

# Autoregressive Strategies

## FRE7241, Fall 2024

Jerzy Pawlowski [jp3900@nyu.edu](mailto:jp3900@nyu.edu)

*NYU Tandon School of Engineering*

November 3, 2024



# The Bid-Ask Spread

The *bid-ask spread* is the difference between the best ask (offer) price minus the best bid price in the market.

The *bid-ask spread* can be estimated from the differences between the execution prices of consecutive buy and sell market orders (roundtrip trades).

Market orders are orders to buy or sell a stock immediately at the best available price in the market.

Market orders guarantee that the trade will be executed, but they do not guarantee the execution price. Market orders are subject to the *bid-ask spread*.

Limit orders are orders to buy or sell a stock at the limit price or better (the investor sets the limit price). Limit orders do not guarantee that the trade will be executed, but they guarantee the execution price. Limit orders are placed only for a certain time when they are "live".

Market orders are executed by matching them with live limit orders through a matching engine at an exchange.

The *bid-ask spread* for many liquid ETFs is about 1 basis point. For example the [XLK ETF](#)

The most liquid [SPY ETF](#) usually trades at a *bid-ask spread* of only one tick (cent=\$0.01, or about 0.2 basis points).

In reality the *bid-ask spread* is not static and depends on many factors, such as market liquidity (trading volume), volatility, and the time of day.

```
> # Load the roundtrip trades
> dtable <- data.table::fread("/Users/jerzy/Develop/lecture_slides/c
> nrows <- NROW(dtable)
> class(dtable$timefill)
> # Sort the trades according to the execution time
> dtable <- dtable[order(dtable$timefill)]
> # Calculate the dollar bid-ask spread
> pricebuy <- dtable$price[dtable$side == "buy"]
> pricesell <- dtable$price[dtable$side == "sell"]
> bidask <- mean(pricebuy-pricesell)
> # Calculate the percentage bid-ask spread
> bidask/mean(pricesell)
```

# Autocorrelations of Stock Returns

Daily stock returns often exhibit negative autocorrelations at one-day lags.

The *VTI* returns appear to have some small, yet significant negative autocorrelations at lag=1, and some positive autocorrelations at larger lags.

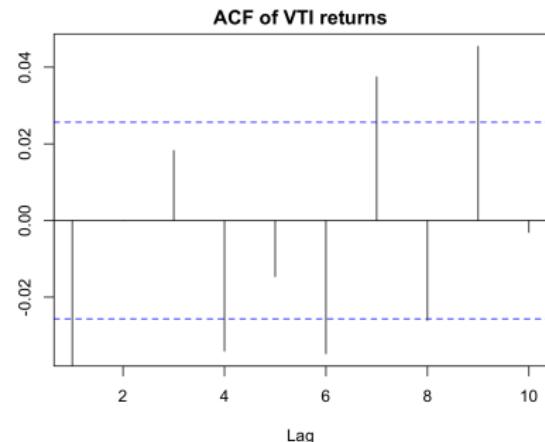
The *autocorrelation* of lag  $k$  of a time series of returns  $r_t$  is equal to:

$$\rho_k = \frac{\sum_{t=k+1}^n (r_t - \bar{r})(r_{t-k} - \bar{r})}{(n - k) \sigma^2}$$

The function `rutils::plot_acf()` calculates and plots the autocorrelations of a time series.

But the visual inspection of the ACF plot alone is not enough to test whether autocorrelations are statistically significant or not.

```
> # Calculate the daily VTI percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> # Calculate the autocorrelations of daily VTI returns
> rutils::plot_acf(retp, lag=10, main="ACF of VTI returns")
```



## Daily Mean Reverting Strategy

The daily mean reverting strategy buys or sells short \$1 of stock at the end of each day (depending on the sign of the previous daily return), and holds the position until the next day.

If the previous daily return was positive, it sells short \$1 of stock. If the previous daily return was negative, it buys \$1 of stock.

The mean reverting strategy has lower returns than *VTI*, but it has a very low correlation to *VTI*, so it has a positive *alpha*.

Thanks to its low correlation to *VTI*, the mean reverting strategy provides diversification of risk, and combined with *VTI*, a higher *Sharpe* ratio than *VTI* alone (it has a positive marginal *alpha*).

Strategies which have low or negative correlations to stocks, can contribute a significant marginal *alpha*, even if they have low returns.

```
> # Simulate the mean reverting strategy
> posv <- -rutils::lagit(sign(retp), lagg=1)
> pnls <- retp*posv
> # Subtract transaction costs from the pnls
> bidask <- 0.0001 # The bid-ask spread is equal to 1 basis point
> costv <- 0.5*bidask*abs(rutils::diffit(posv))
> pnls <- (pnls - costv)
> # Calculate the strategy beta and alpha
> betac <- cov(pnls, retp)/var(retp)
> alphac <- mean(pnls) - betac*mean(retp)
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls, (retp+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="VTI Daily Mean Reverting Strategy") %>%
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dyLegend(show="always", width=300)
```

# Daily Mean Reverting Strategy With a Holding Period

The daily mean reverting strategy can be improved by combining the daily returns from the previous two days. This is equivalent to holding the position for two days, instead of rolling it daily.

The daily mean reverting strategy with a holding period performs better than the simple daily strategy because of risk diversification.

```
> # Simulate mean reverting strategy with two day holding period
> posv <- rutils::roll_sum(sign(retp), lookb=2)/2
> pnls <- retp*rutils::lagit(posv)
```



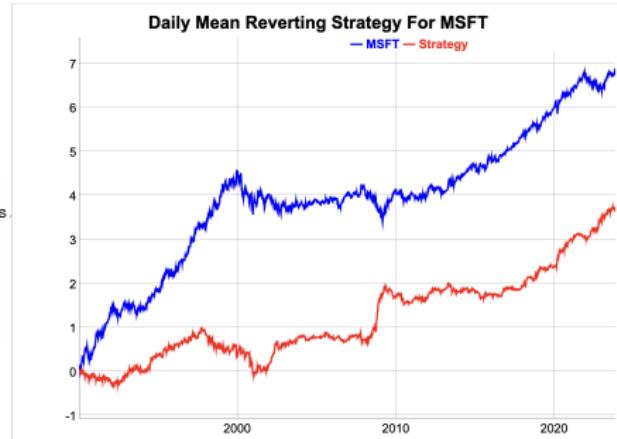
```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Daily Mean Reverting Strategy With Two Day Holding Period",
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Daily Mean Reverting Strategy For Stocks

Some daily stock returns exhibit stronger negative autocorrelations than ETFs.

But the daily mean reverting strategy doesn't perform well for many stocks.

```
> # Load daily S&P500 stock returns
> load(file="/Users/jerzy/Develop/lecture_slides/data/sp500_returns"
> retp <- na.omit(retstock$MSFT)
> rutils::plot_acf(retp)
> # Simulate mean reverting strategy with two day holding period
> posv <- -rutils::roll_sum(sign(retp), lookb=2)/2
> pnls <- retp*rutils::lagit(posv)
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("MSFT", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Daily Mean Reverting Strategy For MSFT") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Daily Mean Reverting Strategy For All Stocks

The combined daily mean reverting strategy for all *S&P500* stocks performed well prior to and during the 2008 financial crisis, but was flat afterwards.

Averaging the stock returns using the function `rowMeans()` with `na.rm=TRUE` is equivalent to rebalancing the portfolio so that stocks with NA returns have zero weight.

```
> # Simulate mean reverting strategy for all S&P500 stocks
> library(parallel) # Load package parallel
> ncores <- detectCores() - 1
> pnll <- mclapply(retstock, function(retp) {
+   retp <- na.omit(retp)
+   posv <- -rutils::roll_sum(sign(retp), lookb=2)/2
+   retp*rutils::lagit(posv)
+ }, mc.cores=ncores) # end mclapply
> pnls <- do.call(cbind, pnll)
> pnls <- rowMeans(pnls, na.rm=TRUE)
> # Calculate the average returns of all S&P500 stocks
> datev <- zoo::index(retstock)
> datev <- datev[-1]
> indeks <- rowMeans(retstock, na.rm=TRUE)
> indeks <- indeks[-1]
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(indeks, pnls)
> wealthv <- xts::xts(wealthv, datev)
> colnames(wealthv) <- c("All Stocks", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Daily Mean Reverting Strategy For All Stocks") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Mean Reverting Strategy For Low and High Volatility Stocks

The daily mean reverting strategy performs better for low volatility stocks than for high volatility stocks.

```
> # Calculate the stock volatilities
> volv <- mclapply(retstock, function(retp) {
+   sd(na.omit(retp))
+ }, mc.cores=ncores) # end mclapply
> volv <- do.call(c, volv)
> # Calculate the median volatility
> medianv <- median(volv)
> # Calculate the pnls for low volatility stocks
> pnlovol <- do.call(cbind, pnll[volv == medianv])
> pnlovol <- rowMeans(pnlovol, na.rm=TRUE)
> # Calculate the pnls for high volatility stocks
> pnlhivol <- do.call(cbind, pnll[volv >= medianv])
> pnlhivol <- rowMeans(pnlhivol, na.rm=TRUE)
```

Mean Reverting Strategy For Low and High Volatility Stocks



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(pnlovol, pnlhivol)
> wealthv <- xts::xts(wealthv, datev)
> colnames(wealthv) <- c("Low Vol", "High Vol")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Mean Reverting Strategy For Low and High Volatility Stocks",
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# The EMA Mean-Reversion Strategy

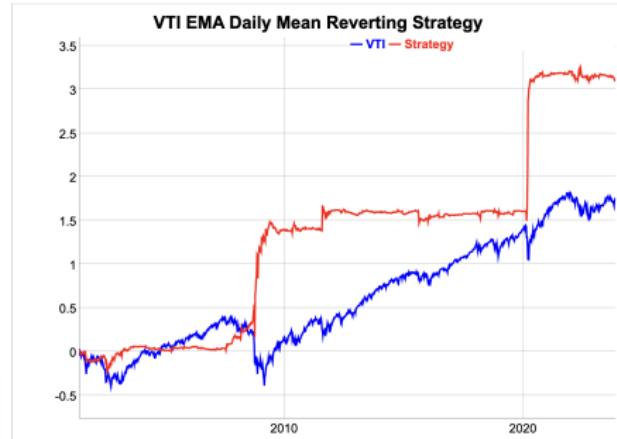
The *EMA* mean-reversion strategy holds either long stock positions or short positions proportional to minus the trailing *EMA* of past returns.

The strategy adjusts its stock position at the end of each day, just before the close of the market.

The strategy takes very large positions in periods of high volatility, when returns are large and highly anti-correlated.

The strategy makes profits mostly in periods of high volatility, but otherwise it's not very profitable.

```
> # Calculate the VTI daily percentage returns
> retp <- na.omit(rutls::etfenv$returns$VTI)
> # Calculate the EMA returns recursively using C++ code
> retma <- HighFreq::run_mean(retp, lambda=0.1)
> # Calculate the positions and PnLs
> posv <- -rutls::lagit(retma, lagg=1)
> pnls <- retp*posv
> # Subtract transaction costs from the pnls
> bidask <- 0.0001 # The bid-ask spread is equal to 1 basis point
> costv <- 0.5*bidask*abs(rutls::diffit(posv))
> pnls <- (pnls - costv)
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutls::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="VTI EMA Daily Mean Reverting Strategy") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# The EMA Mean-Reversion Strategy Scaled By Volatility

Dividing the returns by their trailing volatility reduces the effect of time-dependent volatility.

Scaling the returns by their trailing volatilities reduces the profits in periods of high volatility, but doesn't improve profits in periods of low volatility.

The function `HighFreq::run_var()` calculates the trailing mean and variance of the returns  $r_t$ , by recursively weighting the past variance estimates  $\sigma_{t-1}^2$ , with the squared differences of the returns minus their trailing means  $(r_t - \bar{r}_t)^2$ , using the decay factor  $\lambda$ :

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

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

Where  $\bar{r}_t$  and  $\sigma_t^2$  are the trailing mean and variance at time  $t$ .

The decay factor  $\lambda$  determines how quickly the mean and variance estimates are updated, with smaller values of  $\lambda$  producing faster updating, giving more weight to recent prices, and vice versa.

```
> # Calculate the EMA returns and volatilities
> volv <- HighFreq::run_var(retp, lambda=0.5)
> retma <- volv[, 1]
> volv <- sqrt(volv[, 2])
> # Scale the returns by their trailing volatility
> retsc <- ifelse(volv > 0, retp/volv, 0)
> # Calculate the positions and PnLs
> posv <- -rutils::lagit(retma, lagg=1)
> pnls <- retp*posv
```

VTI EMA Daily Mean Reverting Strategy



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of mean reverting strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="VTI EMA Daily Mean Reverting Strategy") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Autoregressive Model of Stock Returns

The stock returns  $r_t$  can be fitted into an *autoregressive* model  $AR(n)$  with a constant intercept term  $\varphi_0$ :

$$r_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n} + \varepsilon_t$$

The *residuals*  $\varepsilon_t$  are assumed to be normally distributed, independent, and stationary.

The autoregressive model can be written in matrix form as:

$$\mathbf{r} = \boldsymbol{\varphi} \mathbb{P} + \boldsymbol{\varepsilon}$$

Where  $\boldsymbol{\varphi} = \{\varphi_0, \varphi_1, \varphi_2, \dots, \varphi_n\}$  is the vector of autoregressive coefficients.

The *autoregressive* model is equivalent to *multivariate* linear regression, with the *response* equal to the returns  $\mathbf{r}$ , and the columns of the *predictor matrix*  $\mathbb{P}$  equal to the lags of the returns.

```
> # Calculate the VTI daily percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> nrows <- NROW(retp)
> # Define the response and predictor matrices
> respv <- retp
> orderp <- 5
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> # Add constant column for intercept coefficient phi0
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
```

# Forecasting Stock Returns Using Autoregressive Models

The fitted autoregressive coefficients  $\varphi$  are equal to the response  $r$  multiplied by the inverse of the predictor matrix  $P$ :

$$\varphi = P^{-1} r$$

The *in-sample* autoregressive forecasts of the returns are calculated by multiplying the predictor matrix by the fitted AR coefficients:

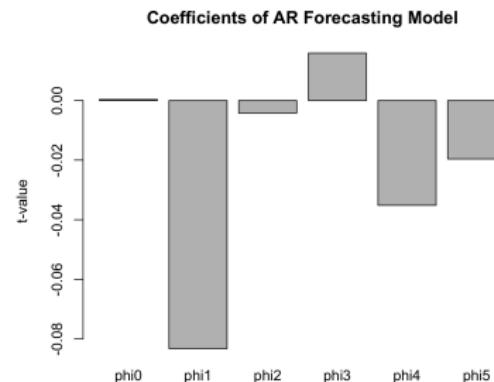
$$f_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n}$$

For *VTI* returns, the intercept coefficient  $\varphi_0$  has a small positive value, while the first autoregressive coefficient  $\varphi_1$  has a small negative value.

This means that the autoregressive forecasting model is a combination of a static long stock position, plus a mean-reverting model which switches its stock position to the reverse of the previous day's return.

