Data Science and Big Data Analysis Notes

```
library(readr)
read csv() # load into dataframe
as.dataframe() # convert to dataframe
str() #structure function
factor() #force column to be Factors
rownames(group_info_df) <- group_info_df$name # set row names to use first names
order(group_info_df$grade) #gives order of rows in ascending grade order, print row num
group info df <- group info df[order(group info df$grade), ] # select the rows in order
group_info_df[3:7,] # selects rows 3,4,5,6,7 all columns
library(readxl)
read_excel("hsb2.xlsx")
library('dplyr')
select(student data, read, write, math, science) #select
select(student_data, female, race, school_type=schtyp) #select and rename schtype to
school type
filter(student data, age>20) #select column that age >20
filter(student data, read>50, write>50)
mutate(student_data,lang=...,) # adding lang to data frame
arrange(data,column) #sort by column by ascending order
arrange(student data, desc(read))
factor(student_data$sex, labels = c("Black", "Asian")) # assign labels to F and M
Bivariate linear regression models—Day3
gdp mean<-mean(world data$gdp, na.rm = TRUE)
n<-length(world_data$gdp[!is.na(world_data$gdp)])
se <- sd(world_data$gdp, na.rm = TRUE) / sqrt(n) # calculate standard error
# lower bound
lb <- gdp mean - 1.96 * se
# upper bound
ub <- gdp mean + 1.96 * se
library(texreg)
names(dat)[which(names(dat) == "..")] <- ".." # change names
 UnemploymentRate ~ NoHighSchool, data = dat,
 xlab = "Adults without High School education (%)",
 ylab = "Unemployment (%)",
 frame.plot = FALSE,
 pch = 16.
 col = rgb(red = 110, green = 200, blue = 110, alpha = 80, maxColorValue = 255)
model1 <- Im(UnemploymentRate ~ NoHighSchool, data = dat)
summary(model1)
```

```
Call: 1 Dependent Variable 2
lm(formula = UnemploymentRate ~ NoHighSchool, data = communities)
                 Difference between the observed values and predicted values of UnemploymentRate 3
   Min 1Q Median 3Q
                                  Max
-0.42347 -0.08499 -0.01189 0.07711 0.56470
Coefficients: UnemploymentRate = 0.078952 + (0.742385 * NoHighSchool)
          Estimate Std. Error t value Pr(>|t|) 🙀
(Intercept) 0.078952 0.006483 12.18 <2e-16 *** means p < 0.01

NoHighSchool 0.742385 0.014955 49.64 <2e-16 *** means p < 0.001

Lvalue = coefficient / std. error 7
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 0.1352 on 1992 degrees of freedom
Multiple R-squared: 0.553, Adjusted R-squared: 0.5527 ←
F-statistic: 2464 on 1 and 1992 DF, p-value: < 2.2e-16
abline(model1, lwd = 3,
     col = rgb(red = 230, green = 150, blue = 0, alpha = 255, maxColorValue = 255))
screenreg(model1)
screenreg(list(model1, model2))
predict(model1, newdata = data.frame(NoHighSchool = 10)) # predict on single value =10
n <- nrow(world_data)
world_data$x <- rnorm(n, mean=0, sd=10
rename(data, new name = old name)
mutate(gender = factor(female, labels=c("M", "F"))) # change 1 and 0 to M and F
```

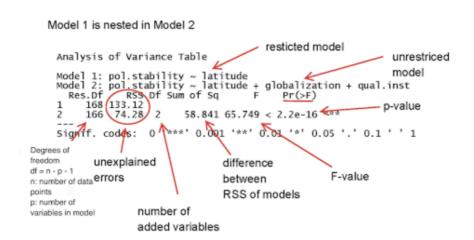
Multiple linear regression models — Day4

```
library(haven) # used to load our data
library(texreg) # used to display fit info
library(dplyr) # used to manipulate data
library(tidyr) # used for the drop_na function
library(ggplot2) # in case we want to make ggplots
```

```
read dta()
```

```
mutate(world_data,democracy = factor(democracy,levels=c(0,1), labels=c("dictatorship","democracy")))
```

