## (2) TWO – Introduction to R

| R code | Description | Example |
|---|---|---|
| c( ) | Used to concatecodnate numbers | x <- c(1, 3, 2, 5)<br><br>> x<br>> 1 3 2 5 |
| seq( ) | Creates a sequence of numbers<br>- Useful for defining grid values to plot with | seq(a, b)<br>makes a vectors of integers (starts with *a* and increments by 1 till *b*)<br>seq(0, 1, length = 10)<br>makes a seq of 10 numbers that are **equally spaced** between *0* and *1*<br>1:10 and seq(1, 10) are equivalent<br>seq(-pi, pi, length = 50)<br>creates something useful for a trigo plot |
| ls( ) | Shows a list of all the objects saved in the workspace | rm( list = ls() )<br><br>is useful in removing all objects from the workspace |
| rm( ) | Used to delete objects from the workspace | |
| matrix( ) | Creates a matrix of numbers<br>**(fills in columns first, by default)** | x <- matrix(data = c(1, 2, 3, 4),<br>                     nrow = 2,<br>                     ncol = 2)<br><br>> x<br>> 1  3<br>   2  4<br>x <- matrix(..., ..., ..., byrow = T) |
| rnorm( ) | Generates normal random variables<br>- get pseudorandom numbers that we can treat as random<br>- almost sure to get a different number each time it's called<br>- **by default, ~N(0,1)** | x <- rnorm(50)<br>generates 50 random *x*-s<br>x <- rnorm(50, mean = 50, sd = 0.1) |
| cor( ) | Compute sample correlation | Creating 2 correlated sets of numbers,<br>x <- rnorm(50)<br>y <- x + rnorm(50, 50, 0.1)<br>cor(x, y) |
| mean( )<br>var( )<br>sd( ) | Compute the mean and sample variance of a vector of numbers | The following are equivalent:<br>sd(y) and sqrt(var(y) |
| set.seed( ) | To reproduce same results, 'resets' randomness to a reproducible state<br>- useful in simulation studies | |
| plot(x, y) | Produces a scatterplot of *y* against *x* | plot(x, y,<br>       xlab = "...", ylab = "...",<br>       main = "...") |
| pdf( )<br>jpeg( ) | Saves the output<br>(alternative to using 'Export' function under "Plots" in RStudio) | pdf("filename.pdf") **or** jpeg("name.jpeg")<br>plot(x, y, ...)<br>dev.off( ) |

| Function | Description | Example |
|---|---|---|
| contour( ) | Produces a contour plot like a **topographical map**<br>- needs *x, y* grid coordinates and a matrix *z* for the height values | y <- x<br>f <- outer(x, y,<br>       function(x, y) cos(y)/(1-x^2))<br>contour(x, y, f)<br>contour(x, y, f, nlevels = 45, add = T) |
| image( ) | Similar to *contour()* but gives a **heatmap** instead of contour lines<br>- lighter colours correspond to higher values | fa <- (f-t(f))/2<br>image(x, y, fa) |
| persp( ) | Generates **3D plots**<br>- *thets*: horizontal rotation<br>- *phi*: vertical rotation | persp(x, y, fa, thets = 30, phi = 40) |
| hist( ) | Plots a histogram | hist(mpg)<br>hist(mpg, col = 2, breaks = 15) |
| pairs( ) | Plots a scatterplot matrix for the entire data set<br>- can also specify just a subset of variables<br>(useful many variables) | pair(Auto)<br>pairs( ~ mpg + displacement, Auto) |
| identify( ) | ** **Label points on plot**<br>- returns the row number of points identified | plot(horsepower, mpg)<br>identify(horsepower, mpg, name)<br>then 1. click on interested points<br>      2. click 'finish' |
| [ ] | Indexes data<br>- Not specifiying *row_idx*: all columns<br>- Not specifying *col_idx*: all rows<br>- Negative indexing to omit/remove elements | > A<br>> 1  5  9 13<br>   2  6 10 14<br>   3  7 11  15<br>   4  8 12  16<br>A[c(1, 3), c(2, 4)]<br>where [*row_idx, col_idx*]<br>> 5 13<br>   7 15 |
| dim( ) | Gives the number of rows and columns of a matrix | |
| setwd( ) | Sets working directory | |
| read.table( ) | Reads in a text file and stores it as a *dataframe*<br>- if it is text data, can use external text editor to peek at data to get things right the first time, i.e. input extra arguments (recommended) | auto <- read.table("Auto.data".<br>      header = T, na.strings = "?") |
| fix( ) | Launches a spreadsheet-like window to view data<br>- not recommended<br>- use *View()* instead | |
| na.omit( ) | Removes rows containing missing data (*NA*) | auto <- na.omit(Auto) |
| names( ) | Displays all the variable names in the data frame | names(Auto) |

| | | |
|---|---|---|
| attach( )<br><br>detach( ) | Makes variables in a data frame available by name<br>- good to detach after you are done<br>- no need to use $ to refer to a variable within a data frame | attach(Auto)<br><br>plot(cylinders, mpg)<br><br>detach(Auto) |
| as.factor( ) | **\*\* Converts quantitative variables into qualitative variables**<br>- using *plot()* will generate a boxplot instead now<br>(where half of the points are in the box) | cylinders <- as.factor(cylinders)<br><br>plot(cylinders, mpg, varwidth = T)<br><br>where *varwidth = T* tell the number of points in a bar by its width/"fatness" |
| summary( ) | Provides a numerical summary of each variable in the data frame<br>- can also use on a single variable | |
| q( ) | To quit R<br>- Recommended to save into a workspace to be loaded instead of default workspace that gets quto-loaded | save.image()<br><br>load() |

