# Project: STAT 525

*KRISHNAN RAMAN, 2nd Year PhD. Statistics*

*12/10/2019*

## An exploration of trading strategies with Dow Jones components

The Dow Jones Index (DIA) is a weighted average of 30 stocks.
These 30 companies are some of the largest industrial giants in the US.
In this analysis, we explore the following questions:

- Does the log-normal assumption apply to Dow stock returns ?
- Predict the DIA using multiple linear regression, as a function of the 30 Dow components.
- Use Model Selection to systematically select the optimal number of covariates.
- Create a high 90% Rsquare Dow Regression model with fewest possible number of covariate stocks.
- Track the performance of the DIA versus the Dow Regression Model over a year
- 1-way Anova: Compare Diversification(Holding all 30 stocks) vs Holding a Single Stock
- Cell Means Model, Means Only Model, Effects Plot
- 1-way Anova: Compare annual returns across 10 different portfolios with 3 stocks in each portfolio
- 1-way Anova: Compare Buy & Hold Trading Strategy vs Buy Highest Sell Lowest vs Contrarian Strategy

**Data**

The Daily Holding Period returns for the Dow Jones Index (DIA) and its 30 stock components were sourced from **Wharton Research Data Services**, a subscription-only source of Financial Time Series. For each of the **30 stocks + DIA = 31 tickers**, we obtained data for all of 2018 ( **There were 251 trading days in 2018** ). That gives us **31x251 = 7781 rows** of daily returns. These are shown below:

```r
rm(list=ls())
library(fitdistrplus)
library(ggplot2)
library(gtools)
library(leaps)
library(effects)


dowfile = "~/Desktop/525/project/dow.csv"
df = read.csv(dowfile)
head(df)
```
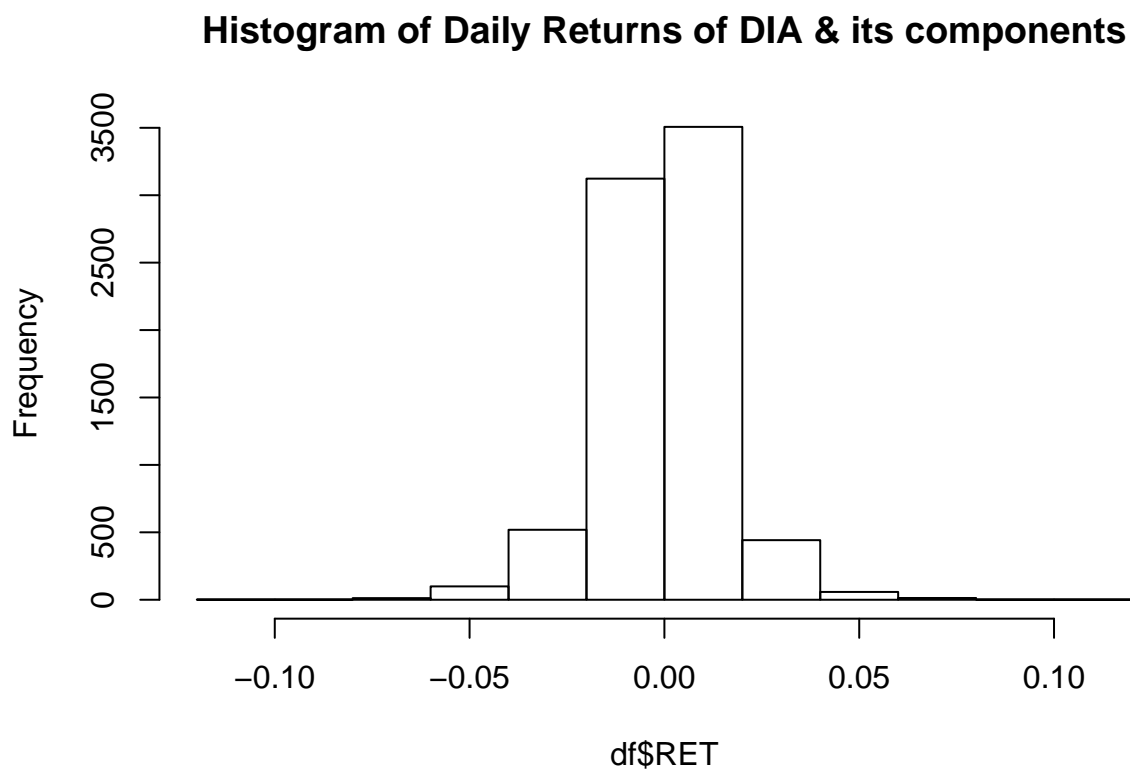
```
##    date TICKER      RET
## 1 20118   MSFT  0.004793
## 2 30118   MSFT  0.004654
## 3 40118   MSFT  0.008801
## 4 50118   MSFT  0.012398
## 5 80118   MSFT  0.001020
## 6 90118   MSFT -0.000680
```

```
tail(df)
```

```
##        date TICKER       RET
## 7776 211218    DIA -0.018280
## 7777 241218    DIA -0.026730
## 7778 261218    DIA  0.048647
## 7779 271218    DIA  0.011149
## 7780 281218    DIA -0.003373
## 7781 311218    DIA  0.011801
```

The histogram of daily returns is firmly centered at 0. On a given day, the Dow components don't move that much.

```
hist(df$RET, main="Histogram of Daily Returns of DIA & its components")
```