The function `MASS::ginv()` calculates the generalized inverse of a matrix.

```
> # Calculate the fitted autoregressive coefficients
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
```



```
> # Plot the AR coefficients
> coeffn <- paste0("phi", 0:(NROW(coeff)-1))
> barplot(coeff ~ coeffn, xlab="", ylab="t-value", col="grey",
+ main="Coefficients of AR Forecasting Model")
```

# The t-values of the Autoregressive Coefficients

The forecast residuals are equal to the differences between the return forecasts minus the actual returns:

$$\varepsilon = f_t - r_t$$

The variance of the autoregressive coefficients  $\sigma_\varphi^2$  is equal to the variance of the forecast residuals  $\sigma_\varepsilon^2$  divided by the squared predictor matrix  $\mathbb{P}$ :

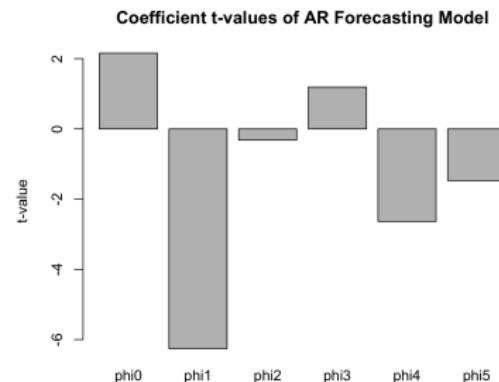
$$\sigma_\varphi^2 = \sigma_\varepsilon^2 (\mathbb{P}^T \mathbb{P})^{-1}$$

The t-values of the autoregressive coefficients are equal to the coefficient values divided by their volatilities:

$$\varphi_{tval} = \frac{\varphi}{\sigma_\varphi}$$

The intercept coefficient  $\varphi_0$  and the first autoregressive coefficient  $\varphi_1$  have statistically significant t-values.

```
> # Calculate the residuals (forecast errors)
> resids <- (fcasts - respv)
> # The residuals are orthogonal to the predictors and the forecasts
> round(cor(resids, fccasts), 6)
> round(sapply(predm[, -1], function(x) cor(resids, x)), 6)
> # Calculate the variance of the residuals
> varv <- sum(resids^2)/(nrows-NROW(coeff))
```



```
> # Calculate the predictor matrix squared
> pred2 <- crossprod(predm)
> # Calculate the covariance matrix of the AR coefficients
> covmat <- varv*MASS::ginv(pred2)
> coefsds <- sqrt(diag(covmat))
> # Calculate the t-values of the AR coefficients
> coefft <- drop(coeff/coefsds)
> coeffn <- paste0("phi", 0:(NROW(coefft)-1))
> # Plot the t-values of the AR coefficients
> barplot(coefft ~ coeffn, xlab="", ylab="t-value", col="grey",
+ main="Coefficient t-values of AR Forecasting Model")
```

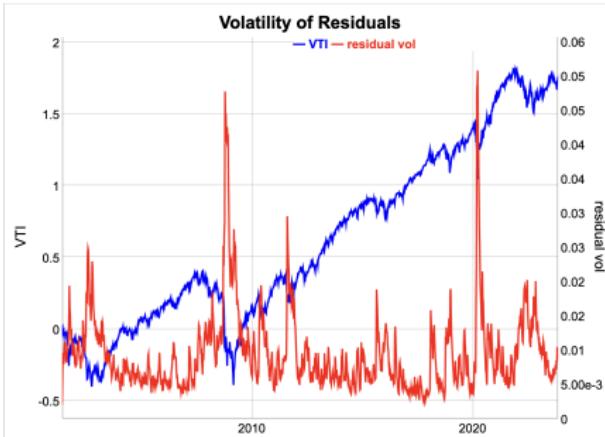
# Residuals of Autoregressive Forecasting Model

The autoregressive model assumes stationary returns and residuals, with similar volatility over time.

In reality stock volatility is highly time dependent, so the volatility of the residuals is also time dependent.

The function `HighFreq::run_var()` calculates the trailing variance of a time series using exponential weights.

```
> # Calculate the trailing volatility of the residuals
> residv <- sqrt(HighFreq::run_var(resids, lambda=0.9)[, 2])
```



```
> # Plot dygraph of volatility of residuals
> datav <- cbind(cumsum(retp), residv)
> colnames(datav) <- c("VTI", "residual vol")
> endd <- rutils::calc_endpoints(datav, interval="weeks")
> dygraphs::dygraph(datav[endd], main="Volatility of Residuals") %>%
+   dyAxis("y", label="VTI", independentTicks=TRUE) %>%
+   dyAxis("y2", label="residual vol", independentTicks=TRUE) %>%
+   dySeries(name="VTI", axis="y", strokeWidth=2, col="blue") %>%
+   dySeries(name="residual vol", axis="y2", strokeWidth=2, col="red") %>%
+   dyLegend(show="always", width=300)
```

## Autoregressive Strategy In-Sample

The first step in strategy development is optimizing it in-sample, even though in practice it can't be implemented. Because a strategy can't perform well out-of-sample if it doesn't perform well in-sample.

The autoregressive strategy invests dollar amounts of *VTI* stock proportional to the in-sample forecasts.

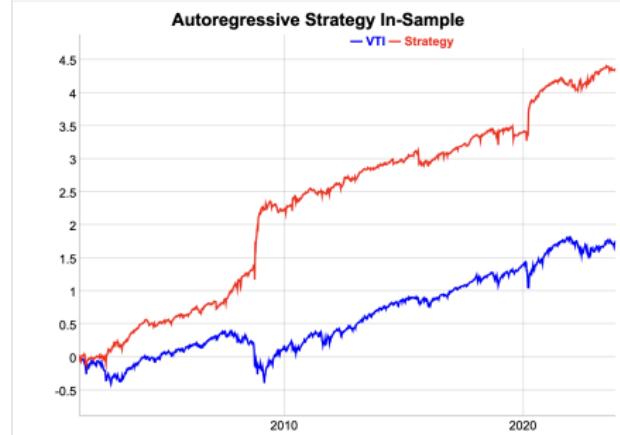
The in-sample autoregressive strategy performs well during periods of high volatility, but not as well in low volatility periods.

The dollar allocations of *VTI* stock are too large in periods of high volatility, which causes over-leverage and very high risk.

The leverage can be reduced by scaling (dividing) the forecasts by their trailing volatility.

The function `HighFreq::run_var()` calculates the trailing variance of a time series using exponential weights.

```
> # Scale the forecasts by their volatility
> fcastv <- sqrt(HighFreq::run_var(fcasts, lambda=0.2)[, 2])
> posv <- ifelse(fcastv > 0, fcasts/fcastv, 0)
> # Calculate the autoregressive strategy PnLs
> pnls <- retp*posv
> costv <- 0.5*bidask*abs(rutils::diffit(posv))
> pnls <- (pnls - costv)
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy In-Sample") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Autoregressive Coefficients in Periods of Low and High Volatility

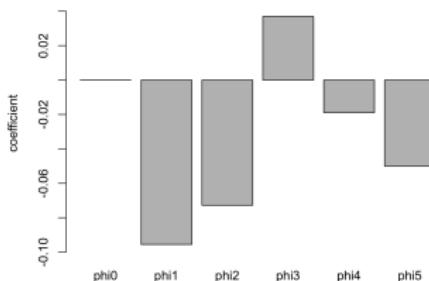
The autoregressive model assumes stationary returns and residuals, with similar volatility over time. In reality stock volatility is highly time dependent.

The autoregressive coefficients in periods of high volatility are very different from those under low volatility.

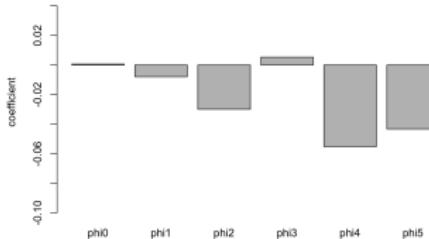
In periods of high volatility, there are larger negative autocorrelations than in low volatility.

```
> # Calculate the high volatility AR coefficients
> respv <- retp["2008/2011"]
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, NROW(predm)), predm)
> predinv <- MASS::ginv(predm)
> coeffh <- drop(predinv %*% respv)
> coeffn <- paste0("phi", 0:(NROW(coeffh)-1))
> barplot(coeffh ~ coeffn, main="High Volatility AR Coefficients",
+   col="grey", xlab="", ylab="coefficient", ylim=c(-0.1, 0.05))
> # Calculate the low volatility AR coefficients
> respv <- retp["2012/2019"]
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, NROW(predm)), predm)
> predinv <- MASS::ginv(predm)
> coeffl <- drop(predinv %*% respv)
> barplot(coeffl ~ coeffn, main="Low Volatility AR Coefficients",
+   xlab="", ylab="coefficient", ylim=c(-0.1, 0.05))
```

High Volatility AR Coefficients



Low Volatility AR Coefficients



# Performance of Low and High Volatility Autoregressive Coefficients

The autoregressive coefficients obtained from periods of high volatility are overfitted and only perform well in periods of high volatility. Similarly the low volatility coefficients.

```
> # Calculate the pnls for the high volatility AR coefficients
> predm <- lapply(1:orderp, rutils::lagit, input=retp)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> fcasts <- predm %*% coeffh
> pnlh <- retp*fcasts
> pnlh <- pnlh*sd(retp[retp<0])/sd(pnlh[pnlh<0])
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnlh)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy High Volatility Coefficients") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
> # Calculate the pnls for the low volatility AR coefficients
> fcasts <- predm %*% coeffl
> pnll <- retp*fcasts
> pnll <- pnll*sd(retp[retp<0])/sd(pnll[pnll<0])
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnll)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy Low Volatility Coefficients") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```



# draft: Regime Switching Autoregressive Strategy

Run two competing autoregressive models for the low and high volatility regimes, with low and high volatility coefficients.

Calculate the Kalman gains from the trailing square forecast errors, and apply the Kalman gains to the forecasts.

Doesn't work because the square forecast errors are similar.

The autoregressive strategy invests dollar amounts of *VTI* stock proportional to the in-sample forecasts.

The in-sample autoregressive strategy performs well during periods of high volatility, but not in low volatility periods.

The dollar allocations of *VTI* stock are too large in periods of high volatility, which causes over-leverage and very high risk.

The autoregressive model assumes stationary returns and residuals, with similar volatility over time. In reality stock volatility is highly time dependent.

```
> # Define the response and predictor matrices
> respv <- retp
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predinv <- MASS::ginv(predm)
> # Simulate strategy with high volatility AR coefficients
> fcasts <- drop(predm %*% coeffh)
>
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy In-Sample") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# The Winsor Function

Some models produce very large dollar allocations, leading to large portfolio leverage (dollars invested divided by the capital).

The *winsor function* maps the *model weight*  $w$  into the dollar amount for investment. The hyperbolic tangent function can serve as a winsor function:

$$W(x) = \frac{\exp(\lambda w) - \exp(-\lambda w)}{\exp(\lambda w) + \exp(-\lambda w)}$$

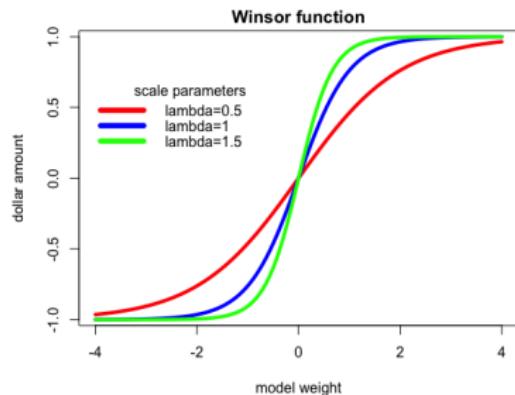
Where  $\lambda$  is the scale parameter.

The hyperbolic tangent is close to linear for small values of the *model weight*  $w$ , and saturates to  $+1\$ / -1\$$  for very large positive and negative values of the *model weight*.

The saturation effect limits (caps) the leverage in the strategy to  $+1\$ / -1\$$ .

For very small values of the scale parameter  $\lambda$ , the invested dollar amount is linear for a wide range of *model weights*. So the strategy is mostly invested in dollar amounts proportional to the *model weights*.

For very large values of the scale parameter  $\lambda$ , the invested dollar amount jumps from  $-1\$$  for negative *model weights* to  $+1\$$  for positive *model weight* values. So the strategy is invested in either  $-1\$$  or  $+1\$$  dollar amounts.



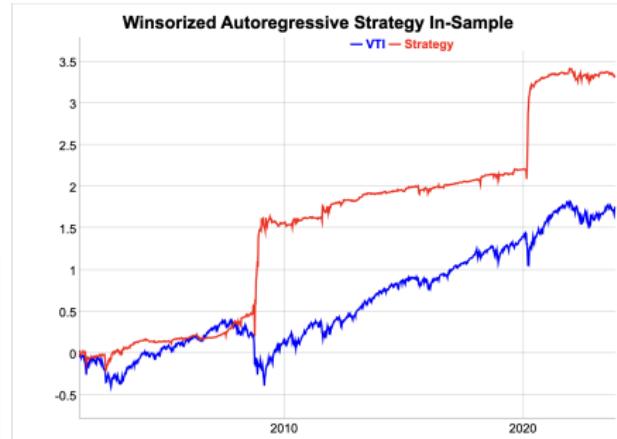
```
> lambdav <- c(0.5, 1, 1.5)
> colovr <- c("red", "blue", "green")
> # Define the winsor function
> winsorfun <- function(retpl, lambdaf) tanh(lambdaf*retpl)
> # Plot three curves in loop
> for (indeks in 1:3) {
+   curve(expr=winsorfun(x, lambda=lambdav[indeks]),
+         xlim=c(-4, 4), type="l", lwd=4,
+         xlab="model weight", ylab="dollar amount",
+         col=colovr[indeks], add=(indeks>1))
+ } # end for
> # Add title and legend
> title(main="Winsor function", line=0.5)
> legend("topleft", title="scale parameters\n",
+        paste("lambdaf", lambdav, sep="="), inset=0.0, cex=1.0,
+        lwd=6, bty="n", y.intersp=0.3, lty=1, col=colovr)
```

# Winsorized Autoregressive Strategy

The performance of the autoregressive strategy can be improved by fitting its coefficients using the *winsorized returns*, to reduce the effect of time-dependent volatility.

The performance can also be improved by *winsorizing* the forecasts, by reducing the leverage due to very large forecasts.

```
> # Winsorize the VTI returns
> retw <- winsorfun(rtp/0.01, lambda=0.1)
> # Define the response and predictor matrices
> predm <- lapply(1:orderp, rutils::lagit, input=retw)
> predm <- rutils::do.call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% retw
> # Calculate the scaled in-sample forecasts of VTI
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fccasts <- fccasts/fcastv
> # Winsorize the forecasts
> # fccasts <- winsorfun(fccasts/mad(fccasts), lambda=1.5)
```



```
> # Calculate the autoregressive strategy PnLs
> pnls <- rtp*fccasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(rtp[rtp<0])/sd(pnls[pnls<0]))
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(rtp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Winsorized Autoregressive Strategy In-Sample") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Autoregressive Strategy With Returns Scaled By Volatility

The performance of the autoregressive strategy can be improved by fitting its coefficients using returns divided by their trailing volatility.

Dividing the returns by their trailing volatility reduces the effect of time-dependent volatility.

The function `HighFreq::run_var()` calculates the trailing variance of a time series using exponential weights.

```
> # Scale the returns by their trailing volatility
> varv <- HighFreq::run_var(retp, lambda=0.99)[, 2]
> retsc <- ifelse(varv > 0, retp/sqrt(varv), 0)
> # Calculate the AR coefficients
> predm <- lapply(1:orderp, rutils::lagit, input=retsc)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% retsc
> # Calculate the scaled in-sample forecasts of VTI
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcاستv[1:100] <- 1
> fcاستs <- fcاستs/fcastv
```



```
> # Calculate the autoregressive strategy PnLs
> pnls <- retp*fcasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy With Returns Scaled By Volatility",
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Autoregressive Strategy in Trading Time

The performance of the autoregressive strategy can be improved by fitting its coefficients using returns divided by the trading volumes (returns in *trading time*).

Dividing the returns by the trading volumes reduces the effect of time-dependent volatility.