# (3A) Three A – Linear Regression (simple, multiple)

| library(MASS) library(ISLR) | - *Mass* contains a large collection of data sets and functions (in base R) - *ISLR* includes data sets from the textbook | |
|---|---|---|
| Simple Linear Regression | | |
| lm(y ~ x) | Runs a simple linear regression | model <- lm(data = Boston, medv ~ lstat) |
| model | Gives the basic coefficient estimates - Just calling the regression model - "*Model*" here is just the name of the variable where output of *lm()* is stored | > model<br>> Call:<br>  lm(formula = medv ~ lstat)<br><br>  Coefficients:<br>  (Intercept)       lstat<br>    34.55         -0.95 |
| **summary(model) | Gives more details - Summary of residuals - Standard errors, t-values, p-values of coefficients - Residual standard error, degree of freedom - Multiple R-squared, adjusted R-squared - F-statistic, p-value - ** In SLR, p-value of coefficient (beta1) should be the same as F-statistic | summary(model) |
| names(model) | Get names of all variables in the list where output of *lm()* is a list - Can use $ to extract out the values of these variables | > names(model)<br>> "coefficients" "residuals" "effects" "rank" "fitted.values" "assign" "qr" "df.residual" "xlevels" "call" "terms" "model" |
| model$coefficients model$coef coef(model) | Returns coefficient estimates - Can use shortcut and type till variable name is unique | |
| **confint(model) | Creates confidence intervals for coefficients | > confint(model)<br>>           2.5 %       97.5%<br>(Intercept)  33.448457   35.6592247<br>lstat       -1.026148   -0.8739505 |
| **predict( ) | Get confidence intervals for specified values of *x* or get prediction intervals for specified values of *x* or for training set - Make sure that data frame for specified values of *x* have the **same column name** as in the model | > predict(model,<br>     data.frame(lstat = c(5, 10, 15)),<br>     interval = "confidence")<br>>    fit    lwr    upr<br>  1    ...     ...     ... |

| | | |
|---|---|---|
| | - These values have to be in a data frame because *predict()* is a general function that works for MLR as well (i.e. more variables) | 2   ...   ...   ...<br>3   ...   ...   ...<br>shows the confidence interval for the respective specified *x* values<br>> predict(model, data.frame(...),<br>       interval = "prediction")<br>> predict(model) |
| abline(model) | Adds a line to an existing plot<br>- Useful for adding the least-squares regression line to a scatterplot of the training points | plot(lstat, medv)<br>abline(model, lwd = 3, col = "red") |
| plot(1:20, 1:20, pch = 1:20) | Shows the different symbols available to plot | |
| plot(model) | Generates multiple plots that can be toggled through<br>1. Residuals against fitted values<br>2. QQ plot<br>3. Scale-location<br>4. Residuals against leverage | par(mfrow = c(2, 2))<br>where default is *par(mfrow = c(1, 1))*<br>plot(model)<br>dev.off() closes the plotting device/resets it |
| residuals(model)<br>rstudent( ) | Get residuals plot directly<br>(i.e. residuals against fitted values) | >plot(predict(model),<br>     residuals(model))<br>>plot(predict(model),<br>     rstudent(model))<br>gives standardised residuals (i.e. dividing by standard error) |
| hatvalues(model) | Computes leverage statistics<br>- Works for multiple linear regression objects as well | > plot(hatvalues(model))<br>> which.max(hatvalues(model))<br>to return the index of the largest element (i.e. the largest leverage point) |
| Multiple Linear Regression | | |
| lm(y ~ x1 + x2 + x3 + ...) | Fits a multiple linear regression model | |
| lm( y ~ . , data = ...) | Regress on **all** variables | |
| summary(model)$... | Access individual components of *summary()* | > names(summary(model))<br>   "call" "terms" "residuals"<br>   "coefficients" "aliased" "sigma"<br>   "df" "r.squared" "adj.r.squared"<br>   "fstatistic" "cov.unscaled"<br>> summary(model)$r.squared<br>> summary(model)$sigma |

| (3B) Three B – Other considerations in linear regression | | |
|---|---|---|
| vif(model) | Computes variance inflation factors<br>- Measures collinearity<br>- Overly high VIF values (>10) | |
| var1:var2<br><br>var1*var2 | Interaction terms<br>- Includes the interaction of *var1* and *var2*<br>- Usually *var1*var2* is used because of hierarchical principle | model <- lm(data = Boston,<br>        medv ~ lstat*age)<br>**includes the main effects *lstat* and *age*terms** |
| lm(y ~ .-var1) | Modelling using all variables except one | |
| update(model, ~.-var1) | Update a model that is already fitted<br>- Either adding or removing a variable | |
| lm(data =...,<br>    y ~ f(var1) ) | Non-linear transformation of predictors<br>- Need *I( )* as a wrapper when using ^, which has a different meaning in *lm( )*<br>(i.e. order of interaction effect) | model2 <- lm(medv ~ lstat + I(lstat^2))<br>model3 <- lm(medv ~ log(rm)) |
| anova(model1,<br>       model2) | Compare 2 models<br>- *Null hypothesis*:<br>2 models fit the data equally well<br>- *Alternative hypothesis*:<br>full model fits better<br>- Can be used to compare nested model when adding multiple predictors | |
| poly(var, order) | Used for higher order polynomial regression<br>- Creates an orthogonal basis<br>- Each *poly(var, order)x* have slightly different coefficients and also different t-statistic values<br>- Generally, high degree polynomials are a bad thing (they have a high increase at the end) | model5 <- lm(medv ~ poly(lstat, 5))<br>**creates a 5$^{th}$ order polynomial fit**<br>> summary(model5)<br>> Call:<br>lm(formula = medv ~ poly(lstat, 5))<br><br>Residuals:<br>Min   1Q   Median   3Q   Max<br> ...   ...   ...   ...   ...<br><br>Coefficients:<br>Estimate   Std.Error   t value   Pr(>\|t\|)<br>(Intercept)<br>poly(lstat, 5)1<br>poly(lstat, 5)2<br>poly(lstat, 5)3<br>poly(lstat, 5)4<br>poly(lstat, 5)5<br>... |

## Qualitative Predictors

| | | |
|---|---|---|
| - *lm( )* will code qualitative predictors automatically | | |

Example:
- For *Carseats* data, we model response *Sales* with a number of predictors
- One of them is the shelving location, *ShelveLoc*, is a qualitative predictor with 3 levels – *Good*, *Medium*, *Bad*

```
model <- lm(data = Carseats,
            Sales ~ . + Income:Advertising + Price:Age) # everything and includes interaction terms
```

- In *summary(model)*, *ShelveLoc* appears as 2 dummy variables – *ShelveLocGood* and *ShelveLocMedium*
- There is no need for *ShelveLocBad* since it can be determined from the 2 dummy variables

| contrasts(var) | Shows the current set of contrasts being used <br> - This coding can be changed according to preference (check *?contrasts*) | > contrasts(ShelveLoc) |
|---|---|---|