## Histogram of Daily Returns of DIA & its components



In fact, the max loss is -10.18 %, and the max gain 11.13 %, on the Dow components. The median daily gain is approx 0%.
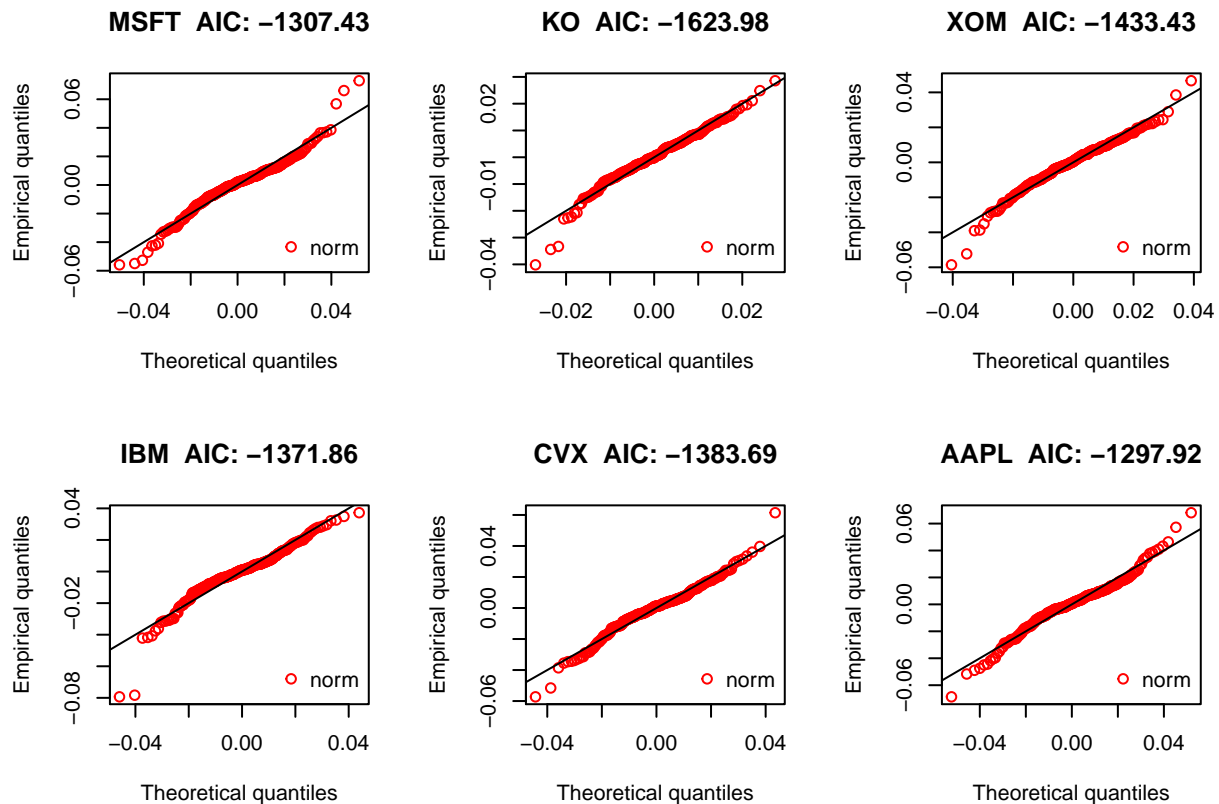
```
summary(df$RET)
```

```
##      Min.   1st Qu.    Median      Mean   3rd Qu.      Max.
## -1.018e-01 -7.184e-03  5.840e-04 5.655e-05  8.283e-03  1.113e-01
```

**LogNormality of Returns**

We split the data by ticker. Financial returns are postulated to be log-normal ie. the log of returns is gaussian. We first fit a Normal distribution upon the log returns. We test the goodness of fit by visualizing the qq plot of six of the fitted distributions below.

```
tickerplot = function(ticker) {
  ret = log(1.0 + tickersplit[[ticker]]$RET )
  fitnorm<-fitdist(ret,"norm") # fit a normal distribution
  qqcomp(fitnorm, main=paste(ticker," AIC:",round(fitnorm$aic, 2)))
}

tickersplit = split(df, df$TICKER)
par(mfrow=c(2,3))
tickers = as.character(unique(df$TICKER))
tickers = tickers[tickers != "DIA"]
sapply(tickers[1:6], tickerplot)
```



```
## $MSFT
## NULL
##
## $KO
## NULL
##
## $XOM
```

3

```
## NULL
##
## $IBM
## NULL
##
## $CVX
## NULL
##
## $AAPL
## NULL
```

```r
par(mfrow=c(1,1))
```

From the QQ plots above, we note that the lognormal assumption is violated in the tails.

**Model Selection**

In order to build a multiple linear regression model with Dow Jones (DIA) as the Response, and each of the 30 Dow stocks as the predictors, it will be convenient to construct a simple matrix with 31 columns & 251 rows. The first column DIA is the Response, the remaining 30 columns are the predictor stocks, and the 251 rows are the daily returns over the 2018 trading year. This process is shown below:

```r
dow = matrix(0,nrow=251,ncol=31)
dow[,1]=tickersplit[["DIA"]]$RET
for(i in 2:31) {
  myticker = tickers[i-1]
  dow[,i] = tickersplit[[myticker]]$RET
}
colnames(dow) <- c("DIA", tickers)
dowdf = data.frame(dow)
dow[1:5,1:7]
```

```
##              DIA      MSFT        KO       XOM      IBM       CVX       AAPL
## [1,]   0.002587  0.004793 -0.007411  0.016619 0.005410  0.019091  0.017905
## [2,]   0.003750  0.004654 -0.002196  0.019640 0.027488  0.007289 -0.000174
## [3,]   0.006628  0.008801  0.014085  0.001384 0.020254 -0.003113  0.004645
## [4,]   0.008460  0.012398 -0.000217 -0.000806 0.004886 -0.001639  0.011385
## [5,]  -0.000514  0.001020 -0.001519  0.004496 0.006031  0.004926 -0.003714
```