```
> # Calculate VTI returns and trading volumes
> ohlc <- rutils::efenv$VTI
> datev <- zoo::index(ohlc)
> nrows <- NROW(ohlc)
> closep <- quantmod::Cl(ohlc)
> retpl <- rutils::diffit(log(closep))
> volumv <- quantmod::Vo(ohlc)
> # Scale the returns using volume clock to trading time
> volumr <- HighFreq::run_mean(volumv, lambda=0.8)
> respv <- retpl*volumr/volumv
> # Calculate the AR coefficients
> orderp <- 5
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> # Calculate the scaled in-sample forecasts of VTI
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fcasts <- fccasts/fcastv
```

Autoregressive Strategy With Returns Scaled By Volume



```
> # Calculate the autoregressive strategy PnLs
> pnls <- retpl*fcasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retpl[retpl<0])/sd(pnls[pnls<0])
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retpl, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv["2010:"], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy With Returns Scaled By Volume") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

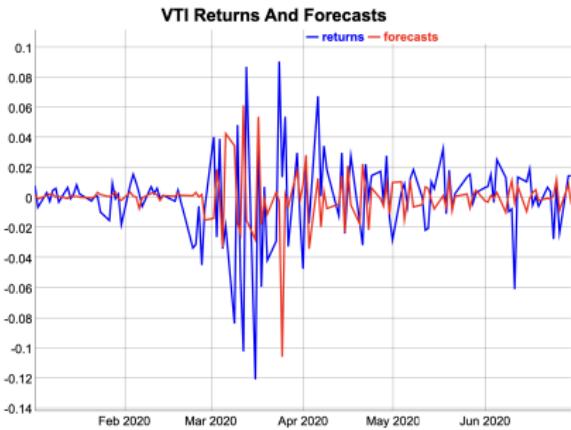
# Mean Squared Error of the Autoregressive Forecasting Model

The accuracy of a forecasting model can be measured using the *mean squared error* and the *correlation*.

The mean squared error (*MSE*) of a forecasting model is the average of the squared forecasting errors  $\varepsilon_i$ , equal to the differences between the *forecasts*  $f_t$  minus the actual values  $r_t$ :  $\varepsilon_i = f_t - r_t$ :

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (r_t - f_t)^2$$

```
> # Define the response and predictor matrices
> respv <- retp
> orderp <- 5
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do.call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> # Calculate the in-sample forecasts of VTI (fitted values)
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> fcasts <- predm %*% coeff
> # Calculate the correlation between forecasts and returns
> cor(fccasts, retp)
> # Calculate the forecasting errors
> errorf <- (fcasts - retp)
> # Mean squared error
> mean(errorf^2)
```



```
> # Plot the forecasts
> datav <- cbind(retp, fcasts)[["2020-01/2020-06"]]
> colnames(datav) <- c("returns", "forecasts")
> dygraphs::dygraph(datav,
+   main="VTI Returns And Forecasts") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# In-sample Order Selection of Autoregressive Forecasting

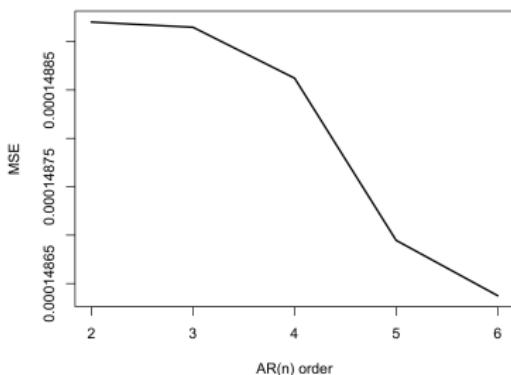
The mean squared errors (*MSE*) of the *in-sample* forecasts decrease steadily with the increasing order parameter  $n$  of the  $AR(n)$  forecasting model.

*In-sample forecasting* consists of first fitting an  $AR(n)$  model to the data, and calculating its coefficients.

The *in-sample* forecasts are calculated by multiplying the predictor matrix by the fitted AR coefficients.

```
> # Calculate the forecasts as a function of the AR order
> fcasts <- lapply(2:NCOL(predm), function(ordern) {
+   # Calculate the fitted AR coefficients
+   predinv <- MASS::ginv(predm[, 1:ordern])
+   coeff <- predinv %*% respv
+   # Calculate the in-sample forecasts of VTI
+   drop(predm[, 1:ordern] %*% coeff)
+ }) # end lapply
> names(fccasts) <- paste0("n=", 2:NCOL(predm))
```

MSE of In-sample  $AR(n)$  Forecasting Model for VTI



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((respv - x)^2), cor=cor(respv, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- names(fccasts)
> # Plot forecasting MSE
> plot(x=2:NCOL(predm), y=mse[, 1],
+       xlab="AR(n) order", ylab="MSE", type="l", lwd=2,
+       main="MSE of In-sample AR(n) Forecasting Model for VTI")
```

# Out-of-Sample Forecasting Using Autoregressive Models

The mean squared errors (*MSE*) of the *out-of-sample* forecasts increase with the increasing order parameter  $n$  of the  $AR(n)$  model.

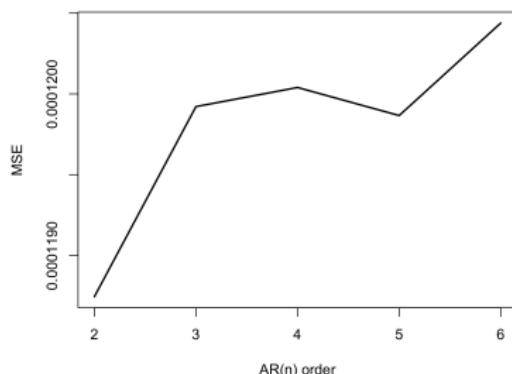
The reason for the increasing out-of-sample MSE is the *overfitting* of the coefficients to the training data for larger order parameters.

*Out-of-sample forecasting* consists of first fitting an  $AR(n)$  model to the training data, and calculating its coefficients.

The *out-of-sample* forecasts are calculated by multiplying the *out-of-sample* predictor matrix by the fitted AR coefficients.

```
> # Define in-sample and out-of-sample intervals
> nrows <- NROW(rtpt)
> insample <- 1:(nrows %/% 2)
> outsample <- (nrows %/% 2 + 1):nrows
> # Calculate the forecasts as a function of the AR order
> fcasts <- lapply(2:NCOL(predm), function(ordern) {
+   # Calculate the fitted AR coefficients
+   predinv <- MASS::ginv(predm[insample, 1:ordern])
+   coeff <- predinv %*% respv[insample]
+   # Calculate the out-of-sample forecasts of VTI
+   drop(predm[outsample, 1:ordern] %*% coeff)
+ }) # end lapply
> names(fccasts) <- paste0("n=", 2:NCOL(predm))
```

MSE of Out-of-sample AR( $n$ ) Forecasting Model for VTI



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((respv[outsample] - x)^2), cor=cor(respv[outsample],
+ })) # end supply
> mse <- t(mse)
> rownames(mse) <- names(fccasts)
> # Plot forecasting MSE
> plot(x=2:NCOL(predm), y=mse[, 1],
+       xlab="AR(n) order", ylab="MSE", type="l", lwd=2,
+       main="MSE of Out-of-sample AR(n) Forecasting Model for VTI")
```

# Out-of-Sample Autoregressive Strategy

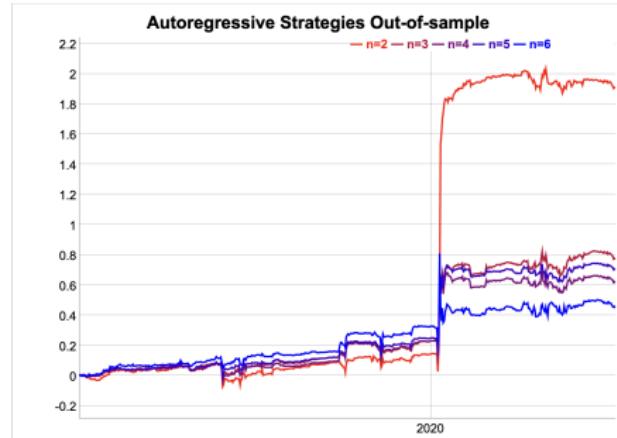
The autoregressive strategy invests a dollar amount of  $VTI$  proportional to the AR forecasts.

The out-of-sample, risk-adjusted performance of the autoregressive strategy is better for a smaller order parameter  $n$  of the  $AR(n)$  model.

The optimal order parameter is  $n = 2$ , with a positive intercept coefficient  $\varphi_0$  (since the average  $VTI$  returns were positive), and a negative coefficient  $\varphi_1$  (because of strong negative autocorrelations in periods of high volatility).

Decreasing the order parameter of the autoregressive model is a form of *shrinkage* because it reduces the number of predictive variables.

```
> # Calculate the optimal AR coefficients
> predinv <- MASS::ginv(predm[insample, 1:2])
> coeff <- drop(predinv %*% respv[insample])
> # Calculate the out-of-sample PnLs
> pnls <- lapply(fcasts, function(fcast) {
+   pnls <- fcast*retp[outsample]
+   pnls<-sd(retp[retp<0])/sd(pnls[pnls<0])
+ }) # end lapply
> pnls <- rutils::do_call(cbind, pnls)
> colnames(pnls) <- names(fcasts)
```



```
> # Plot dygraph of out-of-sample PnLs
> colorv <- colorRampPalette(c("red", "blue"))(NCOL(pnls))
> colv <- colnames(pnls)
> endd <- rutils::calc_endpoints(pnls, interval="weeks")
> dygraphs::dygraph(cumsum(pnls)[endd],
+   main="Autoregressive Strategies Out-of-sample") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyLegend(width=300)
```

# Out-of-Sample Autoregressive Forecasts

The autoregressive coefficients  $\varphi$  are equal to the in-sample response  $r$  times the inverse of the in-sample predictor matrix  $\mathbb{P}$ :

$$\varphi = \mathbb{P}^{-1} r$$

The variance of the autoregressive coefficients  $\sigma_\varphi^2$  is equal to the variance of the in-sample forecast residuals  $\sigma_\varepsilon^2$  divided by the squared predictor matrix  $\mathbb{P}$ :

$$\sigma_\varphi^2 = \sigma_\varepsilon^2 (\mathbb{P}^T \mathbb{P})^{-1}$$

The t-values of the autoregressive coefficients are equal to the coefficient values divided by their volatilities:

$$\varphi_{tval} = \frac{\varphi}{\sigma_\varphi}$$

The *out-of-sample* autoregressive forecast  $f_t$  is equal to the single row of the predictor  $\mathbb{P}_t$  times the fitted AR coefficients  $\varphi$ :

$$f_t = \varphi \mathbb{P}_t$$

The variance  $\sigma_f^2$  of the *forecast value* is equal to the inner product of the predictor  $\mathbb{P}_t$  times the coefficient covariance matrix  $\sigma_\varphi^2$ :

$$\sigma_f^2 = \mathbb{P}_t \sigma_\varphi^2 \mathbb{P}_t^T$$

```
> # Define the look-back range
> lookb <- 100
> tday <- nrow
> startp <- max(1, tday-lookb)
> rangev <- startp:(tday-1)
> # Subset the response and predictors
> respv <- respv[rangev]
> predm <- predm[rangev]
> # Invert the predictor matrix
> predinv <- MASS::ginv(preds)
> # Calculate the fitted AR coefficients
> coeff <- predinv %*% resps
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- preds %*% coeff
> # Calculate the residuals (forecast errors)
> resid <- (fcasts - resps)
> # Calculate the variance of the residuals
> varv <- sum(resid^2)/(NROW(preds)-NROW(coeff))
> # Calculate the predictor matrix squared
> pred2 <- crossprod(preds)
> # Calculate the covariance matrix of the AR coefficients
> covmat <- varv*MASS::ginv(pred2)
> coefsds <- sqrt(diag(covmat))
> # Calculate the t-values of the AR coefficients
> coefft <- drop(coeff/coefsds)
> # Calculate the out-of-sample forecast
> predn <- predm[tday, ]
> fcast <- drop(predn %*% coeff)
> # Calculate the variance of the forecast
> varf <- drop(predn %*% covmat %*% t(predn))
> # Calculate the t-value of the out-of-sample forecast
> fcast/sqrt(varf)
```

# Rolling Autoregressive Forecasting Model

The autoregressive coefficients can be calibrated dynamically over a *rolling* look-back interval, and applied to calculating the *out-of-sample* forecasts.

*Backtesting* is the simulation of a model on historical data to test its forecasting accuracy.

The autoregressive forecasting model can be *backtested* by calculating forecasts over either a *rolling* or an *expanding* look-back interval.

If the start date is fixed at the first row then the look-back interval is *expanding*.

The coefficients of the *AR(n)* process are fitted to past data, and then applied to calculating out-of-sample forecasts.

The *backtesting* procedure allows determining the optimal *meta-parameters* of the forecasting model: the order *n* of the *AR(n)* model and the length of look-back interval (*lookb*).

```
> # Perform rolling forecasting
> lookb <- 500
> fcasts <- parallel::mclapply(1:nrows, function(tday) {
+   if (tday > lookb) {
+     # Define the rolling look-back range
+     startp <- max(1, tday-lookb)
+     # startp <- 1 # Expanding look-back range
+     rangev <- startp:(tday-1) # In-sample range
+     # Subset the response and predictors
+     resps <- respv[rangev]
+     preds <- predm[rangev]
+     # Calculate the fitted AR coefficients
+     predinv <- MASS::ginv(preds)
+     coeff <- predinv %*% resps
+     # Calculate the in-sample forecasts of VTI (fitted values)
+     fcasts <- preds %*% coeff
+     # Calculate the residuals (forecast errors)
+     resid <- (fcasts - resps)
+     # Calculate the variance of the residuals
+     varv <- sum(resid^2)/(NROW(preds)-NROW(coeff))
+     # Calculate the covariance matrix of the AR coefficients
+     pred2 <- crossprod(preds)
+     covmat <- varv*MASS::ginv(pred2)
+     coefs <- sqrt(diag(covmat))
+     coefft <- drop(coeff/coefs) # t-values of the AR coefficients
+     # Calculate the out-of-sample forecast
+     predn <- predm[tday, ]
+     fcast <- drop(predn %*% coeff)
+     # Calculate the variance of the forecast
+     varf <- drop(predn %*% covmat %*% t(predn))
+     return(c(sd(resps), fcast=fcast, fstderr=sqrt(varf), coefft=coefft))
+   } else {
+     return(c(volv=0, fcast=0, fstderr=0, coefft=rep(0, NCOL(predm)))
+   } # end if
+ }) # end sapply
```

# Rolling Autoregressive Strategy Performance

In the rolling autoregressive strategy, the autoregressive coefficients are calibrated on past data from a *rolling* look-back interval, and applied to calculating the *out-of-sample* forecasts.

The rolling autoregressive strategy performance depends on the length of the look-back interval.

```
> # Coerce fcasts to a time series
> fccasts <- t(fccasts)
> ncols <- NCOL(fccasts)
> colnames(fccasts) <- c("volv", "fccasts", "fstderr", colnames(predm))
> fccasts <- xts::xts(fccasts, zoo::index(retp))
> # Calculate the strategy PnLs
> pnls <- retp$fccasts$fccasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+ c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
```



```
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Autoregressive Strategy") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Forecasting Returns Using Rolling Autoregressive Models

Adopt this text to some other slides.

Forecasting using an autoregressive model is performed by first fitting an  $AR(n)$  model to past data, and calculating its coefficients.

The fitted AR coefficients are then applied to calculating the *out-of-sample* forecasts.

The forecasting model depends on two unknown *meta-parameters*: the order  $n$  of the  $AR(n)$  model and the length of the look-back interval (lookb).

# Backtesting the Autoregressive Model

The *meta-parameters* of the *backtesting* function are the order  $n$  of the  $AR(n)$  model and the length of the look-back interval ( $lookb$ ).

The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

*Backtesting* is the simulation of a model on historical data to test its forecasting accuracy.

The autoregressive forecasting model can be *backtested* by calculating forecasts over either a *rolling* or an *expanding* look-back interval.

If the start date is fixed at the first row then the look-back interval is *expanding*.

The coefficients of the  $AR(n)$  process are fitted to past data, and then applied to calculating out-of-sample forecasts.

