightarrow odds probabilities may not make sense

We must model $p(x)$ using a function that gives outputs between 0 to 1.

$$\text{logistic function } p(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

$$\text{odds} := \frac{p(x)}{1-p(x)} = e^{B_0 + B_1 x}, x \in (-\infty, \infty).$$

β_0 & β_1 indicate very low & very high probabilities of response variable.

$$\text{log odds / logit} \log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x.$$

interpretation: increasing x by one unit changes log odds by β_1 , \rightarrow it multiplies them odds by e^{β_1} .

change $p(x)$ due to one unit change in x , depends on the current value of x .

but irrespectively of value of x :

$$\text{if } \begin{cases} (i) \beta_1 > 0 \Rightarrow x \uparrow \text{ then } p(x) \uparrow \\ (ii) \beta_1 < 0 \Rightarrow x \uparrow \text{ then } p(x) \downarrow. \end{cases}$$

Regression coefficients can be estimated by maximising the likelihood function: $L(\beta_0, \beta_1)$

$$L(\beta_0, \beta_1) = \prod p(x_i) \prod (1 - p(x_i))$$

$i: y_i = 1 \quad i: y_i = 0$

prediction: $\hat{p}(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$

increase in one unit of x is associated with an increase in log odds by β_1 .

4.3.4 Multiple logistic regression.

predicting a binary response using multiple predictors (like age, sex, etc.)

$$\log \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p.$$

where $x = (x_1, x_2, \dots, x_p)$ are p predictors.

$$p(x) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$$

• This is done by maximizing the likelihood function to get estimates.

4.4.4: Multinomial logistic regression has ~~two~~ classes
extension of logistic approach to ~~two~~ 

(no. of classes) ≥ 2 ~~case~~ \Rightarrow multinomial logistic regression.

Select a single class to be baseline : k^{th} class.

$$P(Y=k | X=x) = \frac{e^{B_{k0} + B_{k1}x_1 + \dots + B_{kp}x_p}}{1 + \sum_{l=1}^{K-1} e^{B_{l0} + B_{l1}x_1 + \dots + B_{lp}x_p}}$$

for $k=1, \dots, K-1$ &

$$P(Y=k | X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{B_{l0} + B_{l1}x_1 + \dots + B_{lp}x_p}}$$

if $k=1, \dots, K-1$

$$\log \left(\frac{P(Y=k | X=x)}{P(Y=k' | X=x)} \right) = B_{k0} + B_{k1}x_1 + \dots + B_{kp}x_p$$

- it is not necessary to treat the k^{th} class as baseline
- coefficients will be different, but the predictions will be same (the outputs will be same).

4.4. Generative models for classification

Different approach:

- (1) we model the distribution of the predictors x separately in each of the responses class (i.e. each value of y)
- (2) Use Bayes's theorem to find $P(Y=k | X=x)$ estimates.
if $x \sim N$ then model is similar to logistics.

- Q. Why not logistic regression?
- ① If there is substantial separation between the two classes → parameter estimates for the logistic regression model → unstable.
 - ② If distribution of the predictors X is approx normal in each of the classes & sample size is small.
 - ③ easy extension to more than 2 classes.

→ estimated bay fraction of observation belonging to k^{th} class.

π_k : prior probability that a randomly chosen observation comes from the k^{th} class.

$f_k(x) = P(X|Y=k)$ denote the density function of X for an obs that comes from the k^{th} class.

by Baye's thm:

$$P(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)} \rightarrow \text{estimates}$$

Baye's classifier.

$$\text{let } p_k(x) = P(Y=k | X=x)$$

:= This is the posterior probability that an observation $X=x$ belongs to the k^{th} class, given the predictor value for that observation.

Instead of estimating LHS, we estimate RHS.

Three classifiers that use different estimates of $f_k(x)$ to approximate the Bayes classifier:

(1) Linear discriminant analysis.

(2) quadratic discriminant analysis.

(3) naive Bayes.

- (4) Linear Discriminant Analysis for $p=1$
- $p=1 \Rightarrow$ only one predictor, classification rule
 - object classified to κ^{th} class if $P_k(x)$ is max.

