

ISLR book notes: Chapter 3

linear regression: $Y \approx \beta_0 + \beta_1 X + \varepsilon$ w/ $\text{Var}(\varepsilon) = \sigma^2$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

- estimate β_0 and β_1 w/ least squares: $e_i = y_i - \hat{y}_i$

$$\min \sum e_i^2 \rightarrow \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

- unbiased estimators

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

- you estimate the mean of Y as $\hat{\mu} = \bar{y}$

that has $\text{Var}(\hat{\mu}) = \frac{\sigma^2}{n}$

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right] \quad \text{Var}(\hat{\beta}_1) = \sigma^2 \frac{1}{\sum (x_i - \bar{x})^2}$$

↳ variance is smaller if the x_i are more spread out
(easier to fit a line on a longer section)

- estimate for σ : Residual Standard Error RSE

$$\text{RSE} = \sqrt{\frac{\text{RSS}}{n-2}} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}}$$

- If $\varepsilon \sim N(0, \sigma^2)$, then the 95% confidence interval is

$$[\hat{\beta}_1 - 2 \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \text{SE}(\hat{\beta}_1)] \text{ for } \beta_1$$

$$\text{and } [\hat{\beta}_0 - 2 \text{SE}(\hat{\beta}_0), \hat{\beta}_0 + 2 \text{SE}(\hat{\beta}_0)] \text{ for } \beta_0$$

- Hypothesis testing: $H_0: \text{no relationship b/w } X \text{ and } Y \Rightarrow \beta_1 = 0$
 $H_a: \text{some relationship} \Rightarrow \beta_1 \neq 0$

- we want to know how far our estimate $\hat{\beta}_1$ is away from 0.
How far is enough depends on the standard error

→ **t statistic**: $t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)} \quad t \sim T_{n-2}$ distribution
 if $|t|$ is big, we have a steep regression line and no correlation $\approx N$ is $n > 50$

→ **p value**: probability of observing $|t|$ as big if $\beta_1 = 0$

small p value: unlikely that $\beta_1 = 0$

reject null hypothesis

so they are correlated

Assessing how good the regression is:

• RSE (residual error)

Residual standard error:

= average amount the response deviates from the true line

$$RSE = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\frac{1}{n-2} RSS}$$

• measures lack of fit in the context of the data / depends what is an acceptable RSE

• R^2 is similar to RSE but scaled to the data
so is universal

$$RSS = \sum (y_i - \hat{y}_i)^2$$

$$TSS = \sum (y_i - \bar{y})^2$$

$$R^2 = \frac{TSS - RSS}{TSS} = \frac{\sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

↳ total sum of squares

- TSS measures the inherent variance in y regardless of the regression
- RSS the variability unexplained after the regression

$$0 \leq R^2 \leq 1$$

not good good
- model not good
 $\approx \sigma^2$ is high

- depends what a good R^2 -value is

- in simple linear regression, $R^2 = \text{Corr}(X, Y)$, where

$$\text{Corr}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$$

Multilinear regression: $y = \beta_0 + \sum_{i=1}^p \beta_i x_i + \varepsilon$

Q1 Is at least one of the predictors x_i useful?

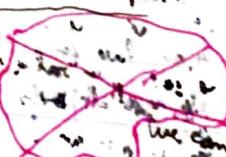
$H_0: \beta_i = 0 \quad \forall i$, no correlation

F-statistic: $\frac{(TSS - RSS)/p}{RSS/(n-p-1)}$

If H_0 is true, F-statistic is close to 1. If H_0 not true, i.e. $(\frac{TSS - RSS}{p}) > \sigma^2$

If H_0 is not true, F-statistic is > 1 . How large it needs to be depends on n, p

