

Time Series Multivariate

FRE6871 & FRE7241, Fall 2024

Jerzy Pawlowski jp3900@nyu.edu

NYU Tandon School of Engineering

October 8, 2024



NYU

**TANDON SCHOOL
OF ENGINEERING**

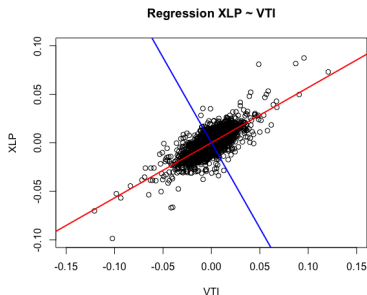
The Alpha and Beta of Stock Returns

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
> retp <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
> raterrf <- 0.03/252
> retp <- (retp - raterrf)
> regmod <- lm(XLP ~ VTI, data=retp)
> regmodsum <- summary(regmod)
> # Get regression coefficients
> coef(regmodsum)
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.27e-05   7.98e-05    0.66   0.509
VTI         5.63e-01   6.58e-03   85.55   0.000
> # Get alpha and beta
> coef(regmodsum)[, 1]
(Intercept)      VTI
 5.27e-05     5.63e-01
```

```
> # Plot scatterplot of returns with aspect ratio 1
> plot(XLP ~ VTI, data=rutils::etfenv$returns, main="Regression XLP
+       xlim=c(-0.1, 0.1), ylim=c(-0.1, 0.1), pch=1, col="blue", asp=1)
> # 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) 5.27e-05   7.98e-05   0.66   0.509
VTI         5.63e-01   6.58e-03  85.55   0.000
> # Calculate regression coefficients from scratch
> betac <- drop(cov(retp$XLP, retp$VTI)/var(retp$VTI))
> alphac <- drop(mean(retp$XLP) - betac*mean(retp$VTI))
> c(alphac, betac)
[1] 5.27e-05 5.63e-01
> # Calculate the residuals
> residuals <- (retp$XLP - (alphac + betac*retp$VTI))
> # Calculate the standard deviation of residuals
> nrows <- NROW(residuals)
> residsd <- sqrt(sum(residuals^2)/(nrows - 2))
> # Calculate the standard errors of beta and alpha
> sum2 <- sum((retp$VTI - mean(retp$VTI))^2)
> betasd <- residsd/sqrt(sum2)
> alphasd <- residsd*sqrt(1/nrows + mean(retp$VTI)^2/sum2)
> c(alphasd, betasd)
[1] 7.98e-05 6.58e-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).