> contrasts(ShelveLoc)

| > | Good | Medium |
|---|---|---|
| Bad | 0 | 0 |
| Good | 1 | 0 |
| Medium | 0 | 1 |

where *Medium* can be changed to (0, 0) to set it as the benchmark

# (4) Four B – Classification

- Train set:
  - Used to train model
  - Size of set depends on number of testing that will be done
- Validation set:
  - Used to **tune parameters** (e.g. *K*) or **select features** (e.g. how many variables)
  - Size of set depends on method
- Test set:
  - Used to make sure the results so far are valid

- Other tips:
  - If working on a **time series**, might not want to randomise because can give glimpse of future
    (fix first part as train and second part as test *(train data comes before test)*, how large each part is can be chosen)

| cor( ) | Gives correlation matrix <br> - for **quantitative** predictors | |

## Logistic Regression

- Tip: consider fitting smaller models by just including predictors that have the smallest p-values in the full model

Note of caution:

**- Use test set sparingly**

- Repeatedly using unseen test set to evaluate error may result in overfitting
- Similar to making use of test set as training set to *select best classifier*
- If keep testing many different classifiers, bound to have one where results are "good" *by chance*
- One way to validate results is to use yet another untouched test set

| glm(y ~ x1 + x2 + x3 +…, <br><br> data = df, <br><br> family = binomial) | Fits generalised linear models <br> - class of models including logistic regression | |
|---|---|---|
| summary(model) | Provides the usual estimated coefficients, standard errors, p-values etc. | - Also gives "*Number of Fisher Scoring Iterations*" <br> - This is how much more times the processing took compared to LS <br> (logistic regression uses iterated related LS to get solutions) |
| coef(model) | Extracts coefficients | |
| preds <- <br> predict(model, <br> type = "response") <br><br><br> predict(model, newdata = …, <br> type = "response") | Get fitted **probabilities** for the *training* set or predicted probabilities for new *test* data | - If *type = "response"* is not specified, the default for logistic regression would be to output the logit <br> - Note that if you get fitted probabilities, later on this will be used in the calculation of the **training error rate** |
| contrasts(y) | To find out what the probabilities are predicting | > contrasts(Response) <br> >      Up <br>   Down  0 <br>   Up     1 |

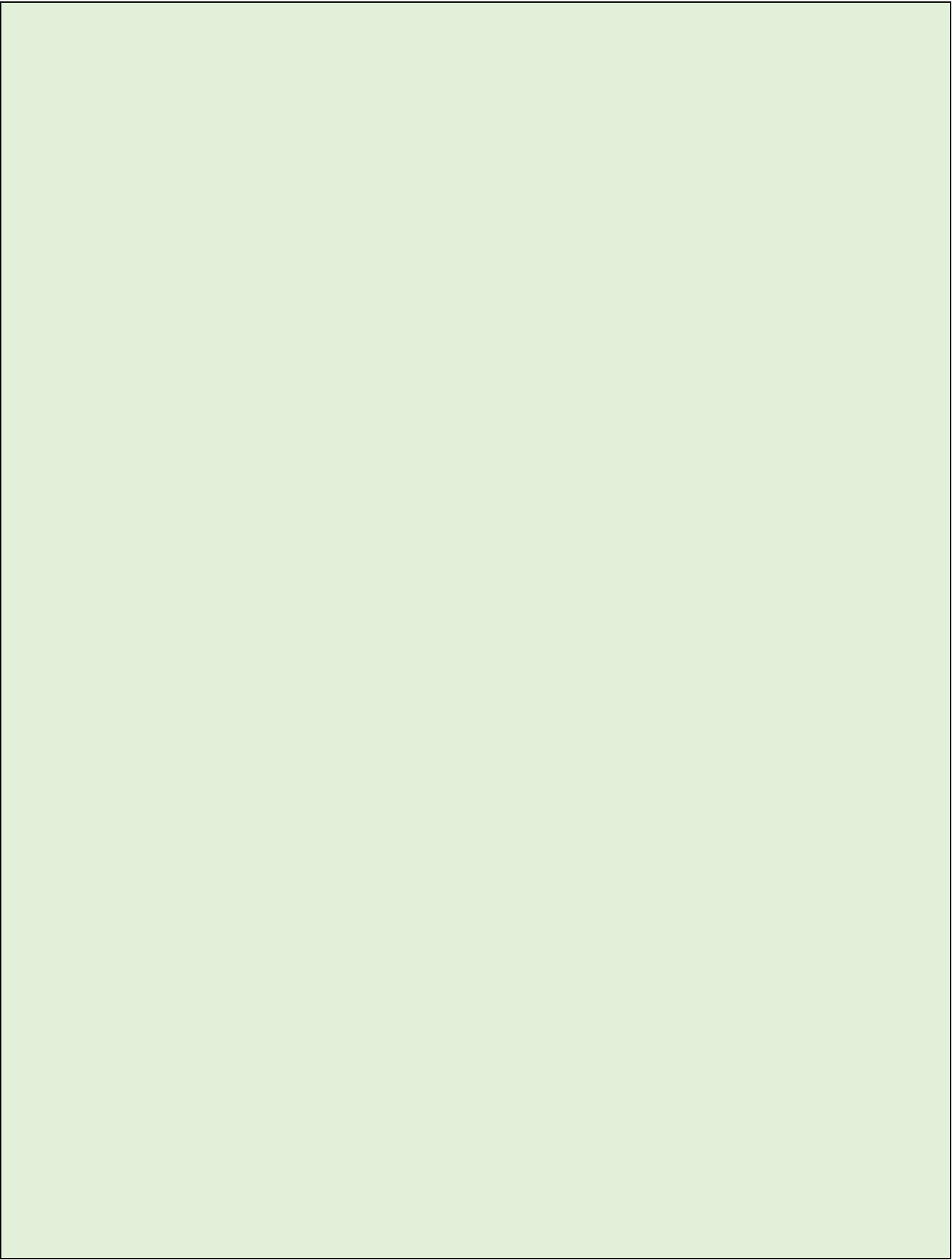| Code | Description | Example |
|---|---|---|
| ifelse(probs > threshold, 1, 0)<br><br>preds <- rep(0, nrow(df))<br>preds[probs > threshold] <- 1 | To get predicted class **labels** by setting a rule | - The second method updates the labels for those whose predicted value is of class '1' |
| con_mat <- table(preds, y) | Creates confusion matrix | > table(pred, Direction)<br>>                Direction<br>  glm.pred   Down     Up<br>  Down        145     141<br>  Up          457     507 |
| mean(preds == y)<br><br>sum(diag(con_mat))/ sum(con_mat) | Gets accuracy rate | |
| df[!train, ] | Subsetting out the test set<br>- This can be done with a vector of indexes, or a **logical** vector | > train <- Year < 2005<br>> Smarket.2005 <- Smarket[!train, ]<br>> Direction.2005 <- Direction[!train, ] |
| glm(y ~ x1 + x2 + x3 +…, data = df, family = binomial, subset = train) | Using logistic regression on a dataframe, separating into train and test sets<br>- Make use of the argument *subset* to specify rows you want for training | |
| Linear Discriminant Analysis | | |
| library(MASS) | *lda()* function is part of the *MASS* library | |
| lda_model <- lda(y ~ x1 + …, data = df, subset = train) | Fit an LDA model | |
| lda_model | Call the model directly to view details about model<br>- Unlike logistic regression, there are no more coefficients<br>- Estimated prior probabilities and means of each class is given directly in output<br><br>- In practice, for computational efficiency, equivalent form of LDA decision rule is computed<br>- For binary case, this is a single vector of coefficients to be multiplied by predictors | > lda_model<br>Call:<br>lda(Direction ~ Lag1 + Lag2, data = Smarket, subset = train)<br><br>Prior probabilities of groups:<br>  Down         Up<br>0.491984     0.508016<br><br>Group means:<br>            Lag1       Lag2<br>Down   0.04279022   0.03389409 |

