

Time Series Multivariate

FRE6871 & FRE7241, Fall 2022

Jerzy Pawlowski jp3900@nyu.edu

NYU Tandon School of Engineering

October 30, 2022



NYU

**TANDON SCHOOL
OF ENGINEERING**

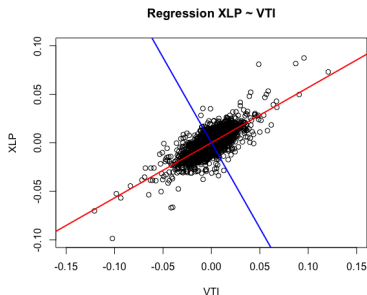
The Alpha and Beta of Stock Returns

The daily stock returns $r_i - r_f$ in excess of the risk-free rate r_f , can be decomposed into *systematic* returns $\beta(r_m - r_f)$ (where $r_m - r_f$ are the excess market returns) plus *idiosyncratic* returns $\alpha + \varepsilon_i$ (which are uncorrelated to the market returns):

$$r_i - r_f = \alpha + \beta(r_m - r_f) + \varepsilon_i$$

The *alpha* α are the abnormal returns in excess of the risk premium, and ε_i are the regression residuals with zero mean.

The *idiosyncratic* risk (equal to ε_i) is uncorrelated to the *systematic* risk, and can be reduced through portfolio diversification.



```
> # Perform regression using formula
> retsp <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
> riskfree <- 0.03/252
> retsp <- (retsp - riskfree)
> regmod <- lm(XLP ~ VTI, data=retsp)
> regmodsum <- summary(regmod)
> # Get regression coefficients
> coef(regmodsum)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	6.71e-05	8.37e-05	0.802	0.423
VTI	5.69e-01	6.80e-03	83.567	0.000

```
> # Get alpha and beta
> coef(regmodsum)[, 1]
(Intercept)      VTI
 6.71e-05      5.69e-01
```

```
> # Plot scatterplot of returns with aspect ratio 1
> plot(XLP ~ VTI, data=rutils::etfenv$returns,
+      xlim=c(-0.1, 0.1), ylim=c(-0.1, 0.1),
+      asp=1, main="Regression XLP ~ VTI")
> # Add regression line and perpendicular line
> abline(regmod, lwd=2, col="red")
> abline(a=0, b=-1/coef(regmodsum)[2, 1], lwd=2, col="blue")
```

The Statistical Significance of α and β

The stock β is independent of the risk-free rate r_f :

$$\beta = \frac{\text{Cov}(r_i, r_m)}{\text{Var}(r_m)}$$

The t -statistic (t -value) is the ratio of the estimated value divided by its standard error.

The p -value is the probability of obtaining values exceeding the t -statistic, assuming the *null hypothesis* is true.

A small p -value means that the regression coefficients are very unlikely to be zero (given the data).

The β values of stock returns are very statistically significant, but the α values are mostly not significant.

The p -value of the *Durbin-Watson* test is large, which indicates that the regression residuals are not autocorrelated.

In practice, the α , β , and the risk-free rate r_f , depend on the time interval of the data, so they're time dependent.

```
> # Get regression coefficients
> coef(regmodsum)
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 6.71e-05   8.37e-05   0.802   0.423
VTI         5.69e-01   6.80e-03  83.567   0.000
> # Calculate regression coefficients from scratch
> betav <- drop(cov(retsp$XLP, retsp$VTI)/var(retsp$VTI))
> alpha <- drop(mean(retsp$XLP) - betav*mean(retsp$VTI))
> c(alpha, betav)
[1] 6.71e-05 5.69e-01
> # Calculate the residuals
> residuals <- (retsp$XLP - (alpha + betav*retsp$VTI))
> # Calculate the standard deviation of residuals
> nrows <- NROW(residuals)
> residsd <- sqrt(sum(residuals^2)/(nrows - 2))
> # Calculate the standard errors of beta and alpha
> sum2 <- sum((retsp$VTI - mean(retsp$VTI))^2)
> betasd <- residsd/sqrt(sum2)
> alphasd <- residsd*sqrt(1/nrows + mean(retsp$VTI)^2/sum2)
> c(alphasd, betasd)
[1] 8.37e-05 6.80e-03
> # Perform the Durbin-Watson test of autocorrelation of residuals
> lmtest::dwtest(regmod)
```

Durbin-Watson test

```
data: regmod
DW = 2, p-value = 1
alternative hypothesis: true autocorrelation is greater than 0
```

The Alpha and Beta of ETF Returns

The *beta* β values of ETF returns are very statistically significant, but the *alpha* α values are mostly not significant.

Some of the ETFs with significant *alpha* α values are the bond ETFs *IEF* and *TLT* (which have performed very well), and the natural resource ETFs *USO* and *DBC* (which have performed very poorly).

```
> retsp <- rutils::etfenv$returns
> symbolv <- colnames(retsp)
> symbolv <- symbolv[symbolv != "VTI"]
> # Perform regressions and collect statistics
> betam <- sapply(symbolv, function(symbol) {
+ # Specify regression formula
+   formulav <- as.formula(paste(symbol, "~ VTI"))
+ # Perform regression
+   regmod <- lm(formulav, data=retsp)
+ # Get regression summary
+   regmodsum <- summary(regmod)
+ # Collect regression statistics
+   with(regmodsum,
+     c(beta=coefficients[2, 1],
+       pbeta=coefficients[2, 4],
+       alpha=coefficients[1, 1],
+       palpha=coefficients[1, 4],
+       pdw=lmtest::dwtest(regmod)$p.value))
+ }) # end sapply
> betam <- t(betam)
> # Sort by palpha
> betam <- betam[order(betam[, "palpha"]), ]
```