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

```
> betam
      beta      pbeta      alpha      palpha      pdw
VXX -2.7525  0.00e+00 -1.50e-03  0.000368  4.29e-01
IEF -0.1130  2.19e-127  1.84e-04  0.000878  5.79e-01
VEU  0.9941  0.00e+00 -2.33e-04  0.013594  1.00e+00
TLT -0.2464  7.71e-138  2.65e-04  0.021975  6.90e-01
USO  0.6988  1.47e-152 -7.02e-04  0.028143  8.14e-02
GLD  0.0577  7.76e-06  3.09e-04  0.048563  8.18e-01
XLF  1.2675  0.00e+00 -2.36e-04  0.054352  1.00e+00
AIEQ 1.3033  5.23e-55 -8.34e-04  0.054665  9.32e-01
VLUE 0.9812  0.00e+00 -1.28e-04  0.157423  9.61e-01
EEM  1.1998  0.00e+00 -1.86e-04  0.165161  9.99e-01
XLP  0.5631  0.00e+00  1.05e-04  0.189820  1.00e+00
USMV 0.7268  0.00e+00  8.30e-05  0.196667  4.68e-01
XLV  0.7396  0.00e+00  9.35e-05  0.267866  8.20e-01
IVE  0.9794  0.00e+00 -5.31e-05  0.270286  1.00e+00
IWD  0.9759  0.00e+00 -4.80e-05  0.305419  1.00e+00
QUAL 0.9776  0.00e+00  4.26e-05  0.307553  9.94e-01
VNQ  1.1703  0.00e+00 -1.69e-04  0.311475  1.00e+00
SVXY 2.1486  3.00e-209 -6.97e-04  0.318576  7.40e-07
DBC  0.4078  1.05e-191 -1.47e-04  0.369352  9.61e-01
QQQ  1.0928  0.00e+00  5.88e-05  0.503593  9.97e-01
IVW  0.9745  0.00e+00  2.63e-05  0.535839  1.00e+00
XLE  1.0936  0.00e+00 -9.64e-05  0.567808  4.81e-01
XLU  0.6466  0.00e+00  6.24e-05  0.609462  9.99e-01
VTV  0.9517  0.00e+00 -2.18e-05  0.693018  1.00e+00
XLK  1.0997  0.00e+00  2.91e-05  0.739335  9.99e-01
IWF  1.0037  0.00e+00  2.04e-05  0.766476  1.00e+00
SPY  0.9854  0.00e+00 -6.24e-06  0.777935  1.00e+00
XLI  1.0012  0.00e+00 -1.92e-05  0.790903  1.00e+00
XLY  1.0386  0.00e+00  1.81e-05  0.817834  1.00e+00
XLB  1.0292  0.00e+00 -2.16e-05  0.831688  1.00e+00
IWB  0.9835  0.00e+00 -3.21e-06  0.872540  1.00e+00
MTUM 1.0018  0.00e+00  1.22e-05  0.902551  3.67e-02
VYM  0.8651  0.00e+00 -2.94e-06  0.966763  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=retp$XLP, Rb=retp$VTI)
[1] 0.563
> # Or
> retxlp <- na.omit(retp[, c("XLP", "VTI")])
> betac <- drop(cov(retxlp$XLP, retxlp$VTI)/var(retxlp$VTI))
> betac
[1] 0.563
> # Calculate XLP alpha
> PerformanceAnalytics::CAPM.alpha(Ra=retp$XLP, Rb=retp$VTI)
[1] 0.000105
> # Or
> mean(retp$XLP - betac*retp$VTI)
[1] NA
> # Calculate XLP bull beta
> PerformanceAnalytics::CAPM.beta.bull(Ra=retp$XLP, Rb=retp$VTI)
[1] 0.578
> # Calculate XLP bear beta
> PerformanceAnalytics::CAPM.beta.bear(Ra=retp$XLP, Rb=retp$VTI)
[1] 0.579
```

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 ($\beta = 0, r = r_f$).

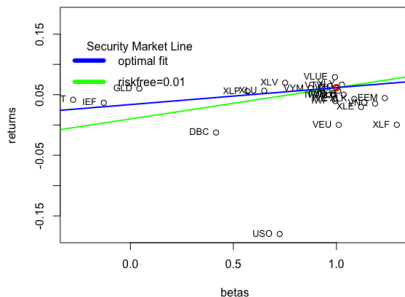
A scatterplot of asset returns versus their β shows which assets earn a positive α , and which don't.

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)
> betac <- betam[-match(c("VXX", "SVXY", "MTUM", "USMV", "QUAL"),
> betac <- c(1, betac)
> names(betac)[1] <- "VTI"
> retsann <- supply(retp[, names(betac)], PerformanceAnalytics::Re
> # Plot scatterplot of returns vs betas
> minrets <- min(retsann)
> plot(retsann ~ betac, 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
> raterf <- 0.01
> abline(a=raterf, b=(retvti-raterf), col="green", lwd=2)
```

Security Market Line for ETFs



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

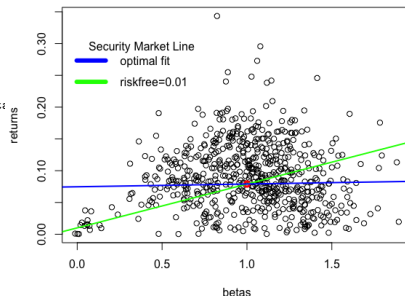
The Security Market Line for Stocks

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)
> retp <- retstock[index(retvti), ]
> nrow <- NROW(retp)
> # Calculate stock betas
> betac <- sapply(retp, function(x) {
+   retp <- na.omit(cbind(x, retvti))
+   drop(cov(retp[, 1], retp[, 2])/var(retp[, 2]))
+ }) # end sapply
> mean(betac)
> # Calculate annual stock returns
> retsann <- retp
> retsann[, 1] <- 0
> retsann <- zoo::na.locf(retsann, na.rm=FALSE)
> retsann <- 252*sapply(retsann, sum)/nrow
> # Remove stocks with zero returns
> sum(retsann == 0)
> betac <- betac[retsann > 0]
> retsann <- retsann[retsann > 0]
> retvti <- 252*mean(retvti)
> # Plot scatterplot of returns vs betas
> plot(retsann ~ betac, 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
> raterf <- 0.01
> abline(a=raterf, b=(retvti-raterf), col="green", lwd=2)
```