| | | |
|---|---|---|
| | - Large values of these linear combination corresponds to a higher posterior for the class coded as '1' in *contrasts()*<br>- Lower values correspond to class '0' | *(for days that are 'Down')*<br>Up       -0.03954635    -0.03132544<br>*(for days that are 'Up')*<br><br>Coefficients of linear discriminants:<br>              LD1<br>Lag1     -0.6420190<br>Lag2     -0.05135293 |
| preds <-<br>predict(lda_model,<br>df_test) | Gives predicted class and also posterior probabilities<br>- Can use other thresholds | |
| names(preds) | There are 3 attributes to the predicted output:<br>   1) **class***<br>   2) posterior<br>   3) x (which is used to compute (2)) | > names(lda.pred)<br>[1]  "class"  "posterior"  "x" |
| preds$posterior | Gives posterior probabilities of each class<br>- Use *colnames(preds$posterior)* to determine what each column is predicting for<br>- Different thresholds can be used to recreate the predicted class labels<br>- Note: Sometimes the thresholds should be changed in small increments, especially if the greatest posterior probability is not so big | For a 50% threshold, find the number of observations predicted as 'Down' (col 1)<br>> sum(lda.pred$posterior[, 1] >= 0.5)<br>[1] 70<br>> sum(lda.pred$posterior[, 1] < 0.5)<br>[1] 182<br><br>**Be wary of setting the threshold too high,**<br>> sum(lda.pred$posterior[, 1] > 0.9)<br>[1] 0<br>> max(lda.pred$posterior[, 1])<br>[1] 0.520235 |
| <td colspan="2" align="center">Quadratic Discriminant Analysis</td> | | |
| library(MASS) | *qda()* is also part of the *MASS* library | |
| qda_model <-<br>qda(y ~ x1 + x2 + ...,<br>data = df, subset = train) | Fit a QDA model | |
| qda_model | Output is similar to *lda()* but missing QDA computation terms (coefficient of discriminants) | |
| qda_pred_class <-<br>predict(qda_model,<br>df_test)$class | | |
| table(qda_pred_class) | Creates confusion matrix (similar to above examples) | |
| <td colspan="3" align="center">K-Nearest Neighbours</td> | | |

| library(class) | *knn()* is part of the *class* library | |
|---|---|---|
| train_X <- cbind(x1, x2)[train, ]<br><br>test_X <- cbind(x1, x2)[!train, ]<br><br>train_y <- y[train, ] | There is no need to train in *knn()*, which predicts in a single command that requires:<br>1. Matrix/dataframe containing **variables** from **train** set<br>2. Matrix/dataframe containing **variables** from **test** set (observations to predict)<br>3. Vector containing **class labels** for train set<br>4. **K**, number of nearest neighbours to look at | |
| set.seed(1) | Set a seed for reproducibility since *knn()* **breaks ties randomly** | A tie happens when<br>- 2 neighbours are equally far apart<br>or<br>- *K* is an even number and there are equal number of votes |
| knn_pred <-<br>knn(train_X,<br>    test_X,<br>    train_y,<br>    k = ...) | | |
| mean(test_y != knn_pred) | Gives the error rate (because classifier, either prediction correct or wrong) | # Comparison to a *trivial classifier*<br># ('No' all the time)<br>> mean(test_y != 'No') |
| *** stdised_x <- scale(df) | *** **When predictors are not all on the same scale**<br>- Standardise data such that all variables have mean = 0 and sd = 1<br><br>- Important since KNN uses Euclidean distance which does take into account the scale<br>- Variables on a large scale will have a larger effect on the distance between observations | # Example 1<br># Difference of $1000 in *salary* is enormous compared to difference of 50 years in *age* (but common sense says otherwise)<br><br># Example 2<br># Variable in column 1 is gonna dominate by 1000<br>> var(Caravan[, 1])<br>[1] 165.0378<br>> var(Caravan[, 2])<br>[1] 0.1647078 |
| *Application: Caravan Insurance Data* | | |

# (5) FIVE – Resampling Methods

- Cross-validation is a very general method that can be used for different methods and performance metrics
- To adapt to other settings, knowing how to code out CV will be useful (using *sample()*, *for()*)