```
> betam
      beta      pbeta      alpha      palpha      pdw
IEF -0.1291 2.26e-170 2.04e-04 0.00020 3.41e-01
VXX -2.7960 0.00e+00 -1.37e-03 0.00201 3.51e-01
TLT -0.2778 2.73e-176 3.13e-04 0.00679 3.95e-01
VEU 1.0115 0.00e+00 -2.52e-04 0.01349 1.00e+00
USO 0.7248 7.48e-147 -7.54e-04 0.02847 1.40e-01
XLF 1.2932 0.00e+00 -2.55e-04 0.05233 1.00e+00
GLD 0.0435 1.46e-03 2.79e-04 0.09847 7.69e-01
XLP 0.5686 0.00e+00 1.18e-04 0.15715 1.00e+00
IVE 0.9906 0.00e+00 -6.66e-05 0.18490 1.00e+00
USMV 0.7424 0.00e+00 8.68e-05 0.20998 6.78e-01
XLV 0.7512 0.00e+00 1.05e-04 0.23998 9.00e-01
EEM 1.2352 0.00e+00 -1.66e-04 0.24228 9.99e-01
SVXY 2.2824 9.94e-177 -8.98e-04 0.27259 5.13e-06
VLUE 0.9930 0.00e+00 -1.05e-04 0.30074 9.89e-01
IWD 0.9860 0.00e+00 -4.20e-05 0.38686 1.00e+00
XLY 1.0273 0.00e+00 5.91e-05 0.46899 1.00e+00
IVW 0.9649 0.00e+00 3.10e-05 0.47553 1.00e+00
XLU 0.6498 0.00e+00 8.98e-05 0.48008 1.00e+00
VNQ 1.1879 0.00e+00 -1.25e-04 0.48820 1.00e+00
DBC 0.4162 3.13e-179 -1.21e-04 0.49141 9.91e-01
XLE 1.1203 0.00e+00 -1.03e-04 0.55764 5.78e-01
VTV 0.9687 0.00e+00 -2.44e-05 0.67054 1.00e+00
QUAL 0.9707 0.00e+00 1.89e-05 0.67869 9.93e-01
XLI 1.0092 0.00e+00 -2.76e-05 0.71849 1.00e+00
XLK 1.0868 0.00e+00 2.39e-05 0.79399 9.99e-01
IWB 0.9832 0.00e+00 -5.57e-06 0.79791 1.00e+00
MTUM 1.0291 0.00e+00 -1.66e-05 0.87321 9.42e-02
IWF 0.9944 0.00e+00 1.15e-05 0.87587 1.00e+00
XLB 1.0354 0.00e+00 -8.31e-06 0.93854 1.00e+00
VYM 0.8761 0.00e+00 5.37e-07 0.99434 1.00e+00
```

Capital Asset Pricing Model (CAPM)

The CAPM model states that the expected return for stock n : $\mathbb{E}[R_n]$ is proportional to its beta β_n times the expected excess return of the market $\mathbb{E}[R_m] - r_f$:

$$\mathbb{E}[R_n] = r_f + \beta_n(\mathbb{E}[R_m] - r_f)$$

The CAPM model states that if a stock has a higher beta then it's also expected to earn higher returns.

According to the CAPM model, assets are on average expected to earn only a *systematic* return proportional to their *systematic* risk.

The CAPM model is not a regression model.

The CAPM model depends on the choice of the risk-free rate r_f .

```
> library(PerformanceAnalytics)
> # Calculate XLP beta
> PerformanceAnalytics::CAPM.beta(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.569
> # Or
> retsxp <- na.omit(retsp[, c("XLP", "VTI")])
> betav <- drop(cov(retsxp$XLP, retsxp$VTI)/var(retsxp$VTI))
> betav
[1] 0.569
> # Calculate XLP alpha
> PerformanceAnalytics::CAPM.alpha(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.000118
> # Or
> mean(retsp$XLP - betav*retsp$VTI)
[1] NA
> # Calculate XLP bull beta
> PerformanceAnalytics::CAPM.beta.bull(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.583
> # Calculate XLP bear beta
> PerformanceAnalytics::CAPM.beta.bear(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.581
```

The Security Market Line for ETFs

The *Security Market Line* (SML) represents the linear relationship between expected stock returns and their systematic risk β .

The SML depends on the choice of the risk-free rate r_f , with a steeper SML line for lower risk-free rates r_f .

All the different SML lines pass through the point ($\beta = 1, r = R_m$) corresponding to the market, and they intersect the y-axis at the risk-free point ($\beta = 0, r = r_f$).

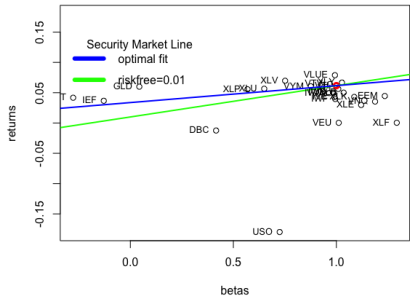
A scatterplot of asset returns versus their β shows which assets earn a positive α , and which don't.

If an asset lies on the SML, then its returns are mostly systematic, and its α is equal to zero.

Assets above the SML have a positive α , and those below have a negative α .

```
> symbolv <- rownames(betam)
> betav <- betam[~match(c("VXX", "SVXY", "MTUM", "USMV", "QUAL"),
> betav <- c(1, betav)
> names(betav)[1] <- "VTI"
> retsann <- sapply(retsp[, names(betav)], PerformanceAnalytics::R
> # Plot scatterplot of returns vs betas
> minrets <- min(retsann)
> plot(retsann ~ betav, xlab="betas", ylab="returns",
+       ylim=c(minrets, -minrets), main="Security Market Line for E
> retvti <- retsann["VTI"]
> points(x=1, y=retvti, col="red", lwd=3, pch=21)
> # Plot Security Market Line
> riskfree <- 0.01
> abline(a=riskfree, b=(retvti-riskfree), col="green", lwd=2)
```

Security Market Line for ETFs