Assumption: $X \sim P_k(x)$ is a Gaussian distribution

$$\textcircled{1} \quad f_k(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{(x-\mu_k)^2}{2\sigma^2}\right\}$$

$$\textcircled{2} \quad \sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2 = \sigma^2$$

$$\therefore P_k(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma^2}\right)$$

$$\text{need to maximize } \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma^2}\right).$$

$$\downarrow \text{maximize: } \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma^2}\right) + \log(\pi_k)$$

equivalent to $\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$

$$\text{boundary: } \delta_1(x) = \delta_2(x) \Rightarrow x \cdot \frac{\mu_1}{\sigma^2} - \frac{\mu_1^2}{2\sigma^2} + \log(\pi_1) = \pi_2.$$

$$x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1^2 + \mu_2^2}{2\sigma^2}$$

$$\text{discriminant function: } \delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

that's: estimates

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$

$$\text{calculus: } \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$$

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i:y_i \neq k} (x_i - \hat{\mu}_k)^2$$

$$\hat{\pi}_k = n_k/n \text{ and } \hat{\mu}_k = \bar{x} \text{ and } \hat{\sigma}^2 = s^2$$

Linear as $\delta_k(x)$ is linear function of x .

Linear discriminant analysis for $p > 1$

(2)

Assumption: $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ drawn $\sim N(\mu_i, \Sigma)$

$$\mu_i = (\mu_{i1}, \mu_{i2}, \dots, \mu_{ip})$$

Σ : common variance covariance matrix

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

observations in the k^{th} class are drawn from a multivariate Gaussian distribution $N(\mu_k, \Sigma)$

→ κ^{th} class s.t. $\delta_k(x)$ is maximum:

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

Boundary: $\delta_k(x) = \delta_l(x)$

$$\delta_k(x) = \delta_l(x)$$

$$x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k = x^T \Sigma^{-1} \mu_l - \frac{1}{2} \mu_l^T \Sigma^{-1} \mu_l$$

sensitivity & specificity are measures of performance of LDA.

4.4.3 Quadratic Discriminant Analysis.

(2) $x \sim N(\mu_k, \Sigma_k)$ (i.e. variance covariance is not common for all classes)
 → observations drawn from each class follow Gaussian distribution.

The Bayes classifier assigns an observation $x = x_i$ to the class for which $\delta_k(x)$ is max where:

$$\phi_{ik}^{(x)} = -\frac{1}{2} (x - \bar{\mu}_k)^T \sum_k^{-1} (x - \bar{\mu}_k) - \frac{1}{2} \log |\Sigma_k| + \log \pi_k.$$



$$= -\frac{1}{2} x^T \sum_k^{-1} x + \frac{1}{2} \sum_k^{-1} \bar{\mu}_k^T \bar{\mu}_k - \frac{1}{2} \log |\Sigma_k| + \log \pi_k.$$

$$-\frac{1}{2} \log |\Sigma_k| + \log \pi_k.$$

→ As first term is quadratic form, QDA gets its name.

* How to decide LDA or QDA?

- No. of parameters to be estimated: for p many predictors.

$$\text{In QDA, no. of parameters} = k \times \left[p + \frac{(p)(p+1)}{2} \right]$$

Classes means. variance
for each x_j Covariance matrix entries.

$$\text{Each class} = k \left[p + \frac{p(p-1)}{2} \right]$$

$$= k \left[\frac{2p + p(p-1)}{2} \right] = \frac{k p (p+1)}{2}$$

$$= k \left[\frac{p(p+1)}{2} \right]$$

$$= k p (p+1) \rightarrow \text{estimate of parameters}$$

In LDA: $k \times p$ per class.

classes, dimension reduced to $p-1$.

⇒ LDA is much less flexible classifier than QDA, and so has substantially lower variance ⇒ improved prediction performance.

- If $\sigma^2 = \sigma_k^2$ assumption is not good then high bias.

- QDA: if dataset is very large. - else LDA.

4.4.

Naive Bayes classifier



Assumption:

within the k^{th} class, the p predictors are independent

mathematically:

$$f_k(x) = f_{k1}(x_1) \times f_{k2}(x_2) \times \dots \times f_{kp}(x_p)$$

where: f_{kj} : is the density function of the j^{th} predictor among observations in the k^{th} class.

- used when n is not very large compared to p .

~~The $P_{kj}(x_j) P_{kj}(x_j)$~~

$$P(Y=k | X=x) =$$

$$\frac{\pi_k \times f_{k1}(x_1) \times f_{k2}(x_2) \times \dots \times f_{kp}(x_p)}{\sum_{j=1}^K \pi_j \times f_{j1}(x_1) \times f_{j2}(x_2) \times \dots \times f_{jp}(x_p)}$$

 $k=1, \dots, K$ To estimate the one-dimensional density functions: f_{kj} (1) x_j is quantitative: $x_j | Y=k \sim N(\mu_{jk}, \sigma_{jk}^2)$ Assumption: within each class, the j^{th} predictor is drawn from a univariate normal distribution.

$$\text{cov}(x_i, x_j) = 0 \quad (\text{if } i \neq j) \quad \sigma_{jk}^2 \neq \sigma_{ik}^2$$

(only different variances).