| R code | Description | | Example |
|---|---|---|---|
| set.seed( ) | - Set a seed for R's random number generator<br>- Use when performing analysis that contains element of randomness for reproducible work | | set.seed(1) |
| library(boot) | Library needed to do CV and bootstrap | | |
| Validation Set Approach | | | |
| sample( ) | Generate random subset<br>- simple random sample | | train <- sample(n, n/2)<br>takes half of indices as training obs |
| lm.fit <- lm(y ~ x,<br>data = df,<br>subset = train) | Fit model only on training set using *subset()* option | | |
| mean(<br>(test_y - predict(model, df)[-train])**2<br>) | | Calculate test MSE<br>- note: this value is random (if re-run, samples change)<br>- Predict/estimate all responses but only take errors corresponding to validation set | |
| fit( ) | For train set | | |
| predict( ) | For validation set | | |
| lm.fit2 <-<br>lm(y ~ poly(x, poly_num),<br>data = df, subset = train) | Higher order polynomial regression | | ...y ~ var1... (linear)<br>...y ~ poly(var1, 2)... (quadratic)<br>...y ~ poly(var1, 3)... (cubic) |
| Leave-One-Out Cross-Validation (LOOCV) | | | |
| glm(y~ x, data = df) | Without *family = "binomial"*, *glm()* is just the usual linear regression and not logistic regression<br>- Use *glm()* for linear regression to do CV | | |
| library(boot) | Library for cross-validation | | |
| cv.error <-<br>cv.glm(df, model)<br><br><br>cv.error$delta | Do LOOCV using fitted model<br>- *cv.glm()* does not take advantage of shortcut for least squares | | glm.fit <- glm(mpg ~ horsepower, data = Auto)<br>cv.glm <- cv.glm(Auto, glm.fit)<br>cv.err$delta<br><br>[ 1 ] 24.23151    24.23114<br>these are the cross-validation results:<br>1. **CV estimate of MSE**<br>2. **bias-corrected version** (only possible if form of bias is known, not possible to correct for bias analytically for different learning methods and error metrics) |
| k-Fold Cross-Validation (k-fold CV) | | | |

| | | |
|---|---|---|
| cv.glm(y~x, data = df, K = num_folds) | Does k-fold cross-validation<br>- common choice for $k = 10$ or $5$<br>- much faster than LOOCV<br>- each model is only fitted $k$ times | |
| Coding k-Fold CV yourself | | |
| sample(n) | Samples n times from 1 to n without replacement<br>- A useful method when coding out k-fold CV by yourself<br>- Used to randomise **indexes** of dataframe in the beginning<br>- **Don't modify** dataframe<br>- Use indexes to **extract out relevant rows** | |
| index = sample(n)<br>foldbreaks = c(0, floor(n/k*1:k))<br>or<br>folds <- sample(rep(1:k, length.out = n), n) | Getting indexes of observations of each fold<br>- Rearrange the index (not the data) | |
| Bootstrap | | |
| - Estimates the accuracy of a statistic of interest | | |
| library(boot) | | |
| fn <- function(df, index) { ... } | To estimate accuracy of a statistic of interest, write a function that computes statistic given data and relevant row indexes<br>- E.g. *mean, median* etc | mean.fn <- function(data, index) {<br>    X = data$X[index]<br>    Y = data$Y[index]<br>    return(mean(X))<br>} |
| fn(df, sample(n, n, replace = T)) *** | To get bootstrap estimate of statistic, generate a **bootstrap sample index** by sampling with replacement<br>- If run this many times using a loop, can store all the bootstrap estimates<br>- With estimates can look at distributions of estimates or do further analysis | |
| boot(df, fn, R = num) *** | Automates the above procedure<br>- Gets all bootstrap estimates for different bootstrap samples<br>- Computes **common bias** and **standard error** estimates | > boot(Portfolio, alpha.fn, R = 1000)<br>where **R** is the number of bootstrap replicates<br>ORDINARY NONPARAMETRIC BOOTSTRAP<br><br>Call:<br>boot(data = Portfolio, statistic = alpha.fn, R = 1000)<br><br>Bootstrap Statistics |

| | | |
|---|---|---|
| | | original       bias       ==std. error==<br>t1*   0.5758321   -7.315422e-05   0.08861826 |
| fn <- function(df, index) { ...} | To ==estimate accuracy of a linear regression model==, write a ==function== that ==computes== the ==coefficient estimates==<br>- Compare bootstrap estimates with those derived from theory (using *summary(model)$coef*)<br>- Can compute standard error estimates analytically because coefficient estimates are unbiased<br>- Note that the estimates from theory and bootstrap may differ because standard error formulas reply on certain theoretical assumptions<br>(e.g. estimate of variance relies on linear model being correct, but true relationship between X and Y may be non-linear)<br>- If true relationship is non-linear then the coefficient estimates should be compared to a different model<br>(e.g. *lm(y~x\*\*2)* in both the bootstrap function and model for full dataset) | ```<br>boot.fn <- function(data, index) {<br>    lm_model =<br>        lm(mpq ~ horsepower,<br>            data = data, subset = index)<br>    return(coef(lm_model))<br>}<br>```<br><br>> boot.fn(Auto, 1:392) # Full dataset<br> (Intercept)       horsepower<br>39.9358610     -0.1578447 |
| boot( ) | Similar to runing function through a loop with bootstrap samples and recording down bootstrap estimates (in a matrix with a column for each coefficient estimate/beta) | > boot(Portfolio, boot.fn, R = 1000)<br><br>ORDINARY NONPARAMETRIC BOOTSTRAP<br><br>Call:<br>boot(data = Portfolio, statistic = boot.fn, R = 1000)<br><br>Bootstrap Statistics<br>     original    bias    ==std. error==<br>t1*  39.9358610  0.02972191  0.860007896<br>t2*  -0.1578447 -0.00030823  0.007404467 |
| colspan Others - Estimates of statistics | | |
| mean(X) | Estimate of population mean/==sample mean== (*mule_hat*) | |
| sqrt(var(X)/n)<br>sd(X)/sqrt(n) | Estimate of standard error of *mule_hat* | |
| median(X) | Estimate of population median value | |
| quantile(X, 0.1) | Estimate of tenth percentile of X in population | |

## (6A) SIX A – Subset Selection Methods

| R code | Description | Example |
|---|---|---|
| - Note that subset of variables chosen by Best Subset Selection, Forward and Backward Stepwise Selection can be different with different models | | |
| Best Subset Selection | | |
| library(leaps) | | |
| regsubsets(y ~ x, data = df) | Identifies best model for a given number of predictors, quantified by $SS_{Res}$<br>- Similar syntax to *lm( )*<br>- * indicates that the variable is included in the corresponding model | > regfit_full <- regsubsets(Salary ~ ., data = Hitters)<br>> summary(regfit_full)<br><br>By default, regsubsets() reports up to 8 best variable model – this can be adjusted:<br>> regsubsets(Salary ~., data = Hitters, nvmax = 19) |
| names(reg.summary) | Summary of best subset model also returns other measures of performance | > names(reg.summary)<br>[ 1 ] "which"  "rsq"  "rss"  "adjr2"  "cp"<br>[ 6 ] "bic"  "outmat"  "obj" |
| reg.summary$rsq | Returns $R^2$<br>- Naturally increases as number of predictors increase, due to greater flexibility | |
| | Plot various measure for all models at once<br>- Visualises results<br><br>Additional:<br>- Argument *type = "l"* connected plotted points with lines<br>- *points()* adds points on a plot already created | > par(mfrow = c(2, 2))<br>1) To plot SSRes:<br>> plot(reg.summary$rss, type = "l")<br>2) To plot adjusted R2:<br>> plot(reg.summary$adjr2, type = "l")<br>To find and label the maximum point:<br>> which.max(reg.summary$adjr2)<br>[ 1 ] 11<br>> points(11, reg.summary$adjr2[11], cex = 2, pch = 20)<br>3) To plot Mallows' Cp:<br>> plot(reg.summary$cp, type = "l")<br>> which.min(reg.summary$cp)<br>[ 1 ] 10<br>> points(10, reg.summary$cp, cex = 2, pch = 20)<br>4) To plot BIC:<br>> plot(reg.summary$bic, type = "l")<br>> which.min(reg.summary$bic)<br>[ 1 ] 6<br>> … |