F-statistic $\sim F_{(p-1, n-p-1)}$ distribution, so we can calculate a p-value



We can prove that $E(RSS/n-p-1) = \sigma^2$

and if H_0 false, $E(\frac{TSS - RSS}{p}) = \sigma^2$

- If we only want to fit test H_0 that only some β_1, \dots, β_q are 0, we fit another model that doesn't use them and calculate RSS_0 (RSS under hypothesis H_0)

then we use the F-statistic

$$F = \frac{(RSS_0 - RSS) / q}{RSS / (n-p-1)}$$

- If $q=1$, this F-statistic is the same as the t-statistic² for that variable in the full model
- So this reports the partial effect of adding that variable to the model
- (If $q=p$, we get $RSS_0 = TSS$ as $\hat{y} = \hat{x}_0 = \bar{y}$)

caveat:

- Why is F-statistic good:
 - If you get many predictors that are not correlated w/ y , by random chance may will have small p-values for their t-scores
 - F-statistic adjusts to number of predictors
- This only works if $p < n$, otherwise can't even fit a line

Q2 Which are the useful variables?

- Various methods to select best model: Mallows' Cp, AIC, BIC, adjusted R²

How to select predictors?

- Forward: start w/ null, then add param that minimises RSS for continue until stopping rule (greedy)

- Backward: start w/ all, remove that w/ largest p-value continue until stopping rule

- Mixed: start w/ null, add until one of the p-values go up (forward) remove largest p-value (backward)

continue until its p-value low enough and add more variables would raise them

Q3 Model fits the data how well?

→ R^2 close to 1. model explains a large portion of the variance

$$R^2 = \text{Cor}(\hat{y}, \hat{y})^2 \text{ in multiple regression}$$

- adding more variables always increases R^2 towards 1, as we can fit better w/ more variables
- this can result in overfitting

→ RSE can increase because of the $n-p-1$ in the denominator when we add a new param

→ plot the data

you can for example spot non-linearity b/w some variables

Q4)

How accurate predictions can you make?

→ we can compute a confidence interval for our reducible error

how close \hat{y} is to $f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$, our best possible linear model

→ there is always some model bias ~~because of noise~~
e.g. linearly approximates reality

→ → even because of E is irreducible error

we can calculate prediction intervals that also incorporate E

- use confidence interval to quantify the uncertainty of average predictions
- use prediction interval for a particular prediction

What to do with qualitative variable?

→ Race has levels

→ use dummy variables e.g. if ~~Race has levels~~ Asian, Caucasian and Afro-American
do $x_{i1} = 1$ if Asian $x_{i2} = 1$ if Caucasian

$$\text{w/ } y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$$

and interpret β_0 as avg for Afro-Americans, β_1 is diff b/w Asians and Afro-Am, and β_2 as diff b/w Caucasians and Afro-Am.

one fewer dummy variable than levels. The last w/o it is the baseline

Extraneous to linear model

additive:

if only x_i changes, the response of y is indep of the values of other x_j .

linear:

if x_i changes one unit, the change it induces in y is constant regardless of the value of x_j .

removing additive assumption:

- There may be synergy b/w variables, an interaction

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

↳ interaction term.

Hierarchical principle:

If we add an interaction term, we should also include the main effect terms even if their p-values are not significant.

removing linear assumption:

polynomial regression - explained later

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

Problems:

1. Non-linearity of data

plot residual plots: $y_i - \hat{y}_i$ vs x_i for single linear regression

$y_i - \hat{y}_i$ vs \hat{y}_i for multilinear

- There should be no pattern identifiable in the plot

• if there is, data might not be linear, try non-linear transformation

2. Correlation of error terms

- If ϵ_i are correlated, our model will do worse

• often a problem for time series

• plot residuals for time series, if no pattern then error terms are not correlated

Problem 3

Non-constant variance of error terms

- heteroscedasticity

- if $\text{Var}(\epsilon_i) = \sigma^2$ is not true \rightarrow sometimes bigger y has bigger error
- can see this from funnel shape on residual plot
- possible solution is to log y or \sqrt{y}