Now we are ready to perform Forward Stepwise Model Selection.

```r
x = dow[,2:31]
y = dow[,1]
forward_varsec = summary(regsubsets(x=x,y=y,method="forward", nbest=1,nvmax=30, all.best=FALSE))
forward_varsec$outmat
```

```
##            MSFT KO  XOM IBM CVX AAPL UTX PG  CAT WBA BA  PFE JNJ MMM MRK
## 1  ( 1 )   " "  " " " " " " " " " "  " " " " " " " " " " " " " " "*" " "
## 2  ( 1 )   " "  " " " " " " " " " "  " " " " " " " " " " " " " " "*" " "
## 3  ( 1 )   " "  " " " " " " " " " "  " " " " " " " " "*" " " " " "*" " "
## 4  ( 1 )   " "  " " " " " " " " " "  " " " " " " " " "*" " " " " "*" " "
```

```
## 5   ( 1 )  " " " " " " " " " " " " " " " " " " " " "*" " " " " " " "*" " "
## 6   ( 1 )  " " " " " " " " " " " " " " " " " " " " "*" " " " " " " "*" " "
## 7   ( 1 )  " " " " " " "*" " " " " " " " " " " " " " "*" " " " " " " "*" " "
## 8   ( 1 )  " " " " " " "*" " " " " " " "*" " " " " " " "*" " " " " " " "*" " "
## 9   ( 1 )  " " " " " " "*" " " " " " " "*" " " " " " " "*" " " " " " " "*" " "
## 10  ( 1 )  " " " " " " "*" " " " " " " "*" " " " " " " "*" " " " " " " "*" " "
## 11  ( 1 )  " " " " " " "*" " " " " " " "*" " " "*" " " " " " " "*" " " " " " " "*" " "
## 12  ( 1 )  " " " " " " "*" " " " " " " "*" " " "*" " " " " " " "*" " " " " " " "*" " "
## 13  ( 1 )  " " " " " " "*" " " " " " " "*" " " "*" " " " " " " "*" " " " " " " "*" " "
## 14  ( 1 )  " " " " " " "*" " " " " " " "*" " " "*" " " " " " " "*" " " " " " " "*" " "
## 15  ( 1 )  " " " " " " "*" "*" " " " " "*" " " "*" " " " " " " "*" " " " " " " "*" " "
## 16  ( 1 )  " " " " " " "*" "*" " " " " "*" " " "*" " " " " " " "*" " " " " "*" "*" " " " "
## 17  ( 1 )  " " " " " " "*" "*" " " " " "*" " " "*" " " " " "*" "*" " " " " "*" "*" " " " "
## 18  ( 1 )  "*" " " " " "*" "*" " " " " "*" " " "*" " " " " "*" "*" " " " " "*" "*" " " " "
## 19  ( 1 )  "*" " " " " "*" "*" " " " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 20  ( 1 )  "*" " " " " "*" "*" " " " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 21  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 22  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 23  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 24  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" " " " "
## 25  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" "*"
## 26  ( 1 )  "*" " " " " "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" "*"
## 27  ( 1 )  "*" " " "*" "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" " " " " "*" "*" "*"
## 28  ( 1 )  "*" " " "*" "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" "*" "*" "*"
## 29  ( 1 )  "*" " " "*" "*" "*" "*" " " "*" " " "*" "*" " " "*" "*" "*" "*" "*"
## 30  ( 1 )  "*" " " "*" "*" "*" "*" " " "*" " " "*" "*" "*" "*" "*" "*" "*" "*"
##            DIS MCD JPM WMT NKE AXP INTC TRV VZ  HD  C   CSCO GS  V   UNH
## 1   ( 1 )  " " " " " " " " " " " " " "  " " " " " " " " " "  " " " " " "
## 2   ( 1 )  " " " " " " " " " " " " " "  " " " " " " " " " "  " " "*" " "
## 3   ( 1 )  " " " " " " " " " " " " " "  " " " " " " " " " "  " " "*" " "
## 4   ( 1 )  " " " " " " "*" " " " " " "  " " " " " " " " " "  " " "*" " "
## 5   ( 1 )  " " " " " " "*" " " " " " "  " " " " " " " " " "  " " "*" "*"
## 6   ( 1 )  " " " " " " "*" " " " " " "  " " " " "*" " " " "  " " "*" "*"
## 7   ( 1 )  " " " " " " "*" " " " " " "  " " " " "*" " " " "  " " "*" "*"
## 8   ( 1 )  " " " " " " "*" " " " " " "  " " " " "*" " " " "  " " "*" "*"
## 9   ( 1 )  " " " " " " "*" " " " " " "  " " "*" "*" " " " "  " " "*" "*"
## 10  ( 1 )  " " " " " " "*" "*" " " " "  " " "*" "*" " " " "  " " "*" "*"
## 11  ( 1 )  " " " " " " "*" "*" " " " "  " " "*" "*" " " " "  " " "*" "*"
## 12  ( 1 )  " " " " "*" "*" "*" " " " "  " " "*" "*" " " " "  " " "*" "*"
## 13  ( 1 )  " " " " "*" "*" "*" " " " "  " " "*" "*" " " " "  "*" "*" "*"
## 14  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 15  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 16  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 17  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 18  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 19  ( 1 )  " " " " "*" "*" "*" "*" " "  " " "*" "*" " " " "  "*" "*" "*"
## 20  ( 1 )  " " " " "*" "*" "*" "*" " "  "*" "*" "*" " " " "  "*" "*" "*"
## 21  ( 1 )  " " " " "*" "*" "*" "*" " "  "*" "*" "*" " " " "  "*" "*" "*"
## 22  ( 1 )  "*" "*" "*" "*" "*" "*" " "  "*" "*" "*" " " " "  "*" "*" "*"
## 23  ( 1 )  "*" "*" "*" "*" "*" "*" " "  "*" "*" "*" " " " "  "*" "*" "*"
## 24  ( 1 )  "*" "*" "*" "*" "*" "*" "*"  "*" "*" "*" " " " "  "*" "*" "*"
## 25  ( 1 )  "*" "*" "*" "*" "*" "*" "*"  "*" "*" "*" " " " "  "*" "*" "*"
## 26  ( 1 )  "*" "*" "*" "*" "*" "*" "*"  "*" "*" "*" " " "*"  "*" "*" "*"
## 27  ( 1 )  "*" "*" "*" "*" "*" "*" "*"  "*" "*" "*" " " "*"  "*" "*" "*"
```

```
## 28  ( 1 ) "*" "*" "*" "*" "*" "*" "*"   "*" "*" "*" " " "*"   "*" "*" "*"
## 29  ( 1 ) "*" "*" "*" "*" "*" "*" "*"   "*" "*" "*" "*" "*"   "*" "*" "*"
## 30  ( 1 ) "*" "*" "*" "*" "*" "*" "*"   "*" "*" "*" "*" "*"   "*" "*" "*"
```

