

FRE6871 R in Finance

Lecture#4, Fall 2023

Jerzy Pawlowski jp3900@nyu.edu

NYU Tandon School of Engineering

October 2, 2023



NYU

**TANDON SCHOOL
OF ENGINEERING**

One-dimensional Optimization Using The Functional optimize()

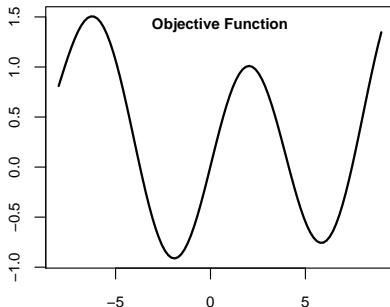
The functional `optimize()` performs *one-dimensional* optimization over a single independent variable.

`optimize()` searches for the minimum of the objective function with respect to its first argument, in the specified interval.

`optimize()` returns a list containing the location of the minimum and the objective function value,

The argument `tol` specifies the numerical accuracy, with smaller values of `tol` requiring more computations.

```
> # Display the structure of optimize()
> str(optimize)
> # Objective function with multiple minima
> objfun <- function(input, param1=0.01) {
+   sin(0.25*pi*input) + param1*(input-1)^2
+ } # end objfun
> opt1ml <- optimize(f=objfun, interval=c(-4, 2))
> class(opt1ml)
> unlist(opt1ml)
> # Find minimum in different interval
> unlist(optimize(f=objfun, interval=c(0, 8)))
> # Find minimum with less accuracy
> accl <- 1e4*.Machine$double.eps^0.25
> unlist(optimize(f=objfun, interval=c(0, 8), tol=accl))
> # Microbenchmark optimize() with less accuracy
> library(microbenchmark)
> summary(microbenchmark(
+   more_accurate = optimize(f=objfun, interval=c(0, 8)),
+   less_accurate = optimize(f=objfun, interval=c(0, 8), tol=accl),
+   times=100))[, c(1, 4, 5)] # end microbenchmark summary
```



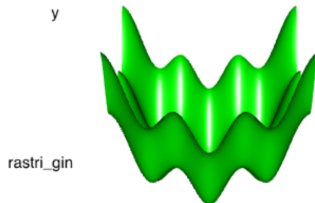
```
> # Plot the objective function
> curve(expr=objfun, type="l", xlim=c(-8, 9),
+ xlab="", ylab="", lwd=2)
> # Add title
> title(main="Objective Function", line=-1)
```

Package *rgl* for Interactive 3d Surface Plots

The package *rgl* creates *interactive* 3d scatter plots and surface plots by calling the [WebGL JavaScript](#) library.

The function `rgl::persp3d()` plots an *interactive* 3d surface plot of a *vectorized* function or a matrix.

```
> # Rastrigin function
> rastrigin <- function(x, y, param=25) {
+   x^2 + y^2 - param*(cos(x) + cos(y))
+ } # end rastrigin
> # Rastrigin function is vectorized!
> rastrigin(c(-10, 5), c(-10, 5))
> # Set rgl options and load package rgl
> library(rgl)
> options(rgl.useNULL=TRUE)
> # Draw 3d surface plot of function
> rgl::persp3d(x=rastrigin, xlim=c(-10, 10), ylim=c(-10, 10),
+   col="green", axes=FALSE, param=15)
> # Render the 3d surface plot of function
> rgl::rglwidget(elementId="plot3drgl", width=400, height=400)
```



Multi-dimensional Optimization Using `optim()`

The function `optim()` performs *multi-dimensional* optimization.

The argument `fn` is the objective function to be minimized.

The argument of `fn` that is to be optimized, must be a vector argument.

The argument `par` is the initial vector argument value.

`optim()` accepts additional parameters bound to the dots `"..."` argument, and passes them to the `fn` objective function.

The arguments `lower` and `upper` specify the search range for the variables of the objective function `fn`.

`method="L-BFGS-B"` specifies the quasi-Newton *gradient* optimization method.

`optim()` returns a list containing the location of the minimum and the objective function value.

The *gradient* methods used by `optim()` can only find the local minimum, not the global minimum.

```
> # Rastrigin function with vector argument for optimization
> rastrigin <- function(vecv, param=25) {
+   sum(vecv^2 - param*cos(vecv))
+ } # end rastrigin
> vecv <- c(pi, pi/4)
> rastrigin(vecv=vecv)
> # Draw 3d surface plot of Rastrigin function
> rgl::persp3d(
+   x=Vectorize(function(x, y) rastrigin(vecv=c(x, y))),
+   xlim=c(-10, 10), ylim=c(-10, 10),
+   col="green", axes=FALSE, zlab="", main="rastrigin")
> # Optimize with respect to vector argument
> optim1 <- optim(par=vecv, fn=rastrigin,
+   method="L-BFGS-B",
+   upper=c(14*pi, 14*pi),
+   lower=c(pi/2, pi/2),
+   param=1)
> # Optimal parameters and value
> optim1$par
> optim1$value
> rastrigin(optim1$par, param=1)
```

The Likelihood Function

The *likelihood* function $\mathcal{L}(\theta|\bar{x})$ is a function of the parameters of a statistical model θ , given a sample of observed values \bar{x} , taken under the model's probability distribution $p(x|\theta)$:

$$\mathcal{L}(\theta|x) = \prod_{i=1}^n p(x_i|\theta)$$

The *likelihood* function measures how *likely* are the parameters of a statistical model, given a sample of observed values \bar{x} .

The *maximum-likelihood* estimate (*MLE*) of the model's parameters are those that maximize the *likelihood* function:

$$\theta_{MLE} = \arg \max_{\theta} \mathcal{L}(\theta|x)$$

In practice the logarithm of the *likelihood* $\log(\mathcal{L})$ is maximized, instead of the *likelihood* itself.

The function `outer()` calculates the *outer* product of two matrices, and by default multiplies the elements of its arguments.

```
> # Sample of normal variables
> datav <- rnorm(1000, mean=4, sd=2)
> # Objective function is log-likelihood
> objfun <- function(parv, datav) {
+   sum(2*log(parv[2])) +
+   ((datav - parv[1])/parv[2])^2)
+ } # end objfun
> # Objective function on parameter grid
> parmean <- seq(1, 6, length=50)
> parsd <- seq(0.5, 3.0, length=50)
> objective_grid <- sapply(parmean, function(m) {
+   sapply(parsd, function(sd) {
+     objfun(c(m, sd), datav)
+   }) # end sapply
+ }) # end sapply
> # Perform grid search for minimum
> objective_min <- which(
+   objective_grid==min(objective_grid),
+   arr.ind=TRUE)
> objective_min
> parmean[objective_min[1]] # mean
> parsd[objective_min[2]] # sd
> objective_grid[objective_min]
> objective_grid[objective_min[, 1] + -1:1,
+   (objective_min[, 2] + -1:1)]
> # Or create parameter grid using function outer()
> objvec <- Vectorize(
+   FUN=function(mean, sd, datav)
+     objfun(c(mean, sd), datav),
+   vectorize.args=c("mean", "sd")
+ ) # end Vectorize
> objective_grid <- outer(parmean, parsd,
+   objvec, datav=datav)
```

Perspective Plot of Likelihood Function

The function `persp()` plots a 3d perspective surface plot of a function specified over a grid of argument values.

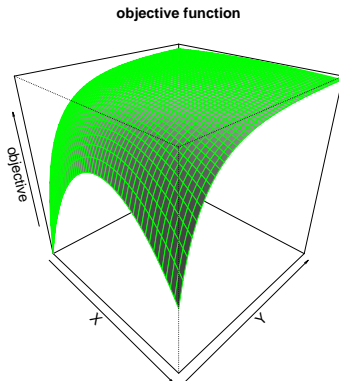
The argument "z" accepts a matrix containing the function values.

`persp()` belongs to the base graphics package, and doesn't create interactive plots.

The function `rgl::persp3d()` plots an *interactive* 3d surface plot of a function or a matrix.

`rgl` is an R package for 3d and perspective plotting, based on the *OpenGL* framework.

```
> # Perspective plot of log-likelihood function
> persp(z=-objective_grid,
+ theta=45, phi=30, shade=0.5,
+ border="green", zlab="objective",
+ main="objective function")
> # Interactive perspective plot of log-likelihood function
> library(rgl) # Load package rgl
> rgl::par3d(cex=2.0) # Scale text by factor of 2
> rgl::persp3d(z=-objective_grid, zlab="objective",
+ col="green", main="objective function")
```

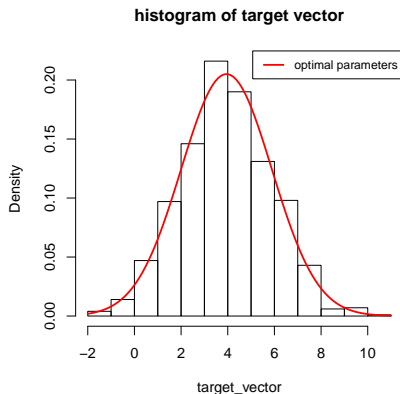


Optimization of Objective Function

The function `optim()` performs optimization of an objective function.

The function `fitdistr()` from package *MASS* fits a univariate distribution to a sample of data, by performing *maximum likelihood* optimization.