world_data <- drop_na(world_data,latitude) # drop na value in latitude and return df anova(latitude_model, inst_model)



```
new data democracy <- data.frame(institutions quality = seg(from = 1.4, to = 9.3, by =
1), globalization = mean(world_data$globalization), latitude = mean(world_data$latitude),
democracy = "democracy")
pred <- predict(inst model, newdata = new data democracy, se.fit = TRUE)
new data democracy$political stability.se <- pred$se.fit
new_data_democracy <- mutate(new_data_democracy,</pre>
political stability.ub = political stability.pred + 1.96*political stability.se.
political stability.lb = political stability.pred - 1.96*political stability.se
# add lines for confidence intervals
# upper bound
lines(x = new data democracy$institutions quality,
   y = new_data_democracy$political_stability.ub,
   Ity = "dashed", Iwd = 1.5)
# lower bound
lines(x = new data democracy$institutions quality,
   y = new data democracy$political stability.lb,
   Ity = "dashed", lwd = 1.5)
abline(a=y0, b=slope, col="blue")
table(world data$former col)
better.model <- lm(human development ~ poly(gdp capita, 2), data = world data)
screenreg(list(bad.model, better.model),
      custom.model.names = c("bad model", "better model"))
Variable Selection — Day5
library(MASS)
Im.1 <- Im(Boston$medv~Boston$lstat,data=Boston)</pre>
Rsquared, (summed squared of residuals) / (variance)
rsq = summary(lm.1)$r.squared
rsq
Adjusted Rsquared, larger better
rsq = summary(lm.1)$adj.r.squared
rsq
MSE, small better
```

mean(residuals(lm.1)^2)

RMSE, small better

sqrt(mean(residuals(lm.1)^2))

MSE_model <- mean((actual data - predict)^2)

AIC(Im.1) # lower AIC means better performance BIC(Im.1) # lower BIC means better performance

```
Im.min <- Im(medv ~ 1, data=data.original)
lm.max <- lm(medv ~ ., data=data.original)</pre>
scp <- list(lower = Im.min, upper = Im.max)</pre>
Im.selected <- stepAIC(Im.min,</pre>
            direction = 'forward',
            scope = scp.
            steps = 1
# Using leaps
library(leaps)
regsubsets.out <- regsubsets( medv ~ .,
                data = data.original,
                nbest = 1.
                nvmax = NULL,
                force.in = NULL, force.out = NULL,
                method = 'exhaustive')
summary(regsubsets.out)
as.data.frame(summary(regsubsets.out)$outmat)
plot(regsubsets.out, scale='adjr2', main='Adjusted Rsg')
library(corrplot)
corrplot(cor(boston housing))
transport <- cut(boston housing$rad, breaks = c(0,15,30),
labels=c("low","high")))
# seperate value 0~15 (low) and 15~30 (high)
credit card$Balance mean <- mean(credit card$Balance)</pre>
credit_card$squared_dev <- (credit_card$Balance</pre>
                                 - credit card$Balance mean)^2
credit card$squared resid <- (credit card$Balance</pre>
                                - credit card$Balance pred)^2
TSS <- sum(credit card$squared dev)
RSS <- sum(credit card$squared resid)
Rsq <- 1 - RSS/TSS
lm.selected <- lm(logSalePrice~OverallQual + OverallCond +</pre>
GrLivArea + YearBuilt, data=house data num)
RSS = sum(lm.selected$residuals)^2
MSE = RSS/nrow(house data)
RMSE = sqrt(MSE)
```

Validation — Day6

```
# create a column with NA data.sample$cv_pred <- NA
```

LOOCV

```
library(cvTools)
cvFit(model, data = data.sample, y = data.sample$Y, K = # of data rows) # This returns sqrt(MSE)
or RMSE value
RMSE = sqrt(sum((model$residuals)^2)/nrow(house_data_num))
RMSE
```

K-fold CV

```
library(cvTools)
cvFit(model, data = data.sample, y = data.sample$Y, K = # of fold, R=..)
#Setting x equal to n yields leave-one-out cross-validation.
# an integer giving the number of replications for repeated
K-fold cross-validation

cv_result = cvFit(model, data = house_data_num, y = house_data_num$logSalePrice, K = 10, R=100)
cv_result$cv #RMSE
cv_result$se #standard error
cv_result$reps
```