Problem 4

Outliers

outliers don't usually affect the fit a lot

help to see them if we plot residuals vs fitted values

to see what is an outlier, we calculate the standardized residuals

$$\text{which is } e_i / \text{SE}(e_i)$$

$\text{SE}(e_i) :$

$$\hat{y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1}}_{= H} X^T y$$

$$e = y - \hat{y} = (I - H)y$$

$$\text{var}(e) = \sigma^2(I - H)$$

$$\text{SE}(e_i) = \sigma \sqrt{1 - h_{ii}}$$

Problem 5

High leverage points

- high leverage points affect the regression a lot
- unusual x values
- difficult to notice a high leverage

h_{ii} = leverage statistic quantifies this \rightarrow it is the same h_{ii} as here

$$\text{for simple linear regression, } h_{ii} = \frac{1}{n} + \frac{(x_{ii} - \bar{x})^2}{\sum (x_{ij} - \bar{x})^2}$$

$$\frac{1}{n} < h_{ii} < 1$$

$$\text{avg}(h_{ii}) \text{ is } \left(\frac{p+1}{n}\right) \text{ for all observations}$$

\therefore if $h_{ii} > \frac{p+1}{n}$, that point x_i has a high leverage

Problem 6 Collinearity

Predictors are closely related to one another

reduces the accuracy of the predictors

which increases the errors and makes the p-values worse
and we might even reject H_0 b/c of it

Looking at correlation matrix doesn't help if 3 or more are correlated

we detect collinearity w/ the

Variance Inflation Factor

$= \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_j) \text{ when full model}}$

no collinearity $\Leftrightarrow 1 \leq VIF < \infty$

$\text{Var}(\hat{\beta}_j) \text{ when only } x_j \text{ fit}$

If VIF > 5 ~ 10, problems w/ collinearity

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_{x_j|x_{-j}}^2}$$

where $R_{x_j|x_{-j}}^2$ is from regression of x_j onto all other predictors
(you fit all other priors to determine x_j , and you look R^2 from it)

If $R_{x_j|x_{-j}}^2 \approx 1$, collinearity, then

solutions:
a, drop one of the predictors

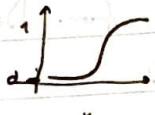
b, combine them into a single value e.g. by taking their mean

Exercises:

Classification → want to predict qualitative variable (0 or 1)
trying to assess the probability of a outcome

Logistic regression:

$$P(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$



$P(X)$ = probability of $y=1$ given X

[odds]: $\frac{P(X)}{1-P(X)} = e^{\beta_0 + \beta_1 X}$ e.g. if $P(X)=0.2 \rightarrow 1 \text{ out of } 5 \rightarrow 1/4 \text{ odds}$

[logistic logit]: $\log\left(\frac{P(X)}{1-P(X)}\right) = \beta_0 + \beta_1 X$

fitting it. we maximize $L(\beta_0, \beta_1) = \prod_{i:y_i=1} P(x_i) \prod_{i:y_i=0} (1 - P(x_i))$ maximize the likelihood fn

assessing if variables are related:

rescale z-statistic for $\hat{\beta}_1$: $\frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$

large |z-statistic|: evidence against $H_0: \beta_1 = 0$

also you can calculate p-value from z-statistic (similar to t-statistic for linear regression)

Multiple logistic regression:

$$\log \frac{P(X)}{1-P(X)} = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

Linear Discriminant analysis

Logistic regression is not always good enough:

- when classes are well separated, LDA is better
- when n is small and predictors X ~ normal, LDA is better
- when we have 2 response classes

π_k = prob. that y is category k \rightarrow prior prob. that an observation belongs to k

$$f_{\pi_k}(x) = P_{\pi_k}(X=x | Y=k)$$

$$\text{Then } P_{\pi_k} = P(Y=k | X=x) = \frac{\pi_k f_{\pi_k}(x)}{\sum_{i=1}^k \pi_i f_{\pi_i}(x)} \text{ from Bayes}$$

$P_{\pi_k}(x) \rightarrow$ posterior prob. that x_{new} belongs to k

We know π_k from our sample (approximately)