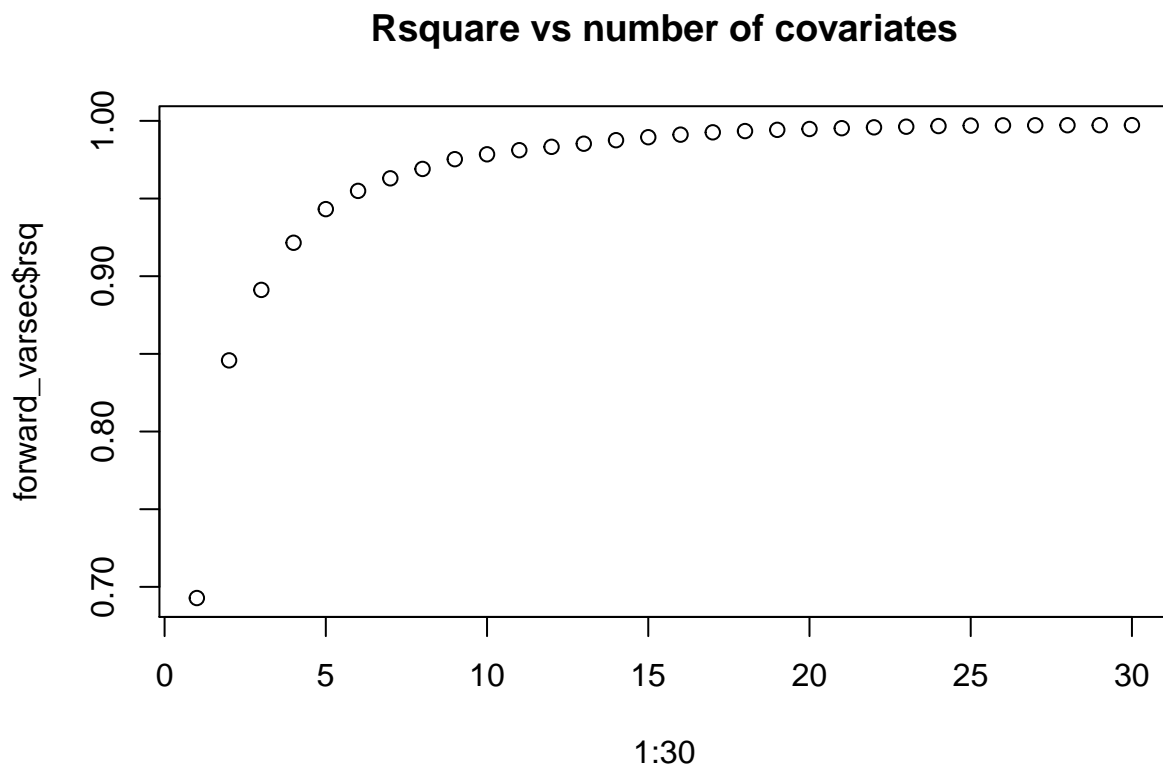It is clear that a single company 3M ( MMM ) alone is a good proxy for the Dow.
From row 5 above, **the most variation in the Dow Jones Index DIA is explained by the top 5 companies: 3M, Boeing(BA), JP Morgan (JPM), Visa (V) & United Healthcare (UNH)** From row 29,30 above, the company that explains the least variation in the Dow is Walgreens Boots Alliance (WBA).
Similar conclusions may be drawn from Backward Selection.

```
#backward_varsec = summary(regsubsets(x=x,y=y,method="backward", nbest=1,nvmax=30,all.best=FALSE))
#backward_varsec$outmat
```

To visualize the explanatory power of individual covariates, let us plot Rsquare as more & more covariates are added to the model.

```
plot(1:30, forward_varsec$rsq, main="Rsquare vs number of covariates")
```