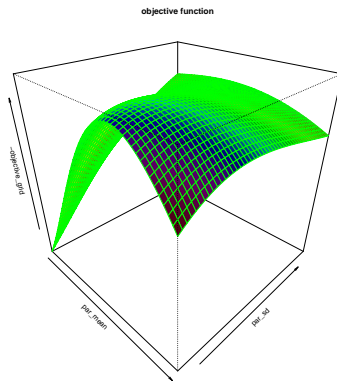
```
> # Initial parameters
> initp <- c(mean=0, sd=1)
> # Perform optimization using optim()
> optim1 <- optim(par=initp,
+   fn=objfun, # Log-likelihood function
+   datav=datav,
+   method="L-BFGS-B", # Quasi-Newton method
+   upper=c(10, 10), # Upper constraint
+   lower=c(-10, 0.1)) # Lower constraint
> # Optimal parameters
> optim1$par
> # Perform optimization using MASS::fitdistr()
> optim1 <- MASS::fitdistr(datav, densfun="normal")
> optim1$estimate
> optim1$sd
> # Plot histogram
> histp <- hist(datav, plot=FALSE)
> plot(histp, freq=FALSE, main="histogram of sample")
> curve(expr=dnorm(x, mean=optim1$par["mean"], sd=optim1$par["sd"]),
+   add=TRUE, type="l", lwd=2, col="red")
> legend("topright", inset=0.0, cex=0.8, title=NULL, y.intersp=0.4,
+   leg="optimal parameters", lwd=2, bg="white", col="red")
```



Mixture Model Likelihood Function

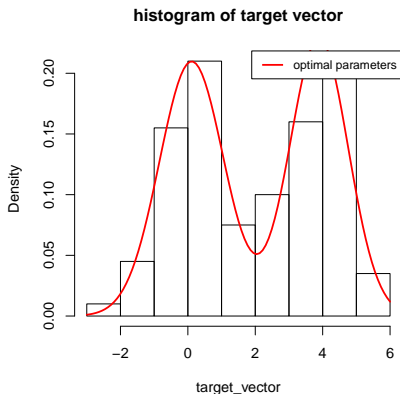
```
> # Sample from mixture of normal distributions
> datav <- c(rnorm(100, sd=1.0),
+           rnorm(100, mean=4, sd=1.0))
> # Objective function is log-likelihood
> objfun <- function(parv, datav) {
+   likew <- parv[1]/parv[3] *
+   dnorm((datav-parv[2])/parv[3]) +
+   (1-parv[1])/parv[5]*dnorm((datav-parv[4])/parv[5])
+   if (any(likew <= 0)) Inf else
+   -sum(log(likew))
+ } # end objfun
> # Vectorize objective function
> objvecive <- Vectorize(
+   FUN=function(mean, sd, w, m1, s1, datav)
+   objfun(c(w, m1, s1, mean, sd), datav),
+   vectorize.args=c("mean", "sd")
+ ) # end Vectorize
> # Objective function on parameter grid
> parmean <- seq(3, 5, length=50)
> parsd <- seq(0.5, 1.5, length=50)
> objective_grid <- outer(parmean, parsd,
+   objvecive, datav=datav,
+   w=0.5, m1=2.0, s1=2.0)
> rownames(objective_grid) <- round(parmean, 2)
> colnames(objective_grid) <- round(parsd, 2)
> objective_min <- which(objective_grid==
+   min(objective_grid), arr.ind=TRUE)
> objective_min
> objective_grid[objective_min]
> objective_grid[(objective_min[, 1] + -1:1),
+   (objective_min[, 2] + -1:1)]
```

```
> # Perspective plot of objective function
> persp(parmean, parsd, -objective_grid,
+   theta=45, phi=30,
+   shade=0.5,
+   col=rainbow(50),
+   border="green",
+   main="objective function")
```



Optimization of Mixture Model

```
> # Initial parameters
> initp <- c(weight=0.5, m1=0, s1=1, m2=2, s2=1)
> # Perform optimization
> optim1 <- optim(par=initp,
+   fn=objfun,
+   datav=datav,
+   method="L-BFGS-B",
+   upper=c(1,10,10,10,10),
+   lower=c(0,-10,0.2,-10,0.2))
> optim1$par
> # Plot histogram
> histp <- hist(datav, plot=FALSE)
> plot(histp, freq=FALSE,
+   main="histogram of sample")
> fitfun <- function(x, parv) {
+   parv["weight"]*dnorm(x, mean=parv["m1"], sd=parv["s1"]) +
+   (1-parv["weight"])*dnorm(x, mean=parv["m2"], sd=parv["s2"])
+ } # end fitfun
> curve(expr=fitfun(x, parv=optim1$par), add=TRUE,
+ type="l", lwd=2, col="red")
> legend("topright", inset=0.0, cex=0.8, title=NULL,
+   leg="optimal parameters", y.intersp=0.4,
+   lwd=2, bg="white", col="red")
```



Package *DEoptim* for Global Optimization

The function `DEoptim()` from package *DEoptim* performs *global* optimization using the *Differential Evolution* algorithm.

Differential Evolution is a genetic algorithm which evolves a population of solutions over several generations:

<https://link.springer.com/content/pdf/10.1023/A:1008202821328.pdf>

The first generation of solutions is selected randomly.

Each new generation is obtained by combining the best solutions from the previous generation.

The *Differential Evolution* algorithm is well suited for very large multi-dimensional optimization problems, such as portfolio optimization.

Gradient optimization methods are more efficient than *Differential Evolution* for smooth objective functions with no local minima.

```
> # Rastrigin function with vector argument for optimization
> rastrigin <- function(vecv, param=25) {
+   sum(vecv^2 - param*cos(vecv))
+ } # end rastrigin
> vecv <- c(pi/6, pi/6)
> rastrigin(vecv=vecv)
> library(DEoptim)
> # Optimize rastrigin using DEoptim
> optim1 <- DEoptim(rastrigin,
+   upper=c(6, 6), lower=c(-6, -6),
+   DEoptim.control(trace=FALSE, itermax=50))
> # Optimal parameters and value
> optim1$optim$bestmem
> rastrigin(optim1$optim$bestmem)
> summary(optim1)
> plot(optim1)
```

Downloading Treasury Bond Rates from FRED

The constant maturity Treasury rates are yields of hypothetical fixed-maturity bonds, interpolated from the market yields of actual Treasury bonds.

The *FRED* database contains current and historical constant maturity Treasury rates,

<https://fred.stlouisfed.org/series/DGS5>

`quantmod::getSymbols()` creates objects in the specified *environment* from the input strings (names).

It then assigns the data to those objects, without returning them as a function value, as a *side effect*.

```
> # Symbols for constant maturity Treasury rates
> symbolv <- c("DGS1", "DGS2", "DGS5", "DGS10", "DGS20", "DGS30")
> # Create new environment for time series
> ratesenv <- new.env()
> # Download time series for symbolv into ratesenv
> quantmod::getSymbols(symbolv, env=ratesenv, src="FRED")
> # List files in ratesenv
> ls(ratesenv)
> # Get class of all objects in ratesenv
> sapply(ratesenv, class)
> # Get class of all objects in R workspace
> sapply(ls(), function(name) class(get(name)))
> # Save the time series environment into a binary .RData file
> save(ratesenv, file="/Users/jerzy/Develop/lecture_slides/data/ra
```



```
> # Get class of time series object DGS10
> class(get(x="DGS10", env=ratesenv))
> # Another way
> class(ratesenv$DGS10)
> # Get first 6 rows of time series
> head(ratesenv$DGS10)
> # Plot dygraphs of 10-year Treasury rate
> dygraphs::dygraph(ratesenv$DGS10, main="10-year Treasury Rate") %>%
+   dyOptions(colors="blue", strokeWidth=2)
> # Plot 10-year constant maturity Treasury rate
> x11(width=6, height=5)
> par(mar=c(2, 2, 0, 0), oma=c(0, 0, 0, 0))
> chart_Series(ratesenv$DGS10["1990/"], name="10-year Treasury Rate")
```

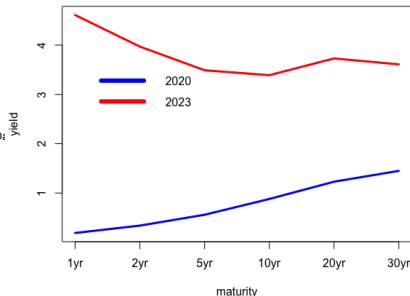
Treasury Yield Curve

The *yield curve* is a vector of interest rates at different maturities, on a given date.

The *yield curve* shape changes depending on the economic conditions: in recessions rates drop and the curve flattens, while in expansions rates rise and the curve steepens.

```
> # Load constant maturity Treasury rates
> load(file="/Users/jerzy/Develop/lecture_slides/data/rates_data.RData")
> # Get most recent yield curve
> ycnow <- eapply(ratesenv, xts::last)
> class(ycnow)
> ycnow <- do.call(cbind, ycnow)
> # Check if 2020-03-25 is not a holiday
> date2020 <- as.Date("2020-03-25")
> weekdays(date2020)
> # Get yield curve from 2020-03-25
> yc2020 <- eapply(ratesenv, function(x) x[date2020])
> yc2020 <- do.call(cbind, yc2020)
> # Combine the yield curves
> ycurves <- c(yc2020, ycnow)
> # Rename columns and rows, sort columns, and transpose into matrix
> colnames(ycurves) <- substr(colnames(ycurves), start=4, stop=11)
> ycurves <- ycurves[, order(as.numeric(colnames(ycurves)))]
> colnames(ycurves) <- paste0(colnames(ycurves), "yr")
> ycurves <- t(ycurves)
> colnames(ycurves) <- substr(colnames(ycurves), start=1, stop=4)
```

Yield Curves in 2020 and 2023



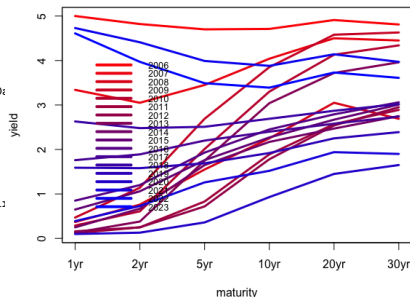
```
> # Plot using matplot()
> colorv <- c("blue", "red")
> matplot(ycurves, main="Yield Curves in 2020 and 2023", xaxt="n",
+ type="l", xlab="maturity", ylab="yield", col=colorv)
> # Add x-axis
> axis(1, seq_along(rownames(ycurves)), rownames(ycurves))
> # Add legend
> legend("topleft", legend=colnames(ycurves), y.intersp=0.1,
+ bty="n", col=colorv, lty=1, lwd=6, inset=0.05, cex=1.0)
```