The *backtesting* procedure allows determining the optimal *meta-parameters* of the forecasting model: the order  $n$  of the  $AR(n)$  model and the length of look-back interval ( $lookb$ ).

```
> # Define backtesting function
> sim_fcasts <- function(lookb=100, ordern=5, fixedlb=TRUE) {
+   # Perform rolling forecasting
+   fcasts <- sapply((lookb+1):nrows, function(tday) {
+     # Rolling look-back range
+     startp <- max(1, tday-lookb)
+     # Expanding look-back range
+     if (!fixedlb) {startp <- 1}
+     startp <- max(1, tday-lookb)
+     rangev <- startp:(tday-1) # In-sample range
+     # Subset the response and predictors
+     respv <- respv[rangev]
+     predm <- predm[rangev, 1:ordern]
+     # Invert the predictor matrix
+     predinv <- MASS::ginv(preds)
+     # Calculate the fitted AR coefficients
+     coeff <- predinv %*% respv
+     # Calculate the out-of-sample forecast
+     drop(predm[tday, 1:ordern] %*% coeff)
+   }) # end apply
+   # Add warmup period
+   fcasts <- c(rep(0, lookb), fccasts)
+ } # end sim_fcasts
> # Simulate the rolling autoregressive forecasts
> fcasts <- sim_fcasts(lookb=100, ordern=5)
> c(mse=mean((fcasts - retp)^2), cor=cor(retp, fcasts))
```

# Forecasting Dependence On the Look-back Interval

The *backtesting* function can be used to find the optimal *meta-parameters* of the autoregressive forecasting model.

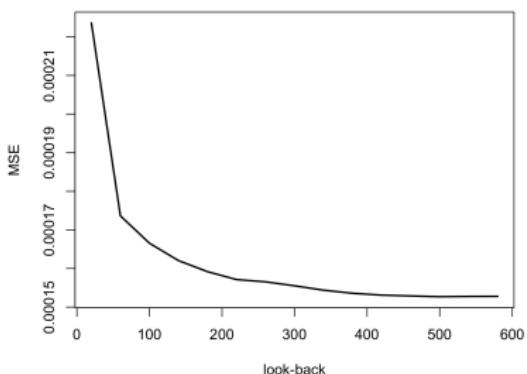
The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model and on the length of the look-back interval ( $lookb$ ).

The accuracy of the forecasting model increases with longer look-back intervals ( $lookb$ ), because more data improves the estimates of the autoregressive coefficients.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> lookbv <- seq(20, 600, 40)
> fcasts <- parLapply(compclust, lookbv, sim_fccasts, ordern=6)
> # Perform parallel bootstrap under Mac-OSX or Linux
> fccasts <- mclapply(lookbv, sim_fccasts, ordern=6, mc.cores=ncores)
```

MSE of AR Forecasting Model As Function of Look-back



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retp - x)^2), cor=cor(retp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- lookbv
> # Select optimal lookb interval
> lookb <- lookbv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=lookbv, y=mse[, 1],
+       xlab="look-back", ylab="MSE", type="l", lwd=2,
+       main="MSE of AR Forecasting Model As Function of Look-back")
```

# Order Dependence With Fixed Look-back

The *backtesting* function can be used to find the optimal *meta-parameters* of the autoregressive forecasting model.

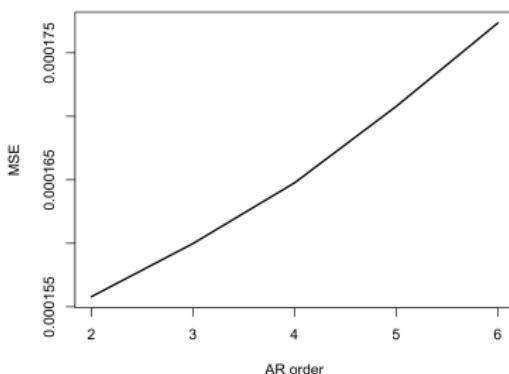
The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model and on the length of the look-back interval ( $lookb$ ).

The accuracy of the forecasting model decreases for larger AR order parameters, because of overfitting in-sample.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> orderv <- 2:6
> fcasts <- parLapply(compclust, orderv, sim_fcasts, lookb=lookb)
> stopCluster(compclust) # Stop R processes over cluster under Wi
> # Perform parallel bootstrap under Mac-OSX or Linux
> fcasts <- mclapply(orderv, sim_fcasts,
+   lookb=lookb, mc.cores=ncores)
```

MSE of Forecasting Model As Function of AR Order



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retp - x)^2), cor=cor(retp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- orderv
> # Select optimal order parameter
> ordern <- orderv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=orderv, y=mse[, 1],
+   xlab="AR order", ylab="MSE", type="l", lwd=2,
+   main="MSE of Forecasting Model As Function of AR Order")
```

# Autoregressive Strategy With Fixed Look-back

The return forecasts are calculated just before the close of the markets, so that trades can be executed before the close.

The autoregressive strategy returns are large in periods of high volatility, but much smaller in periods of low volatility. This because the forecasts are bigger in periods of high volatility, and also because the forecasts are more accurate, because the autocorrelations of stock returns are much higher in periods of high volatility.

Using the return forecasts as portfolio weights produces very large weights in periods of high volatility, and creates excessive risk.

To reduce excessive risk, a binary strategy can be used, with portfolio weights equal to the sign of the forecasts.

```
> # Simulate the rolling autoregressive forecasts
> fcasts <- sim_fcasts(lookb=lookb, ordern=ordern)
> # Calculate the strategy PnLs
> pnls <- fcasts*retp
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls, (retp+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
```



```
> # Annualized Sharpe ratios of VTI and AR strategy
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of AR strategy combined with VTI
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy Fixed Look-back") %>%
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dySeries(name="Combined", strokeWidth=3) %>%
+   dyLegend(show="always", width=300)
```

# Order Dependence With Expanding Look-back

The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

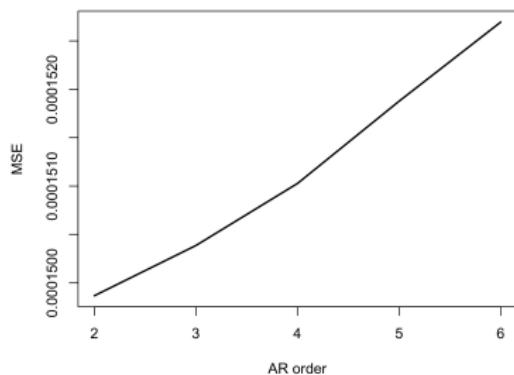
The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model.

Longer look-back intervals (`lookb`) are usually better for the autoregressive forecasting model.

The return forecasts are calculated just before the close of the markets, so that trades can be executed before the close.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> orderv <- 2:6
> fcasts <- parLapply(compclust, orderv, sim_fccasts,
+   lookb=lookb, fixedlb=FALSE)
> stopCluster(compclust) # Stop R processes over cluster under Win
> # Perform parallel bootstrap under Mac-OSX or Linux
> fcasts <- mclapply(orderv, sim_fccasts,
+   lookb=lookb, fixedlb=FALSE, mc.cores=ncores)
```

MSE With Expanding Look-back As Function of AR Order



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retlp - x)^2), cor=cor(retlp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- orderv
> # Select optimal order parameter
> ordern <- orderv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=orderv, y=mse[, 1],
+   xlab="AR order", ylab="MSE", type="l", lwd=2,
+   main="MSE With Expanding Look-back As Function of AR Order")
```

# Autoregressive Strategy With Expanding Look-back

The model with an *expanding* look-back interval has better performance compared to the *fixed* look-back interval.

The autoregressive strategy returns are large in periods of high volatility, but much smaller in periods of low volatility. This because the forecasts are bigger in periods of high volatility, and also because the forecasts are more accurate, because the autocorrelations of stock returns are much higher in periods of high volatility.

```
> # Simulate the autoregressive forecasts with expanding look-back
> fcasts <- sim_fcasts(lookb=lookb, ordern=ordern, fixedlb=FALSE)
> # Calculate the strategy PnLs
> pnls <- fcasts*retp
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls, (retp+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
```



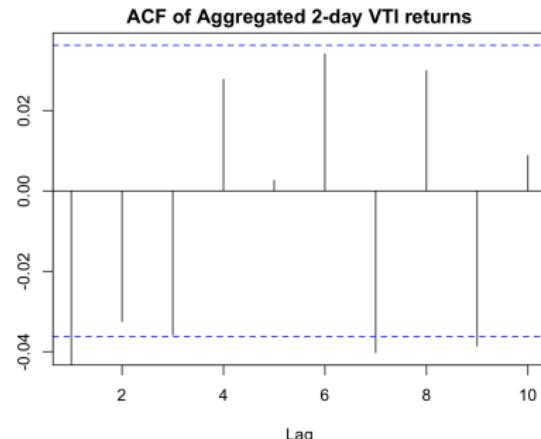
```
> # Annualized Sharpe ratios of VTI and AR strategy
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of AR strategy combined with VTI
> dygraphs::dygraph(cumsum(wealthv)[endd]),
+   main="Autoregressive Strategy Expanding Look-back") %>%
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dySeries(name="Combined", strokeWidth=3) %>%
+   dyLegend(show="always", width=300)
```

# Autocorrelations of Aggregated Stock Returns

Stock returns aggregated over several days have much smaller autocorrelations.

But aggregating stock returns can reduce their noise and improve the performance of autoregressive strategies.

```
> # Calculate VTI returns over non-overlapping 2-day intervals  
> pricev <- na.omit(rutils::etfenv$prices$VTI)  
> reta <- rutils:::diffit(log(pricev), lag=2)  
> reta <- reta[2*(1:(NROW(pricev) %% 2))]  
> # Calculate the autocorrelations of daily VTI returns  
> rutils:::plot_acf(reta, lag=10, main="ACF of Aggregated 2-day VTI")
```



# Autoregressive Strategy With Aggregated Stock Returns

The autoregressive strategy with aggregated stock returns calculates the AR coefficients and the forecasts using stock returns aggregated over several days.

The strategy holds the stock position for several days, and rebalances the position at the end of the aggregation interval.

In this example, the returns are calculated over 2-day intervals, and the stock is held for 2 days. Effectively, there are two strategies, rebalancing every other day.

The aggregated autoregressive strategy has good performance both in periods of high volatility and in low volatility.

```
> # Define the response and predictor matrices
> reta <- rutils::diffit(log(pricev), lag=2)/2
> orderp <- 5
> predm <- lapply(2*(1:orderp), rutils::lagit, input=reta)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, NROW(reta)), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> # Calculate the AR coefficients
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% reta
> coeffn <- paste0("phi", 0:(NROW(coeff)-1))
> barplot(coeff ~ coeffn, xlab="", ylab="t-value", col="grey",
+   main="Coefficients of AR Forecasting Model")
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fcasts <- fccasts/fcastv
```

## Autoregressive Strategy With Aggregated Stock Returns



```
> # Calculate the autoregressive strategy PnLs
> pnls <- reta*fcasts
> pnls <- pnls*sd(reta[reta<0])/sd(pnls[pnls<0])
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(reta, pnls, (reta+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv["2010/"], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy With Aggregated Stock Returns") %
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dyLegend(show="always", width=300)
```

# Daytime and Overnight Stock Strategies

The overnight stock strategy consists of holding a long position only overnight (buying at the market close and selling at the open the next day).

The daytime stock strategy consists of holding a long position only during the daytime (buying at the market open and selling at the close the same day).

The *Overnight Market Anomaly* is the consistent outperformance of overnight returns relative to the daytime returns.

The *Overnight Market Anomaly* has been observed for many decades for most stock market indices, but not always for all stock sectors.

The *Overnight Market Anomaly* is not as pronounced after the 2008–2009 financial crisis.

```
> # Calculate the log of OHLC VTI prices
> ohlc <- log(rutils::etfenv$VTI)
> nrow <- NROW(ohlc)
> openp <- quantmod::Op(ohlc)
> highp <- quantmod::Hi(ohlc)
> lowp <- quantmod::Lo(ohlc)
> closep <- quantmod::Cl(ohlc)
> # Calculate the close-to-close log returns,
> # the daytime open-to-close returns
> # and the overnight close-to-open returns.
> retp <- rutils::diffit(closep)
> colnames(retp) <- "daily"
> retd <- (closep - openp)
> colnames(retd) <- "daytime"
> reton <- (openp - rutils::lagit(closep, lagg=1, pad_zeros=FALSE))
> colnames(reton) <- "overnight"
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retp, reton, retd)
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the Daytime and Overnight strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Wealth of Close-to-Close, Overnight, and Daytime Strategies",
+   dySeries(name="daily", strokeWidth=2, col="blue") %>%
+   dySeries(name="overnight", strokeWidth=2, col="red") %>%
+   dySeries(name="daytime", strokeWidth=2, col="green") %>%
+   dyLegend(width=600)
```

# EMA Mean-Reversion Strategy For Daytime Returns

After the 2008–2009 financial crisis, the cumulative daytime stock index returns have been range-bound. So the daytime returns have more significant negative autocorrelations than overnight returns.

The *EMA* mean-reversion strategy holds stock positions equal to the *sign* of the trailing *EMA* of past returns.

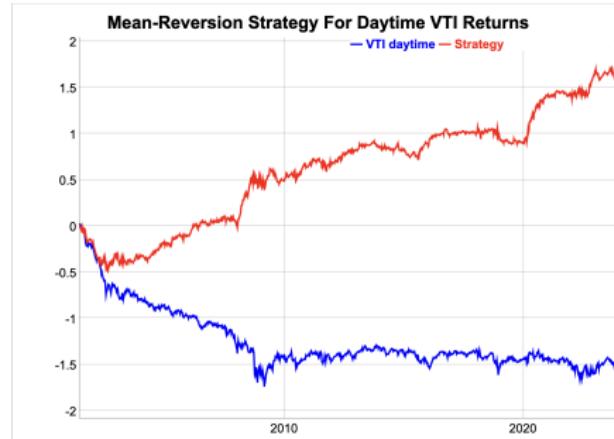
The strategy adjusts its stock position at the end of each day, just before the close of the market.

An alternative strategy holds positions proportional to minus of the sign of the trailing *EMA* of past returns.

The strategy takes very large positions in periods of high volatility, when returns are large and highly anti-correlated.

The limitation is that this strategy makes most of its profits in periods of high volatility, but otherwise it's profits are small.

```
> # Calculate the autocorrelations of daytime and overnight returns
> pacfl <- pacf(retnd, lag.max=10, plot=FALSE)
> sum(pacfl$acf)
> pacfl <- pacf(retov, lag.max=10, plot=FALSE)
> sum(pacfl$acf)
> # Calculate the EMA returns recursively using C++ code
> retma <- HighFreq::run_mean(retnd, lambda=0.4)
> # Calculate the positions and PnLs
> posv <- -rutils::lagit(sign(retma), lagg=1)
> pnls <- retnd*posv
```



```
> # Calculate the pnls and the transaction costs
> bidask <- 0.0001 # The bid-ask spread is equal to 1 basis point
> costv <- 0.5*bidask*abs(rutils::diffit(posv))
> pnls <- (pnls - costv)
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retd, pnls)
> colnames(wealthv) <- c("VTI daytime", "Strategy")
> sqrt(252)*apply(wealthv, function(x)
+ + c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of crossover strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+ + main="Mean-Reversion Strategy For Daytime VTI Returns") %>%
+ + dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+ + dyLegend(show="always", width=300)
```

# Bollinger Strategy For Daytime Returns

The Bollinger z-score for daytime returns  $z_t$ , is equal to the difference between the cumulative returns  $p_t = \sum r_t$  minus their trailing mean  $\bar{p}_t$ , divided by their volatility  $\sigma_t$ :

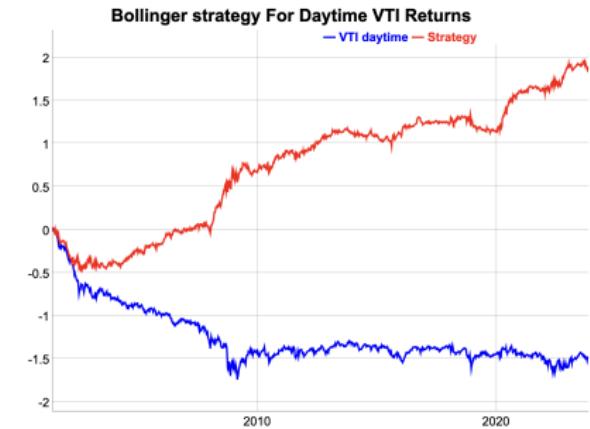
$$z_t = \frac{p_t - \bar{p}_t}{\sigma_t}$$

The Bollinger strategy switches to \$1 dollar long stock if the z-score drops below the threshold of -1 (indicating the prices are cheap), and switches to -\$1 dollar short if the z-score exceeds the threshold of 1 (indicating the prices are rich - expensive).