**Rsquare vs number of covariates**

**Create a high 90% Rsquare Dow Regression model with fewest possible number of covariate stocks.**

From the above plot, we obtain about 90% Rsq from just the top 3 stocks. Lets now build a multiple linear regression model to predict the Dow Jones Returns (DIA), using just the top 3 stocks: 3M (MMM), Boeing (BA) & Visa (V)

```
# from the 3 best (forward selection) predictors
mod1 = lm(dow[,"DIA"]~dow[,"MMM"]+dow[,"BA"]+dow[,"V"])
summary(mod1)
```

```
##
## Call:
## lm(formula = dow[, "DIA"] ~ dow[, "MMM"] + dow[, "BA"] + dow[,
##      "V"])
##
## Residuals:
##         Min         1Q     Median         3Q        Max
## -0.0089642 -0.0020556 -0.0000545  0.0017604  0.0162088
##
## Coefficients:
##                Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.0001877  0.0002379  -0.789    0.431
## dow[, "MMM"]  0.3198295  0.0213530  14.978   <2e-16 ***
## dow[, "BA"]   0.1689630  0.0166476  10.149   <2e-16 ***
## dow[, "V"]    0.2624178  0.0206945  12.681   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.003748 on 247 degrees of freedom
## Multiple R-squared:  0.8912, Adjusted R-squared:  0.8899
## F-statistic: 674.3 on 3 and 247 DF,  p-value: < 2.2e-16
```

```
anova(mod1)
```

```
## Analysis of Variance Table
##
## Response: dow[, "DIA"]
##               Df    Sum Sq    Mean Sq F value    Pr(>F)
## dow[, "MMM"]   1 0.0220972 0.0220972 1572.67 < 2.2e-16 ***
## dow[, "BA"]    1 0.0040676 0.0040676  289.49 < 2.2e-16 ***
## dow[, "V"]     1 0.0022593 0.0022593  160.80 < 2.2e-16 ***
## Residuals    247 0.0034705 0.0000141
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
#residual plot & diagnostics ( lecture notes Chp 11)
plot(mod1$residuals)
```

```
library(car)
qqPlot(mod1)
```

```
## [1]  78 205
```

As expected, we obtain 89% Rsquare from a linear model with just 3 covariates. The linear model's F statistic is highly significant so the model is a good fit. Also, each of the 3 covariates have a highly significant t statistic.Finally, the residual plot does not show any pattern or overdispersion. While the QQ plot does show 2 outliers (78,205) that may have outsize leverage, the normality assummption of the residual errors holds.

**Track the performance of the DIA versus the Dow Regression Model over a year**

We construct a portfolio with similar returns as the Dow, using the betas (coefficients) from the regression above. We track the performance of this portfolio versus the Dow over 1 year.

```
myportfolio = 0.32*dow[, "MMM"] + 0.17*dow[, "BA"]  + 0.26*dow[, "V"]
dowdf$myportfolio = myportfolio
ggplot(dowdf, aes(x=1:251, y=dowdf$DIA)) + geom_smooth(method="loess", span = 0.5, color="blue") + geom_
```

## DIA vs Regression Portfolio



So we see the Dow in blue tracks the returns of myportfolio in red very closely. We are able to visualize them apart only because the 2 loess curves have different spans (In fact, if we match the span of the 2 loess curves, we cannot distinguish beween the 2 curves! )

**1-way Anova: Compare Diversification(Holding all 30 stocks) vs Holding a Single Stock**

Assume the first investor buys just 1 stock, IBM. Whereas the second investor diversifies i.e. distributes money among all 30 Dow stocks ie. buys the DIA index. They both hold the instrument for a week. Lets compare their returns to see if there is a statistically significant difference.

```r
#Nondiversification vs Diversification
n = 7
returns = matrix(0, nrow=n*2, ncol=2)
colnames(returns) <- c("Returns", "Portfolio")
returns[,1] = c(dow[, "IBM"][1:n], dow[, "DIA"][1:n])
returns[,2] = c(rep(1, n),rep(2, n))
returns = data.frame(returns)
returns$Returns = as.double(as.character( returns$Returns ))
returns$Portfolio = factor( returns$Portfolio )
res.aov <- aov(Returns~Portfolio, data = returns)
summary(res.aov)
```

```
##             Df    Sum Sq   Mean Sq F value Pr(>F)
## Portfolio    1 0.0001346 1.346e-04   2.419  0.146
## Residuals   12 0.0006679 5.566e-05
```

```
#plot(res.aov)
TukeyHSD(res.aov)
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = Returns ~ Portfolio, data = returns)
##
## $Portfolio
##          diff         lwr         upr     p adj
## 2-1 -0.006202 -0.01489037 0.002486369 0.1458423
```

```
boxplot(Returns~Portfolio, data = returns)
```



```
# Tukey vs Welch's t test since homogenaity of variance is not met
```

From the p-value on the Tukey Test, we conclude there is no significant difference in the returns of the Dow versus individual Dow components. This test yields similar results when performed against each of the 30 tickers, not just IBM. The Dow components perform mostly alike over a weekly duration. But it is clear from the Boxplots that the homogenaiety of variance is not met. So we can perform a Welch's T test to check the Tukey Results.

```
t.test(returns$Returns[returns$Portfolio=="1"],returns$Returns[returns$Portfolio=="2"])
```

```
##
##  Welch Two Sample t-test
##
## data:  returns$Returns[returns$Portfolio == "1"] and returns$Returns[returns$Portfolio == "2"]
## t = 1.5553, df = 7.421, p-value = 0.1614
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.003119957  0.015523957
## sample estimates:
##   mean of x   mean of y
## 0.009772429 0.003570429
```

Once again, the T test p-value (0.16 vs 0.14 in Tukey, since Tukey uses a common variance) shows no significant difference in the mean returns (Equivalently, the confidence interval of mean differences includes zero)