Security Market Line for Stocks



```
> # Find optimal risk-free rate by minimizing residuals
> retsadj <- (retsann - retvti*betac)
> betadj <- (1-betac)
> raterf <- sum(retsadj*betadj)/sum(betadj^2)
> abline(a=raterf, b=(retvti-raterf), col="blue", lwd=2)
> legend(x="topleft", bty="n", title="Security Market Line",
+   legend=c("optimal fit", "raterf=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 rate) 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=retp$XLP, Rb=retp$VTI)
[1] 0.101
> # Calculate XLP Information ratio
> InformationRatio(Ra=retp$XLP, Rb=retp$VTI)
[1] -0.0265
```


CAPM Summary Statistics

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

```
> PerformanceAnalytics::table.CAPM(Ra=retp[, c("XLP", "XLF")],
+                                Rb=retp$VTI, scale=252)
                                XLP to VTI XLF to VTI
Alpha                          0.0001   -0.0002
Beta                           0.5631   1.2675
Beta+                           0.5784   1.3433
Beta-                           0.5794   1.3425
R-squared                       0.5559   0.7284
Annualized Alpha                0.0267  -0.0578
Correlation                     0.7456   0.8535
Correlation p-value             0.0000   0.0000
Tracking Error                  0.1283   0.1576
Active Premium                  -0.0034  -0.0613
Information Ratio                -0.0265  -0.3887
Treynor Ratio                   0.1010   0.0100

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

```
> rutils::etfenv$capmstats[, c("Beta", "Alpha", "Information", "Treynor")
Beta Alpha Information Treynor
GLD 0.0508 0.0818 -0.0489 1.3873
TLT -0.2500 0.0690 -0.2160 -0.1226
IEF -0.1147 0.0475 -0.2584 -0.2841
XLV 0.7442 0.0286 0.0932 0.0968
XLP 0.5570 0.0242 -0.0275 0.1020
XLU 0.6212 0.0238 -0.0383 0.0872
USMV 0.7490 0.0158 -0.2299 0.1544
XLY 1.0073 0.0109 0.0387 0.0647
XLI 0.9691 0.0102 0.0404 0.0668
XLB 0.9501 0.0076 -0.0411 0.0574
QQQ 1.1892 0.0039 -0.0002 0.0492
VTI 0.9945 0.0033 0.1038 0.0701
QUAL 1.0023 0.0031 0.1015 0.1211
IWB 0.9789 0.0021 0.0180 0.0602
XLE 0.9901 0.0016 -0.1097 0.0372
IVW 0.9872 0.0004 -0.0328 0.0582
SPY 1.0000 0.0000 NaN 0.0851
IWD 0.9517 0.0000 -0.0662 0.0579
VYM 0.8722 -0.0020 -0.1698 0.0771
IVE 0.9655 -0.0024 -0.0961 0.0552
XLK 1.1514 -0.0032 -0.0469 0.0473
MTUM 1.0227 -0.0036 -0.0645 0.1166
IWF 1.0253 -0.0039 -0.0832 0.0505
VTV 0.9591 -0.0055 -0.1695 0.0743
DBC 0.4080 -0.0368 -0.4691 -0.0389
VLUE 0.9950 -0.0372 -0.5773 0.0792
VNQ 1.1671 -0.0405 -0.2720 0.0281
XLF 1.2322 -0.0410 -0.2909 0.0103
EEM 1.2092 -0.0444 -0.2892 0.0372
VEU 1.0066 -0.0590 -0.6956 0.0132
USO 0.7064 -0.1633 -0.7062 -0.2321
SVXY 2.1851 -0.1706 NaN NaN
AIEQ 1.2978 -0.1997 -1.8897 0.0601
VXX -2.8192 -0.3074 -0.9509 0.2212
```

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

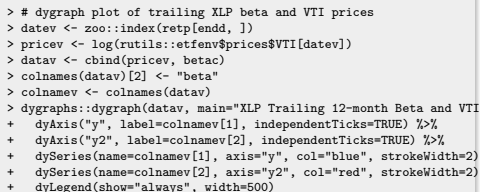
```
> # Calculate XLP and VTI returns
> retp <- na.omit(rutils::etfenv$returns[, c("XLP", "VTI")])
> # Calculate monthly end points
> endd <- xts::endpoints(retp, on="months")[-1]
> # Calculate start points from look-back interval
> lookb <- 12 # Look back 12 months
> startp <- c(retp(1, lookb), endd[1:(NROW(endd)-lookb)])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



```
> # dygraph plot of trailing XLP beta and VTI prices
> datev <- zoo::index(retp[emdd, 1])
> pricev <- rutils::etfenv$prices$VTI[datev]
> datav <- cbind(pricev, betac)
> colnames(datav)[2] <- "beta"
> colnamev <- colnames(datav)
> dygraphs::dygraph(datav, main="XLP Trailing 12-month Beta and VTI
+ 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=2)
+ dyLegend(show="always", width=500)
```