| | | |
|---|---|---|
| plot(regfit.full, scale = "...") | Show models ranked by performance measure using *regsubsets()*'s built-in *plot()* function<br>- black boxes in the top row show variables in the optimal model according to the statistic | > ... scale = "bic" ...<br>> ... scale = "cp" ... |
| coeff(regfit.full, n) | Extract coefficients for each predictor size<br>- Argument *n* is the number of coefficients to extract<br>- Gives coefficient estimates for intercept and each of the selected variables in subset | |

## Forward and Backward Stepwise Selection

- Use same function *regsubsets()* but with argument *method* specified

| | | |
|---|---|---|
| regsubsets( y ~ x, data = df, nvmax = n, method = "forward") | Forward stepwise | |
| regsubsets( y ~ x, data = df, nvmax = n, method = "backward") | Backward stepwise | |

## Choosing models

- Important to use **ONLY TRAINING OBSERVATIONS** for model fitting and **variable selection**

### Validation Set Approach

| | | |
|---|---|---|
| set.seed(1)<br>train = sample(c(T, F),<br>        nrow(df), rep = T)<br>test = !train | Split observations into training and test set | |
| regsubsets(y ~ x, data = df[train,],<br>       nvmax = n) | Run *regsubsets()* on **training** set<br>- Best subsets for training data only | |
| test.mat =<br>    model.matrix(y ~ .,<br>       data = df[test,]) | Do something similar to *predict()* which is not available<br>  1.  Create a test "X" matrix<br>    - *model.matrix()* adds a colum of "1"s for intercept term, beta_0<br>    - Creates dummy variables for qualitative variables | |

| | | |
|---|---|---|
| val.errors = rep(NA, n)<br>for (i in 1:n) {<br>    coefi = coef(regfit.best, id = i)<br>    pred =<br>       test.mat[, names(coefi)]<br>       %*% coefi<br>    val.errors[i] =<br>       mean(df$y[test] − pred)**2)<br>}<br>which.min(val.errors) | 2. Compute validation error for each of the *n* subset models (specified by *nvmax* argument)<br>- Validation error is MSE of all observations in test set<br>- Choose number of variables with minimum validation error | Own predict function that returns predicted values for specifief number of variables in model, *id*:<br>> function(object, newdata, id, ...) {<br>    # object is a *regsubsets()* model<br>    form =<br>    as.formula(object$call[[2]])<br>    mat = model.matrix(form, newdata)<br>    # id is number of variables<br>    coefi = coef(object, id = id)<br>    xvars = names(coefi)<br>    return(mat[, xvars] %*% coefi)<br>} |
| regsubsets(y ~ ., data = df, nvmax = n) | 3. **Refit with full data set** | |
| Cross-validation Approach | | |

- Each fold: 1) fit
        2) variable selection
- **DO NOT** do variable selection before k-fold CV ***

| | | |
|---|---|---|
| | Do k-fold cross validation:<br>1. Assign data to different folds | > fold_index <- sample(rep(1:k, nrow(df), nrow(df)) |
| | 2. Set up matrix to store CV error estimates | # Each fold has n errors for different number of variables<br>> cv.errors <- matrix(NA, k, n) |
| | 3. Put Validation Set Approach code through a loop across the k folds | > for (j in 1:k) {<br>    best.fit = regsubsets(Salary ~ .,<br>        data = Hitters[folds != j,]<br>        nvmax = 19)<br>    for (i in 1:19) {<br>        pred = predict(best.fit,<br>           Hitters[folds == j,], id = 1)<br>        cv.errors[j, i] =<br>        mean((Hitters$Salary[folds == j,] −<br>           pred)**2)<br>    }<br>} |

| | | |
|---|---|---|
| | 4. Aggregate across all folds to see mean CV errors for models of different number of variables | # Apply by column<br><br>> mean.cv.errors = apply(cv.erros, 2, mean)<br><br>> which.min(mean.cv.errors)<br><br>[ 1 ] 11 |
| | 5. Refit model with full data set and get coefficients for number of variables determined in Step 4. | > regsubsets(Salary ~ ., data = Hitters, nvmax = 19)<br><br>> coef(reg.best, 11) |

## (6B) SIX B – Shrinkage Methods

| R code | Description | Example |
|---|---|---|
| library(glmnet) | - Package written to perform elastic net on generalised linear models where *elastic net* is a hybrid version of ridge regression and lasso<br>- Takes in quantitative response *y* and quantitative predictor matrix *x* | |
| x = model.matrix(<br>    y ~ ., df)[, -1]<br>y = df$y | Prepare data for *glmnet()* function as a matrix | |
| Ridge Regression | | |
| grid =<br>    10^seq(10, -2,<br>    length = 100) | Vector of lambda values used to fit the model<br>- Argument for *lambda* in *glmnet()*<br>- If argument left empty, *glmnet()* will automatically select a range<br><br>- Example: range of lambda values from $10^{10}$ to $10^{-2}$<br>- To cover full range of scenarios for data | |
| glmnet(x, y,<br>    alpha = 0,<br>    lambda = grid) | Fit a ridge regression model<br>- By default, variables are standardised before fitting and returns coefficients in the original scale<br>- Turn off by setting argument *standardize = F*<br>- Especially when variables are already in same units | |
| ridge.mod$lambda<br>[index]<br><br>coef(ridge.mod)[,<br>index] | Returns vector of RR coefficient estimates for each value of lambda<br>- First, find lambda at specified index<br>- Next, find coefficients for this value of lambda<br><br>- Larger values of lambda, have smaller coefficients, in terms of *l2* norm | In this example, there are 100 lambda values and each value has 20 coefficient estimates for each predictor + intercept:<br>> dim(coef(ridge.mod))<br>[ 1 ] 20   100<br>> ridge.mod$lambda[50]<br>[ 1 ] 11497.57<br>> coef(ridge.mod)[, 50]<br>  (Intercept)        AtBat    ...<br>407.356050200    0.036957182   ... |
| predict(ridge.mod,<br>s = *new_lambda_value*,<br>type = "coefficients") | Predict coefficients for a new value of lambda | |
| ridge.pred =<br>   predict(ridge.mod,<br>    s = *lambda*,<br>    newx = *test_set*) | Returns predictions for a new data set<br>- Usually for test data<br>- For lambda values (*s*) that were not specified during fitting of the model, *predict()* interpolates/extrapolates | |