**Cell Means Model, Means Only Model, Effects Plot**

We obtain the Cell Means model. We also obtain the regular linear model with intercept & compare it versus the mean-only model. Further, we examine an effects plot for the 2 strategies.

```
lm1 = lm(returns$Returns~returns$Portfolio - 1) # Cell means model
summary(lm1)
```

```
##
## Call:
## lm(formula = returns$Returns ~ returns$Portfolio - 1)
##
## Residuals:
##      Min       1Q    Median       3Q      Max
## -0.007636 -0.004361 -0.002362  0.002618  0.017716
##
## Coefficients:
##                   Estimate Std. Error t value Pr(>|t|)
## returns$Portfolio1 0.009772    0.002820   3.466  0.00467 **
## returns$Portfolio2 0.003570    0.002820   1.266  0.22946
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.00746 on 12 degrees of freedom
## Multiple R-squared:  0.5315, Adjusted R-squared:  0.4534
## F-statistic: 6.807 on 2 and 12 DF,  p-value: 0.01057
```

```
lm2 = lm(Returns~Portfolio,data=returns) # Regular linear model
plot(allEffects(lm2))
```

**Portfolio effect plot**



```r
lm3 <- lm(Returns~1,data=returns) # mean only model
summary(lm3)
```

```
##
## Call:
## lm(formula = Returns ~ 1, data = returns)
##
## Residuals:
##        Min         1Q     Median         3Q        Max
## -0.0074594 -0.0043732 -0.0017934 -0.0001927  0.0208166
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.006671   0.002100   3.177  0.00728 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.007857 on 13 degrees of freedom
```

```r
anova(lm2,lm3)
```

```
## Analysis of Variance Table
##
## Model 1: Returns ~ Portfolio
## Model 2: Returns ~ 1
```

```
##   Res.Df      RSS Df   Sum of Sq      F Pr(>F)
## 1     12 0.00066786
## 2     13 0.00080249 -1 -0.00013463 2.4189 0.1458
```