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

```
time=10))[-1:4,5]] # and midpoint of each column
```



Recursive Trailing Stock Beta

The trailing beta β of a stock with returns r_t with respect to a stock index with returns R_t can be updated using these recursive formulas with the decay factor λ :

$$\bar{r}_t = \lambda \bar{r}_{t-1} + (1 - \lambda) r_t$$

$$\bar{R}_t = \lambda \bar{R}_{t-1} + (1 - \lambda) R_t$$

$$\sigma_t^2 = \lambda \sigma_{t-1}^2 + (1 - \lambda) (R_t - \bar{R}_t)^2$$

$$\text{cov}_t = \lambda \text{cov}_{t-1} + (1 - \lambda) (r_t - \bar{r}_t) (R_t - \bar{R}_t)$$

$$\beta_t = \frac{\text{cov}_t}{\sigma_t^2}$$

The parameter λ determines the rate of decay of the weight of past returns. If λ is close to 1 then the decay is weak and past returns have a greater weight, and the trailing mean values have a stronger dependence on past returns. This is equivalent to a long look-back interval. And vice versa if λ is close to 0.

The function `HighFreq::run_covar()` calculates the trailing variances, covariances, and means of two *time series*.

```
> # Calculate the trailing betas
> lambdaf <- 0.99
> covarv <- HighFreq::run_covar(retp, lambdaf)
> betac <- covarv[, 1]/covarv[, 3]
```



```
> # dygraph plot of trailing XLP beta and VTI prices
> datav <- cbind(pricev, betac[endsd])[-(1:11)] # Remove warmup period
> colnames(datav)[2] <- "beta"
> colnamev <- colnames(datav)
> dygraphs::dygraph(datav, main="XLP Trailing 12-month Beta and VTI Prices")
+ dyAxis("y", label=colnamev[1], independentTicks=TRUE) %>%
+ dyAxis("y2", label=colnamev[2], independentTicks=TRUE) %>%
+ dySeries(name=colnamev[1], axis="y", col="blue", strokeWidth=2)
+ dySeries(name=colnamev[2], axis="y2", col="red", strokeWidth=2)
+ dyLegend(show="always", width=500)
```

Principal Components of S&P500 Stock Constituents

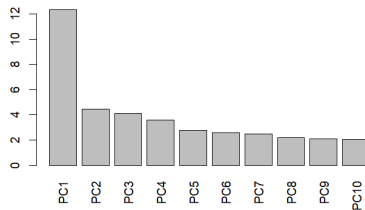
The *PCA* standard deviations are the volatilities of the *principal component* time series.