```
> # Add labels
> text(x=betav, y=retsann, labels=names(betav), pos=2, cex=0.8)
> # Find optimal risk-free rate by minimizing residuals
> rss <- function(riskfree) {
+   sum((retsann - riskfree - betav*(retvti-riskfree))^2)
+ } # end rss
> optimrss <- optimize(rss, c(-1, 1))
> riskfree <- optimrss$minimum
> # Or simply
> retsadj <- (retsann - retvti*betav)
> betadj <- (1-betav)
> riskfree <- sum(retsadj*betadj)/sum(betadj^2)
> abline(a=riskfree, b=(retvti-riskfree), col="blue", lwd=2)
> legend(x="topleft", bty="n", title="Security Market Line",
+       legend=c("optimal fit", "riskfree=0.01"),
+       y.intersp=0.5, cex=1.0, lwd=6, lty=1, col=c("blue", "green"))
```

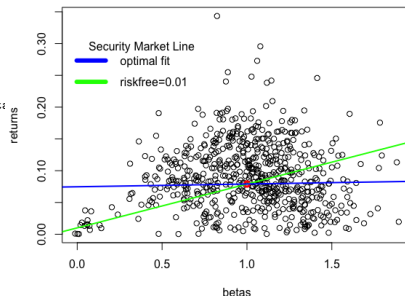
The Security Market Line for Stocks

The best fitting *Security Market Line* (SML) for stocks is almost flat, which shows that stocks with higher β don't earn higher returns.

This is called the *low beta anomaly*.

```
> # Load S&P500 constituent stock returns
> load("/Users/jerzy/Develop/lecture_slides/data/sp500_returns.RData")
> retvti <- na.omit(rutils::etfenv$returns$VTI)
> retsp <- returns[index(retvti), ]
> nrow <- NROW(retsp)
> # Calculate stock betas
> betav <- sapply(retsp, function(x) {
+   retsp <- na.omit(cbind(x, retvti))
+   drop(cov(retsp[, 1], retsp[, 2])/var(retsp[, 2]))
+ }) # end sapply
> mean(betav)
> # Calculate annual stock returns
> retsann <- retsp
> retsann[, ] <- 0
> retsann <- zoo::na.locf(retsann, na.rm=FALSE)
> retsann <- 252*sapply(retsann, sum)/nrow
> # Remove stocks with zero returns
> sum(retsann == 0)
> betav <- betav[retsann > 0]
> retsann <- retsann[retsann > 0]
> retvti <- 252*mean(retvti)
> # Plot scatterplot of returns vs betas
> plot(retsann ~ betav, xlab="betas", ylab="returns",
+   main="Security Market Line for Stocks")
> points(x=1, y=retvti, col="red", lwd=3, pch=21)
> # Plot Security Market Line
> riskfree <- 0.01
> abline(a=riskfree, b=(retvti-riskfree), col="green", lwd=2)
```

Security Market Line for Stocks



```
> # Find optimal risk-free rate by minimizing residuals
> retsadj <- (retsann - retvti*betav)
> betadj <- (1-betav)
> riskfree <- sum(retsadj*betadj)/sum(betadj^2)
> abline(a=riskfree, b=(retvti-riskfree), col="blue", lwd=2)
> legend(x="topleft", bty="n", title="Security Market Line",
+   legend=c("optimal fit", "riskfree=0.01"),
+   y.intersp=0.5, cex=1.0, lwd=6, lty=1, col=c("blue", "green"))
```

Beta-adjusted Performance Measurement

The *Treynor* ratio measures the excess returns per unit of the *systematic* risk β , and is equal to the excess returns (over a risk-free return) divided by the β :

$$T_r = \frac{E[R - r_f]}{\beta}$$

The *Treynor* ratio is similar to the *Sharpe* ratio, with the difference that its denominator represents only *systematic* risk, not total risk.

The *Information* ratio is equal to the excess returns (over a benchmark) divided by the *tracking error* (standard deviation of excess returns):

$$I_r = \frac{E[R - R_b]}{\sqrt{\sum_{i=1}^n (R_i - R_{i,b})^2}}$$

The *Information* ratio measures the amount of outperformance versus the benchmark, and the consistency of outperformance.

```
> library(PerformanceAnalytics)
> # Calculate XLP Treynor ratio
> TreynorRatio(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.098
> # Calculate XLP Information ratio
> InformationRatio(Ra=retsp$XLP, Rb=retsp$VTI)
[1] 0.0334
```


CAPM Summary Statistics

PerformanceAnalytics::table.CAPM() calculates the β and α values, the Treynor ratio, and other performance statistics.

```
> PerformanceAnalytics::table.CAPM(Ra=retsp[, c("XLP", "XLF")],
+ Rb=retsp$VTI, scale=252)
```

	XLP to VTI	XLF to VTI
Alpha	0.0001	-0.0003
Beta	0.5686	1.2932
Beta+	0.5828	1.3695
Beta-	0.5812	1.3578
R-squared	0.5673	0.7344
Annualized Alpha	0.0303	-0.0621
Correlation	0.7532	0.8570
Correlation p-value	0.0000	0.0000
Tracking Error	0.1285	0.1623
Active Premium	0.0043	-0.0673
Information Ratio	0.0334	-0.4144
Treynor Ratio	0.0980	0.0003

```
> capmstats <- table.CAPM(Ra=retsp[, symbolv],
+ Rb=retsp$VTI, scale=252)
> colnamev <- strsplit(colnames(capmstats), split=" ")
> colnamev <- do.call(cbind, colnamev[1, ])
> colnames(capmstats) <- colnamev
> capmstats <- t(capmstats)
> capmstats[, -1]
> colnamev <- colnames(capmstats)
> whichv <- match(c("Annualized Alpha", "Information Ratio", "Treynor Rat
> colnamev[whichv] <- c("Alpha", "Information", "Treynor")
> colnames(capmstats) <- colnamev
> capmstats <- capmstats[order(capmstats[, "Alpha"], decreasing=TRUE), ]
> # Copy capmstats into etfenv and save to .RData file
> etfenv <- rutils::etfenv
> etfenv$capmstats <- capmstats
> save(etfenv, file="/Users/jerzy/Develop/lecture_slides/data/etf_data.RData")
```