The Bollinger strategy is a *mean reverting* (contrarian) strategy because it bets on the cumulative returns reverting to their mean value.

The Bollinger strategy has performed well for daytime VTI returns because they exhibit significant mean-reversion.

```
> # Calculate the z-scores of the cumulative daytime returns
> retc <- cumsum(retd)
> lambdaf <- 0.24
> retm <- rutils::lagit(HighFreq::run_mean(retc, lambda=lambdadaf))
> retv <- sqrt(rutils::lagit(HighFreq::run_var(retc, lambda=lambdadaf)))
> zscores <- ifelse(retv > 0, (retc - retm)/retv, 0)
> # Calculate the positions from the Bollinger z-scores
> posv <- rep(NA_integer_, nrowz)
> posv[1] <- 0
> posv <- ifelse(zscores > 1, -1, posv)
> posv <- ifelse(zscores < -1, 1, posv)
> posv <- zoo::na.locf(posv)
```



```
> # Calculate the pnls and the transaction costs
> pnls <- retd*posv
> costv <- 0.5*bidask*abs(rutils::diffit(posv))
> pnls <- (pnls - costv)
> # Calculate the Sharpe ratios
> wealthv <- cbind(retd, pnls)
> colnames(wealthv) <- c("VTI daytime", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of daytime Bollinger strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Bollinger strategy For Daytime VTI Returns") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Overnight Trend Strategy

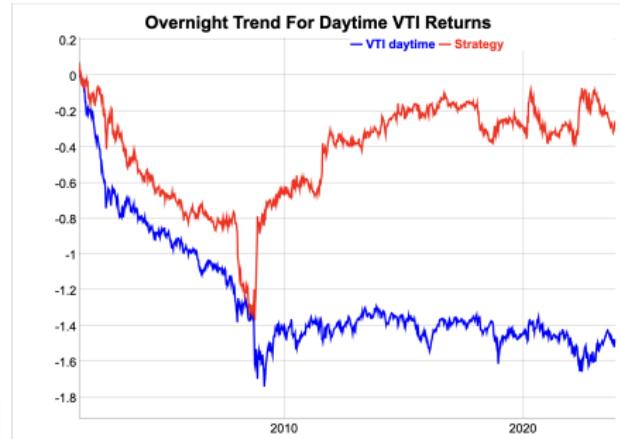
Analysts at *JPMorgan* and at the *Federal Reserve* have observed that there is a trend in the overnight returns.

Positive overnight returns are often followed by positive daytime returns, and vice versa.

If the overnight returns were positive, then the strategy buys \$1 dollar of stock at the market open and sells it at the market close, or if the overnight returns were negative then it shorts -\$1 dollar of stock.

The strategy has performed well immediately after the 2008–2009 financial crisis, but it has waned in recent years.

```
> # Calculate the pnls and the transaction costs
> posv <- sign(ret0n)
> pnls <- posv*retd
> costv <- 0.5*bidask*abs(rutils::diffit(posv))
> pnls <- (pnls - costv)
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retd, pnls)
> colnames(wealthv) <- c("VTI daytime", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of crossover strategy
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Overnight Trend For Daytime VTI Returns") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Multifactor Autoregressive Strategy

Multifactor autoregressive strategies can have a large number of predictors, and are often called *kitchen sink* strategies.

Multifactor strategies are overfit to the in-sample data because they have a large number of parameters.

The out-of-sample performance of the multifactor strategy is much worse than its in-sample performance, because the strategy is overfit to the in-sample data.

```
> # Calculate the VTI daily percentage returns
> retp <- na.omit(rutils::etfenv$returns$VTI)
> datev <- index(retp)
> nrows <- NROW(retp)
> # Define in-sample and out-of-sample intervals
> insample <- 1:(nrows %/% 2)
> outsample <- (nrows %/% 2 + 1):nrows
> cutoff <- nrows %/% 2
> # Define the response and predictor matrices
> respv <- retp
> orderp <- 8 # 9 predictors!!!
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> # Calculate the in-sample fitted autoregressive coefficients
> predinv <- MASS::ginv(predm[insample, ])
> coeff <- predinv %*% respv[insample, ]
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> fcstv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcstv[1:100] <- 1
> fccasts <- fccasts/fcstv
```



```
> # Calculate the autoregressive strategy PnLs
> pnls <- retp*fcasts
> pnls <- pnls*sd(retp<0))/sd(pnls[pnls<0])
> # Calculate the in-sample and out-of-sample Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "AR_multifact")
> sqrt(252)*sapply(wealthv[insample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[outsample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Multifactor Autoregressive Strategy") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# draft: Multifactor Autoregressive Strategy Using t-Values

The out-of-sample performance of the multifactor AR strategy can be improved by using the *t-values* of the AR coefficients.

The forecasts oscillate around 1. Why?

```
> # Calculate the t-values of the AR coefficients
> residss <- (fcasts[insample, ] - respv[insample, ])
> varv <- sum(residss^2)/(nrows-NROW(coeff))
> pred2 <- crossprod(predm[insample, ])
> covmat <- varv*MASS::ginv(pred2)
> coefsds <- sqrt(diag(covmat))
> coefft <- drop(coeff/coefsds)
> coeffn <- paste0("phi", 0:(NROW(coefft)-1))
> # Plot the t-values of the AR coefficients
> barplot(coefft ~ coeffn, xlab="", ylab="t-value", col="grey",
+ main="Coefficient t-values of AR Forecasting Model")
> # Calculate the autoregressive strategy PnLs
> fcasts <- predm %*% coefft
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcavt[1:100] <- 1
> fcavt <- fccasts/fcastv
> pnls <- repv*fcavt
> pnls <- pnls*sd(rep[rep<0])/sd(pnls[pnls<0])
```



```
> # Calculate the in-sample and out-of-sample Sharpe and Sortino rates
> wealthv <- cbind(repv, pnls)
> colnames(wealthv) <- c("VTI", "AR_multifact")
> sqrt(252)*sapply(wealthv[insample, ], function(x)
+ c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[outsample, ], function(x)
+ c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv)[endd],
+ main="Multifactor Autoregressive Strategy Using t-Values") %>%
+ dyOptions(colors=colorv, strokeWidth=2) %>%
+ dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+ dyLegend(show="always", width=300)
```

# Regularization of the Inverse Predictor Matrix

The *SVD* of a rectangular matrix  $\mathbb{A}$  is defined as the factorization:

$$\mathbb{A} = \mathbb{U}\Sigma\mathbb{V}^T$$

Where  $\mathbb{U}$  and  $\mathbb{V}$  are the *singular matrices*, and  $\Sigma$  is a diagonal matrix of *singular values*.

The *generalized inverse* matrix  $\mathbb{A}^{-1}$  satisfies the inverse equation:  $\mathbb{A}\mathbb{A}^{-1}\mathbb{A} = \mathbb{A}$ , and it can be expressed as a product of the *SVD* matrices as follows:

$$\mathbb{A}^{-1} = \mathbb{V}\Sigma^{-1}\mathbb{U}^T$$

If any of the *singular values* are zero then the *generalized inverse* does not exist.

*Regularization* is the removal of zero singular values, to make calculating the inverse matrix possible.

The *generalized inverse* is obtained by removing the zero *singular values*:

$$\mathbb{A}^{-1} = \mathbb{V}_n\Sigma_n^{-1}\mathbb{U}_n^T$$

Where  $\mathbb{U}_n$ ,  $\mathbb{V}_n$  and  $\Sigma_n$  are the *SVD* matrices without the zero *singular values*.

The generalized inverse satisfies the inverse matrix equation:  $\mathbb{A}\mathbb{A}^{-1}\mathbb{A} = \mathbb{A}$ .

```
> # Calculate singular value decomposition of the predictor matrix
> svdec <- svd(predm)
> barplot(svdec$d, main="Singular Values of Predictor Matrix")
> # Calculate generalized inverse from SVD
> invsvd <- svdec$v %*% (t(svdec$u) / svdec$d)
> # Verify inverse property of the inverse
> all.equal(zoo::coredata(predm), predm %*% invsvd %*% predm)
> # Compare with the generalized inverse using MASS::ginv()
> invreg <- MASS::ginv(predm)
> all.equal(invreg, invsvd)
> # Set tolerance for determining zero singular values
> precv <- sqrt(.Machine$double.eps)
> # Check for zero singular values
> round(svdec$d, 12)
> nonzero <- (svdec$d > (precv*svdec$d[1]))
> # Calculate generalized inverse from SVD
> invsvd <- svdec$v[, nonzero] %*%
+ (t(svdec$u[, nonzero]) / svdec$d[nonzero])
> # Verify inverse property of invsvd
> all.equal(zoo::coredata(predm), predm %*% invsvd %*% predm)
> all.equal(invsvd, invreg)
```

# Reduced Inverse of the Predictor Matrix

*Regularization* is the removal of zero singular values, to make calculating the inverse matrix possible.

If the higher order singular values are very small then the inverse matrix will amplify the noise in the response matrix.

*Dimension reduction* is achieved by the removal of small singular values, to improve the out-of-sample performance of the inverse matrix.

The *reduced inverse* is obtained by removing the very small *singular values*.

$$\mathbb{A}^{-1} = \mathbb{V}_n \Sigma_n^{-1} \mathbb{U}_n^T$$

This effectively reduces the number of parameters in the model.

The *reduced inverse* satisfies the inverse equation only approximately (it is *biased*), but it's often used in machine learning because it produces a lower *variance* of the forecasts than the exact inverse.

```
> # Calculate reduced inverse from SVD
> dimax <- 3 # Number of dimensions to keep
> invred <- svdec$v[, 1:dimax] %*%
+   (t(svdec$u[, 1:dimax]) / svdec$d[1:dimax])
> # Inverse property fails for invred
> all.equal(zoo::coredata(predm), predm %*% invred %*% predm)
> # Calculate reduced inverse using RcppArmadillo
> invrcpp <- HighFreq::calc_invsvd(predm, dimax=dimax)
> all.equal(invred, invrcpp, check.attributes=FALSE)
```

# Autoregressive Strategy With Dimension Reduction

Dimension reduction improves the out-of-sample, risk-adjusted performance of the multifactor autoregressive strategy.

The best performance is achieved with the smallest order parameter (strongest reduction)  $\text{dimax} = 2$

```
> # Calculate the in-sample SVD
> svdec <- svd(predm[insample, ])
> # Calculate the in-sample fitted AR coefficients for different dimensions
> dimv <- 2:5
> # dimv <- c(2, 5, 10, NCOL(predm))
> coeffm <- sapply(dimv, function(dimax) {
+   predinv <- svdec$v[, 1:dimax] %*%
+     (t(svdec$u[, 1:dimax]) / svdec$d[1:dimax])
+   predinv %*% respv[insample]
+ }) # end lapply
> colnames(coeffm) <- paste0("dimax=", dimv)
> colorv <- colorRampPalette(c("red", "blue"))(NCOL(coeffm))
> matplot(y=coeffm, type="l", lty="solid", lwd=1, col=colorv,
+   xlab="predictor", ylab="coefficient",
+   main="AR Coefficients For Different Dimensions")
> # Calculate the forecasts of VTI
> fcasts <- predm %*% coeffm
> fccasts <- apply(fccasts, 2, function(x) {
+   fcavt <- sqrt(HighFreq::run_var(matrix(x), lambda=0.8)[, 2])
+   fcavt[1:100] <- 1
+   x/fcavt
+ }) # end apply
> # Simulate the autoregressive strategies
> retn <- coredata(retp)
> pnls <- apply(fccasts, 2, function(x) (x*retn))
> pnls <- xts(pnls, datev)
> # Scale the PnL volatility to that of VTI
> pnls <- lapply(pnls, function(x) x/sd(x))
> pnls <- sd(retp)*do.call(cbind, pnls)
```

## Kitchen Sink Strategies With Dimension Reduction



```
> # Calculate the in-sample and out-of-sample Sharpe and Sortino ratios
> wealthv <- cbind(retp, pnls)
> sqrt(252)*sapply(wealthv[insample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[outsample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategies With Dimension Reduction") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=500)
```

# Shrinkage of the Autoregressive Coefficients

The AR coefficients can be found by minimizing the sum of the squared errors of the in-sample forecasts.

The objective function is the sum of the squared errors of the in-sample forecasts, plus a penalty term proportional to the square of the AR coefficients:

$$\text{ObjFunc} = \sum_{i=1}^n (r_t - f_t)^2 + \sum_{i=1}^k \varphi_i^2 \lambda^i$$

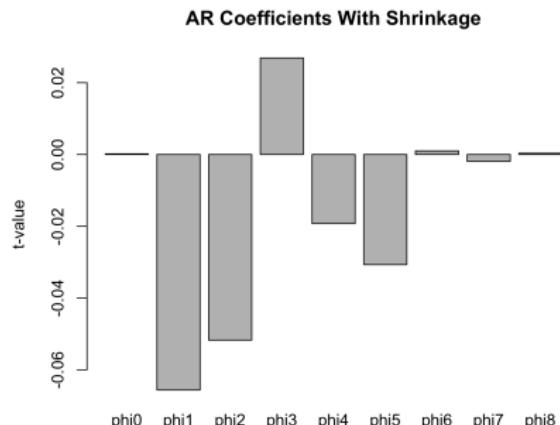
As the objective function is minimized, the *shrinkage* penalty shrinks the AR coefficients closer to zero.

A larger *shrinkage factor*  $\lambda$  applies more *shrinkage* to the AR coefficients, to shrink them closer to zero.

The *shrinkage* penalty is greater for the higher order coefficients, to reflect the fact that returns from the more distant past should be less important for forecasting the current returns.

The coefficient shrinkage also produces a dimension reduction effect, because the higher order coefficients are shrunk closer to zero.

```
> # Objective function for the in-sample AR coefficients
> objfun <- function(coef, lambdaf) {
+   fcasts <- predm[insample, ] %*% coef
+   lambdav <- lambdaf^(0:orderp)
+   10000*sum((respv[insample, ] - fcasts)^2) + sum(lambdav*coef^2)
+ } # end objfun
```



```
> # Perform optimization using the quasi-Newton method
> optiml <- optim(par=numeric(orderp+1),
+   fn=objfun, lambdaf=5.0, method="L-BFGS-B",
+   upper=rep(10000, orderp+1),
+   lower=rep(-10000, orderp+1))
> # Extract the AR coefficients
> coeff <- optiml$par
> coeffn <- paste0("phi", 0:orderp)
> names(coeff) <- coeffn
> barplot(coeff ~ coeffn, xlab="", ylab="t-value", col="grey",
+   main="AR Coefficients With Shrinkage")
```

# Multifactor Autoregressive Strategy With Shrinkage

The out-of-sample performance of the multifactor autoregressive strategy can be improved by applying *shrinkage* to the AR coefficients.

The value of the *shrinkage factor*  $\lambda$  can be chosen to maximize the out-of-sample performance.

```
> # Calculate the forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fcasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fcasts <- fcasts/fcastv
> # Calculate the autoregressive strategy PnLs
> pnls <- retpl*fcasts
> pnls <- pnls*sd(retpl[retpl<0])/sd(pnls[pnls<0])
```



```
> # Calculate the in-sample and out-of-sample Sharpe and Sortino rates
> wealthv <- cbind(retpl, pnls)
> colnames(wealthv) <- c("VTI", "AR_multifact")
> sqrt(252)*sapply(wealthv[,1], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[,2], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv[,endd]),
+   main="Multifactor Autoregressive Strategy With Shrinkage") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# Kitchen Sink Model of Stock Returns

The "kitchen sink" model uses many possible predictor variables, so the predictor matrix has a very large number of columns.