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)
> retp <- rutils::diffit(log(pricets))
> # Standardize (center and scale) the returns
> retp <- lapply(retp, function(x) {(x - mean(x))/sd(x)})
> retp <- rutils::do_call(cbind, retp)
> # Perform principal component analysis PCA
> pcad <- prcomp(retp, scale=TRUE)
> # Find number of components with variance greater than 2
> ncomp <- which(pcad$sdev^2 < 2)[1]
```

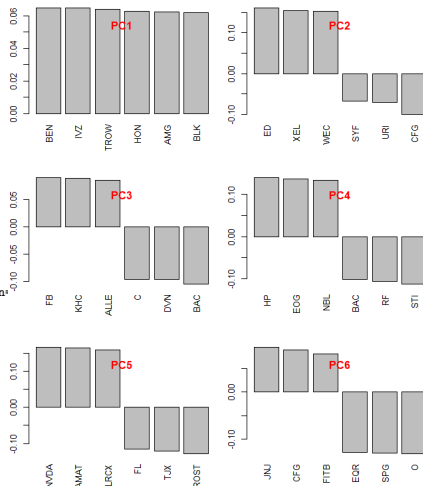
```
> # Plot standard deviations of principal components
> barplot(pcad$sdev[1:ncomp],
+   names.arg=colnames(pcad$rotation[, 1:ncomp]),
+   las=3, xlab="", ylab="",
+   main="Volatilities of S&P500 Principal Components")
```

S&P500 Principal Component Loadings

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
> weightv <- sort(pcad$rotation[, 1], decreasing=TRUE)
> barplot(weightv[1:6], las=3, xlab="", ylab="", main="")
> title(paste0("PC", 1), line=-2.0, col.main="red")
> for (ordern in 2:ncomps) {
+   weightv <- sort(pcad$rotation[, ordern], decreasing=TRUE)
+   barplot(weightv[c(1:3, 498:500)], las=3, xlab="", ylab="", main="")
+   title(paste0("PC", ordern), line=-2.0, col.main="red")
+ } # end for
```

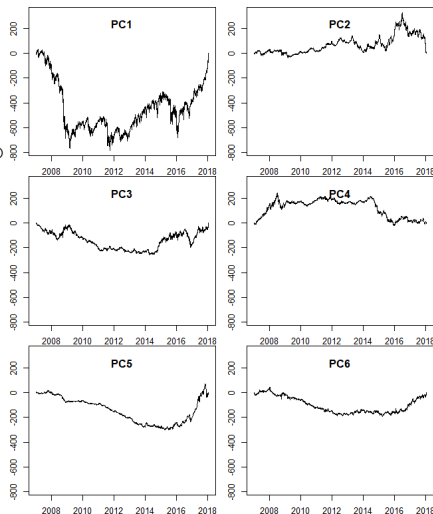


S&P500 Principal Component Time Series

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
> retpc <- xts(retp %>% pcad$rotation[, 1:ncomps], order.by=datev)
> round(cov(retpc), 3)
> retpcac <- cumsum(retpc)
> # 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(retpcac)
> for (ordern in 1:ncomps) {
+   plot.zoo(retpcac[, 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*:

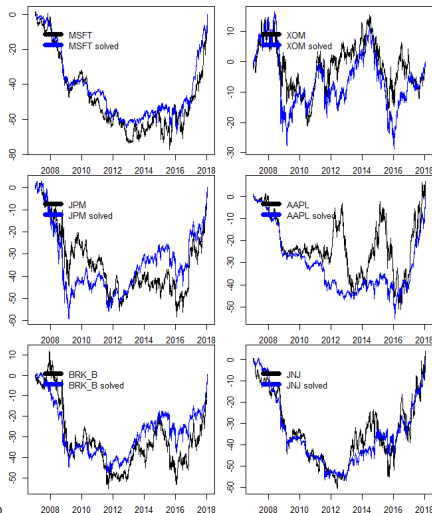
$$\mathbf{r}_i = \alpha_i + \sum_{j=1}^k \beta_{ji} \mathbf{F}_j + \varepsilon_i$$