```
> rutils::etfenv$capmstats[, c("Beta", "Alpha", "Information", "Treynor")]
```

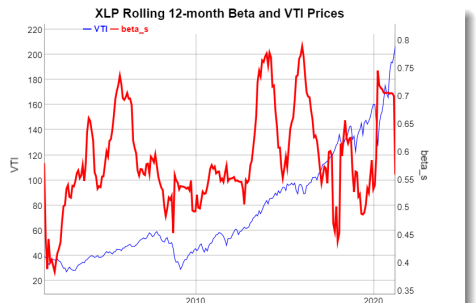
	Beta	Alpha	Information	Treynor
TLT	-0.2778	0.0820	-0.1550	-0.1503
GLD	0.0435	0.0729	-0.0567	1.3797
IEF	-0.1291	0.0526	-0.2127	-0.2844
XLP	0.5686	0.0303	0.0334	0.0980
XLV	0.7512	0.0267	0.0913	0.0928
XLU	0.6498	0.0229	-0.0313	0.0870
USMV	0.7424	0.0221	-0.1017	0.1548
XLY	1.0273	0.0150	0.1294	0.0649
IVW	0.9649	0.0078	0.1067	0.0512
XLK	1.0868	0.0061	0.0378	0.0396
QUAL	0.9707	0.0048	0.0566	0.1079
IWF	0.9944	0.0029	-0.0142	0.0398
VYM	0.8761	0.0001	-0.1183	0.0692
VTI	1.0000	0.0000	NaN	0.0617
IWB	0.9832	-0.0014	-0.1009	0.0510
XLB	1.0354	-0.0021	-0.0723	0.0484
MTUM	1.0291	-0.0042	-0.0661	0.1035
VTV	0.9687	-0.0061	-0.1680	0.0663
XLI	1.0092	-0.0069	-0.1254	0.0557
IWD	0.9860	-0.0105	-0.2393	0.0505
IVE	0.9906	-0.0167	-0.3412	0.0434
XLE	1.1203	-0.0256	-0.2150	0.0265
VLUE	0.9930	-0.0260	-0.4243	0.0795
DBC	0.4162	-0.0301	-0.3980	-0.0297
VNQ	1.1879	-0.0311	-0.2188	0.0298
EEM	1.2352	-0.0410	-0.2600	0.0361
VEU	1.0115	-0.0614	-0.6992	0.0004
XLF	1.2932	-0.0621	-0.4144	0.0003
USO	0.7248	-0.1730	-0.7117	-0.2481
SVXY	2.2824	-0.2025	NaN	NaN
VXX	-2.7960	-0.2921	-0.9119	0.2191

Rolling Beta Regressions Over Time

The rolling beta of *XLP* versus *VTI* changes over time, with lower beta in periods of *VTI* selloffs.

The function `roll_reg()` from package *HighFreq* performs rolling regressions in C++ (*RcppArmadillo*), so it's therefore much faster than equivalent R code.

```
> # Calculate XLP and VTI returns
> retsp <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
> # Calculate monthly end points
> endp <- xts::endpoints(retsp, on="months")[-1]
> # Calculate start points from look-back interval
> look_back <- 12 # Look back 12 months
> startp <- c(rep(1, look_back), endp[1:(NROW(endp)-look_back)])
> head(cbind(endp, startp), look_back+2)
> # Calculate rolling beta regressions every month in R
> formulav <- XLP ~ VTI # Specify regression formula
> betar <- sapply(1:NROW(endp), FUN=function(ep) {
+   datav <- retsp[startp[ep]:endp[ep], ]
+   # coef(lm(formulav, data=datav))[2]
+   drop(cov(datav$XLP, datav$VTI)/var(datav$VTI))
+ }) # end sapply
> # Calculate rolling betas using RcppArmadillo
> reg_stats <- HighFreq::roll_reg(response=retsp$XLP, retsp=retsp$
> betas <- reg_stats$VTI
> all.equal(betas, betar)
> # Compare the speed of RcppArmadillo with R code
> library(microbenchmark)
> summary(microbenchmark(
+   Rcpp=HighFreq::roll_reg(response=retsp$XLP, retsp=retsp$VTI, et
+   Rcode=sapply(1:NROW(endp), FUN=function(ep) {
+     datav <- retsp[startp[ep]:endp[ep], ]
+     drop(cov(datav$XLP, datav$VTI)/var(datav$VTI))
+   }),
+   times=10))[, c(1, 4, 5)] # end microbenchmark summary
```



```
> # dygraph plot of rolling XLP beta and VTI prices
> dates <- zoo::index(retsp[endp, ])
> pricev <- rutils::etfenv$prices$VTI[dates]
> datav <- cbind(pricev, betas)
> colnamev <- colnames(datav)
> dygraphs::dygraph(datav, main="XLP Rolling 12-month Beta and VTI P
+   dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
+   dySeries(name=colnamev[1], axis="y", col="blue") %>%
+   dySeries(name=colnamev[2], axis="y2", col="red", strokeWidth=3)
+   dyLegend(show="always", width=500)
```

draft: Engle-Granger Two-step Procedure for Cointegration

Use log prices?

The *ADF* test can be applied to test for the cointegration of time series of prices.

The Engle-Granger two-step procedure for two time series consists of:

- Performing a regression to calculate the cointegrating factor β ,
- Applying the *ADF* test to the residuals of the regression to determine that they don't have a unit root (they are mean reverting).

The regression of prices is not statistically valid because they are not stationary and not normally distributed.