**1-way Anova: Compare annual returns across 10 different portfolios with 3 stocks in each portfolio**

Lets construct 10 combinations of 3 Dow stocks, and see if the mean returns are statistically different over a 1 year duration. Each portfolio is a linear combination of 3 different Dow stocks.

```
# 1-way anova lect notes 12

threecomb = combinations(5,3,tickers)
nc3 = 10 # 5 choose 3 = 10
returns = matrix(0, nrow=251*(1+nc3), ncol=2)
colnames(returns) <- c("Returns", "Portfolio")

s = 1
for(i in 1:nc3) {
  threetickers = threecomb[i,]
  ticker1 = threetickers[1]
  ticker2 = threetickers[2]
  ticker3 = threetickers[3]
  # portfolio is simple linear combination of 3 tickers
  portfolio = dow[, ticker1] + dow[, ticker2] + dow[, ticker3]
  returns[s:(s+250),1] = portfolio
  returns[s:(s+250),2] = rep(i, 251)
  s = s + 251
}
returns[s:(s+250),1] = dow[, "DIA"]
returns[s:(s+250),2] = rep(11, 251)

returns = data.frame(returns)
returns$Returns = as.double(as.character( returns$Returns ))
returns$Portfolio = factor( returns$Portfolio )
res.aov <- aov(Returns~Portfolio, data = returns)
summary(res.aov)
```

```
##               Df Sum Sq   Mean Sq F value Pr(>F)
## Portfolio     10  0.002 0.0001532   0.082      1
## Residuals   2750  5.111 0.0018584
```

```
TukeyHSD(res.aov)
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = Returns ~ Portfolio, data = returns)
##
## $Portfolio
##              diff         lwr         upr      p adj
## 2-1   -1.850562e-03 -0.01424685 0.010545728 0.9999942
```

```
## 3-1    -1.186096e-03 -0.01358239 0.011210194 0.9999999
## 4-1    -1.230466e-03 -0.01362676 0.011165824 0.9999999
## 5-1    -5.660000e-04 -0.01296229 0.011830290 1.0000000
## 6-1    -2.416562e-03 -0.01481285 0.009979728 0.9999291
## 7-1    -1.163845e-03 -0.01356013 0.011232445 0.9999999
## 8-1    -4.993785e-04 -0.01289567 0.011896911 1.0000000
## 9-1    -2.349940e-03 -0.01474623 0.010046350 0.9999453
## 10-1   -1.729845e-03 -0.01412613 0.010666445 0.9999970
## 11-1   -6.694940e-04 -0.01306578 0.011726796 1.0000000
## 3-2     6.644661e-04 -0.01173182 0.013060756 1.0000000
## 4-2     6.200956e-04 -0.01177619 0.013016386 1.0000000
## 5-2     1.284562e-03 -0.01111173 0.013680852 0.9999998
## 6-2    -5.660000e-04 -0.01296229 0.011830290 1.0000000
## 7-2     6.867171e-04 -0.01170957 0.013083007 1.0000000
## 8-2     1.351183e-03 -0.01104511 0.013747473 0.9999997
## 9-2    -4.993785e-04 -0.01289567 0.011896911 1.0000000
## 10-2    1.207171e-04 -0.01227557 0.012517007 1.0000000
## 11-2    1.181068e-03 -0.01121522 0.013577358 0.9999999
## 4-3    -4.437052e-05 -0.01244066 0.012351919 1.0000000
## 5-3     6.200956e-04 -0.01177619 0.013016386 1.0000000
## 6-3    -1.230466e-03 -0.01362676 0.011165824 0.9999999
## 7-3     2.225100e-05 -0.01237404 0.012418541 1.0000000
## 8-3     6.867171e-04 -0.01170957 0.013083007 1.0000000
## 9-3    -1.163845e-03 -0.01356013 0.011232445 0.9999999
## 10-3   -5.437490e-04 -0.01294004 0.011852541 1.0000000
## 11-3    5.166016e-04 -0.01187969 0.012912892 1.0000000
## 5-4     6.644661e-04 -0.01173182 0.013060756 1.0000000
## 6-4    -1.186096e-03 -0.01358239 0.011210194 0.9999999
## 7-4     6.662151e-05 -0.01232967 0.012462911 1.0000000
## 8-4     7.310876e-04 -0.01166520 0.013127378 1.0000000
## 9-4    -1.119474e-03 -0.01351576 0.011276816 1.0000000
## 10-4   -4.993785e-04 -0.01289567 0.011896911 1.0000000
## 11-4    5.609721e-04 -0.01183532 0.012957262 1.0000000
## 6-5    -1.850562e-03 -0.01424685 0.010545728 0.9999942
## 7-5    -5.978446e-04 -0.01299413 0.011798445 1.0000000
## 8-5     6.662151e-05 -0.01232967 0.012462911 1.0000000
## 9-5    -1.783940e-03 -0.01418023 0.010612350 0.9999959
## 10-5   -1.163845e-03 -0.01356013 0.011232445 0.9999999
## 11-5   -1.034940e-04 -0.01249978 0.012292796 1.0000000
## 7-6     1.252717e-03 -0.01114357 0.013649007 0.9999999
## 8-6     1.917183e-03 -0.01047911 0.014313473 0.9999919
## 9-6     6.662151e-05 -0.01232967 0.012462911 1.0000000
## 10-6    6.867171e-04 -0.01170957 0.013083007 1.0000000
## 11-6    1.747068e-03 -0.01064922 0.014143358 0.9999967
## 8-7     6.644661e-04 -0.01173182 0.013060756 1.0000000
## 9-7    -1.186096e-03 -0.01358239 0.011210194 0.9999999
## 10-7   -5.660000e-04 -0.01296229 0.011830290 1.0000000
## 11-7    4.943506e-04 -0.01190194 0.012890641 1.0000000
## 9-8    -1.850562e-03 -0.01424685 0.010545728 0.9999942
## 10-8   -1.230466e-03 -0.01362676 0.011165824 0.9999999
## 11-8   -1.701155e-04 -0.01256641 0.012226174 1.0000000
## 10-9    6.200956e-04 -0.01177619 0.013016386 1.0000000
## 11-9    1.680446e-03 -0.01071584 0.014076736 0.9999977
## 11-10   1.060351e-03 -0.01133594 0.013456641 1.0000000
```
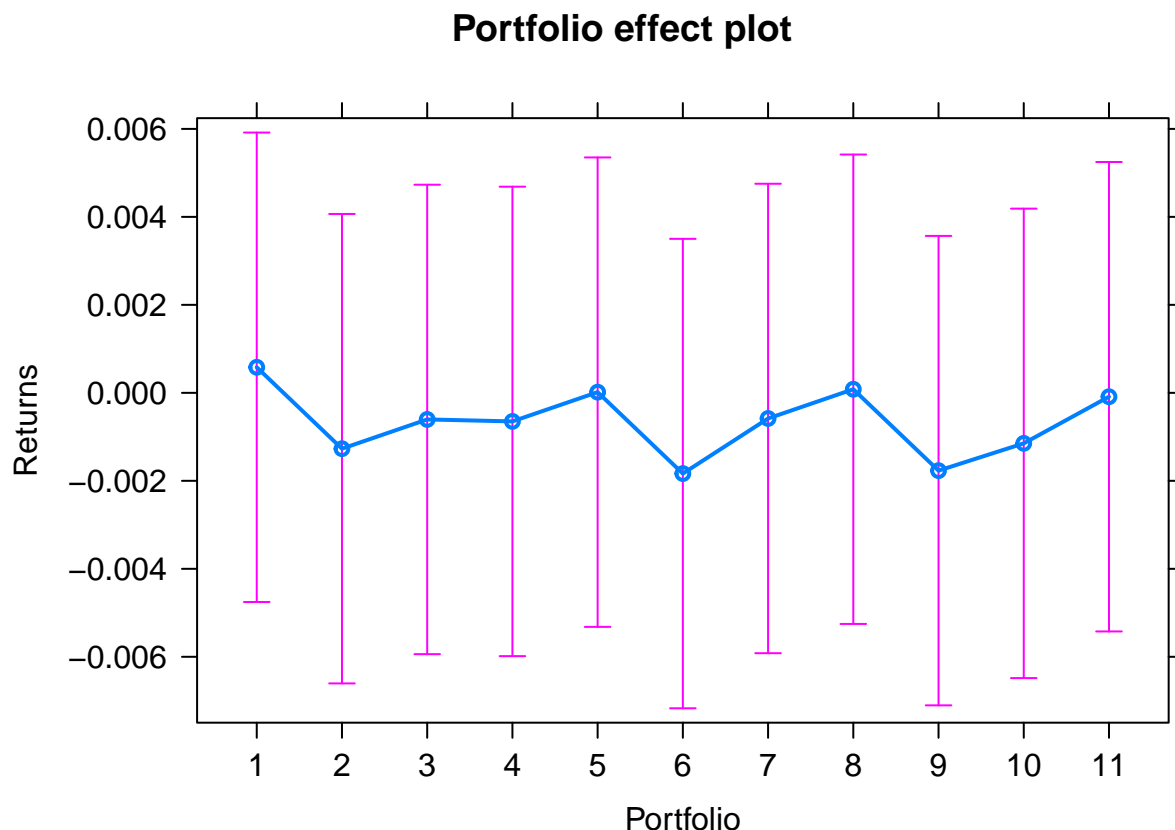
```r
boxplot(Returns~Portfolio, data = returns)
```



```r
plot(allEffects(lm(Returns~Portfolio,data=returns)))
```

## Portfolio effect plot



Amazingly, none of the 10 portfolios are significantly different over the annual duration! Further, the 11th portfolio is the Dow itself! This tells us that it is incredibly difficult to find an effective Buy and Hold Strategy with linear combination of Dow stocks, that outperforms the Dow Jones Index. No matter how we construct a portfolio, it performs as good as holding the index. This is why mutual funds which attempt to beat the Dow by performing manual selection of Dow stocks intelligently, fail to do so. From the Bpx plot we see that the 11th box ie. the Dow has the same returns but the lowest variance. Thus, holding the index is the best Buy & Hold Strategy.