(2) Estimating $f_{kj}(x_j)$ as the fraction of the training observations in the k^{th} class that belong to the same histogram bin as x_j .

OR Kernel density estimator.

(3) x_j is qualitative: proportion of training observations for the j^{th} predictor corresponding to each class.

Comparison of classification methods



4.5

- 4.5.1 An Analytical Comparison
- ↳ generalizes that
 - ↳ a class that needs to be assigned to new observation
if the condition is satisfied

K: baseline

maximizes:

$$\log \left[\frac{P(Y=k|x=x)}{P(Y=K|x=x)} \right] \quad \text{for } k = 1, 2, \dots, K.$$

③ LDA: $x \sim N_k(\mu_k, \Sigma)$.

$$\begin{aligned} \log \left[\frac{P(Y=k|x=x)}{P(Y=K|x=x)} \right] &= \log \left(\frac{\pi_k f_k(x)}{\pi_K f_K(x)} \right) \\ &= \log \left(\frac{\pi_k \exp(-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k))}{\pi_K \exp(-\frac{1}{2} (x - \mu_K)^T \Sigma^{-1} (x - \mu_K))} \right) \\ &= \log \left(\frac{\pi_k}{\pi_K} \right) - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \frac{1}{2} (x - \mu_K)^T \Sigma^{-1} (x - \mu_K) \end{aligned}$$

$$= a_k + \sum_{j=1}^p b_{kj} x_j$$

$$\therefore \text{LDA: } a_k + \sum_{j=1}^p b_{kj} x_j$$

$$\text{Q.DA: } a_k + \sum_{j=1}^p b_{kj} x_j + \sum_{j=1}^p \sum_{l=1}^p c_{kjl} x_j x_l$$

Naive Bayes: $\pi_k + \sum_{j=1}^p g_{kj}(x_j)$



→ Form of generalized linear models.

- (1) LDA is special case of QDA with $C_{kj}=0 \forall j, k$.
- (2) LDA is special case of Naive Bayes with ~~b_{kj}~~ $g_{kj}(x_j) = b_{kj}x_j$.
- (3) ~~Opposed LDA~~, $f_{kj}(x_j) = N(\mu_{kj}, \sigma_j^2)$
then $g_{kj}(x_j) = b_{kj}x_j$

$$\text{where } b_{kj} = \frac{(\mu_{kj} - \bar{\mu}_k)}{\sigma_j^2}$$

Then Naive Bayes is special case of LDA with Σ restricted to diagonal matrix with j th diagonal element equal to σ_j^2 .

- (1) Neither QDA nor ~~Bayes~~ Naive Bayes is a special case of the other.
As QDA has $C_{kj}x_jx_j$ terms : it has potential to be more accurate in settings where interactions among the predictors are important in discriminating between classes.
- Choice of method will depend on (1) the distributions of the predictors
(2) values of π_k .

For multi-nomial logistic regression:

$$\text{eg } \left[\frac{P(Y=k | X=x)}{P(Y=1 | X=x)} \right] = \beta_{k0} + \sum_{j=1}^p \beta_{kj}x_j$$

→ identical to LDA

- KNN is completely non-parametric approach:
no assumptions are made about the shape of decision boundary



- (1) Dominates LDA & logistic when the decision boundary is highly non-linear, provided $n \gg p$.
- (2) $n \gg p$, tends to reduce the bias while incurring a lot of variance.
- (3) n is not very large, p is not very small then QDA.
- (4) KNN does not tell us which predictors are important.

~~Empirical Comparison~~

Poisson regression: $y \sim \text{Poi}(\lambda)$ if y : counts.

neither quantitative nor qualitative

$$\ln(\lambda) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\rightarrow \lambda(x) = e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}$$

Estimation:

$$L(\beta) \text{ with } \prod_{i=1}^n \frac{\lambda(x_i)^{y_i}}{y_i!} \text{ called } (\lambda)$$

where $\lambda(x_i) = e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}$. Generalized Linear Model.

$$f \rightarrow y \rightarrow E(y|x) \rightarrow \eta(y|x) = \beta_0 + \sum \beta_j x_j$$

* Link function (η) applies transformation to

$E(y|x_1, \dots, x_p)$ so that the transformed mean is a linear function of predictors.

$$\eta(x) = x \rightarrow \text{linear}$$

$$\eta(u) = \log\left(\frac{u}{1-u}\right) \rightarrow \text{logistic}$$

$$\eta(x) = \log x \rightarrow \text{poisson}$$