```
> # Calculate XLB and XLE prices
> pricev <- na.omit(rutils::etfenv$prices[, c("XLB", "XLE")])
> cor(rutils::diffit(log(pricev)))
> xlb <- drop(zoo::coredata(pricev$XLB))
> xle <- drop(zoo::coredata(pricev$XLE))
> # Calculate regression coefficients of  $XLB \sim XLE$ 
> betav <- cov(xlb, xle)/var(xle)
> alpha <- (mean(xlb) - betav*mean(xle))
> # Calculate regression residuals
> fittedv <- (alpha + betav*xle)
> residuals <- (xlb - fittedv)
> # Perform ADF test on residuals
> tseries::adf.test(residuals, k=1)
```



```
> # Plot prices
> dygraphs::dygraph(pricev, main="XLB and XLE Prices") %>%
+   dyOptions(colors=c("blue", "red"))
> # Plot cointegration residuals
> residuals <- xts::xts(residuals, zoo::index(pricev))
> dygraphs::dygraph(residuals, main="XLB and XLE Cointegration Residuals")
```

Principal Components of S&P500 Stock Constituents

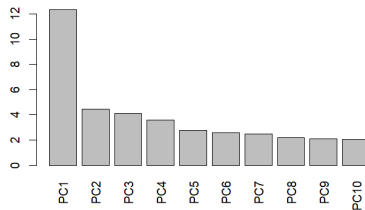
The *PCA* standard deviations are the volatilities of the *principal component* time series.

The original time series of returns can be calculated approximately from the first few *principal components* with the largest standard deviations.

The *Kaiser-Guttman* rule uses only *principal components* with *variance* greater than 1.

Another rule of thumb is to use the *principal components* with the largest standard deviations which sum up to 80% of the total variance of returns.

Volatilities of S&P500 Principal Components



```
> # Load S&P500 constituent stock prices
> load("/Users/jerzy/Develop/lecture_slides/data/sp500_prices.RData")
> # Calculate stock prices and percentage returns
> pricets <- zoo::na.locf(pricets, na.rm=FALSE)
> pricets <- zoo::na.locf(pricets, fromLast=TRUE)
> retsp <- rutils::diffit(log(pricetv))
> # Standardize (de-mean and scale) the returns
> retsp <- lapply(retsp, function(x) {(x - mean(x))/sd(x)})
> retsp <- rutils::do_call(cbind, retsp)
> # Perform principal component analysis PCA
> pcad <- prcomp(retsp, scale=TRUE)
> # Find number of components with variance greater than 2
> ncomp <- which(pcad$sdev^2 < 2)[1]
```

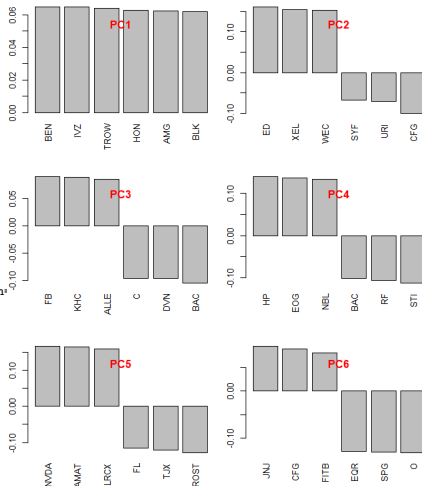
```
> # Plot standard deviations of principal components
> barplot(pcad$sdev[1:ncomp],
+   names.arg=colnames(pcad$rotation[, 1:ncomp]),
+   las=3, xlab="", ylab="",
+   main="Volatilities of S&P500 Principal Components")
```

S&P500 Principal Component Loadings (Weights)

Principal component loadings are the weights of *principal component* portfolios.

The *principal component* portfolios have mutually orthogonal returns represent the different orthogonal modes of the return variance.

```
> # Calculate principal component loadings (weights)
> # Plot barplots with PCA weights in multiple panels
> ncomps <- 6
> par(mfrow=c(ncomps/2, 2))
> par(mar=c(4, 2, 2, 1), oma=c(0, 0, 0, 0))
> # First principal component weights
> weights <- sort(pcad$rotation[, 1], decreasing=TRUE)
> barplot(weights[1:6], las=3, xlab="", ylab="", main="")
> title(paste0("PC", 1), line=-2.0, col.main="red")
> for (ordern in 2:ncomps) {
+   weights <- sort(pcad$rotation[, ordern], decreasing=TRUE)
+   barplot(weights[c(1:3, 498:500)], las=3, xlab="", ylab="", main="")
+   title(paste0("PC", ordern), line=-2.0, col.main="red")
+ } # end for
```

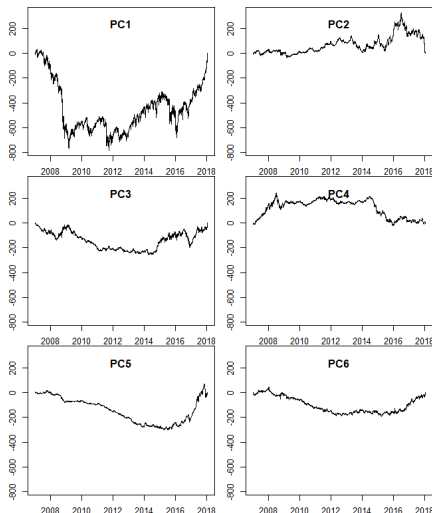


S&P500 Principal Component Time Series

The time series of the *principal components* can be calculated by multiplying the loadings (weights) times the original data.

Higher order *principal components* are gradually less volatile.