| | | |
|---|---|---|
| ```
predict(ridge.mod,
    s = lambda,
    newx = test_set,
    exact = T,
    x = train_set,
    y = train_set$y,
    thresh = 1e-16)
``` | - Set argument *exact = T* to recompute exact coefficients, original *x* and *y* need to be supplied as well<br><br>- Set lower threshold (algorithm runs longer) for closer approximation to coefficient estimates because *glmnet()* is a numerical fitting method | |
| ```
cv.out =
cv.glmnet(x_train,
    y_train,
    alpha = 0,
    nfolds = k)


plot(cv.out)
bestlam =
log(cv.out$lambda.min)
``` | Choose optimal lambda using cross-validation | |
| **LASSO** | | |
| ```
glmnet(x[train,], y[train],
    alpha = 1,
    lambda = grid)
``` | - Fit model using train set<br>- Find best lambda using CV on train set<br>- Predict test set<br>- Find test MSE | ```
> lasso.mod = glmnet(x[train,], y[train,],
alpha = 1, lambda = grid)
> plot(lasso.mod)
```<br><br>```
> set.seed(1)
> cv.out = cv.glmnet(x[train,], y[train],
alpha = 1)
> plot(cv.out)
> bestlam = cv.out$lambda.min
```<br><br>```
> lasso.pred = predict(lasso.mod, s =
bestlam, newx = x[test,])
> mean((lasso.pred – y.test)**2)
``` |
| ```
out = glmnet(x, y,
    alpha = 1,
    lambda = grid)
predict(out,
    type = "coefficients",
    s = bestlam)
``` | Returns coefficient estimates<br>- Note that some coefficients can be exactly zero in lasso | |

# (6C) SIX C – Dimension Reduction Methods

| R code | Description | Example |
|---|---|---|
| library(pls) | | |

## Principal Components Regression

- Attempts to maximise the amount of variance explained in predictors

| R code | Description | Example |
|---|---|---|
| set.seed(2)<br><br>pcr(y ~ .,<br>    data = df,<br>    subset = train,<br>    scale = T,<br>    validation = "CV") | Fit a PCR on train set<br>- PCR is closely related to Ridge Regression<br>- *scale = T* to standardise each predictor<br>- *validation = "CV"* to compute ten-fold CV error for each possible value of M | |
| summary(model) | View resulting fit of PCR model<br>- CV score is sqrt(MSE)<br>- So MSE = CV²<br><br>- Cumulative PVE in predictors and response also provided<br>- Amount of information about the predictors or the response that is captured using M principal components<br>- Using all M = p will increase cumulative PVE to 100% | > summary(pcr.fit)<br>Data: X dimension: 263  19<br>Y dimension: 263  1<br>Fit method: svdpc<br>Number of components considered: 12<br><br>VALIDATION: RMSEP<br>Cross-validated using 10 random segments<br>(Intercept)   1 comps   2 comps   ...<br>CV          452      348.9   ...<br>adjCV      452      348.7   ...<br>...<br><br>TRAINING: % variance explained<br>       1 comps  2 comps  3 comps  ...<br>X     38.31    60.16    70.84   ...<br>Salary  40.63    41.58    42.17   ... |
| validationplot(model,<br>    val.type = "MSEP") | Plot CV scores<br>- MSEP against number of components<br><br>- CV errors are more variable/higher variance because:<br>    1. train set choice has much higher variance across different CV-folds<br>    2. smaller dataset to train | |

| | | |
|---|---|---|
| predict(model,<br>        x[test,],<br>ncomp = M) | Estimate test error using M that gave lowest CV error | |
| full_model <- pcr( y ~ x, scale = T, ncomp = M)<br>summary(full_model) | Fit PCR on **full** dataset to get final PCR coefficient estimates | |
| Partial Least Squares | | |
| - No guarantee that *x, y* will be optimal (especially when n is small)<br>- Uses correlation between *x* and *y* to find approximate direction<br>- Searches for directions that explain variance in both predictors and response<br>- Variance of *y* plays a huge part | | |
| plsr(y ~ ., data = df,<br>      subset = train,<br>      scale = T) | Fit on training set | |
| summary(model) | Find best number of components, M | |
| predict(model, test_set, ncomp = M) | Find test MSE<br>- Useful when comparing models<br>(E.g. PCR model) | |
| plsr(y ~ ., data = df,<br>scale = T,<br>ncomp = 2) | Fit on full data set for best M components found | |

# (10A) TEN A – PCA, K-means

| R code | Description | Example |
|---|---|---|
| Principal Components Analysis | | |
| prcomp(df, scale = T) | Perform PCA<br>- By default, variables centered to have mean = 0<br>- *scale = T* for sd = 1 | |
| names(pca_output) | Attributes of the output of an *procomp()* object | > names(pr.out)<br>[ 1 ] "sdev" "rotation" "center" "scale" "x" |
| **Means**<br>pca_output$center<br>**S.d.**<br>pca_output$scale | Means and standard deviations of variables used for scaling prior to implementing PCA | |
| pca_output$rotation | Gives rotation matrix of size p by M/corresponding principal component **loading** vectors<br>- If multiplied with X of size n by p, returns principal component scores/coordinates in rotated coordinate system | |
| pca_output$x | Returns principal component **score** vectors of size n by M<br>- kth column is for kth PC | |
| pca_score <-<br>      pca_output$x<br>plot(<br>pca_score[, c(PC_num1,<br>      PC_num2)],<br>col = Cols(df$y)) | Plot first few PC scores to visualise data | # To assign colours to corresponding points instead of words<br>Cols = function(vec) {<br>    cols = rainbow(length(unique(vec))<br>    return(cols[as.numeric(as.factor(vec))])<br>}<br># Plot PC2 against PC1<br>> plot(pca_output$x[, 1:2])<br># Plot PC3 against PC1<br>> plot(pca_output$x[,c(1,3)]) |