\Rightarrow we only need to estimate $f_{\pi_k}(x) \Pr(X=x | Y=k)$

$p=1$ only one predictor

• assume $f_{\pi_k}(x) \sim \text{normal}$

$$f_{\pi_k}(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right)$$

• assume $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$

$$P_{\pi_k}(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)}{\sum \pi_i \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_i)^2\right)}$$

$\Rightarrow P_{\pi_k}(x)$ is maximal when

$$\hat{\sigma}_k^2(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k) \text{ is maximal}$$

LDA estimates μ_k , π_k and σ^2 to help maximizing $P_{\pi_k}(x)$:

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

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

$$\hat{\pi}_k = \frac{n_k}{n}$$

so we get $x=x$ assignment to the class $k=1$.

$$\text{if } \hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k) \text{ is the largest}$$

$\hat{\delta}_k$ discriminant function

linear in x

$p > 1$:

Assume (x_1, \dots, x_p) drawn from a multi-variate Gaussian

$\sim N(\mu_k, \Sigma)$ for k th class

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

: algebra

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

linear in x

$\hat{\delta}_k(x)$ is maximal when $\hat{\delta}_k(x)$ is maximal

so decision boundaries are when $\hat{\delta}_k(x) = \hat{\delta}_l(x) \rightarrow$ Bayes decision boundary
"the beneficial test"

① estimate $\hat{\mu}_1, \dots, \hat{\mu}_k$ with some formulas

② plug in into $\hat{\delta}_1(x), \dots, \hat{\delta}_k(x) \rightarrow \hat{\delta}_k(x) = \hat{\delta}_l(x) \rightarrow$ LDA decision boundary
trying to approximate

③ classify x s.t. $\hat{\delta}_k(x)$ is the highest

$$\text{argmax}_{k=1}^K \text{ if } P(Y=k | X=x) \geq P(Y=i | X=x) \forall i$$

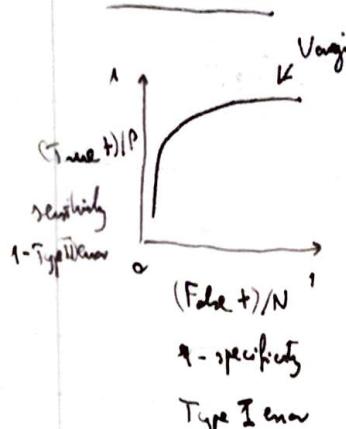
LDA approximates the Bayes classifier, e.g. trying to minimize the misclassification error

we might want to do something else, for example minimize the false positives, or the false negatives.

\rightarrow If $K=2$, Bayes classifier picks a threshold 0.5 for $P(Y=0 | X=x) > 0.5 \Rightarrow y=0$
 $< 0.5 \Rightarrow y=1$

we can make this smaller for less false negatives
 = max true positives

Roc curve:



Area under the curve measures the quality of the classifier

perfect: 1

random: 0.5

what threshold you want to pick depends on
what you need the classifier for

Quadratic discriminant analysis:

- same as LDA, but assumes that $X \sim N(\mu_2, \Sigma_2)$ \rightarrow to each class

we get $\bar{D}_2(x) = -\frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} \log |\Sigma_2| + \log \bar{\pi}_2$ from might be covariance matrix

$$\rightarrow = -\frac{1}{2} x^T \Sigma_2^{-1} x + x^T \Sigma_2^{-1} \mu_2 - \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2 - \frac{1}{2} \log |\Sigma_2| + \log \bar{\pi}_2$$

- we have formulas for $\hat{\mu}_2, \bar{\pi}_2$ and $\hat{\Sigma}_2$

Why use QDA and LDA? when is one better?