```
> # Calculate principal component time series
> retspca <- xts(retsp %*% pcad$rotation[, 1:ncomps],
+   order.by=dates)
> round(cov(retspca), 3)
> pcacum <- cumsum(retspca)
> # Plot principal component time series in multiple panels
> par(mfrow=c(ncomps/2, 2))
> par(mar=c(2, 2, 0, 1), oma=c(0, 0, 0, 0))
> rangev <- range(pcacum)
> for (ordern in 1:ncomps) {
+   plot.zoo(pcacum[, ordern], ylim=rangev, xlab="", ylab="")
+   title(paste0("PC", ordern), line=-2.0)
+ } # end for
```



S&P500 Factor Model From Principal Components

By inverting the *PCA* analysis, the *S&P500* constituent returns can be calculated from the first k *principal components* under a *factor model*:

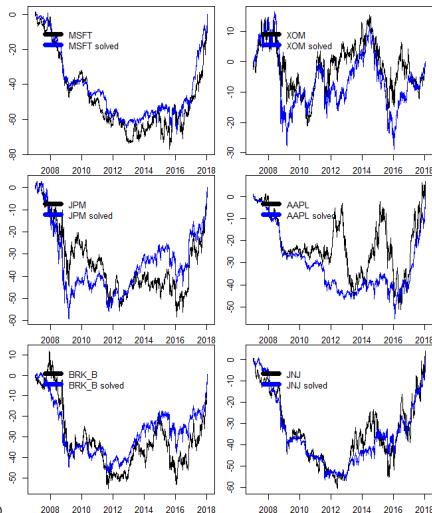
$$\mathbf{r}_i = \alpha_i + \sum_{j=1}^k \beta_{ji} \mathbf{F}_j + \varepsilon_i$$

The *principal components* are interpreted as *market factors*: $\mathbf{F}_j = \mathbf{pc}_j$.

The *market betas* are the inverse of the *principal component loadings*: $\beta_{ji} = w_{ij}$.

The ε_i are the *idiosyncratic* returns, which should be mutually independent and uncorrelated to the *market factor* returns.

```
> # Invert principal component time series
> invmat <- solve(pcad$rotation)
> all.equal(invmat, t(pcad$rotation))
> solved <- retspca %*% invmat[1:ncomps, ]
> solved <- xts::xts(solved, dates)
> solved <- cumsum(solved)
> retc <- cumsum(retsp)
> # Plot the solved returns
> symbolv <- c("MSFT", "XOM", "JPM", "AAPL", "BRK_B", "JNJ")
> for (symbol in symbolv) {
+   plot.zoo(cbind(retc[, symbol], solved[, symbol]),
+   plot.type="single", col=c("black", "blue"), xlab="", ylab="")
+   legend(x="topleft", bty="n",
+   legend=paste0(symbol, c("", " solved")),
+   title=NULL, inset=0.05, cex=1.0, lwd=6,
+   lty=1, col=c("black", "blue"))
+ } # end for
```



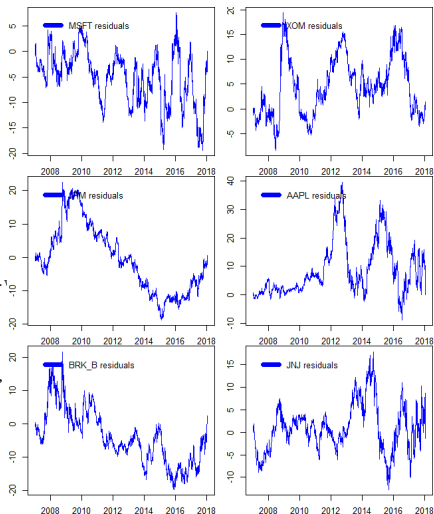
S&P500 Factor Model Residuals

The original time series of returns can be calculated exactly from the time series of all the *principal components*, by inverting the loadings matrix.

The original time series of returns can be calculated approximately from just the first few *principal components*, which demonstrates that *PCA* is a form of *dimension reduction*.

The function `solve()` solves systems of linear equations, and also inverts square matrices.

```
> # Perform ADF unit root tests on original series and residuals
> sapply(symbolv, function(symbol) {
+   c(series=tseries::adf.test(retc[, symbol])$p.value,
+     resid=tseries::adf.test(retc[, symbol] - solved[, symbol])$p.,
+   }) # end sapply
> # Plot the residuals
> for (symbol in symbolv) {
+   plot.zoo(retc[, symbol] - solved[, symbol],
+     plot.type="single", col="blue", xlab="", ylab="")
+   legend(x="topleft", bty="n", legend=paste(symbol, "residuals"),
+     title=NULL, inset=0.05, cex=1.0, lwd=6, lty=1, col="blue")
+ } # end for
> # Perform ADF unit root test on principal component time series
> retspca <- xts(retsp %*% pcad$rotation, order.by=dates)
> pcacum <- cumsum(retspca)
> adf_pvalues <- sapply(1:NCOL(pcacum), function(ordern)
+   tseries::adf.test(pcacum[, ordern])$p.value)
> # AdF unit root test on stationary time series
> tseries::adf.test(rnorm(1e5))
```

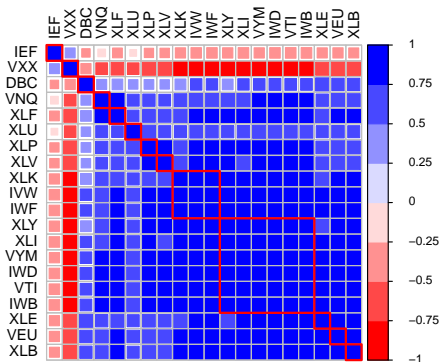


Correlation and Factor Analysis