| biplot(pca_output,<br>    scale = 0)<br><br>pca_output$rotation =<br>  -pca_output$rotation<br>pca_output$x =<br>  -pca_output$x | Plots first 2 principal components in a **biplot**<br>- *scale = 0* ensures that loadings are in the right scale<br>- PCs unique up to a sign change<br>- Can get flipped version by flipping signs in loadings and scores | |
| pca_output$sdev | Returns standard deviation of each principal component<br>- With respect to loadings | |

| | | |
|---|---|---|
| pr_var =<br><br>    pca_output$sdev**2<br><br><br>pve =<br><br>    pr_var/sum(pr_var) | Compute PVE of each component<br>- Squared of *$sdev* divided by total variance | |
| plot(pve,<br><br>    xlab = "PC",<br><br>    lab = "PVE",<br><br>    ylim = c(0, 1),<br><br>    type = 'b') | Plot PVE<br>- *type* = 'b' for points to be represented as a circle and connected by lines | |
| plot(cumsum(pve),<br><br>    ...) | Plot cumulative PVE | |
| **K-Means** | | |
| - If data has more than 2 variables, 1) Perform PCA first<br>                                         2) Plot the first 2 PC score vectors | | |
| kmeans(x, k,<br><br>        nstart = n) | Perform K-means clustering<br>- Argument *nstart* can be used to specify how many initial cluster assignments to run algorithm through<br>- Generally a good idea to run with a large value of *nstart* (e.g. 20, 50) to discard undesirable local optima<br>- Algorithm is very efficient | |
| km_out$cluster | View cluster assignments for each observation by index | > km.out = kmeans(x, 2, nstart = 20)<br>> km.out$cluster<br>[ 1 ] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2<br>2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 … |

| | | | |
|---|---|---|---|
| plot(x,<br><br>    col= (km_out$cluster + 1),<br><br>    pch = 20, cex = 2) | | Plot data, with each observation coloured according to its cluster assignment | |

| | | |
|---|---|---|
| km_out$tot.withinss | Returns local minimum | |
| km_out | Returns *kmeans()* object | > km_out<br><br>…<br><br>Available components:<br>[ 1 ] "cluster"    "centers"       "totss"<br>[ 4 ] "withinss" "tot.withinss" "betweenss"<br>[ 7 ] "size"      "iter"          "ifault" |

## (10B) TEN B – Hierarchical Clustering

| R code | Description | | Example |
|---|---|---|---|
| hclust(dist(x),<br><br>method = "...")<br><br><br>hclust(dist(scale(x)), ...) | Hierarchical clustering<br>**- EUCLIDEAN-BASED DISTANCE**<br>- *dist()* outputs a distance output/distance matrix<br>- Can scale variables before performing clustering | | |
| | colspan="2" | method = | |
| | "complete" | Complete linkage<br>- Default | |
| | "average" | Average linkage | |
| | "single" | Single linkage<br>- Commonly gives a trailing cluster phenomenon | |
| | "centroid" | Centroid linkage | |
| dd <- as.dist(1-cor(t(x)))<br>hclust(dd, ...) | Hierarchical clustering<br>**- CORRELATION-BASED DISTANCE**<br>- convert arbitrary square symmetric correlation matrix into a form that *hclust()* can recognise as a distance matrix<br>- Highly correlated/ correlation value close to 1, smaller distance<br>- Correlation value close to -1, larger distance (to consider close to -1 as close, take \|*cor(x)*\|)<br><br>- Dissimilarity function here is **(1-r)** where **r** is sample correlation<br><br>- Other choices:<br>   1. **sqrt(1-r)**<br>   2. **1 - \|r\|**<br>   3. **sqrt(1 - \|r\|)** | | |
| plot(hclust(...),<br>   labels = df$y,<br>   cex = 0.9) | Plot dendrograme directly from *hclust()* output | | |
| abline(<br>   h = height_of_cut,<br>   col = "red") | Plot cut | | |
| cutree(hclust(...),<br>   c) | Cut tree for a specified number of clusters<br>- *c* is number of clusters | | |

| | | |
|---|---|---|
| hc_clusters <-<br>cutree(...)<br>table(hc_clusters, df$y) | View how observations were clustered against their true labels<br>- Something similar to a confusion matrix | hc_clusters BRE CNS COL K5A K5B LEU<br><br>| hc_clusters | BRE | CNS | COL | K5A | K5B | LEU |<br>|---|---|---|---|---|---|---|<br>| 1 | 2 | 3 | 4 | 0 | 0 | 0 |<br>| 2 | 3 | 2 | 0 | 0 | 0 | 0 |<br>| 3 | 0 | 0 | 0 | 1 | 1 | 6 |<br>| 4 | 2 | 0 | 5 | 0 | 0 | 0 |<br><br>- *BRE* and *CNS* are very spread out across different clusters<br>- *LEU* all clustered together |
| km_cluster <-<br>        km_out$cluster<br>hc_cluster <-<br>        cutree(..., k)<br><br>table(km_cluster,<br>        hc_cluster) | Compare hierarchical cluster assignments with k-means cluster assignments | hc_clusters<br><br>| km_clusters | 1 | 2 | 3 | 4 |<br>|---|---|---|---|---|<br>| 1 | 11 | 0 | 0 | 9 |<br>| 2 | 0 | 0 | 8 | 0 |<br>| 3 | 9 | 0 | 0 | 0 |<br>| 4 | 20 | 7 | 0 | 0 |<br><br>- Cluster 1 for *hc_cluster* has observations spread out across different clusters in *km_clusters*<br>- Cluster 3 confirmed by both clustering methods |
| | Perform clustering based on first few PC score vectors<br>- Results can be different from using full dataset<br>- Sometimes PCA yields better results if most of signals captured in first few components<br>- **Using PCA as a denoising pre-processing step** (get a cleaner result)<br>- **Looking at scree-plot decides where to make the cut** | # Do for first 5 PCs<br>> hc_out <- hclust(dist(pr_out$x[1:5]))<br>> plot(hc_out,<br>        labels = y,<br>        main = "Hier. Clust. on First Five<br>                Score Vectors")<br>> table(cutree(hc_out, 4), y) |