Treasury Yield Curve Over Time

The *yield curve* has changed shape dramatically depending on the economic conditions: in recessions rates drop and the curve flattens, while in expansions rates rise and the curve steepens.

```
> # Load constant maturity Treasury rates
> load(file="/Users/jerzy/Develop/lecture_slides/data/rates_data.RD")
> # Get end-of-year dates since 2006
> datev <- xts::endpoints(ratesenv$DGS1["2006/"], on="years")
> datev <- zoo::index(ratesenv$DGS1["2006/"][datev])
> # Create time series of end-of-year rates
> ycurves <- eapply(ratesenv, function(ratev) ratev[datev])
> ycurves <- rutils::do_call(cbind, ycurves)
> # Rename columns and rows, sort columns, and transpose into matrix
> colnames(ycurves) <- substr(colnames(ycurves), start=4, stop=11)
> ycurves <- ycurves[, order(as.numeric(colnames(ycurves))))]
> colnames(ycurves) <- paste0(colnames(ycurves), "yr")
> ycurves <- t(ycurves)
> colnames(ycurves) <- substr(colnames(ycurves), start=1, stop=4)
> # Plot matrix using plot.zoo()
> colorv <- colorRampPalette(c("red", "blue"))(NCOL(ycurves))
> plot.zoo(ycurves, main="Yield curve since 2006", lwd=3, xaxt="n"
+   plot.type="single", xlab="maturity", ylab="yield", col=colorv)
> # Add x-axis
> axis(1, seq_along(rownames(ycurves)), rownames(ycurves))
> # Add legend
> legend("topleft", legend=colnames(ycurves), y.intersp=0.1,
+   bty="n", col=colorv, lty=1, lwd=4, inset=0.05, cex=0.8)
```

Yield curve since 2006



```
> # Alternative plot using matplot()
> matplot(ycurves, main="Yield curve since 2006", xaxt="n", lwd=3,
+   type="l", xlab="maturity", ylab="yield", col=colorv)
> # Add x-axis
> axis(1, seq_along(rownames(ycurves)), rownames(ycurves))
> # Add legend
> legend("topleft", legend=colnames(ycurves), y.intersp=0.1,
+   bty="n", col=colorv, lty=1, lwd=4, inset=0.05, cex=0.8)
```

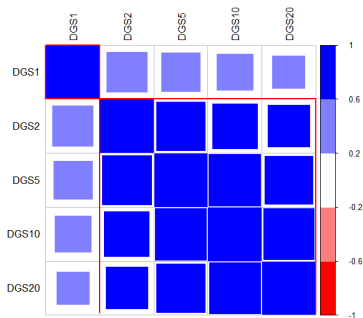
Covariance Matrix of Interest Rates

The covariance matrix \mathbb{C} , of the interest rate matrix \mathbf{r} is given by:

$$\mathbb{C} = \frac{(\mathbf{r} - \bar{\mathbf{r}})^T (\mathbf{r} - \bar{\mathbf{r}})}{n - 1}$$

```
> # Extract rates from ratesenv
> symbolv <- c("DGS1", "DGS2", "DGS5", "DGS10", "DGS20")
> ratem <- mget(symbolv, envir=ratesenv)
> ratem <- rutils::do_call(cbind, ratem)
> ratem <- zoo::na.locf(ratem, na.rm=FALSE)
> ratem <- zoo::na.locf(ratem, fromLast=TRUE)
> # Calculate daily percentage rates changes
> retp <- rutils::diffit(log(ratem))
> # Center (de-mean) the returns
> retp <- lapply(retp, function(x) {x - mean(x)})
> retp <- rutils::do_call(cbind, retp)
> sapply(retp, mean)
> # Covariance and Correlation matrices of Treasury rates
> covmat <- cov(retp)
> cormat <- cor(retp)
> # Reorder correlation matrix based on clusters
> library(corrplot)
> ordern <- corrMatOrder(cormat, order="hclust",
+   hclust.method="complete")
> cormat <- cormat[ordern, ordern]
```

Correlation of Treasury Rates



```
> # Plot the correlation matrix
> x11(width=6, height=6)
> colorv <- colorRampPalette(c("red", "white", "blue"))
> corrplot(cormat, title=NA, tl.col="black",
+   method="square", col=colorv(NCOL(cormat)), tl.cex=0.8,
+   cl.offset=0.75, cl.cex=0.7, cl.align.text="l", cl.ratio=0.25)
> title("Correlation of Treasury Rates", line=1)
> # Draw rectangles on the correlation matrix plot
> corrRect.hclust(cormat, k=NROW(cormat) %/% 2,
+   method="complete", col="red")
```

Principal Component Vectors

Principal components are linear combinations of the k return vectors \mathbf{r}_i :

$$\mathbf{pc}_j = \sum_{i=1}^k w_{ij} \mathbf{r}_i$$

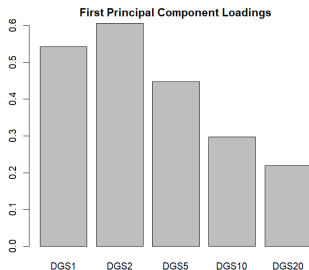
Where \mathbf{w}_j is a vector of weights (loadings) of the *principal component* j , with $\mathbf{w}_j^T \mathbf{w}_j = 1$.

The weights \mathbf{w}_j are chosen to maximize the variance of the *principal components*, under the condition that they are orthogonal to:

$$\mathbf{w}_j = \arg \max \left\{ \mathbf{pc}_j^T \mathbf{pc}_j \right\}$$

$$\mathbf{pc}_i^T \mathbf{pc}_j = 0 \quad (i \neq j)$$

```
> # Create initial vector of portfolio weights
> nweights <- NROW(symbolv)
> weightv <- rep(1/sqrt(nweights), nweights)
> names(weightv) <- symbolv
> # Objective function equal to minus portfolio variance
> objfun <- function(weightv, retp) {
+   retp <- retp %*% weightv
+   -1e7*var(retp) + 1e7*(1 - sum(weightv*weightv))^2
+ } # end objfun
> # Objective function for equal weight portfolio
> objfun(weightv, retp)
> # Compare speed of vector multiplication methods
> library(microbenchmark)
> summary(microbenchmark(
+   transp=t(retp) %*% retp,
+   sumv=sum(retp*retp),
```



```
> # Find weights with maximum variance
> optim1 <- optim(par=weightv,
+   fn=objfun,
+   retp=retp,
+   method="L-BFGS-B",
+   upper=rep(5.0, nweights),
+   lower=rep(-5.0, nweights))
> # Optimal weights and maximum variance
> weights1 <- optim1$par
> objfun(weights1, retp)
> # Plot first principal component loadings
> x11(width=6, height=5)
> par(mar=c(3, 3, 2, 1), oma=c(0, 0, 0, 0), mgp=c(2, 1, 0))
> barplot(weights1, names.arg=names(weights1),
+   xlab="", ylab="", main="First Principal Component Loadings")
```

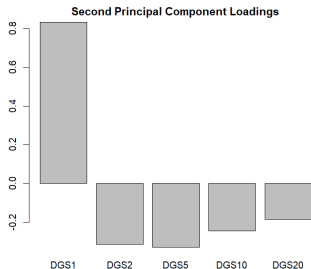
Higher Order Principal Components

The *second principal component* can be calculated by maximizing its variance, under the constraint that it must be orthogonal to the *first principal component*.

Similarly, higher order *principal components* can be calculated by maximizing their variances, under the constraint that they must be orthogonal to all the previous *principal components*.

The number of principal components is equal to the dimension of the covariance matrix.

```
> # pc1 weights and returns
> pc1 <- drop(retp %*% weights1)
> # Redefine objective function
> objfun <- function(weightv, retp) {
+   retp <- retp %*% weightv
+   -1e7*var(retp) + 1e7*(1 - sum(weightv^2))^2 +
+   1e7*sum(weights1*weightv)^2
+ } # end objfun
> # Find second principal component weights
> optim1 <- optim(par=weightv,
+   fn=objfun,
+   retp=retp,
+   method="L-BFGS-B",
+   upper=rep(5.0, nweights),
+   lower=rep(-5.0, nweights))
```



```
> # pc2 weights and returns
> weights2 <- optim1$par
> pc2 <- drop(retp %*% weights2)
> sum(pc1*pc2)
> # Plot second principal component loadings
> barplot(weights2, names.arg=names(weights2),
+   xlab="", ylab="", main="Second Principal Component Loadings")
```


Eigenvalues of the Covariance Matrix

The portfolio variance: $\mathbf{w}^T \mathbb{C} \mathbf{w}$ can be maximized under the *quadratic* weights constraint $\mathbf{w}^T \mathbf{w} = 1$, by maximizing the *Lagrangian* \mathcal{L} :

$$\mathcal{L} = \mathbf{w}^T \mathbb{C} \mathbf{w} - \lambda (\mathbf{w}^T \mathbf{w} - 1)$$

Where λ is a *Lagrange multiplier*.

The maximum variance portfolio weights can be found by differentiating \mathcal{L} with respect to \mathbf{w} and setting it to zero:

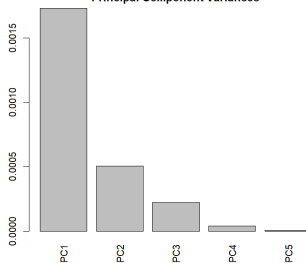
$$\mathbb{C} \mathbf{w} = \lambda \mathbf{w}$$

The above is the *eigenvalue* equation of the covariance matrix \mathbb{C} , with the optimal weights \mathbf{w} forming an *eigenvector*, and λ is the *eigenvalue* corresponding to the *eigenvector* \mathbf{w} .