**1-way Anova: Compare Buy & Hold Trading Strategy vs Buy Highest Sell Lowest vs Contrarian Strategy**

Trading Strategy 1. Sort yesterday's returns. Sell the best performing stock & buy the worst performing. Trading Strategy 2. Sort yesterday's returns. BUY the best performing stock & SELL the worst performing. Trading Strategy 3. Just buy the index (DIA).

```r
components = dow[,2:31]
trading = matrix(0,nrow=250,ncol=3)
for (i in 2:251) {
  yesterday = components[i-1,]
  today = components[i,]
  res = sort(yesterday, decreasing=TRUE, index.return=TRUE)
  best = res$ix[1]
  worst = res$ix[30]
  trading[i-1,1] = today[worst]-today[best]
  trading[i-1,2] = today[best]-today[worst]
```

```r
  trading[i-1,3] = dow[i,1]
}
tradingreturns = matrix(0, nrow=250*3, ncol=2)
tradingreturns[1:250,1] = trading[,1]
tradingreturns[251:500,1] = trading[,2]
tradingreturns[501:750,1] = trading[,3]
tradingreturns[1:250,2] = rep(1,250)
tradingreturns[251:500,2] = rep(2,250)
tradingreturns[501:750,2] = rep(3,250)

colnames(tradingreturns) <- c("Returns", "Portfolio")
tradingreturns = data.frame(tradingreturns)
tradingreturns$Returns = as.double(as.character( tradingreturns$Returns ))
tradingreturns$Portfolio = factor( tradingreturns$Portfolio )
res.aov <- aov(Returns~Portfolio, data = tradingreturns)
summary(res.aov)
```

```
##               Df  Sum Sq   Mean Sq F value Pr(>F)
## Portfolio      2 0.00103 0.0005139   1.873  0.154
## Residuals    747 0.20493 0.0002743
```
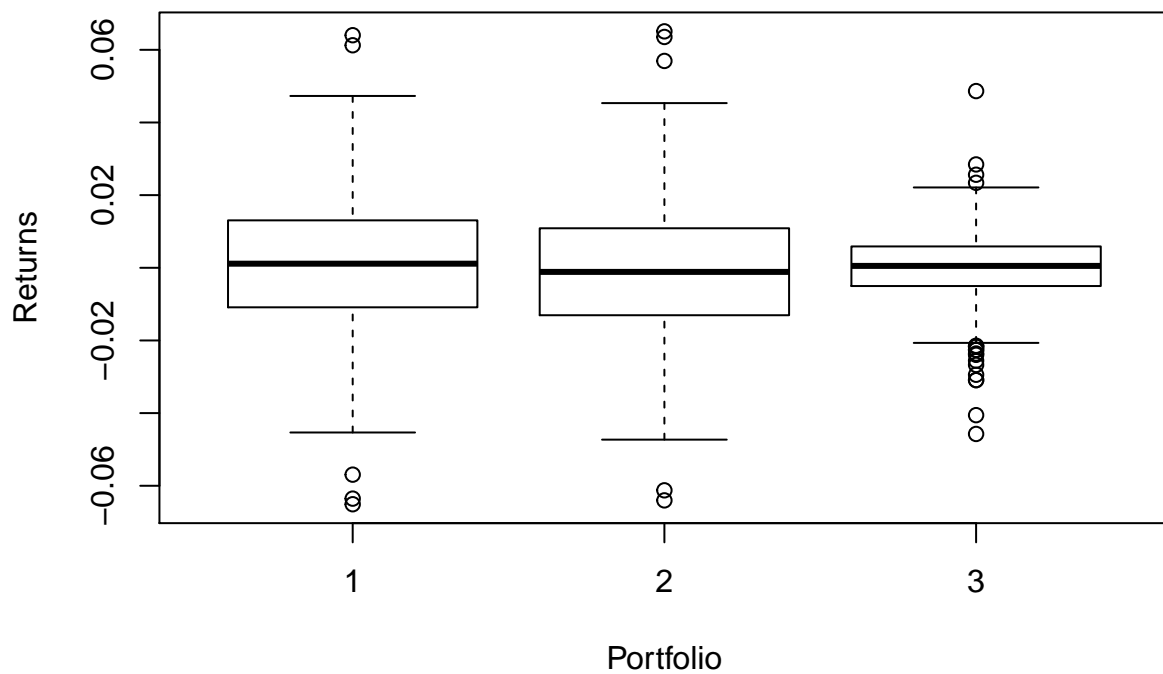
```r
TukeyHSD(res.aov)
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = Returns ~ Portfolio, data = tradingreturns)
##
## $Portfolio
##             diff          lwr          upr       p adj
## 2-1 -0.002865104 -0.006344091 0.0006138832 0.1298846
## 3-1 -0.001531376 -0.005010363 0.0019476112 0.5557975
## 3-2  0.001333728 -0.002145259 0.0048127152 0.6402838
```

```r
boxplot(Returns~Portfolio, data = tradingreturns)
```

Once again, we notice that even our moderately sophisticated trading strategy has failed to yield a profit indistinguishable from simply Buying & Holding the Dow Index! When traders complain about "lack of alpha", this is exactly what they mean: the trading strategy returns are not significantly different from zero.