```

> ### Perform pair-wise correlation analysis
> # Calculate correlation matrix
> cormat <- cor(retsp)
> colnames(cormat) <- colnames(retsp)
> rownames(cormat) <- colnames(retsp)
> # Reorder correlation matrix based on clusters
> # Calculate permutation vector
> library(corrplot)
> ordern <- corrMatOrder(cormat, order="hclust",
+   hclust.method="complete")
> # Apply permutation vector
> cormat <- cormat[ordern, ordern]
> # Plot the correlation matrix
> colorv <- colorRampPalette(c("red", "white", "blue"))
> corrplot(cormat, tl.col="black", tl.cex=0.8,
+   method="square", col=colorv(8),
+   cl.offset=0.75, cl.cex=0.7,
+   cl.align.text="l", cl.ratio=0.25)
> # draw rectangles on the correlation matrix plot
> corrRect.hclust(cormat, k=NROW(cormat) %/% 2,
+   method="complete", col="red")

```



Hierarchical Clustering Analysis

The function `as.dist()` converts a matrix representing the *distance* (dissimilarity) between elements, into a list of class "dist".

For example, `as.dist()` converts $(1 - \text{correlation})$ to distance.

The function `hclust()` recursively combines elements into clusters based on their mutual *distance*.

First `hclust()` combines individual elements that are closest to each other.

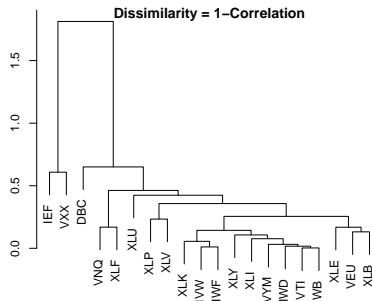
Then it combines elements to the closest clusters, then clusters with other clusters, until all elements are combined into one cluster.

This process of recursive clustering can be represented as a *dendrogram* (tree diagram).

Branches of a *dendrogram* represent clusters.

Neighboring branches contain elements that are close to each other (have small distance).

Neighboring branches combine into larger branches, that then combine with their closest branches, etc.

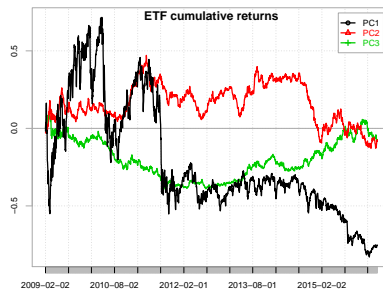


```
> # Convert correlation matrix into distance object
> distancev <- as.dist(1-cormat)
> # Perform hierarchical clustering analysis
> cluster <- hclust(distancev)
> plot(cluster, ann=FALSE, xlab="", ylab="")
> title("Dendrogram representing hierarchical clustering
+ \nwith dissimilarity = 1-correlation", line=-0.5)
```

depr: Principal Component Returns Time Series

```
> # PC returns from rotation and scaled returns
> retscaled <- apply(retsp, 2, scale)
> retspca <- retscaled %*% pcad$rotation
> # "x" matrix contains time series of PC returns
> dim(pcad$x)
> class(pcad$x)
> head(pcad$x[, 1:3], 3)
> # Convert PC matrix to xts and rescale to decimals
> retspca <- xts(pcad$x/100, order.by=zoo::index(retsp))
```

```
> chart.CumReturns(
+   retspca[, 1:3], lwd=2, ylab="",
+   legend.loc="topright", main="")
> # Add title
> title(main="ETF cumulative returns", line=-1)
```



depr: *Principal Component Returns Analysis*

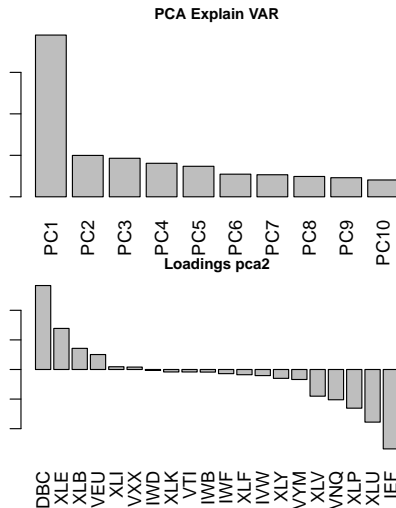
```
> # Calculate PC correlation matrix
> cormat <- cor(retspca)
> colnames(cormat) <- colnames(retspca)
> rownames(cormat) <- colnames(retspca)
> cormat[1:3, 1:3]
> table.CAPM(Ra=retspca[, 1:3], Rb=retsp$VTI, scale=252)
```

depr: Principal Component Analysis

```

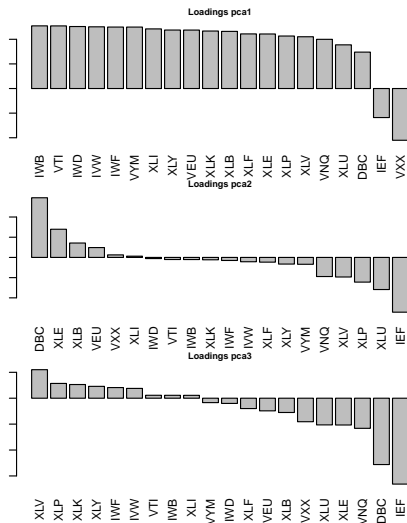
> ### Perform principal component analysis PCA
> retsp <- na.omit(rutils::etfenv$returns)
> pcad <- prcomp(retsp, center=TRUE, scale=TRUE)
> barplot(pcad$sdev[1:10],
+   names.arg=colnames(pcad$rotation)[1:10],
+   las=3, ylab="STDEV", xlab="PCVec",
+   main="PCA Explain VAR")
> # Show first three principal component loadings
> head(pcad$rotation[,1:3], 3)
> # Permute second principal component loadings by size
> pca2 <- as.matrix(
+   pcad$rotation[order(pcad$rotation[, 2],
+     decreasing=TRUE), 2])
> colnames(pca2) <- "pca2"
> head(pca2, 3)
> # The option las=3 rotates the names.arg labels
> barplot(as.vector(pca2),
+   names.arg=rownames(pca2),
+   las=3, ylab="Loadings",
+   xlab="Symbol", main="Loadings pca2")

```



depr: Principal Component Vectors

```
> # Get list of principal component vectors
> pca_vecs <- lapply(1:3, function(orden) {
+   pca_vec <- as.matrix(
+     pcad$rotation[
+       order(pcad$rotation[, orden],
+         decreasing=TRUE), orden])
+   colnames(pca_vec) <- paste0("pca", orden)
+   pca_vec
+ }) # end lapply
> names(pca_vecs) <- c("pca1", "pca2", "pca3")
> # The option las=3 rotates the names.arg labels
> for (orden in 1:3) {
+   barplot(as.vector(pca_vecs[[orden]]),
+     names.arg=rownames(pca_vecs[[orden]]),
+     las=3, xlab="", ylab="",
+     main=paste("Loadings",
+       colnames(pca_vecs[[orden]])))
+ } # end for
```



depr: Package *factorAnalytics*

The package *factorAnalytics* performs estimation and risk analysis of linear factor models for portfolio asset returns.