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

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
> pcinv <- solve(pcad$rotation)
> all.equal(pcinv, t(pcad$rotation))
> solved <- retpca %*% pcinv[1:ncomps, ]
> solved <- xts::xts(solved, datev)
> solved <- cumsum(solved)
> retc <- cumsum(retp)
> # 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
```



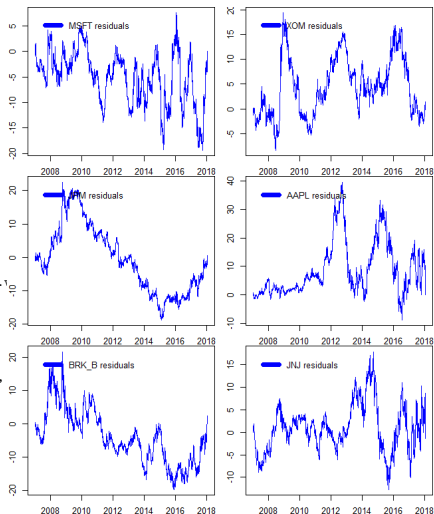
S&P500 Factor Model Residuals

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
> retpcac <- xts(retp %>% pcad$rotation, order.by=datev)
> retpcac <- cumsum(retpcac)
> adf_pvalues <- sapply(1:NCOL(retpcac), function(ordern)
+   tseries::adf.test(retpcac[, ordern])$p.value)
> # AdF unit root test on stationary time series
> tseries::adf.test(rnorm(1e5))
```

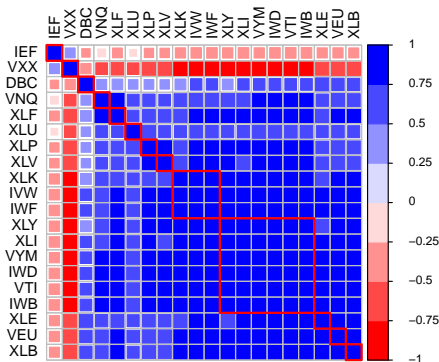


Correlation and Factor Analysis

```

> ### Perform pair-wise correlation analysis
> # Calculate correlation matrix
> cormat <- cor(retp)
> colnames(cormat) <- colnames(retp)
> rownames(cormat) <- colnames(retp)
> # 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.

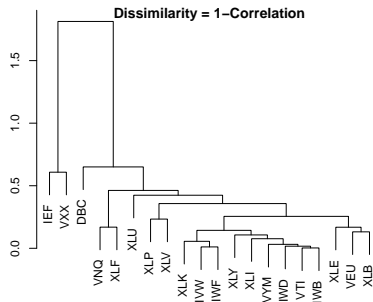
Then it combines elements to the closest clusters, then clusters with other clusters, until all elements are combined into one cluster.

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.

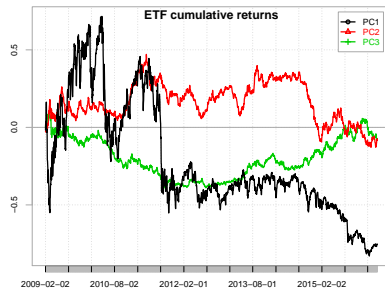


```
> # Convert correlation matrix into distance object
> distancev <- as.dist(1-cormat)
> # Perform hierarchical clustering analysis
> compclust <- hclust(distancev)
> plot(compclust, ann=FALSE, xlab="", ylab="")
> title("Dendrogram representing hierarchical clustering
+ \nwith dissimilarity = 1-correlation",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

depr: Principal Component Returns Time Series

```
> # PC returns from rotation and scaled returns
> retsc <- apply(retp, 2, scale)
> retpca <- retsc %*% 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
> retpca <- xts(pcad$x/100, order.by=zoo::index(retp))
```

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



depr: *Principal Component* Returns Analysis

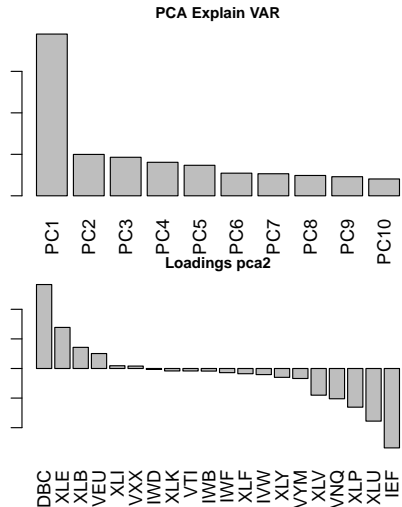
```
> # Calculate PC correlation matrix
> cormat <- cor(retpca)
> colnames(cormat) <- colnames(retpca)
> rownames(cormat) <- colnames(retpca)
> cormat[1:3, 1:3]
> table.CAPM(Ra=retpca[, 1:3], Rb=retpca$VTI, scale=252)
```