bootstrapping

library(qvlma)

```
R = 1000
coef_out \leftarrow rep(NA,R)
MSE_out \leftarrow rep(NA,R)
for (i in 1:R)
 resampled_rows <- sample(1:n,n,replace=TRUE)</pre>
 data.resampled <- data.sample[resampled_rows,]
 model <- lm(Y \sim X, data = data.resampled)
 coef out[i] <- coef(model)[[2]]
 MSE_out[i] <- mean(resid(model)^2)
}
hist(coef_out)
hist(MSE_out)
mean(coef out)
mean(MSE_out)
plot(house_data$logSalePrice~lm.selected$fitted.values) # stepwise, actual~fittedvaule
plot(lm.selected$residuals~lm.selected$fitted.values)
#install.packages("gvlma")
```

gvlma(lm.selected) # assess linear model assumption

```
library(corrplot)
corrplot(cor(predictor_set))
res <- cor(my_data)
round(res, 2)
cor(y = world_data$hdi, x = world_data$corruption, use =
"complete.obs") #correlation coefficient
cor.test(~ hdi+ gdp, data=world_data, ) #estimated
correlation # find evidence to reject the null hypothesis
at the 5% significance level, and accept the hypothesis
that the variables are correlated.

residuals <- actual data - pred
RMSE <- sqrt(sum((residuals)^2)/length(residuals))
RMSE</pre>
```

bootstrapping, non-linear fit methods and shrinkage methods—Day7

```
Non-linear fit methods
library(ISLR)
library(splines)
# can specify knot points
fit spline1=lm(wage~bs(age,knots=c(25,40,60)),data=Wage)
pred <- predict(fit spline1,newdata=age grid df)</pre>
plot(wage~age, pch='.', data=Wage)
lines(age grid,pred,lwd=2,col="red")
# can specify degrees of freedom (more=more flexible)
fit spline2=lm(wage~bs(age,df=6),data=Wage)
pred <- predict(fit spline2,newdata=age grid df)</pre>
plot(wage~age, pch='.', data=Wage)
lines(age grid,pred,lwd=2,lty=2,col="red")
# loess to build a LOcal regrESSion
fit loess1=loess(wage~age,span=.2,data=Wage)
fit loess2=loess(wage~age,span=1,data=Wage)
plot(wage~age, pch='.', data=Wage)
lines(age grid,predict(fit loess1,newdata=age grid df),col="red",l
wd=2)
lines(age grid,predict(fit loess2,newdata=age grid df),col="blue",
lwd=2)
cut points <- quantile(Wage$age,c(0,0.25,0.5,0.75,1.0))
Wage$age groups <- cut(Wage$age, breaks = cut points)</pre>
```

assess spline with different degree of freedom

```
library(cvTools)
df vals < c(4,,8,12,16,20,24,28,32,36,40)
n vals <- length(df vals)</pre>
results <- 0
results <- data.frame(df = rep(0,n_vals),
                       rmse = rep(0, n vals),
                       se = rep(0, n vals))
for (i in 1:n vals){
  df i <- df vals[i]</pre>
  fit spline2=lm(wage~bs(age,df=df i),data=Wage)
  cv result <- cvFit(fit spline2, data = Wage, y = Wage$wage, K =
10, R = 10)
 results$df[i] <- df_i
  results$rmse[i] <- cv_result$cv</pre>
  results$se[i] <- cv result$se
}
head(results)
plot(results$rmse~results$df, pch=1, ylim=c(39.8,40.9))
points(results$df, results$rmse + 1.96*results$se, pch=2)
points(results$df, results$rmse - 1.96*results$se, pch=2)
glmnet Ridge Regression and Lasso method
fit lasso <- glmnet(housing data.x, housing data.y,
```

```
alpha = 1, lambda = 1)
fit ridge <- glmnet(housing data.x, housing data.y,
                alpha = 0, lambda = 1)
# view fit coefficients
coef(fit lasso)
# to predict new points based on model
pred <- predict(fit lasso,newx=housing data.x)</pre>
# if we do not specify lambda glmnet will test out
# a range of penalties and we can view how coefficents
# behave by plotting the result
fit lasso <- glmnet(housing data.x, housing data.y,
                alpha = 1)
plot(fit lasso, xvar = "lambda") # lambda is on log-scale
# cross validation for best choice of lambda
results <- cv.glmnet(housing_data.x, housing_data.y,
         alpha = 0, nfolds=10
plot(results)
```