```
> library(factorAnalytics) # Load package "factorAnalytics"
> # Get documentation for package "factorAnalytics"
> packageDescription("factorAnalytics") # Get short description
> help(package="factorAnalytics") # Load help page
```

```
> # List all objects in "factorAnalytics"
> ls("package:factorAnalytics")
>
> # List all datasets in "factorAnalytics"
> # data(package="factorAnalytics")
>
> # Remove factorAnalytics from search path
> detach("package:factorAnalytics")
```

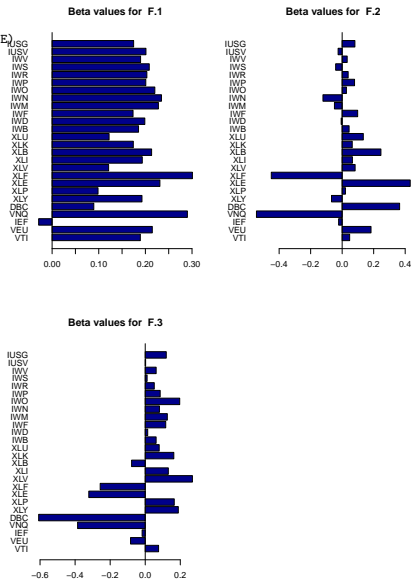
depr: Fitting Factor Models Using PCA

```
> library(factorAnalytics)
> # Fit a three-factor model using PCA
> factpca <- fitSfm(rutils::etfenv$returns, k=3)
> head(factpca$loadings, 3) # Factor loadings
> # Factor realizations (time series)
> head(factpca$factors)
> # Residuals from regression
> factpca$residuals[1:3, 1:3]
```

```
> factpca$alpha # Estimated alphas
> factpca$r2 # R-squared regression
> # Covariance matrix estimated by factor model
> factpca$Omega[1:3, 4:6]
```

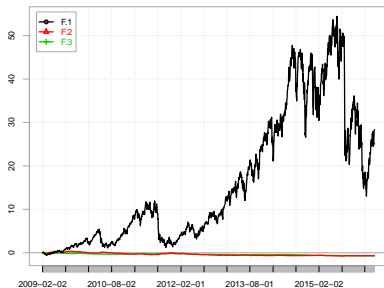

depr: Factor Loadings

```
> plot(factpca, which.plot.group=3, n.max=30, loop=FALSE)
> # ?plot.sfm
```



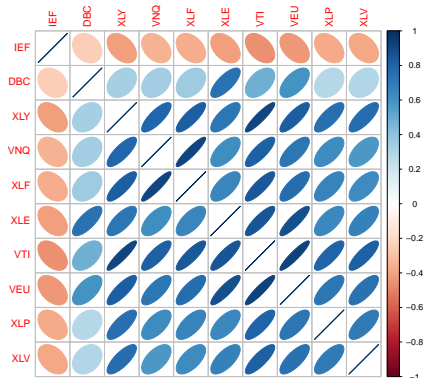
depr: Time Series of Factors

```
> library(PortfolioAnalytics)
> # Plot factor cumulative returns
> chart.CumReturns(factpca$factors,
+   lwd=2, ylab="", legend.loc="topleft", main="")
>
> # Plot time series of factor returns
> # Plot(factpca, which.plot.group=2,
> #   loop=FALSE)
```



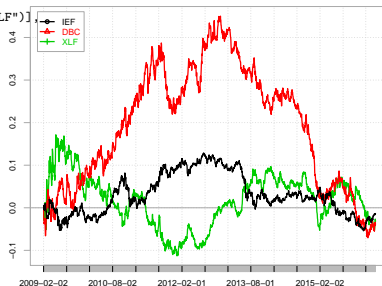
depr: Asset Correlations

```
> # Asset correlations "hclust" hierarchical clustering
> plot(factpca, which.plot.group=7, loop=FALSE,
+      order="hclust", method="ellipse")
```



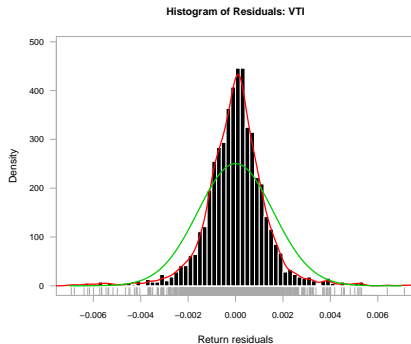
depr: Time Series of Residuals

```
> library(PortfolioAnalytics)
> # Plot residual cumulative returns
> chart.CumReturns(factpca$residuals[, c("IEF", "DBC", "XLF")],
+   lwd=2, ylab="", legend.loc="topleft", main="")
```



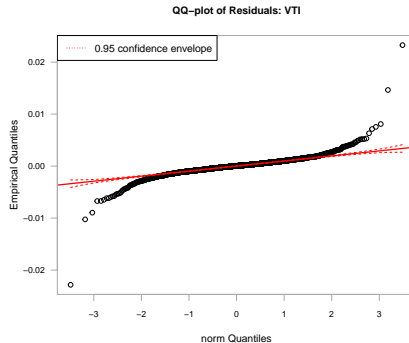
depr: Residual Returns Histogram

```
> library(PortfolioAnalytics)
> # Plot residual histogram with normal curve
> plot(factpca, asset.name="VTI",
+      which.plot.single=8,
+      plot.single=TRUE, loop=FALSE,
+      xlim=c(-0.007, 0.007))
```



depr: Residual Returns and the Q-Q Plot

```
> # Residual Q-Q plot  
> plot(factpca, asset.name="VTI",  
+       which.plot.single=9,  
+       plot.single=TRUE, loop=FALSE)
```



depr: Autocorrelation of Residuals

```
> # SACF and PACF of residuals  
> plot(factpca, asset.name="VTI",  
+      which.plot.single=5,  
+      plot.single=TRUE, loop=FALSE)
```