The *eigenvalues* are the variances of the *eigenvectors*, and their sum is equal to the sum of the return variances:

$$\sum_{i=1}^k \lambda_i = \frac{1}{1-k} \sum_{i=1}^k \mathbf{r}_i^T \mathbf{r}_i$$

Principal Component Variances



```
> eigend <- eigen(covmat)
> eigend$eigenvectors
> # Compare with optimization
> all.equal(sum(diag(covmat)), sum(eigend$values))
> all.equal(abs(eigend$eigenvectors[, 1]), abs(weights1), check.attributes=FALSE)
> all.equal(abs(eigend$eigenvectors[, 2]), abs(weights2), check.attributes=FALSE)
> all.equal(eigend$values[1], var(pc1), check.attributes=FALSE)
> all.equal(eigend$values[2], var(pc2), check.attributes=FALSE)
> # Eigenvalue equations are satisfied approximately
> (covmat %*% weights1) / weights1 / var(pc1)
> (covmat %*% weights2) / weights2 / var(pc2)
> # Plot eigenvalues
> barplot(eigend$values, names.arg=paste0("PC", 1:nweights),
+ las=3, xlab="", ylab="", main="Principal Component Variances")
```

Principal Component Analysis Versus Eigen Decomposition

Principal Component Analysis (PCA) is equivalent to the *eigen decomposition* of either the correlation or the covariance matrix.

If the input time series *are* scaled, then *PCA* is equivalent to the eigen decomposition of the *correlation matrix*.

If the input time series *are not* scaled, then *PCA* is equivalent to the eigen decomposition of the *covariance matrix*.

Scaling the input time series improves the accuracy of the *PCA dimension reduction*, allowing a smaller number of *principal components* to more accurately capture the data contained in the input time series.

The function `prcomp()` performs *Principal Component Analysis* on a matrix of data (with the time series as columns), and returns the results as a list of class `prcomp`.

The `prcomp()` argument `scale=TRUE` specifies that the input time series should be scaled by their standard deviations.

```
> # Eigen decomposition of correlation matrix
> eigend <- eigen(cormat)
> # Perform PCA with scaling
> pcad <- prcomp(retp, scale=TRUE)
> # Compare outputs
> all.equal(eigend$values, pcad$sdev^2)
> all.equal(abs(eigend$vectors), abs(pcad$rotation),
+   check.attributes=FALSE)
> # Eigen decomposition of covariance matrix
> eigend <- eigen(covmat)
> # Perform PCA without scaling
> pcad <- prcomp(retp, scale=FALSE)
> # Compare outputs
> all.equal(eigend$values, pcad$sdev^2)
> all.equal(abs(eigend$vectors), abs(pcad$rotation),
+   check.attributes=FALSE)
```

Principal Component Analysis of the Yield Curve

Principal Component Analysis (PCA) is a *dimension reduction* technique, that explains the returns of a large number of correlated time series as linear combinations of a smaller number of principal component time series.

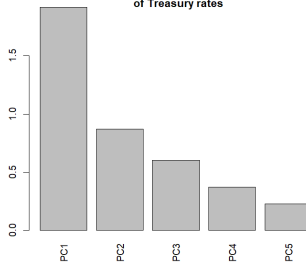
The input time series are often scaled by their standard deviations, to improve the accuracy of *PCA dimension reduction*, so that more information is retained by the first few *principal component* time series.

If the input time series are not scaled, then *PCA* analysis is equivalent to the *eigen decomposition* of the covariance matrix, and if they are scaled, then *PCA* analysis is equivalent to the *eigen decomposition* of the correlation matrix.

The function `prcomp()` performs *Principal Component Analysis* on a matrix of data (with the time series as columns), and returns the results as a list of class `prcomp`.

The `prcomp()` argument `scale=TRUE` specifies that the input time series should be scaled by their standard deviations.

Scree Plot: Volatilities of Principal Components of Treasury rates



A *scree plot* is a bar plot of the volatilities of the *principal components*.

```
> # Perform principal component analysis PCA
> pcad <- prcomp(retp, scale=TRUE)
> # Plot standard deviations
> barplot(pcad$sdev, names.arg=colnames(pcad$rotation),
+   las=3, xlab="", ylab="",
+   main="Scree Plot: Volatilities of Principal Components
+   of Treasury rates")
```

Yield Curve Principal Component Loadings (Weights)

Principal component loadings are the weights of portfolios which have mutually orthogonal returns.

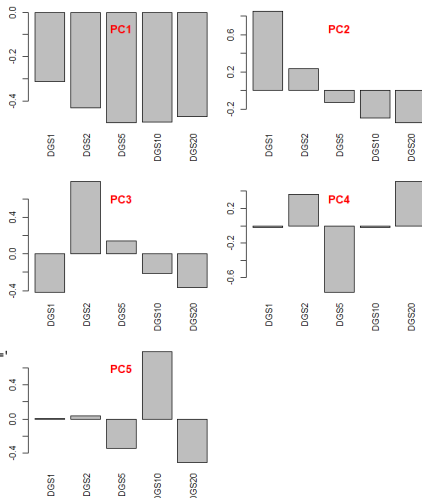
The *principal component* portfolios represent the different orthogonal modes of the data variance.

The first *principal component* of the *yield curve* is the correlated movement of all rates up and down.

The second *principal component* is *yield curve* steepening and flattening.

The third *principal component* is the *yield curve* butterfly movement.

```
> # Calculate principal component loadings (weights)
> pcad$rotation
> # Plot loading barplots in multiple panels
> par(mfrow=c(3,2))
> par(mar=c(3.5, 2, 2, 1), oma=c(0, 0, 0, 0))
> for (ordern in 1:NCOL(pcad$rotation)) {
+   barplot(pcad$rotation[, ordern], las=3, xlab="", ylab="", main=
+ title(paste0("PC", ordern), line=-2.0, col.main="red")
+ } # end for
```



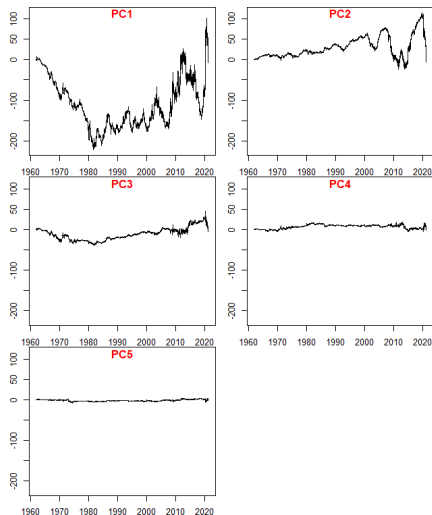
Yield Curve Principal Component Time Series

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

The *principal component* time series have mutually orthogonal returns.

Higher order *principal components* are gradually less volatile.

```
> # Standardize (center and scale) the returns
> retp <- lapply(retp, function(x) {(x - mean(x))/sd(x)})
> retp <- rutils::do_call(cbind, retp)
> sapply(retp, mean)
> sapply(retp, sd)
> # Calculate principal component time series
> pcacum <- retp %*% pcad$rotation
> all.equal(pcad$x, pcacum, check.attributes=FALSE)
> # Calculate products of principal component time series
> round(t(pcacum) %*% pcacum, 2)
> # Coerce to xts time series
> pcacum <- xts(pcacum, order.by=zoo::index(retp))
> pcacum <- cumsum(pcacum)
> # Plot principal component time series in multiple panels
> par(mfrow=c(3,2))
> par(mar=c(2, 2, 0, 1), oma=c(0, 0, 0, 0))
> rangev <- range(pcacum)
> for (ordern in 1:NCOL(pcacum)) {
+   plot.zoo(pcacum[, ordern], ylim=rangev, xlab="", ylab="")
+   title(paste0("PC", ordern), line=-1, col.main="red")
+ } # end for
```



Inverting Principal Component Analysis

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

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

```
> # Invert all the principal component time series
> retpca <- retp %*% pcad$rotation
> solved <- retpca %*% solve(pcad$rotation)
> all.equal(coredata(retp), solved)
```

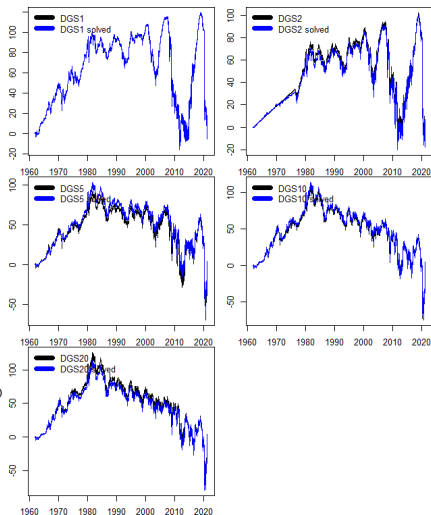
Dimension Reduction Using Principal Component Analysis

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

A popular rule of thumb is to use the *principal components* with the largest variances, which sum up to 80% of the total variance of returns.

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

```
> # Invert first 3 principal component time series
> solved <- retpcal[, 1:3] %*% solve(pcad$rotation)[1:3, ]
> solved <- xts::xts(solved, zoo::index(retp))
> solved <- cumsum(solved)
> retc <- cumsum(retp)
> # Plot the solved returns
> par(mfrow=c(3,2))
> par(mar=c(2, 2, 0, 1), oma=c(0, 0, 0, 0))
> 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", y.intersp=0.1,
+   legend=paste0(symbol, c("", " solved")),
+   title=NULL, inset=0.0, cex=1.0, lwd=6,
+   lty=1, col=c("black", "blue"))
+ } # end for
```



Calibrating Yield Curve Using Package *RQuantLib*

The package *RQuantLib* is an interface to the *QuantLib* open source C/C++ library for quantitative finance, mostly designed for pricing fixed-income instruments and options.

The function `DiscountCurve()` calibrates a *zero coupon yield curve* from *money market rates*, *Eurodollar futures*, and *swap rates*.