The predictor matrix includes the lagged and scaled daily returns and the squared returns.

stock returns  $r_t$  can be fitted into an *autoregressive* model  $AR(n)$  with a constant intercept term  $\varphi_0$ :

$$r_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n} + \varepsilon_t$$

The *residuals*  $\varepsilon_t$  are assumed to be normally distributed, independent, and stationary.

The autoregressive model can be written in matrix form as:

$$\mathbf{r} = \boldsymbol{\varphi} \mathbb{P} + \boldsymbol{\varepsilon}$$

Where  $\boldsymbol{\varphi} = \{\varphi_0, \varphi_1, \varphi_2, \dots, \varphi_n\}$  is the vector of autoregressive coefficients.

The *autoregressive* model is equivalent to *multivariate* linear regression, with the *response* equal to the returns  $\mathbf{r}$ , and the columns of the *predictor matrix*  $\mathbb{P}$  equal to the lags of the returns.

```
> # Calculate the returns of VTI, TLT, and VXX
> retp <- na.omit(rutils::etfenv$returns[, c("VTI", "TLT", "VXX")])
> datev <- zoo::index(retp)
> nrows <- NROW(retp)
> # Define the response and the VTI predictor matrix
> respv <- retp$VTI
> orderp <- 5
> predm <- lapply(1:orderp, rutils::lagit, input=respv)
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, nrows), predm)
> colnames(predm) <- c("phi0", paste0("lag", 1:orderp))
> # Add the TLT predictor matrix
> predx <- lapply(1:orderp, rutils::lagit, input=retp$TLT)
> predx <- rutils::do_call(cbind, predx)
> colnames(predx) <- paste0("TLT", 1:orderp)
> predm <- cbind(predm, predx)
> # Add the VXX predictor matrix
> predx <- lapply(1:orderp, rutils::lagit, input=retp$VXX)
> predx <- rutils::do_call(cbind, predx)
> colnames(predx) <- paste0("VXX", 1:orderp)
> predm <- cbind(predm, predx)
> # Perform the multivariate linear regression
> regmod <- lm(respv ~ predm - 1)
> summary(regmod)
```

# Kitchen Sink Autoregressive Strategy

Multifactor autoregressive strategies can have a large number of predictors, and are often called *kitchen sink* strategies.

Multifactor strategies are overfit to the in-sample data because they have a large number of parameters.

The out-of-sample performance of the multifactor strategy is much worse than its in-sample performance, because the strategy is overfit to the in-sample data.

```
> # Define in-sample and out-of-sample intervals
> insample <- 1:(nrows %/% 2)
> outsample <- (nrows %/% 2 + 1):nrows
> cutoff <- nrows %/% 2
> # Calculate the in-sample fitted autoregressive coefficients
> predinv <- MASS::ginv(predm[insample, ])
> coeff <- predinv %*% respv[insample, ]
> coeffn <- colnames(predm)
> barplot(coeff ~ coeffn, xlab="", ylab="t-value", col="grey",
+   main="Coefficients of Kitchen Sink Autoregressive Model")
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fccasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fcasts <- fccasts/fcastv
```



```
> # Calculate the autoregressive strategy PnLs
> pnls <- respv*fcasts
> pnls <- pnls*sd(respv[respv<0])/sd(pnls[pnls<0])
> # Calculate the in-sample and out-of-sample Sharpe and Sortino ratios
> wealthv <- cbind(respv, pnls)
> colnames(wealthv) <- c("VTI", "Kitchen sink")
> sqrt(252)*sapply(wealthv[insample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[outsample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Kitchen Sink Autoregressive Strategy") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datenv[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# Kitchen Sink Strategy With Dimension Reduction

Dimension reduction improves the out-of-sample, risk-adjusted performance of the multifactor autoregressive strategy.

The best performance is achieved with the smallest order parameter (strongest reduction)  $\text{dimax} = 2$

```
> # Calculate the in-sample SVD
> svdec <- svd(predm[insample, ])
> # Calculate the in-sample fitted AR coefficients for different dimensions
> dimv <- 2:5
> # dimv <- c(2, 5, 10, NCOL(predm))
> coeffm <- sapply(dimv, function(dimax) {
+   predinv <- svdec$v[, 1:dimax] %*%
+     (t(svdec$u[, 1:dimax]) / svdec$d[1:dimax])
+   predinv %*% respv[insample]
+ }) # end lapply
> colnames(coeffm) <- paste0("dimax=", dimv)
> colorv <- colorRampPalette(c("red", "blue"))(NCOL(coeffm))
> matplot(y=coeffm, type="l", lty="solid", lwd=1, col=colorv,
+   xlab="predictor", ylab="coeff",
+   main="AR Coefficients For Different Dimensions")
> # Calculate the forecasts of VTI
> fcasts <- predm %*% coeffm
> fccasts <- apply(fccasts, 2, function(x) {
+   fcavt <- sqrt(HighFreq::run_var(matrix(x), lambda=0.8)[, 2])
+   fcavt[1:100] <- 1
+   x/fcavt
+ }) # end apply
> # Simulate the autoregressive strategies
> retn <- coredata(respv)
> pnls <- apply(fccasts, 2, function(x) (x*retn))
> pnls <- xts(pnls, datev)
> # Scale the PnL volatility to that of VTI
> pnls <- lapply(pnls, function(x) x/sd(x))
> pnls <- sd(respv)*do.call(cbind, pnls)
```

**Kitchen Sink Strategies With Dimension Reduction**



```
> # Calculate the in-sample and out-of-sample Sharpe and Sortino ratios
> wealthv <- cbind(respv, pnls)
> sqrt(252)*sapply(wealthv[insample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[outsample, ], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Kitchen Sink Strategies With Dimension Reduction") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=500)
```

# Shrinkage of the Kitchen Sink Coefficients

The AR coefficients can be found by minimizing the sum of the squared errors of the in-sample forecasts.

The objective function is the sum of the squared errors of the in-sample forecasts, plus a penalty term proportional to the square of the AR coefficients:

$$\text{ObjFunc} = \sum_{i=1}^n (r_t - f_t)^2 + \sum_{i=1}^k \varphi_i^2 \lambda^i$$

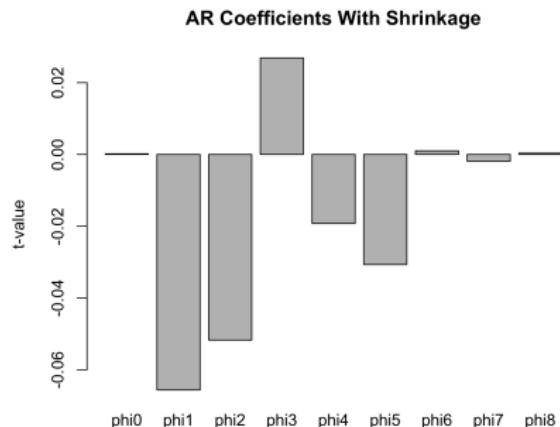
As the objective function is minimized, the *shrinkage* penalty shrinks the AR coefficients closer to zero.

A larger *shrinkage factor*  $\lambda$  applies more *shrinkage* to the AR coefficients, to shrink them closer to zero.

The *shrinkage* penalty is greater for the higher order coefficients, to reflect the fact that returns from the more distant past should be less important for forecasting the current returns.

The coefficient shrinkage also produces a dimension reduction effect, because the higher order coefficients are shrunk closer to zero.

```
> # Objective function for the in-sample AR coefficients
> objfun <- function(coeff, respv, predm, lambdaf) {
+   fcasts <- predm %*% coeff
+   10000*sum((respv - fcasts)^2) + sum(lambdaf*coeff^2)
+ } # end objfun
```



```
> # Perform optimization using the quasi-Newton method
> ncoeff <- NROW(coeff)
> optiml <- optim(par=numeric(ncoeff),
+   fn=objfun,
+   respv=respv[insample, ],
+   predm=predm[insample, ],
+   lambdaf=2.0,
+   method="L-BFGS-B",
+   upper=rep(10000, ncoeff),
+   lower=rep(-10000, ncoeff))
> # Extract the AR coefficients
> coeff <- optiml$par
> names(coeff) <- coeffn
> barplot(coeff ~ coeffn, xlab="", ylab="t-value", col="grey",
+   main="AR Coefficients With Shrinkage")
```

# Kitchen Sink Strategy With Shrinkage

The out-of-sample performance of the multifactor autoregressive strategy can be improved by applying *shrinkage* to the AR coefficients.

The value of the *shrinkage factor*  $\lambda$  can be chosen to maximize the out-of-sample performance.

```
> # Calculate the forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> fcastv <- sqrt(HighFreq::run_var(fcasts, lambda=0.8)[, 2])
> fcastv[1:100] <- 1
> fcasts <- fcasts/fcastv
> # Calculate the autoregressive strategy PnLs
> pnls <- respv*fcasts
> pnls <- pnls*sd(respv[respv<0])/sd(pnls[pnls<0])
```



```
> # Calculate the in-sample and out-of-sample Sharpe and Sortino rates
> wealthv <- cbind(respv, pnls)
> colnames(wealthv) <- c("VTI", "AR_multifact")
> sqrt(252)*sapply(wealthv[,1], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> sqrt(252)*sapply(wealthv[,2], function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategies
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(wealthv))
> dygraphs::dygraph(cumsum(wealthv[,endd]),
+   main="Multifactor Autoregressive Strategy With Shrinkage") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# draft: Multifactor Autoregressive Strategy With Features

When  $VXX$  and  $SVXY$  returns are used as predictors of  $VTI$  returns, then the coefficients are overfitted to profit from single events, like the 2018 flash crash and the 2020 pandemic crash. But the model doesn't perform well outside of those single events.

The stock returns  $r_t$  can be fitted into an *autoregressive* model  $AR(n)$  with a constant intercept term  $\varphi_0$ :

$$r_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n} + \varepsilon_t$$

The *residuals*  $\varepsilon_t$  are assumed to be normally distributed, independent, and stationary.

```
> # Calculate the VTI daily percentage returns
> retp <- na.omit(rutils::etfenv$returns[, c("VTI", "VXX", "SVXY")])
> nrows <- NROW(retp)
> # Define the response and predictor matrices
> respv <- retp["/2019", "VTI"]
> orderp <- 3
> predm <- lapply(1:orderp, rutils::lagit, input=retp["/2019", c("VXX", "SVXY")])
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, NROW(predm)), predm)
> colnames(predm) <- c("phi0", paste0(c("VXX", "SVXY"), rep(1:orderp, each=2)))
> # Calculate the fitted autoregressive coefficients
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> # Calculate the residuals (forecast errors)
> resids <- (fcasts - respv)
> # The residuals are orthogonal to the predictors and the forecasts
> round(cor(resids, fcasts), 6)
> round(sapply(predm[, -1], function(x) cor(resids, x))), 6)
```

# draft: Forecasting Stock Returns Using Autoregressive Models

Explain why the sum of the AR coefficients is negative, even though the stock returns are positive on average.

The fitted autoregressive coefficients  $\varphi$  are equal to the response  $r$  multiplied by the inverse of the predictor matrix  $P$ :

$$\varphi = P^{-1} r$$

The *in-sample* autoregressive forecasts of the returns are calculated by multiplying the predictor matrix by the fitted AR coefficients:

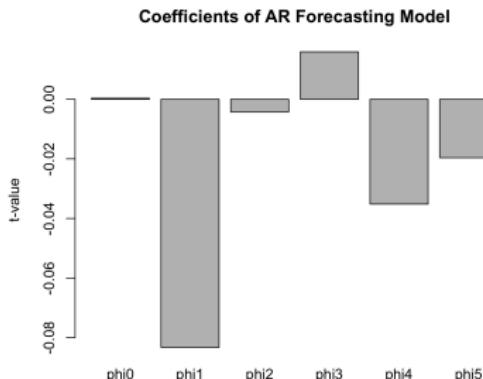
$$f_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n}$$

For *VTI* returns, the intercept coefficient  $\varphi_0$  has a small positive value, while the first autoregressive coefficient  $\varphi_1$  has a small negative value.

This means that the autoregressive forecasting model is a combination of a static long stock position, plus a mean-reverting model which switches its stock position to the reverse of the previous day's return.

The function `MASS::ginv()` calculates the generalized inverse of a matrix.

```
> # Calculate the fitted autoregressive coefficients
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> sum(coeff[1:orderp])
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
```



```
> # Plot the AR coefficients
> coefffn <- paste0("phi", 0:(NROW(coeff)-1))
> barplot(coeff ~ coefffn, xlab="", ylab="t-value", col="grey",
+ main="Coefficients of AR Forecasting Model")
```

# draft: The t-values of the Autoregressive Coefficients

The forecast residuals are equal to the differences between the return forecasts minus the actual returns:

$$\varepsilon = f_t - r_t$$

The variance of the autoregressive coefficients  $\sigma_\varphi^2$  is equal to the variance of the forecast residuals  $\sigma_\varepsilon^2$  divided by the squared predictor matrix  $\mathbb{P}$ :

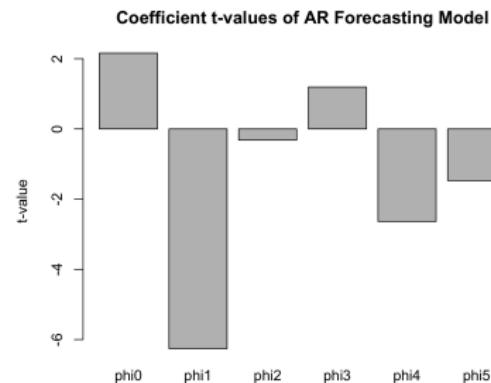
$$\sigma_\varphi^2 = \sigma_\varepsilon^2 (\mathbb{P}^T \mathbb{P})^{-1}$$

The t-values of the autoregressive coefficients are equal to the coefficient values divided by their volatilities:

$$\varphi_{tval} = \frac{\varphi}{\sigma_\varphi}$$

The intercept coefficient  $\varphi_0$  and the first autoregressive coefficient  $\varphi_1$  have statistically significant t-values.

```
> # Calculate the residuals (forecast errors)
> resids <- (fcasts - respv)
> # The residuals are orthogonal to the predictors and the forecasts
> round(cor(resids, fccasts), 6)
> round(sapply(predm[, -1], function(x) cor(resids, x)), 6)
> # Calculate the variance of the residuals
> varv <- sum(resids^2)/(nrows-NROW(coeff))
```



```
> # Calculate the predictor matrix squared
> pred2 <- crossprod(predm)
> # Calculate the covariance matrix of the AR coefficients
> covmat <- varv*MASS::ginv(pred2)
> coefsds <- sqrt(diag(covmat))
> # Calculate the t-values of the AR coefficients
> coefft <- drop(coeff/coefsds)
> coeffn <- paste0("phi", 0:(NROW(coeff)-1))
> # Plot the t-values of the AR coefficients
> barplot(coefft ~ coeffn, xlab="", ylab="t-value", col="grey",
+ main="Coefficient t-values of AR Forecasting Model")
```

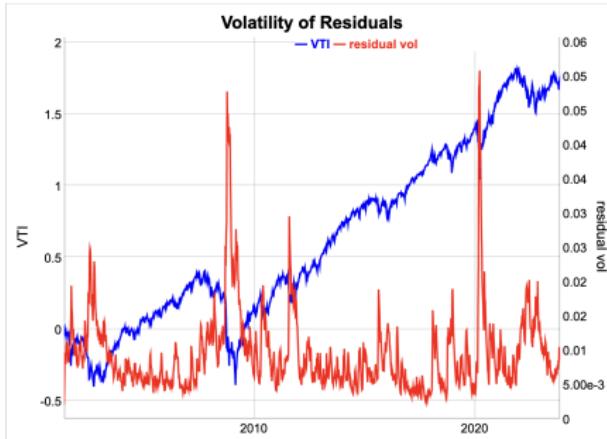
# draft: Residuals of Autoregressive Forecasting Model

The autoregressive model assumes stationary returns and residuals, with similar volatility over time.

In reality stock volatility is highly time dependent, so the volatility of the residuals is also time dependent.