logistic regression, LDA and QDA classification -Day8

```
par(mfrow=c(2,2))
```

```
# the following command creates a 2x2 plot but
# adjusts margins (plot spacing) to be smaller
par(mfcol=c(2,2), mar=c(4,4,0.5,0.5), oma=c(1.5,2,1,1))
# with command lets you run a
# command using a given data set
# similar to using the data=... argument
with(mf_training, scatter.smooth(height ~ hair,
                    pch=20,
                     col = ifelse(gender==1,'red','blue')
   )
plot(height ~ eye, data=mf_training)
plot(height ~ as.factor(gender), data=mf training)
plot(height ~ as.factor(glasses), data=mf_training)
par(mfrow=c(1,2))
plot(regsubsets.out, scale="adjr2")
adj_rsquare_values <- summary(regsubsets.out)$adjr2
n predictor vals <- 1:5
plot(adj_rsquare_values~n_predictor_vals)
nfemales <- length(which(mf_training$gender==1)) # number of female in data
rows_glasses = which( mf_training$glasses==1)
result <- mean(mf_training[rows_glasses,]$gender) #P(femalelglasses), if female if they wear
glasses
result <- mean(mf_training[-rows_glasses,]$gender) #P(femalelno-glasses)
mean(mf_training[which(mf_training$hair>10),]$gender) #P(femalelhair>10)
lm(gender ~ height,mf training)
predict(gender_height.lm, data.frame(height = c(150, 180))) # probability for female when height
is 150 and 180
b0 = coef(gender_height.lm)[["(Intercept)"]]
b1 = coef(gender_height.lm)[["height"]]
result < (0.5 - b0)/b1 # take x out of y = b0 + b1*x
total <- nrow(mf_training)
misclassified <- length(which(mf_training$mod1_prediction != mf_training$gender))
rate <- misclassified/total
hist(subset(mf_training,gender==1)$height, col='blue',
   xlab='height',main='Male/Females height distributions')
hist(subset(mf_training,gender==0)$height, col=rgb(1,0,0,0.5), add=T)
```

```
Logistic
```

```
#predicts the probability given height, is female
gender_height.glm=glm(gender~height,data=mf_training,family=binomial)
plot(gender~height,data=mf_training)
curve(predict(gender_height.glm,data.frame(height=x),type="resp"),add=TRUE)
logistic fit formula is P(femalelheight=H)/y=1/(1 + \exp(-B0 - B1*X))
mf trainingprediction \leftarrow factor(mf training\\mod1 prediction, levels=<math>c(0,1),
labels=c("male","female"))
plot(iris$Petal.Length, iris$Petal.Width, pch=24,
   bg=c("red","green","blue")[unclass(iris$Species)],
   main="Iris Species\n R: Setosa, G: Versicolor, B: Virginica")
unclass() #converts a factor to its integer representation
levels(iris$Species) # print factor
subset(iris,Species=="virginica" | Species=="versicolor") #subset rows with only species
iris_simple[,c("Species","Sepal.Length","Sepal.Width")]
correct = sum(logistic.model.prediction == iris simple$Species)
correct
LDA
library(MASS)
Ida.model = Ida(Species ~ Sepal.Length+Sepal.Width, data = iris_simple)
Ida.model.prediction <- predict(Ida.model)$class
Ida.model = Ida(Species ~ ., data = iris, CV=TRUE)
QDA
qda.model = qda(Species ~ ., data = iris)
qda.model.prediction <- predict(qda.model)$class
plot(iris$Sepal.Length, iris$Sepal.Width, pch=24,
   bg=c("red", "green", "blue") [unclass(iris$Species)],
   main="Iris Species\n R: Setosa, G: Versicolor, B: Virginica")
misclassified = sum(iris test$pred lda!=iris test$Species)
  misclassification rate lda = misclassified/nrow(iris test)
  misclassification rate lda
```

KNN - Day9

normalise <- function(x){</pre>

```
norm_x <- (x-min(x))/(max(x)-min(x))
return(norm_x)
}

Set.seed(xx)
sample(1:nrow(iris),size=floor(nrow(iris)/2))

knn.pred=knn(iris_X_train,iris_X_test,iris_Y_train,k=30)

knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)</pre>
```

Arguments

train matrix or data frame of training set cases.

test matrix or data frame of test set cases. A vector will be interpreted as a row

vector for a single case.

cl factor of true classifications of training set

k number of neighbours considered.

table(predicted=knn.pred,actual=iris_Y_test)
knn.cv() performs LOOCV
plot(X2~X1,col=Y,pch=2, data=data.cols)