The function `DiscountCurve()` interpolates the *zero coupon rates* into a vector of dates specified by the `times` argument.

```
> library(RQuantLib) # Load RQuantLib
> # Specify curve parameters
> curve_params <- list(tradeDate=as.Date("2018-01-17"),
+                      settleDate=as.Date("2018-01-19"),
+                      dt=0.25,
+                      interpWhat="discount",
+                      interpHow="loglinear")
> # Specify market data: prices of FI instruments
> market_data <- list(d3m=0.0363,
+                     fut1=96.2875,
+                     fut2=96.7875,
+                     fut3=96.9875,
+                     fut4=96.6875,
+                     s5y=0.0443,
+                     s10y=0.05165,
+                     s15y=0.055175)
> # Specify dates for calculating the zero rates
> disc_dates <- seq(0, 10, 0.25)
> # Specify the evaluation (as of) date
> setEvaluationDate(as.Date("2018-01-17"))
> # Calculate the zero rates
> disc_curves <- DiscountCurve(params=curve_params,
+                              tsQuotes=market_data,
+                              times=disc_dates)
> # Plot the zero rates
> x11()
> plot(x=disc_curves$zerorates, t="l", main="zerorates")
```


Vector and Matrix Calculus

Let \mathbf{v} and \mathbf{w} be vectors, with $\mathbf{v} = \{v_i\}_{i=1}^{i=n}$, and let $\mathbf{1}$ be the unit vector, with $\mathbf{1} = \{1\}_{i=1}^{i=n}$.

Then the inner product of \mathbf{v} and \mathbf{w} can be written as $\mathbf{v}^T \mathbf{w} = \mathbf{w}^T \mathbf{v} = \sum_{i=1}^n v_i w_i$.

We can then express the sum of the elements of \mathbf{v} as the inner product: $\mathbf{v}^T \mathbf{1} = \mathbf{1}^T \mathbf{v} = \sum_{i=1}^n v_i$.

And the sum of squares of \mathbf{v} as the inner product: $\mathbf{v}^T \mathbf{v} = \sum_{i=1}^n v_i^2$.

Let \mathbb{A} be a matrix, with $\mathbb{A} = \{A_{ij}\}_{i,j=1}^{i,j=n}$.

Then the inner product of matrix \mathbb{A} with vectors \mathbf{v} and \mathbf{w} can be written as:

$$\mathbf{v}^T \mathbb{A} \mathbf{w} = \mathbf{w}^T \mathbb{A}^T \mathbf{v} = \sum_{i,j=1}^n A_{ij} v_i w_j$$

The derivative of a scalar variable with respect to a vector variable is a vector, for example:

$$\frac{d(\mathbf{v}^T \mathbf{1})}{d\mathbf{v}} = d_v[\mathbf{v}^T \mathbf{1}] = d_v[\mathbf{1}^T \mathbf{v}] = \mathbf{1}^T$$

$$d_v[\mathbf{v}^T \mathbf{w}] = d_v[\mathbf{w}^T \mathbf{v}] = \mathbf{w}^T$$

$$d_v[\mathbf{v}^T \mathbb{A} \mathbf{w}] = \mathbf{w}^T \mathbb{A}^T$$

$$d_v[\mathbf{v}^T \mathbb{A} \mathbf{v}] = \mathbf{v}^T \mathbb{A} + \mathbf{v}^T \mathbb{A}^T$$

Formula Objects

Formulas in R are defined using the "~" operator followed by a series of terms separated by the "+" operator.

Formulas can be defined as separate objects, manipulated, and passed to functions.

The formula " $z \sim x$ " means the *response vector* z is explained by the *predictor* x (also called the *explanatory variable* or *independent variable*).

The formula " $z \sim x + y$ " represents a linear model: $z = ax + by + c$.

The formula " $z \sim x - 1$ " or " $z \sim x + 0$ " represents a linear model with zero intercept: $z = ax$.

The function `update()` modifies existing formulas.

The "." symbol represents either all the remaining data, or the variable that was in this part of the formula.

```
> # Formula of linear model with zero intercept
> formulav <- z ~ x + y - 1
> formulav
>
> # Collapse vector of strings into single text string
> paste0("x", 1:5)
> paste(paste0("x", 1:5), collapse="+")
>
> # Create formula from text string
> formulav <- as.formula(
+   # Coerce text strings to formula
+   paste("z ~ ",
+   paste(paste0("x", 1:5), collapse="+")
+   ) # end paste
+ ) # end as.formula
> class(formulav)
> formulav
> # Modify the formula using "update"
> update(formulav, log(.) ~ . + beta)
```

Simple Linear Regression

A Simple Linear Regression is a linear model between a *response vector* y and a single *predictor* x , defined by the formula:

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

α and β are the unknown *regression coefficients*.

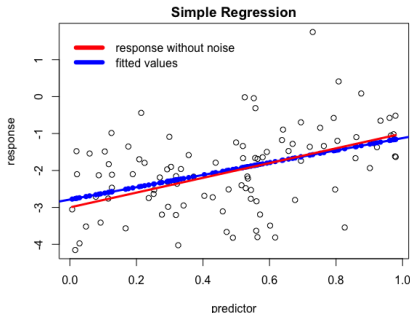
ε_i are the *residuals*, which are usually assumed to be standard normally distributed $\phi(0, \sigma_\varepsilon)$, independent, and stationary.

In the Ordinary Least Squares method (OLS), the regression parameters are estimated by minimizing the *Residual Sum of Squares (RSS)*:

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \\ &= (y - \alpha \mathbf{1} - \beta x)^T (y - \alpha \mathbf{1} - \beta x) \end{aligned}$$

Where $\mathbf{1}$ is the unit vector, with $\mathbf{1}^T \mathbf{1} = n$ and $\mathbf{1}^T x = x^T \mathbf{1} = \sum_{i=1}^n x_i$

The data consists of n pairs of observations (x_i, y_i) of the response and predictor variables, with the index i ranging from 1 to n .



```
> # Define explanatory (predm) variable
> nrows <- 100
> set.seed(1121) # Initialize random number generator
> predm <- runif(nrows)
> noisev <- rnorm(nrows)
> # Response equals linear form plus random noise
> respv <- (-3 + 2*predm + noisev)
```

The *response vector* and the *predictor matrix* don't have to be normally distributed.

Solution of Linear Regression

The *OLS* solution for the *regression coefficients* is found by equating the *RSS* derivatives to zero:

$$RSS_{\alpha} = -2(y - \alpha \mathbf{1} - \beta x)^T \mathbf{1} = 0$$

$$RSS_{\beta} = -2(y - \alpha \mathbf{1} - \beta x)^T x = 0$$

The solution for α is given by:

$$\alpha = \bar{y} - \beta \bar{x}$$

The solution for β can be obtained by manipulating the equation for RSS_{β} as follows:

$$(y - (\bar{y} - \beta \bar{x})\mathbf{1} - \beta x)^T (x - \bar{x}\mathbf{1}) =$$

$$((y - \bar{y}\mathbf{1}) - \beta(x - \bar{x}\mathbf{1}))^T (x - \bar{x}\mathbf{1}) =$$

$$(\hat{y} - \beta \hat{x})^T \hat{x} = \hat{y}^T \hat{x} - \beta \hat{x}^T \hat{x} = 0$$

Where $\hat{x} = x - \bar{x}\mathbf{1}$ and $\hat{y} = y - \bar{y}\mathbf{1}$ are the centered (de-meanned) variables. Then β is given by:

$$\beta = \frac{\hat{y}^T \hat{x}}{\hat{x}^T \hat{x}} = \frac{\sigma_y}{\sigma_x} \rho_{xy}$$

β is proportional to the correlation coefficient ρ_{xy} between the response and predictor variables.

If the response and predictor variables have zero mean, then $\alpha = 0$ and $\beta = \frac{y^T x}{x^T x}$.

The *residuals* $\varepsilon = y - \alpha \mathbf{1} - \beta x$ have zero mean: $RSS_{\alpha} = -2\varepsilon^T \mathbf{1} = 0$.

The *residuals* ε are orthogonal to the *predictor* x : $RSS_{\beta} = -2\varepsilon^T x = 0$.

The expected value of the *RSS* is equal to the *degrees of freedom* $(n - 2)$ times the variance σ_{ε}^2 of the *residuals* ε_i : $\mathbb{E}[RSS] = (n - 2)\sigma_{\varepsilon}^2$.

```
> # Calculate centered (de-meanned) predictor and response vectors
> predc <- predm - mean(predm)
> respc <- respv - mean(respv)
> # Calculate the regression beta
> betav <- cov(predm, respv)/var(predm)
> # Calculate the regression alpha
> alpha <- mean(respv) - betav*mean(predm)
```

Linear Regression Using Function lm()

Let the data generating process for the response variable be given as: $z = \alpha_{lat} + \beta_{lat}x + \varepsilon_{lat}$

Where α_{lat} and β_{lat} are latent (unknown) coefficients, and ε_{lat} is an unknown vector of random noise (error terms).

The error terms are the difference between the measured values of the response minus the (unknown) actual response values.

The function `lm()` fits a linear model into a set of data, and returns an object of class "lm", which is a list containing the results of fitting the model:

- call - the model formula,
- coefficients - the fitted model coefficients (α , β_j),
- residuals - the model residuals (respv minus fitted values),

The regression *residuals* are not the same as the error terms, because the regression coefficients are not equal to the coefficients of the data generating process.

```
> # Specify regression formula
> formulav <- respv ~ predm
> regmod <- lm(formulav) # Perform regression
> class(regmod) # Regressions have class lm
[1] "lm"
> attributes(regmod)
$names
  [1] "coefficients" "residuals" "effects" "rank"
  [5] "fitted.values" "assign" "qr" "df.residual"
  [9] "xlevels" "call" "terms" "model"

$class
[1] "lm"
> eval(regmod$call$formula) # Regression formula
respv ~ predm
> regmod$coeff # Regression coefficients
(Intercept)      predm
    -2.79      1.67
> all.equal(coef(regmod), c(alpha, betav),
+           check.attributes=FALSE)
[1] TRUE
```

The Fitted Values of Linear Regression

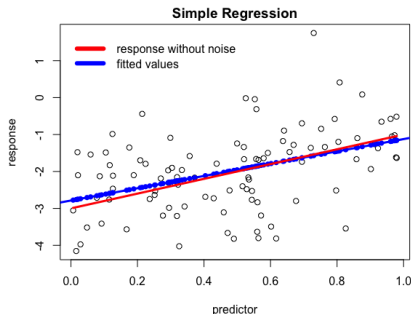
The *fitted values* y_{fit} are the estimates of the *response vector* obtained from the regression model:

$$y_{fit} = \alpha + \beta x$$

The *generic function* `plot()` produces a scatterplot when it's called on the regression formula.