depr: Principal Component Analysis

```

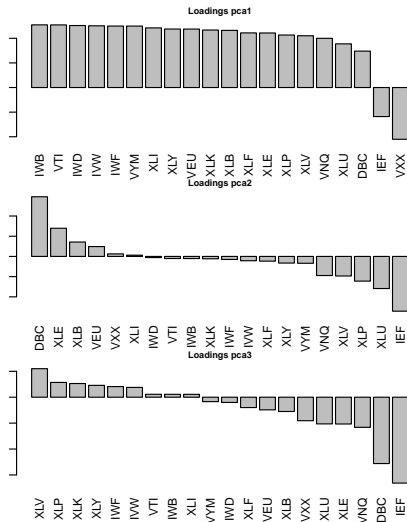
> ### Perform principal component analysis PCA
> retp <- na.omit(rutils::etfenv$returns)
> pcad <- prcomp(retp, center=TRUE, scale=TRUE)
> barplot(pcad$sdev[1:10],
+   names.arg=colnames(pcad$rotation)[1:10],
+   las=3, ylab="STDEV", xlab="PCVec",
+   main="PCA Explain VAR")
> # Show first three principal component loadings
> head(pcad$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")
```

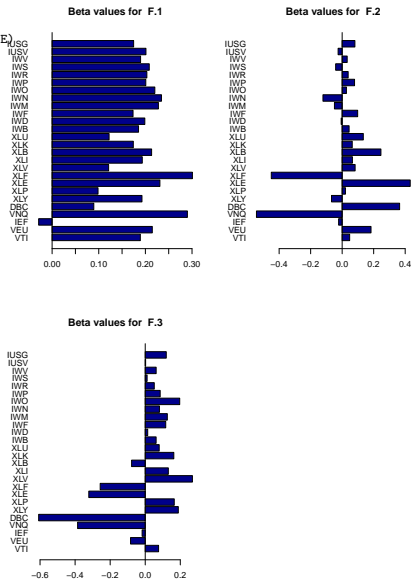

depr: Fitting Factor Models Using PCA

```
> 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
```



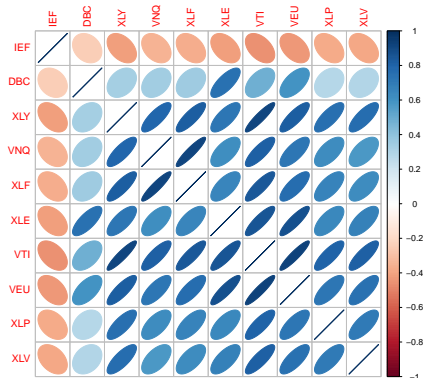
depr: Time Series of Factors

```
> 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)
```



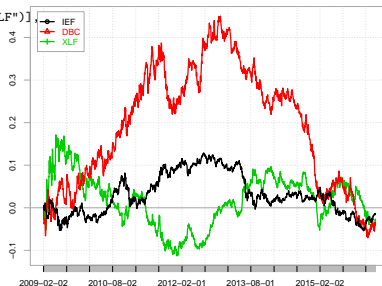
depr: Asset Correlations

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



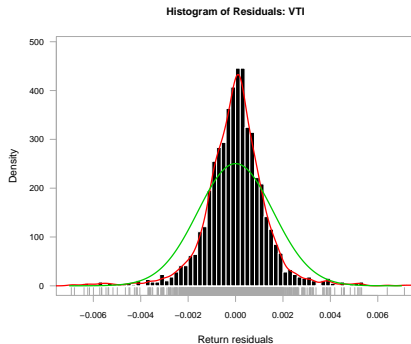
depr: Time Series of Residuals

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



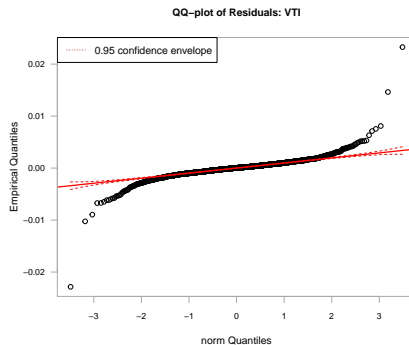
depr: Residual Returns Histogram

```
> 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)
```