QDA estimates $K \cdot \frac{P(p+1)}{2}$ pairs for $\Sigma_1, \dots, \Sigma_K \rightarrow$ quadratic in no. of pairs

In QDA, the model is linear in $x \rightarrow$ K_p linear coeff to estimate

\rightarrow LDA has less pairs \rightarrow less flexible \rightarrow lower variance \rightarrow less overfit
 \rightarrow if covar Σ is wrong \rightarrow can lead to high bias

LDA: when fewer observations and don't want to overfit

QDA: when large training set

or when $\Sigma_1 = \Sigma_2 = \Sigma_K$ is clearly very wrong

Comparison of LDA, QDA, Logistic regression and KNN:

logistic regression and LDA are essentially the same model:

both fit linear decision boundary:

$$\log \left(\frac{P_1}{1-P_1} \right) = \beta_0 + \beta_1 x$$

$$\log \left(\frac{P_1(x)}{1-P_1(x)} \right) = \log \left(\frac{P_1(x)}{P_2(x)} \right) = C_0 + C_1 x$$

where C_0 and C_1 are a function of μ_1, μ_2 and σ^2

the difference is how you fit the params:

w/ maximum likelihood	w/ assuming a normal distribution w/ a common σ^2
better workflow if model assumption is not true	better workflow if this is true

KNN is very different: non-parametric

- better when decision boundary is highly non-linear
- does not tell us which predictors are important

QDA between these above: more flexible boundary than linear, but not as flexible as KNN

- better than KNN of fewer observations

When to use which?

if decision boundary is linear (e.g. predictors are not correlated differently for different classes):
normally distributed: LDA
not normally distributed: Logistic regression

if normally distributed: LDA

if not normally distributed: Logistic regression

if decision boundary is non-linear (e.g. predictors have different correlation for each class):
or depend polynomially on the predictors e.g. x_1^2, x_1x_2

if normally distributed: QDA

if not normally distributed: KNN (cross-validated)

and very complex boundary:

for example

$$\text{Exercises: } 1, \quad p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$\text{log-req.} \quad 1 - p(x) = \frac{1 + e^{\beta_0 + \beta_1 x} - e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

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

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

$$\text{LDA } 2, \quad p_g(x) = \frac{\pi_g \frac{1}{2\sqrt{\pi}\sigma} \exp(-\frac{1}{2\sigma^2}(x-\mu_g)^2)}{\sum \pi_i \frac{1}{2\sqrt{\pi}\sigma} \exp(-\frac{1}{2\sigma^2}(x-\mu_i)^2)} = P(Y=g | X=x)$$

when $P(X=x | Y=g) \sim N(\mu_g, \sigma^2)$

The Bayes classifier is assigning s.t. $p_g(x)$ the largest

This means $\pi_g \exp(-\frac{1}{2\sigma^2}(x-\mu_g)^2)$ is the largest

$\log \pi_g - \frac{1}{2\sigma^2}(x-\mu_g)^2$ is the largest

which means $\log \pi_g + \frac{x\mu_g - \mu_g^2}{\sigma^2}$ is the largest.

This is indeed $\sigma_g^2(x)$.

$$\text{QDA } p=1, \quad 3, \quad \text{QDA } p=1 \quad x \sim N(\mu_g, \sigma_g^2) \text{ in } g\text{th class}$$

$$f_g(x) = \frac{1}{\sqrt{2\pi\sigma_g^2}} \exp\left(-\frac{1}{2\sigma_g^2}(x-\mu_g)^2\right)$$

$$P_g(x) = \frac{\pi_g f_g(x)}{\sum \pi_i f_i(x)}, \quad \text{so Bayes maximises } \pi_g f_g(x)$$

$$\text{so want to maximise } \frac{\pi_g}{\sqrt{\sigma_g^2}} \exp\left(-\frac{1}{2\sigma_g^2}(x-\mu_g)^2\right)$$

$$\text{maximise } \log \pi_g - \frac{1}{2} \log \sigma_g^2 \leftrightarrow -\frac{1}{2\sigma_g^2}(x-\mu_g)^2$$

$$\text{maximise } \log \pi_g - \frac{1}{2} \log \sigma_g^2 \left(-\frac{1}{2\sigma_g^2} x^2 + \frac{1}{\sigma_g^2} x \mu_g - \frac{1}{2\sigma_g^2} \mu_g^2 \right)$$

decision boundary is not linear in x
because of this term

4, $X \sim U(0,1)$

KNN

- a, we want to predict x using all the observations within 10% of the range of x
 closest to the test observation e.g. for $x=0.6$, use $[0.55, 0.65]$

on average, what fraction of the observations will we use?