`abline()` plots a straight line corresponding to the regression coefficients, when it's called on the regression object.

```
> fitv <- (alpha + betav*predm)
> all.equal(fitv, regmod$fitted.values, check.attributes=FALSE)
> # Plot scatterplot using formula
> plot(formulav, xlab="predictor", ylab="response")
> title(main="Simple Regression", line=0.5)
> # Add regression line
> abline(regmod, lwd=3, col="blue")
> # Plot fitted (forecast) response values
> points(x=predm, y=regmod$fitted.values, pch=16, col="blue")
```



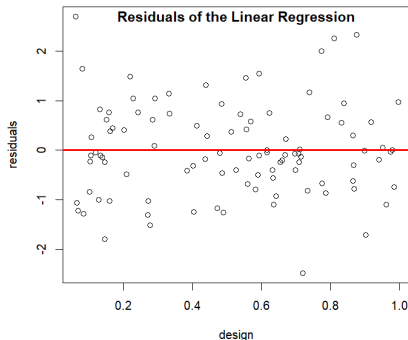
```
> # Plot response without noise
> lines(x=predm, y=(respv-noisev), col="red", lwd=3)
> legend(x="topleft", # Add legend
+       legend=c("response without noise", "fitted values"),
+       title=NULL, inset=0.0, cex=1.0, y.intersp=0.3,
+       bty="n", lwd=6, lty=1, col=c("red", "blue"))
```

Linear Regression Residuals

The residuals ε_i of a linear regression are defined as the response vector minus the fitted values:

$$\varepsilon_i = y_i - y_{fit}$$

```
> # Calculate the residuals
> fitv <- (alpha + betav*predm)
> resids <- (respv - fitv)
> all.equal(resids, regmod$residuals, check.attributes=FALSE)
[1] TRUE
> # Residuals are orthogonal to the predictor
> all.equal(sum(resids*predm), target=0)
[1] TRUE
> # Residuals are orthogonal to the fitted values
> all.equal(sum(resids*fitv), target=0)
[1] TRUE
> # Sum of residuals is equal to zero
> all.equal(mean(resids), target=0)
[1] TRUE
```



```
> x11(width=6, height=5) # Open x11 for plotting
> # Set plot parameters to reduce whitespace around plot
> par(mar=c(5, 5, 1, 1), oma=c(0, 0, 0, 0))
> # Extract residuals
> datav <- cbind(predm, regmod$residuals)
> colnames(datav) <- c("predictor", "residuals")
> # Plot residuals
> plot(datav)
> title(main="Residuals of the Linear Regression", line=-1)
> abline(h=0, lwd=3, col="red")
```

Standard Errors of Regression Coefficients

The *residuals* are the source of error in the regression model, producing uncertainty in the *response vector* y and in the regression coefficients: $y_i = \alpha + \beta x_i + \varepsilon_i$.

The standard errors of the regression coefficients are equal to their standard deviations, given the *residuals* as the source of error.

Since $\beta = \frac{\hat{y}^T \hat{x}}{\hat{x}^T \hat{x}}$, then its variance is equal to:

$$\sigma_\beta^2 = \frac{1}{(n-2)} \frac{E[(\varepsilon^T \hat{x})^2]}{(\hat{x}^T \hat{x})^2} = \frac{1}{(n-2)} \frac{E[\varepsilon^2]}{\hat{x}^T \hat{x}} = \frac{\sigma_\varepsilon^2}{\hat{x}^T \hat{x}}$$

Since $\alpha = \bar{y} - \beta \bar{x}$, then its variance is equal to:

$$\sigma_\alpha^2 = \frac{\sigma_\varepsilon^2}{n} + \sigma_\beta^2 \bar{x}^2 = \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\hat{x}^T \hat{x}} \right)$$

```
> # Degrees of freedom of residuals
> degf <- regmod$df.residual
> # Standard deviation of residuals
> residsd <- sqrt(sum(resids^2)/degf)
> # Standard error of beta
> betasd <- residsd/sqrt(sum(predc^2))
> # Standard error of alpha
> alphasd <- residsd*sqrt(1/nrows + mean(predm)^2/sum(predc^2))
```


Linear Regression Summary

The function `summary.lm()` produces a list of regression model diagnostic statistics:

- `coefficients`: matrix with estimated coefficients, their t -statistics, and p -values,
- `r.squared`: fraction of response variance explained by the model,
- `adj.r.squared`: `r.squared` adjusted for higher model complexity,
- `fstatistic`: ratio of variance explained by the model divided by unexplained variance,

The regression `summary` is a list, and its elements can be accessed individually.

```
> regsum <- summary(regmod) # Copy regression summary
> regsum # Print the summary to console
```

```
Call:
lm(formula = formulav)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-2.133 -0.649  0.106  0.590  3.321
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   -2.787       0.196  -14.20 < 2e-16 ***
predm           1.665       0.357   4.67 9.8e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.988 on 98 degrees of freedom
Multiple R-squared:  0.182, Adjusted R-squared:  0.173
F-statistic: 21.8 on 1 and 98 DF,  p-value: 9.75e-06
```

```
> attributes(regsum)$names # get summary elements
[1] "call"          "terms"         "residuals"     "coefficients"
[5] "aliased"       "sigma"         "df"            "r.squared"
[9] "adj.r.squared" "fstatistic"    "cov.unscaled"
```

Regression Model Diagnostic Statistics

The *null hypothesis* for regression is that the coefficients are zero.

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 key assumption in the formula for the standard error is that the *residuals* are normally distributed, independent, and stationary.

If they are not, then the standard error and the *p*-value may be much bigger than reported by `summary.lm()`, and therefore the regression may not be statistically significant.

Asset returns are very far from normal, so the small *p*-values shouldn't be automatically interpreted as meaning that the regression is statistically significant.

```
> regsum$coeff
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   -2.79      0.196   -14.20 1.61e-25
predm          1.67      0.357    4.67 9.75e-06
> # Standard errors
> regsum$coefficients[2, "Std. Error"]
[1] 0.357
> all.equal(c(alphasd, betasd), regsum$coefficients[, "Std. Error"],
+   check.attributes=FALSE)
[1] TRUE
> # R-squared
> regsum$r.squared
[1] 0.182
> regsum$adj.r.squared
[1] 0.173
> # F-statistic and ANOVA
> regsum$fstatistic
value numdf den df
21.8    1.0   98.0
> anova(regmod)
Analysis of Variance Table

Response: resp
      Df Sum Sq Mean Sq F value    Pr(>F)
predm   1   21.3   21.25    21.8 9.8e-06 ***
Residuals 98   95.7    0.98
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Weak Regression

If the relationship between the response and predictor variables is weak compared to the error terms (noise), then the regression will have low statistical significance.

```
> # High noise compared to coefficient
> respv <- (-3 + 2*predm + rnorm(nrows, sd=8))
> regmod <- lm(formulav) # Perform regression
> # Values of regression coefficients are not
> # Statistically significant
> summary(regmod)
```

```
Call:
lm(formula = formulav)
```

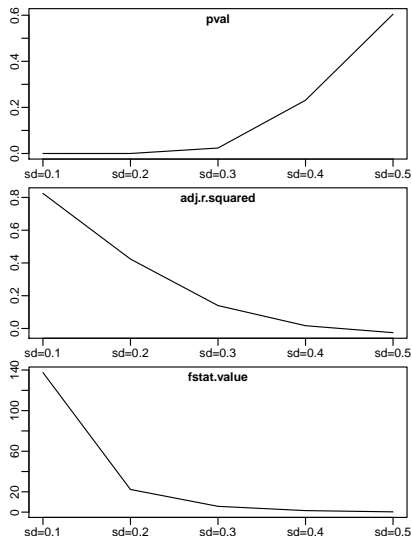
```
Residuals:
    Min       1Q   Median       3Q      Max
-16.430  -4.325   0.735   4.365  16.720
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    -1.65      1.44    -1.14   0.26
predm          -1.70      2.62    -0.65   0.52
```

```
Residual standard error: 7.25 on 98 degrees of freedom
Multiple R-squared:  0.0043, Adjusted R-squared:  -0.00586
F-statistic: 0.423 on 1 and 98 DF, p-value: 0.517
```

Influence of Noise on Regression

```
> regstats <- function(stdev) { # Noisy regression
+   set.seed(1121) # initialize number generator
+   # Define explanatory (predm) and response variables
+   predm <- rnorm(100, mean=2)
+   respv <- (1 + 0.2*predm + rnorm(nrows, sd=stdev))
+   # Specify regression formula
+   formulav <- respv ~ predm
+   # Perform regression and get summary
+   regsum <- summary(lm(formulav))
+   # Extract regression statistics
+   with(regsum, c(pval=coefficients[2, 4],
+     adj_rsquared=adj.r.squared,
+     fstat=fstatistic[1]))
+ } # end regstats
> # Apply regstats() to vector of stdev dev values
> vecsd <- seq(from=0.1, to=0.5, by=0.1)
> names(vecsd) <- paste0("sd=", vecsd)
> statsmat <- t(sapply(vecsd, regstats))
> # Plot in loop
> par(mfrow=c(NCOL(statsmat), 1))
> for (it in 1:NROW(statsmat)) {
+   plot(statsmat[, it], type="l",
+     xaxt="n", xlab="", ylab="", main="")
+   title(main=colnames(statsmat)[it], line=-1.0)
+   axis(1, at=1:(NROW(statsmat)), labels=rownames(statsmat))
+ } # end for
```

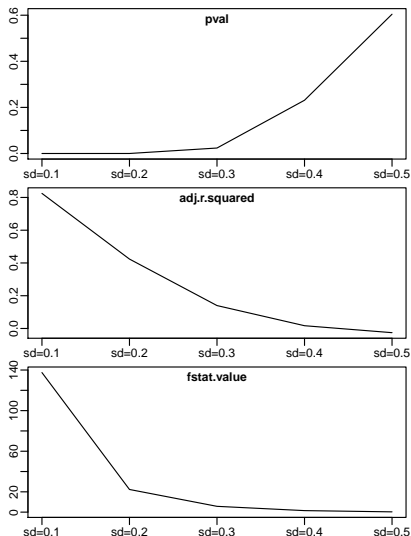


Influence of Noise on Regression Another Method