The function `HighFreq::run_var()` calculates the trailing variance of a time series using exponential weights.

```
> # Calculate the trailing volatility of the residuals
> residv <- sqrt(HighFreq::run_var(resids, lambda=0.9)[, 2])
```



```
> # Plot dygraph of volatility of residuals
> datav <- cbind(cumsum(retp), residv)
> colnames(datav) <- c("VTI", "residual vol")
> endd <- rutils::calc_endpoints(datav, interval="weeks")
> dygraphs::dygraph(datav[endd], main="Volatility of Residuals") %>%
+   dyAxis("y", label="VTI", independentTicks=TRUE) %>%
+   dyAxis("y2", label="residual vol", independentTicks=TRUE) %>%
+   dySeries(name="VTI", axis="y", strokeWidth=2, col="blue") %>%
+   dySeries(name="residual vol", axis="y2", strokeWidth=2, col="red") %>%
+   dyLegend(show="always", width=300)
```

## draft: Autoregressive Strategy In-Sample

The kitchen sink strategy has more predictor variables (and more degrees of freedom), so in-sample it can forecast the returns more closely.

The first step in strategy development is optimizing it in-sample, even though in practice it can't be implemented. Because a strategy can't perform well out-of-sample if it doesn't perform well in-sample.

The autoregressive strategy invests dollar amounts of *VTI* stock proportional to the in-sample forecasts.

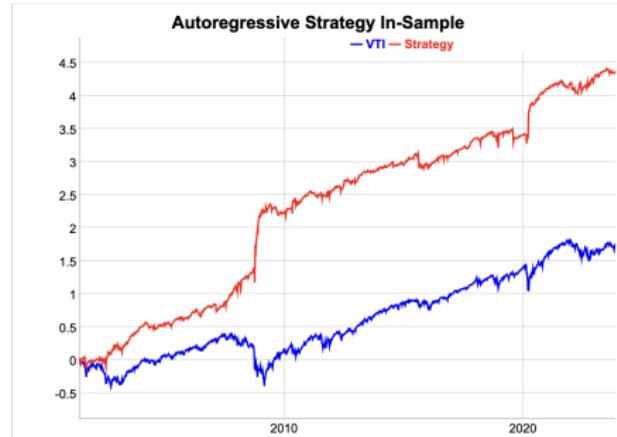
The in-sample autoregressive strategy performs well during periods of high volatility, but not as well in low volatility periods.

The dollar allocations of *VTI* stock are too large in periods of high volatility, which causes over-leverage and very high risk.

The leverage can be reduced by scaling (dividing) the forecasts by their trailing volatility.

The function `HighFreq::run_var()` calculates the trailing variance of a time series using exponential weights.

```
> # Calculate the kitchen sink strategy in-sample
> pnls <- retpl*fcasts
> pnls <- pnls*sd(retpl[retpl<0])/sd(pnls[pnls<0])
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(retpl, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutilss::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy In-Sample") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Out-of-Sample Performance of the Kitchen Sink Strategy

The out-of-sample performance of the kitchen sink strategy is much worse than its in-sample performance, because the strategy is overfit to the in-sample data.

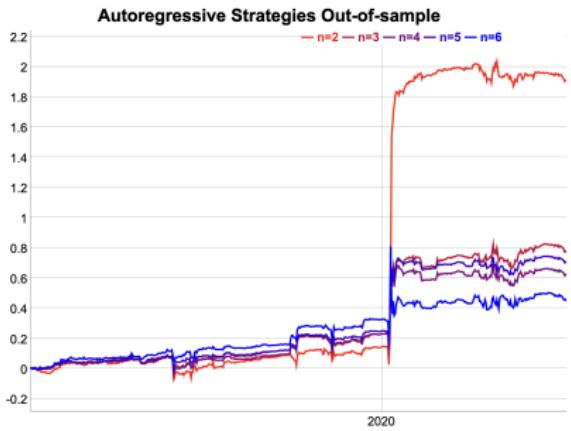
The autoregressive strategy invests a dollar amount of  $VTI$  proportional to the AR forecasts.

The out-of-sample, risk-adjusted performance of the autoregressive strategy is better for a smaller order parameter  $n$  of the  $AR(n)$  model.

The optimal order parameter is  $n = 2$ , with a positive intercept coefficient  $\varphi_0$  (since the average  $VTI$  returns were positive), and a negative coefficient  $\varphi_1$  (because of strong negative autocorrelations in periods of high volatility).

Decreasing the order parameter of the autoregressive model is a form of *shrinkage* because it reduces the number of predictive variables.

```
> # Define in-sample and out-of-sample intervals
> nrows <- NROW(retpt)
> cutoff <- nrows %/% 2
> datev[cutoff]
> insample <- 1:cutoff
> outsample <- (cutoff + 1):nrows
> # Calculate the optimal AR coefficients
> predinv <- MASS::ginv(predm[insample, ])
> coeff <- drop(predinv %*% respv[insample])
> # Calculate the strategy PnLs
> fcasts <- predm %*% coeff
> pnls <- sign(fcasts)*retpt
> wealthv <- cbind(retpt, pnls)
```



```
> # Calculate the in-sample Sharpe and Sortino ratios
> sqrt(252)*sapply(wealthv[insample, ],
+   function(x) c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Calculate the out-of-sample Sharpe and Sortino ratios
> sqrt(252)*sapply(wealthv[outsample, ],
+   function(x) c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Kitchen Sink Strategy Out-of-sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# draft: Out-of-Sample Performance With Dimension Reduction

The out-of-sample performance of the kitchen sink strategy is much worse than its in-sample performance, because the strategy is overfit to the in-sample data.

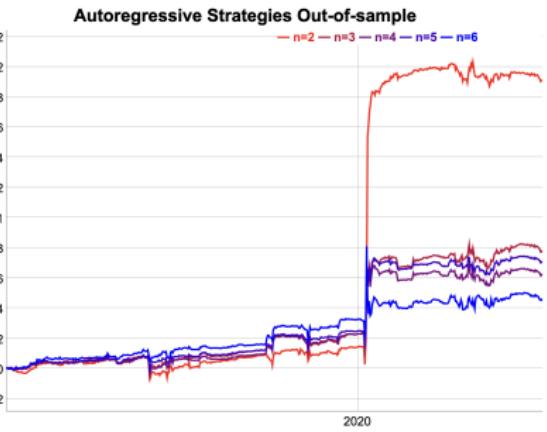
The autoregressive strategy invests a dollar amount of  $VTI$  proportional to the AR forecasts.

The out-of-sample, risk-adjusted performance of the autoregressive strategy is better for a smaller order parameter  $n$  of the  $AR(n)$  model.

The optimal order parameter is  $n = 2$ , with a positive intercept coefficient  $\varphi_0$  (since the average  $VTI$  returns were positive), and a negative coefficient  $\varphi_1$  (because of strong negative autocorrelations in periods of high volatility).

Decreasing the order parameter of the autoregressive model is a form of *shrinkage* because it reduces the number of predictive variables.

```
> # Define in-sample and out-of-sample intervals
> nrows <- NROW(rtpt)
> cutoff <- nrows %/% 2
> datev[cutoff]
> insample <- 1:cutoff
> outsample <- (cutoff + 1):nrows
> # Calculate reduced inverse of the predictor matrix from SVD
> svdec <- svd(predm[insample, ])
> dimax <- 2
> predinv <- svdec$v[, 1:dimax] %*%
+   (t(svdec$u[, 1:dimax]) / svdec$d[1:dimax])
> coeff <- drop(predinv %*% respv[insample])
> # Calculate the strategy PnLs
```



```
> # Calculate the in-sample Sharpe and Sortino ratios
> sqrt(252)*sapply(wealthv[insample, ], 
+   function(x) c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Calculate the out-of-sample Sharpe and Sortino ratios
> sqrt(252)*sapply(wealthv[outsample, ], 
+   function(x) c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Kitchen Sink Strategy Out-of-sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyEvent(datev[cutoff], label="cutoff", strokePattern="solid",
+   dyLegend(show="always", width=300)
```

# draft: Out-of-Sample Autoregressive Forecasts

The autoregressive coefficients  $\varphi$  are equal to the in-sample *response*  $r$  times the inverse of the in-sample *predictor matrix*  $\mathbb{P}$ :

$$\varphi = \mathbb{P}^{-1} r$$

The variance of the autoregressive coefficients  $\sigma_\varphi^2$  is equal to the variance of the in-sample forecast residuals  $\sigma_\varepsilon^2$  divided by the squared *predictor matrix*  $\mathbb{P}$ :

$$\sigma_\varphi^2 = \sigma_\varepsilon^2 (\mathbb{P}^T \mathbb{P})^{-1}$$

The t-values of the autoregressive coefficients are equal to the coefficient values divided by their volatilities:

$$\varphi_{tval} = \frac{\varphi}{\sigma_\varphi}$$

The *out-of-sample* autoregressive forecast  $f_t$  is equal to the single row of the predictor  $\mathbb{P}_t$  times the fitted AR coefficients  $\varphi$ :

$$f_t = \varphi \mathbb{P}_t$$

The variance  $\sigma_f^2$  of the *forecast value* is equal to the inner product of the predictor  $\mathbb{P}_t$  times the coefficient covariance matrix  $\sigma_\varphi^2$ :

$$\sigma_f^2 = \mathbb{P}_t \sigma_\varphi^2 \mathbb{P}_t^T$$

```
> # Define the look-back range
> lookb <- 100
> tday <- nrow
> startp <- max(1, tday-lookb)
> rangev <- startp:(tday-1)
> # Subset the response and predictors
> respv <- respv[rangev]
> predm <- predm[rangev]
> # Invert the predictor matrix
> predinv <- MASS::ginv(preds)
> # Calculate the fitted AR coefficients
> coeff <- predinv %*% respv
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- preds %*% coeff
> # Calculate the residuals (forecast errors)
> resids <- (fcasts - respv)
> # Calculate the variance of the residuals
> varv <- sum(resids^2)/(NROW(preds)-NROW(coeff))
> # Calculate the predictor matrix squared
> pred2 <- crossprod(preds)
> # Calculate the covariance matrix of the AR coefficients
> covmat <- varv*MASS::ginv(pred2)
> coefsds <- sqrt(diag(covmat))
> # Calculate the t-values of the AR coefficients
> coefft <- drop(coeff/coefsds)
> # Calculate the out-of-sample forecast
> predn <- predm[tday, ]
> fcast <- drop(predn %*% coeff)
> # Calculate the variance of the forecast
> varf <- drop(predn %*% covmat %*% t(predn))
> # Calculate the t-value of the out-of-sample forecast
> fcast/sqrt(varf)
```

# draft: Rolling Autoregressive Forecasting Model

The autoregressive coefficients can be calibrated dynamically over a *rolling* look-back interval, and applied to calculating the *out-of-sample* forecasts.

*Backtesting* is the simulation of a model on historical data to test its forecasting accuracy.

The autoregressive forecasting model can be *backtested* by calculating forecasts over either a *rolling* or an *expanding* look-back interval.

If the start date is fixed at the first row then the look-back interval is *expanding*.

The coefficients of the *AR(n)* process are fitted to past data, and then applied to calculating out-of-sample forecasts.

The *backtesting* procedure allows determining the optimal *meta-parameters* of the forecasting model: the order *n* of the *AR(n)* model and the length of look-back interval (*lookb*).

```
> # Perform rolling forecasting
> lookb <- 100
> fcasts <- sapply(1:nrows, function(tday) {
+   if (tday > lookb) {
+     # Define the rolling look-back range
+     startp <- max(1, tday-lookb)
+     # startp <- 1 # Expanding look-back range
+     rangev <- startp:(tday-1) # In-sample range
+     # Subset the response and predictors
+     resps <- respv[rangev]
+     preds <- predm[rangev]
+     # Calculate the fitted AR coefficients
+     predinv <- MASS::ginv(preds)
+     coeff <- predinv %*% resps
+     # Calculate the in-sample forecasts of VTI (fitted values)
+     fcasts <- preds %*% coeff
+     # Calculate the residuals (forecast errors)
+     resid <- (fcasts - resps)
+     # Calculate the variance of the residuals
+     varv <- sum(resid^2)/(NROW(preds)-NROW(coeff))
+     # Calculate the covariance matrix of the AR coefficients
+     pred2 <- crossprod(preds)
+     covmat <- varv*MASS::ginv(pred2)
+     coefs <- sqrt(diag(covmat))
+     coefft <- drop(coeff/coefs) # t-values of the AR coefficients
+     # Calculate the out-of-sample forecast
+     predn <- predm[tday, ]
+     fcast <- drop(predn %*% coeff)
+     # Calculate the variance of the forecast
+     varf <- drop(predn %*% covmat %*% t(predn))
+     return(c(sd(resps), fcast=fcast, fstderr=sqrt(varf), coefft=coefft))
+   } else {
+     return(c(volv=0, fcast=0, fstderr=0, coefft=rep(0, NCOL(predm)))
+   } # end if
+ }) # end sapply
```

# draft: Rolling Autoregressive Strategy Performance

In the rolling autoregressive strategy, the autoregressive coefficients are calibrated on past data from a *rolling* look-back interval, and applied to calculating the *out-of-sample* forecasts.

The rolling autoregressive strategy performance depends on the length of the look-back interval.

```
> # Coerce fcasts to a time series
> fccasts <- t(fccasts)
> ncols <- NCOL(fccasts)
> colnames(fccasts) <- c("volv", "fccasts", "fstderr", colnames(predm))
> fccasts <- xts::xts(fccasts, zoo::index(retp))
> # Calculate the strategy PnLs
> pnls <- retp$fccasts$fccasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+ c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
```



```
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Autoregressive Strategy") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Forecasting Returns Using Rolling Autoregressive Models

Adopt this text to some other slides.

Forecasting using an autoregressive model is performed by first fitting an  $AR(n)$  model to past data, and calculating its coefficients.

The fitted AR coefficients are then applied to calculating the *out-of-sample* forecasts.

The forecasting model depends on two unknown *meta-parameters*: the order  $n$  of the  $AR(n)$  model and the length of the look-back interval (lookb).

# draft: Backtesting the Autoregressive Model

The *meta-parameters* of the *backtesting* function are the order  $n$  of the  $AR(n)$  model and the length of the look-back interval (`lookb`).

The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

*Backtesting* is the simulation of a model on historical data to test its forecasting accuracy.

The autoregressive forecasting model can be *backtested* by calculating forecasts over either a *rolling* or an *expanding* look-back interval.

If the start date is fixed at the first row then the look-back interval is *expanding*.

The coefficients of the  $AR(n)$  process are fitted to past data, and then applied to calculating out-of-sample forecasts.

The *backtesting* procedure allows determining the optimal *meta-parameters* of the forecasting model: the order  $n$  of the  $AR(n)$  model and the length of look-back interval (`lookb`).