KNN

$$\rightarrow 10\% = \frac{1}{10}$$

b, $x_1, x_2 \sim U[0,1] \times [0,1]$

within 10% of x_1 and within 10% of x_2

e.g. for $(0.6, 0.35)$, we use $[0.25, 0.65] \times [0.3, 0.4]$

$$\rightarrow 1\% \text{ of the observations} = \frac{1}{100}$$

c, $p=100$

$$\frac{1}{100} \text{ fraction of the observations}$$

d, \rightarrow when p is large, KNN gets weaker because fewer observations will be nearer the point. The space gets bigger and the points get further apart.

e, size of hypercube that contains 10% of the training observations:

$$p=1: \frac{1}{10} \rightarrow \text{length} = 0.1$$

$$p=2: x_1^2 = \frac{1}{10} \rightarrow x_1 = \sqrt{\frac{1}{10}} = \frac{1}{\sqrt{10}} \quad x_1 = \frac{1}{\sqrt{10}} \text{ length} = 0.316$$

$$p=100: x_1^{100} = \frac{1}{10} \rightarrow x_1 = \frac{1}{10^{1/100}} \quad x_1 = \frac{1}{10} \text{ length} = 0.997$$

5/ a, if Bayes boundary is linear: LDA better on test set

QDA might be better on training set as more flexible

LDA/QDA

b, if non-linear: QDA better on both

c, if sample size bigger: we expect QDA to be better than LDA as training sample size increases
 QDA includes LDA within it, so wouldn't be worse, but

d, False: if we don't have a huge sample size, QDA will overfit and be worse than LDA
 i.e. if Bayes is linear

9) x_1 = hours studied

x_2 = undergrad GPA

γ = get an A

logistic regn: $\hat{\beta}_0 = -6$

$\hat{\beta}_1 = 0.05$

$\hat{\beta}_2 = 1$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$$

log. reg.

$$x_1 = 40$$

$$x_2 = 3.5$$

$$\text{prob of getting an A} = e^{-6 + 0.05 \cdot 40 + 3.5}$$

$$\text{logit } f(x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}$$

$$a, \quad p(\text{get an A}) = \frac{e^{-6 + 3.5}}{1 + e^{-6 + 3.5}} = 0.3775$$

$$b, \quad p(\text{get an A}) = 0.5$$

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\log(0) = -6 + 0.05 \cdot 40 + 3.5$$

LDA

7) predict if stock gives dividend

$$\bar{x} = 10 \quad \text{if yes (1)} \quad x: \text{last year profit}$$

$$\bar{x} = 0 \quad \text{if no (0)}$$

$$\sigma^2 = 36$$

80% of companies gave dividends.

$$x_1 = \frac{2.5}{0.05} = 50$$

$$\mu_0 = 0 \quad \sigma^2 = 36$$

$$\mu_1 = 10$$

$$\pi_0 = 0.2$$

$$\pi_1 = 0.8$$

for Bayes:

$$P(\text{yes} | x=4) \approx \frac{\pi_1 \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(4-\mu_1)^2)}{\sum \pi_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(4-\mu_i)^2)}$$

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)}$$

$$P(\text{good } x=4) = \frac{0.4852}{0.4852 + 0.16} = 0.752$$

log reg

8) log reg: 20% on training, 30 on test having me = left side

KNN

KNN: 18% over both

which is better could be log reg. KNN on training, 30 on test having me = left side

KNN less 0% having me, so

an KNN on training and 18 on test is really good

not have 36% test set error \rightarrow log reg is better

9, odd of 0.37 to defaults
on average what fraction defaults?
odds

$$\text{odd: } \frac{p(x)}{1-p(x)} = 0.37 \rightarrow p(x) = \frac{0.37}{1+0.37} = 0.27$$

1, 16% of default.

$$\text{odd: } \frac{0.16}{1-0.16} = 0.19$$

Resampling methods - Chapter 5