```

> regstats <- function(datav) { # get regression
+ # Perform regression and get summary
+   colnamev <- colnames(datav)
+   formulav <- paste(colnamev[2], colnamev[1], sep="~")
+   regsum <- summary(lm(formulav, data=datav))
+ # Extract regression statistics
+   with(regsum, c(pval=coefficients[2, 4],
+     adj.rsquared=adj.r.squared,
+     fstat=fstatistic[1]))
+ } # end regstats
> # Apply regstats() to vector of stdev dev values
> vecsd <- seq(from=0.1, to=0.5, by=0.1)
> names(vecsd) <- paste0("sd=", vecsd)
> statsmat <- t(sapply(vecsd, function(stdev) {
+   set.seed(1121) # initialize number generator
+ # Define explanatory (predm) and response variables
+   predm <- rnorm(100, mean=2)
+   respv <- (1 + 0.2*predm + rnorm(nrows, sd=stdev))
+   regstats(data.frame(predm, respv))
+ }))
> # Plot in loop
> par(mfrow=c(NCOL(statsmat), 1))
> for (it in 1:NCOL(statsmat)) {
+   plot(statsmat[, it], type="l",
+     xaxt="n", xlab="", ylab="", main="")
+   title(main=colnames(statsmat)[it], line=-1.0)
+   axis(1, at=1:(NROW(statsmat)),
+     labels=rownames(statsmat))
+ } # end for

```



Linear Regression Diagnostic Plots

`plot()` produces diagnostic scatterplots for the *residuals*, when called on the regression object.

The diagnostic scatterplots allow for visual inspection to determine the quality of the regression fit.

"Residuals vs Fitted" is a scatterplot of the residuals vs. the forecast responses.

"Scale-Location" is a scatterplot of the square root of the standardized residuals vs. the forecast responses.

The residuals should be randomly distributed around the horizontal line representing zero residual error.

A pattern in the residuals indicates that the model was not able to capture the relationship between the variables, or that the variables don't follow the statistical assumptions of the regression model.

"Normal Q-Q" is the standard Q-Q plot, and the points should fall on the diagonal line, indicating that the residuals are normally distributed.

"Residuals vs Leverage" is a scatterplot of the residuals vs. their leverage.

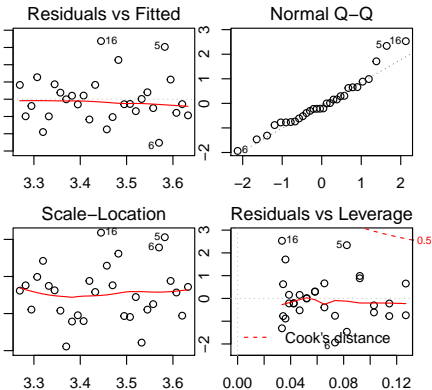
Leverage measures the amount by which the fitted values would change if the response values were shifted by a small amount.

Cook's distance measures the influence of a single observation on the fitted values, and is proportional to the sum of the squared differences between forecasts made with all observations and forecasts made without the observation.

Points with large leverage, or a Cook's distance greater than 1 suggest the presence of an outlier or a poor model,

```
> par(mfrow=c(2, 2)) # Plot 2x2 panels
> plot(regmod) # Plot diagnostic scatterplots
> plot(regmod, which=2) # Plot just Q-Q
```

lm(reg_formula)



Durbin-Watson Test of Autocorrelation of Residuals

The *Durbin-Watson* test is designed to test the *null hypothesis* that the autocorrelations of regression *residuals* are equal to zero.

The test statistic is equal to:

$$DW = \frac{\sum_{i=2}^n (\varepsilon_i - \varepsilon_{i-1})^2}{\sum_{i=1}^n \varepsilon_i^2}$$

Where ε_i are the regression *residuals*.

The value of the *Durbin-Watson* statistic *DW* is close to zero for large positive autocorrelations, and close to four for large negative autocorrelations.

The *DW* is close to two for autocorrelations close to zero.

The *p*-value for the *reg_model* regression is large, and we conclude that the *null hypothesis* is TRUE, and the regression *residuals* are uncorrelated.

```
> library(lmtest) # Load lmtest
> # Perform Durbin-Watson test
> lmtest::dwtest(regmod)
```

Durbin-Watson test

```
data: regmod
DW = 2, p-value = 0.7
alternative hypothesis: true autocorrelation is greater than 0
```

The Leverage for Univariate Regression

We can add an extra unit column to the *predictor matrix* \mathbb{X} so that the univariate regression can be written in *homogeneous form* as:

$$y = \mathbb{X}\beta + \varepsilon$$

With two *regression coefficients*: $\beta = (\alpha, \beta_1)$, and a *predictor matrix* \mathbb{X} with two columns, with the first column equal to a unit vector.

After the second column of the *predictor matrix* \mathbb{X} is centered (de-meanned), its *covariance matrix* is given by:

$$\mathbb{X}^T \mathbb{X} = \begin{pmatrix} n & 0 \\ 0 & \sum_{i=1}^n (x_i - \bar{x})^2 \end{pmatrix}$$

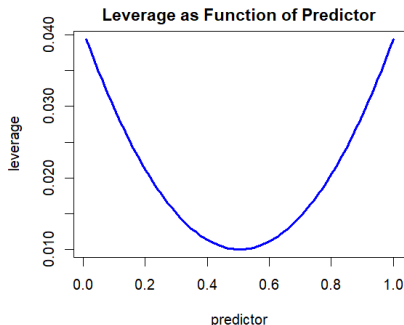
And the *influence matrix* \mathbb{H} is given by:

$$\mathbb{H}_{ij} = [\mathbb{X}(\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T]_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

The first term above is due to the influence of the regression intercept α , and the second term is due to the influence of the regression slope β_1 .

The diagonal elements of the *influence matrix* \mathbb{H}_{ii} form the *leverage vector*.

```
> # Define linear regression data
> set.seed(1121) # Initialize random number generator
> nrows <- 100
> predm <- runif(nrows)
> noisew <- rnorm(nrows)
> respy <- (-3 + 2*predm + noisew)
```



```
> # Add unit column to the predictor matrix
> predm <- cbind(rep(1, nrows), predm)
> # Calculate generalized inverse of the predictor matrix
> predinv <- MASS::ginv(predm)
> # Calculate the influence matrix
> infmat <- predm %*% predinv
> # Plot the leverage vector
> ordern <- order(predm[, 2])
> plot(x=predm[ordern, 2], y=diag(infmat)[ordern],
+      type="l", lwd=3, col="blue",
+      xlab="predictor", ylab="leverage",
+      main="Leverage as Function of Predictor")
```


Covariance Matrix of Fitted Values in Univariate Regression

The *fitted values* y_{fit} can be considered to be *random variables* \hat{y}_{fit} :

$$\hat{y}_{fit} = \mathbb{H}\hat{y} = \mathbb{H}(y_{fit} + \hat{\epsilon}) = y_{fit} + \mathbb{H}\hat{\epsilon}$$

The *covariance matrix* of the *fitted values* \hat{y}_{fit} is:

$$\sigma_{fit}^2 = \frac{\mathbb{E}[\mathbb{H}\hat{\epsilon}(\mathbb{H}\hat{\epsilon})^T]}{d_{free}} = \frac{\mathbb{E}[\mathbb{H}\hat{\epsilon}\hat{\epsilon}^T\mathbb{H}^T]}{d_{free}} = \frac{\mathbb{H}\mathbb{E}[\hat{\epsilon}\hat{\epsilon}^T]\mathbb{H}^T}{d_{free}} = \sigma_{\epsilon}^2\mathbb{H} = \sigma_{\epsilon}^2\mathbb{X}(\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T$$

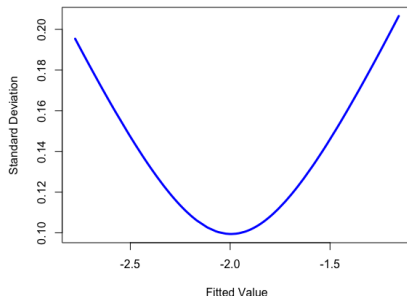
The square of the *influence matrix* \mathbb{H} is equal to itself (it's idempotent): $\mathbb{H}\mathbb{H}^T = \mathbb{H}$.

The variance of the *fitted values* σ_{fit}^2 increases with the distance of the *predictors* from their mean values.

This is because the *fitted values* farther from their mean are more sensitive to the variance of the regression slope.

```
> # Calculate the influence matrix
> infmat <- predm %*% predinv
> # The influence matrix is idempotent
> all.equal(infmat, infmat %*% infmat)
```

Standard Deviations of Fitted Values
in Univariate Regression



```
> # Calculate covariance and standard deviations of fitted values
> betav <- predinv %*% respv
> fitv <- drop(predm %*% betav)
> residv <- drop(respv - fitv)
> degf <- (NROW(predm) - NCOL(predm))
> residstd <- sqrt(sum(residv^2)/degf)
> fitcovar <- residstd*infmat
> fitsd <- sqrt(diag(fitcovar))
> # Plot the standard deviations
> fitdata <- cbind(fitted=fitv, stdev=fitsd)
> fitdata <- fitdata[order(fitv), ]
> plot(fitdata, type="l", lwd=3, col="blue",
+      xlab="Fitted Value", ylab="Standard Deviation",
+      main="Standard Deviations of Fitted Values\nin Univariate Regression")
```

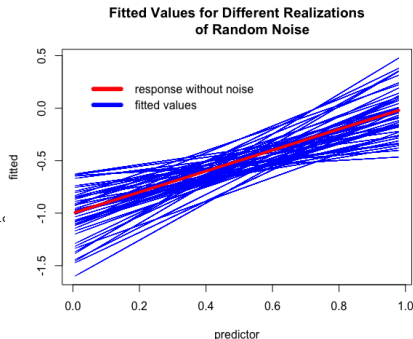
Fitted Values for Different Realizations of Random Noise

The fitted values are more volatile for *predictor* values that are further away from their mean, because those points have higher *leverage*.

The higher *leverage* of points further away from the mean of the *predictor* is due to their greater sensitivity to changes in the slope of the regression.

The fitted values for different realizations of random noise can be calculated using the influence matrix.

```
> # Calculate response without random noise for univariate regression
> # equal to weighted sum over columns of predictor.
> respn <- predm %*% c(-1, 1)
> # Perform loop over different realizations of random noise
> fitm <- lapply(1:50, function(it) {
+   # Add random noise to response
+   respv <- respn + rnorm(nrows, sd=1.0)
+   # Calculate fitted values using influence matrix
+   infmat %*% respv
+ }) # end lapply
> fitm <- rutils::do_call(cbind, fitm)
```



```
> # Plot fitted values
> matplot(x=predm[, 2], y=fitm,
+ type="l", lty="solid", lwd=1, col="blue",
+ xlab="predictor", ylab="fitted",
+ main="Fitted Values for Different Realizations
+ of Random Noise")
> lines(x=predm[, 2], y=respn, col="red", lwd=4)
> legend(x="topleft", # Add legend
+ legend=c("response without noise", "fitted values"),
+ title=NULL, inset=0.05, cex=1.0, lwd=6, y.intersp=0.4,
+ bty="n", lty=1, col=c("red", "blue"))
```

Forecasts From *Univariate Regression Models*

The forecast $y_{forecast}$ from a regression model is equal to the *response value* corresponding to the *predictor* vector with the new data \mathbb{X}_{new} :

$$y_{forecast} = \mathbb{X}_{new} \beta$$

The variance $\sigma_{forecast}^2$ of the *forecast value* is equal to the *predictor* vector multiplied by the *covariance matrix* of the *regression coefficients* σ_{β}^2 :

$$\sigma_{forecast}^2 = \frac{\mathbb{E}[\mathbb{X}_{new} \mathbb{X}_{inv} \hat{\hat{\epsilon}} (\mathbb{X}_{new} \mathbb{X}_{inv} \hat{\hat{\epsilon}})^T]}{d_{free}} =$$

$$\frac{\mathbb{E}[\mathbb{X}_{new} \mathbb{X}_{inv} \hat{\hat{\epsilon}} \hat{\hat{\epsilon}}^T \mathbb{X}_{inv}^T \mathbb{X}_{new}^T]}{d_{free}} = \sigma_{\epsilon}^2 \mathbb{X}_{new} \mathbb{X}_{inv} \mathbb{X}_{inv}^T \mathbb{X}_{new}^T =$$

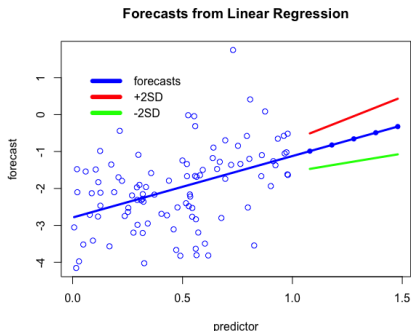
$$\sigma_{\epsilon}^2 \mathbb{X}_{new} (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}_{new}^T = \mathbb{X}_{new} \sigma_{\beta}^2 \mathbb{X}_{new}^T$$

```
> # Define new predictor
> newdata <- (max(predm[, 2]) + 10*(1:5)/nrows)
> predn <- cbind(rep(1, NROW(newdata)), newdata)
> # Calculate the forecast values and standard errors
> predm2 <- MASS::ginv(crossprod(predm)) # Inverse of predictor matrix
> preds2 <- residm*sqrt(predn %*% predm2 %*% t(predn))
> fcast <- cbind(forecast=drop(predn %*% betav),
+   stdev=diag(preds2))
```

Confidence Intervals of Regression Forecasts

The variables σ_ε^2 and σ_y^2 follow the *chi-squared* distribution with $d_{\text{free}} = (n - k - 1)$ degrees of freedom, so the *forecast value* y_{forecast} follows the *t-distribution*.

```
> # Prepare plot data
> xdata <- c(predm[, 2], newdata)
> ydata <- c(fitv, fcast[, 1])
> # Calculate t-quantile
> tquant <- qt(pnorm(2), df=degf)
> fcastl <- fcast[, 1] - tquant*fcast[, 2]
> fcasth <- fcast[, 1] + tquant*fcast[, 2]
> # Plot the regression forecasts
> xlim <- range(xdata)
> ylim <- range(c(respv, ydata, fcastl, fcasth))
> plot(x=xdata, y=ydata, xlim=xlim, ylim=ylim,
+      type="l", lwd=3, col="blue",
+      xlab="predictor", ylab="forecast",
+      main="Forecasts from Linear Regression")
> points(x=predm[, 2], y=respv, col="blue")
> points(x=newdata, y=fcast[, 1], pch=16, col="blue")
> lines(x=newdata, y=fcasth, lwd=3, col="red")
> lines(x=newdata, y=fcastl, lwd=3, col="green")
> legend(x="topleft", # Add legend
+       legend=c("forecasts", "+2SD", "-2SD"),
+       title=NULL, inset=0.05, cex=1.0, lwd=6, y.intersp=0.4,
+       bty="n", lty=1, col=c("blue", "red", "green"))
```



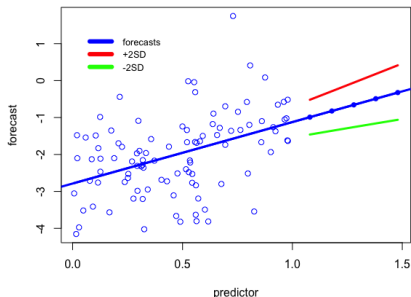
Forecasts of Linear Regression Using predict.lm()

The function `predict()` is a *generic function* for forecasting based on a given model.

`predict.lm()` is the forecasting method for linear models (regressions) produced by the function `lm()`.

```
> # Perform univariate regression
> dframe <- data.frame(resp=respv, pred=predm[, 2])
> regmod <- lm(resp ~ pred, data=dframe)
> # Calculate forecasts from regression
> newdf <- data.frame(pred=predn[, 2]) # Same column name
> fcastlm <- predict.lm(object=regmod,
+   newdata=newdf, confl=1-2*(1-pnorm(2)),
+   interval="confidence")
> rownames(fcastlm) <- NULL
> all.equal(fcastlm[, "fit"], fcast[, 1])
> all.equal(fcastlm[, "lwr"], fcast[, 1])
> all.equal(fcastlm[, "upr"], fcast[, 1])
> plot(x=xdata, y=ydata, xlim=xlim, ylim=ylim,
+   type="l", lwd=3, col="blue",
+   xlab="predictor", ylab="forecast",
+   main="Forecasts from lm() Regression")
> points(x=predm[, 2], y=respv, col="blue")
```

Forecasts from lm() Regression



```
> abline(regmod, col="blue", lwd=3)
> points(x=newdata, y=fcastlm[, "fit"], pch=16, col="blue")
> lines(x=newdata, y=fcastlm[, "lwr"], lwd=3, col="green")
> lines(x=newdata, y=fcastlm[, "upr"], lwd=3, col="red")
> legend(x="topleft", # Add legend
+   legend=c("forecasts", "+2SD", "-2SD"),
+   title=NULL, inset=0.05, cex=0.8, lwd=6, y.intersp=0.4,
+   bty="n", lty=1, col=c("blue", "red", "green"))
```

Spurious Time Series Regression

Regression of non-stationary time series creates *spurious* regressions.

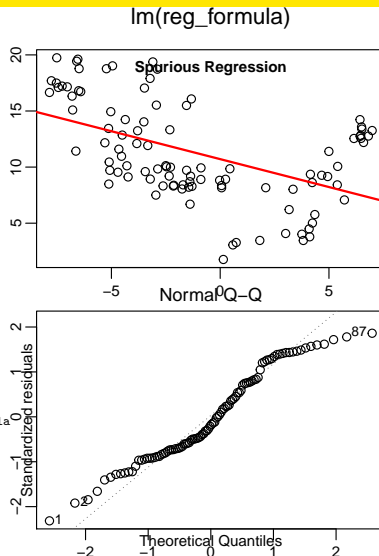
The t -statistics, p -values, and R -squared all indicate a statistically significant regression.

But the Durbin-Watson test shows residuals are autocorrelated, which invalidates the other tests.

The Q-Q plot also shows that residuals are *not* normally distributed.

```
> predm <- cumsum(rnorm(100)) # Unit root time series
> respv <- cumsum(rnorm(100))
> formulav <- respv ~ predm
> regmod <- lm(formulav) # Perform regression
> # Summary indicates statistically significant regression
> regsum <- summary(regmod)
> regsum$coeff
> regsum$r.squared
> # Durbin-Watson test shows residuals are autocorrelated
> dwtest <- lmtest::dwtest(regmod)
> c(dwtest$statistic[[1]], dwtest$p.value)

> plot(formulav, xlab="", ylab="") # Plot scatterplot using formula
> title(main="Spurious Regression", line=-1)
> # Add regression line
> abline(regmod, lwd=2, col="red")
> plot(regmod, which=2, ask=FALSE) # Plot just Q-Q
```



Homework Assignment

Required

- Study all the lecture slides in *FRE6871_Lecture_4.pdf*, and run all the code in *FRE6871_Lecture_4.R*