```
> # Define backtesting function
> sim_fcasts <- function(lookb=100, ordern=5, fixedlb=TRUE) {
+   # Perform rolling forecasting
+   fcasts <- sapply((lookb+1):nrows, function(tday) {
+     # Rolling look-back range
+     startp <- max(1, tday-lookb)
+     # Expanding look-back range
+     if (!fixedlb) {startp <- 1}
+     startp <- max(1, tday-lookb)
+     rangev <- startp:(tday-1) # In-sample range
+     # Subset the response and predictors
+     respv <- respv[rangev]
+     predm <- predm[rangev, 1:ordern]
+     # Invert the predictor matrix
+     predinv <- MASS::ginv(preds)
+     # Calculate the fitted AR coefficients
+     coeff <- predinv %*% respv
+     # Calculate the out-of-sample forecast
+     drop(predm[tday, 1:ordern] %*% coeff)
+   }) # end apply
+   # Add warmup period
+   fccasts <- c(rep(0, lookb), fccasts)
+ } # end sim_fcasts
> # Simulate the rolling autoregressive forecasts
> fccasts <- sim_fccasts(lookb=100, ordern=5)
> c(mse=mean((fccasts - retp)^2), cor=cor(retp, fccasts))
```

# draft: Forecasting Dependence On the Look-back Interval

The *backtesting* function can be used to find the optimal *meta-parameters* of the autoregressive forecasting model.

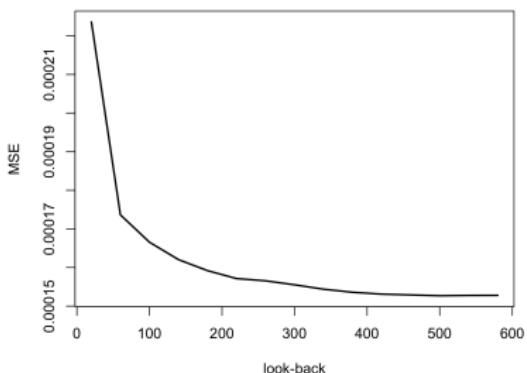
The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model and on the length of the look-back interval (`lookb`).

The accuracy of the forecasting model increases with longer look-back intervals (`lookb`), because more data improves the estimates of the autoregressive coefficients.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> lookbv <- seq(20, 600, 40)
> fcasts <- parLapply(compclust, lookbv, sim_fccasts, ordern=6)
> # Perform parallel bootstrap under Mac-OSX or Linux
> fccasts <- mclapply(lookbv, sim_fccasts, ordern=6, mc.cores=ncores)
```

**MSE of AR Forecasting Model As Function of Look-back**



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retp - x)^2), cor=cor(retp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- lookbv
> # Select optimal lookb interval
> lookb <- lookbv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=lookbv, y=mse[, 1],
+       xlab="look-back", ylab="MSE", type="l", lwd=2,
+       main="MSE of AR Forecasting Model As Function of Look-back")
```

# draft: Order Dependence With Fixed Look-back

The *backtesting* function can be used to find the optimal *meta-parameters* of the autoregressive forecasting model.

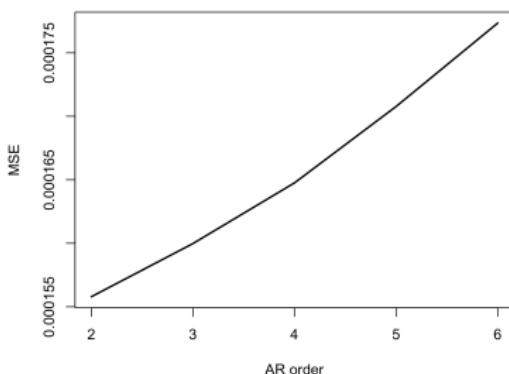
The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model and on the length of the look-back interval ( $lookb$ ).

The accuracy of the forecasting model decreases for larger AR order parameters, because of overfitting in-sample.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> orderv <- 2:6
> fcasts <- parLapply(compclust, orderv, sim_fcasts, lookb=lookb)
> stopCluster(compclust) # Stop R processes over cluster under Wi
> # Perform parallel bootstrap under Mac-OSX or Linux
> fcasts <- mclapply(orderv, sim_fcasts,
+   lookb=lookb, mc.cores=ncores)
```

**MSE of Forecasting Model As Function of AR Order**



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retp - x)^2), cor=cor(retp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- orderv
> # Select optimal order parameter
> ordern <- orderv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=orderv, y=mse[, 1],
+   xlab="AR order", ylab="MSE", type="l", lwd=2,
+   main="MSE of Forecasting Model As Function of AR Order")
```

# draft: Autoregressive Strategy With Fixed Look-back

The return forecasts are calculated just before the close of the markets, so that trades can be executed before the close.

The autoregressive strategy returns are large in periods of high volatility, but much smaller in periods of low volatility. This because the forecasts are bigger in periods of high volatility, and also because the forecasts are more accurate, because the autocorrelations of stock returns are much higher in periods of high volatility.

Using the return forecasts as portfolio weights produces very large weights in periods of high volatility, and creates excessive risk.

To reduce excessive risk, a binary strategy can be used, with portfolio weights equal to the sign of the forecasts.

```
> # Simulate the rolling autoregressive forecasts
> fcasts <- sim_fcasts(lookb=lookb, ordern=ordern)
> # Calculate the strategy PnLs
> pnls <- fcasts*retp
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls, (retp+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
```



```
> # Annualized Sharpe ratios of VTI and AR strategy
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of AR strategy combined with VTI
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy Fixed Look-back") %>%
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dySeries(name="Combined", strokeWidth=3) %>%
+   dyLegend(show="always", width=300)
```

# draft: Order Dependence With Expanding Look-back

The two *meta-parameters* can be chosen by minimizing the *MSE* of the model forecasts in a *backtest* simulation.

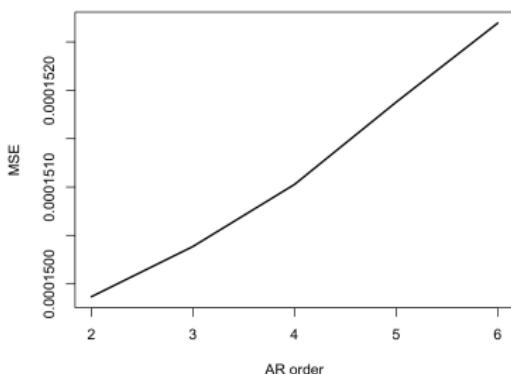
The accuracy of the forecasting model depends on the order  $n$  of the  $AR(n)$  model.

Longer look-back intervals (`lookb`) are usually better for the autoregressive forecasting model.

The return forecasts are calculated just before the close of the markets, so that trades can be executed before the close.

```
> library(parallel) # Load package parallel
> # Calculate the number of available cores
> ncores <- detectCores() - 1
> # Initialize compute cluster under Windows
> compclust <- makeCluster(ncores)
> # Perform parallel loop under Windows
> orderv <- 2:6
> fcasts <- parLapply(compclust, orderv, sim_fcasts,
+   lookb=lookb, fixedlb=FALSE)
> stopCluster(compclust) # Stop R processes over cluster under Wi
> # Perform parallel bootstrap under Mac-OSX or Linux
> fcasts <- mclapply(orderv, sim_fcasts,
+   lookb=lookb, fixedlb=FALSE, mc.cores=ncores)
```

**MSE With Expanding Look-back As Function of AR Order**



```
> # Calculate the mean squared errors
> mse <- sapply(fccasts, function(x) {
+   c(mse=mean((retlp - x)^2), cor=cor(retlp, x))
+ }) # end sapply
> mse <- t(mse)
> rownames(mse) <- orderv
> # Select optimal order parameter
> ordern <- orderv[which.min(mse[, 1])]
> # Plot forecasting MSE
> plot(x=orderv, y=mse[, 1],
+   xlab="AR order", ylab="MSE", type="l", lwd=2,
+   main="MSE With Expanding Look-back As Function of AR Order")
```

# draft: Autoregressive Strategy With Expanding Look-back

The model with an *expanding* look-back interval has better performance compared to the *fixed* look-back interval.

The autoregressive strategy returns are large in periods of high volatility, but much smaller in periods of low volatility. This because the forecasts are bigger in periods of high volatility, and also because the forecasts are more accurate, because the autocorrelations of stock returns are much higher in periods of high volatility.

```
> # Simulate the autoregressive forecasts with expanding look-back
> fcasts <- sim_fcasts(lookb=lookb, ordern=ordern, fixedlb=FALSE)
> # Calculate the strategy PnLs
> pnls <- fcasts*retp
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(retp[retp<0])/sd(pnls[pnls<0])
> wealthv <- cbind(retp, pnls, (retp+pnls)/2)
> colnames(wealthv) <- c("VTI", "AR_Strategy", "Combined")
> cor(wealthv)
```



```
> # Annualized Sharpe ratios of VTI and AR strategy
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of AR strategy combined with VTI
> dygraphs::dygraph(cumsum(wealthv)[endd]),
+   main="Autoregressive Strategy Expanding Look-back") %>%
+   dyOptions(colors=c("blue", "red", "green"), strokeWidth=1) %>%
+   dySeries(name="Combined", strokeWidth=3) %>%
+   dyLegend(show="always", width=300)
```

# draft: Multifactor Autoregressive Strategy of Stock Returns

When  $VXX$  and  $SVXY$  returns are used as predictors of  $VTI$  returns, then the coefficients are overfitted to profit from single events, like the 2018 flash crash and the 2020 pandemic crash. But the model doesn't perform well outside of those single events.

The stock returns  $r_t$  can be fitted into an *autoregressive* model  $AR(n)$  with a constant intercept term  $\varphi_0$ :

$$r_t = \varphi_0 + \varphi_1 r_{t-1} + \varphi_2 r_{t-2} + \dots + \varphi_n r_{t-n} + \varepsilon_t$$

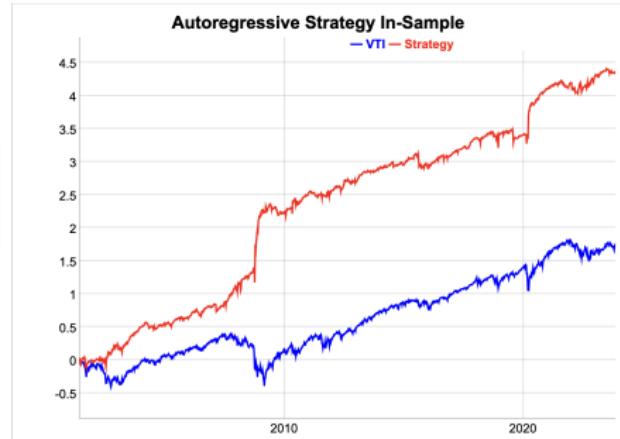
The *residuals*  $\varepsilon_t$  are assumed to be normally distributed, independent, and stationary.

```
> # Calculate the VTI daily percentage returns
> retp <- na.omit(rutils::etfenv$returns[, c("VTI", "VXX", "SVXY")])
> nrows <- NROW(retp)
> # Define the response and predictor matrices
> respv <- retp["/2019", "VTI"]
> orderp <- 3
> predm <- lapply(1:orderp, rutils::lagit, input=retp["/2019", c("VXX", "SVXY")])
> predm <- rutils::do_call(cbind, predm)
> predm <- cbind(rep(1, NROW(predm)), predm)
> colnames(predm) <- c("phi0", paste0(c("VXX", "SVXY"), rep(1:orderp, each=2)))
> # Calculate the fitted autoregressive coefficients
> predinv <- MASS::ginv(predm)
> coeff <- predinv %*% respv
> # Calculate the in-sample forecasts of VTI (fitted values)
> fcasts <- predm %*% coeff
> # Calculate the residuals (forecast errors)
> resids <- (fcasts - respv)
> # The residuals are orthogonal to the predictors and the forecasts
> round(cor(resids, fcasts), 6)
> round(sapply(predm[, -1], function(x) cor(resids, x))), 6)
```

# draft: Multifactor Autoregressive Strategy In-Sample

The coefficients are calibrated with returns before 2020, but the model still makes profit during the 2020 stock crash.

```
> # Calculate the autoregressive strategy PnLs
> pnls <- respv*fcasts
> # Scale the PnL volatility to that of VTI
> pnls <- pnls*sd(respv)/sd(pnls)
```



```
> # Calculate the Sharpe and Sortino ratios
> wealthv <- cbind(respv, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> sqrt(252)*sapply(wealthv, function(x)
+   c(Sharpe=mean(x)/sd(x), Sortino=mean(x)/sd(x[x<0])))
> # Plot dygraph of the autoregressive strategy
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Autoregressive Strategy In-Sample") %>%
+   dyOptions(colors=c("blue", "red"), strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Interest Rate Yield Curve and Stock Returns

Daily stock returns have insignificant correlations with the daily changes in interest rates, with the possible exception of the 10-year bond yield.

And these correlations change significantly over time.

```
> # Load constant maturity Treasury rates
> load(file="/Users/jerzy/Develop/lecture_slides/data/rates_data.R"
> # Combine rates into single xts series
> ratev <- do.call(cbind, as.list(ratesenv))
> # Sort the columns of rates according bond maturity
> namev <- colnames(ratev)
> namev <- substr(namev, start=4, stop=10)
> namev <- as.numeric(namev)
> indeks <- order(namev)
> ratev <- ratev[, indeks]
> # Align rates dates with VTI prices
> closep <- log(quantmod::Cl(rutils::etfenv$VTI))
> colnames(closep) <- "VTI"
> nrows <- NROW(closep)
> datev <- zoo::index(closep)
> ratev <- na.omit(ratev[datev])
> closep <- closep[zoo::index(ratev)]
> datev <- zoo::index(closep)
```

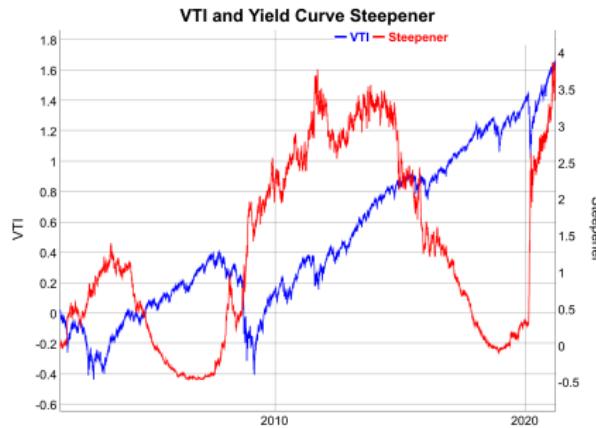
```
> # Calculate VTI returns and IR changes
> retp <- rutils::diffit(log(closep))
> retr <- rutils::diffit(log(ratev))
> # Regress VTI returns versus the lagged rate differences
> predm <- rutils::lagit(retr)
> regmod <- lm(retp ~ predm)
> summary(regmod)
> # Regress VTI returns before and after 2012
> summary(lm(retp["/2012"] ~ predm["/2012"]))
> summary(lm(retp["2012/] ~ predm["2012/"]))
```

# Yield Curve Principal Components and Stock Returns

The principal components of the interest rate yield curve can also be used as predictors of stock indices.

The second principal component describes the steepening and flattening of the yield curve, and it's an indicator of investor risk appetite. So it's also related to bullish and bearish market periods.

```
> # Calculate PCA of rates correlation matrix
> eigend <- eigen(cor(retr))
> pcar <- -(retr %*% eigend$vectors)
> colnames(pcar) <- paste0("PC", 1:6)
> # Define predictor as the YC PCAs
> predm <- rutils::lagit(pcar)
> regmod <- lm(retp ~ predm)
> summary(regmod)
```



```
> # Plot YC steepener principal component with VTI
> datav <- cbind(retp, pcar[, 2])
> colnames(datav) <- c("VTI", "Steepener")
> colv <- colnames(datav)
> dygraphs::dygraph(cumsum(datav),
+   main="VTI and Yield Curve Steepener") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", strokeWidth=2, col="blue") %>%
+   dySeries(name=colv[2], axis="y2", strokeWidth=2, col="red") %>%
+   dyLegend(show="always", width=300)
```

# Yield Curve Strategy In-Sample

For in-sample forecasts, the training set and the test set are the same. The model is calibrated on the data that is used for forecasting.

## Yield Curve Strategy

Although it's not realistic to achieve the in-sample performance, it's useful because it provides insights into how the model works.

The in-sample strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define predictor with intercept term
> predm <- rutils::lagit(retr)
> predm <- cbind(rep(1, NROW(predm)), predm)
> colnames(predm)[1] <- "intercept"
> # Calculate inverse of predictor
> invreg <- MASS::ginv(predm)
> # Calculate coefficients from response and inverse of predictor
> respv <- retr
> coeff <- drop(invreg %*% respv)
> # Calculate forecasts and PnLs in-sample
> fcasts <- (predm %*% coeff)
> pnls <- sign(fccasts)*respv
> # Calculate in-sample factors
> factv <- (predm*coeff)
> apply(factv, 2, sd)
```



```
> # Plot dygraph of in-sample IR strategy
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> endd <- rutils::calc_endpoints(wealthv, interval="weeks")
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Yield Curve Strategy In-sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Yield Curve Strategy Out-of-Sample

For out-of-sample forecasts, the training set and the test set are separate. The model is calibrated on the training data, and forecasts are calculated using the test data.

The out-of-sample strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define in-sample and out-of-sample intervals
> insample <- (datev < as.Date("2020-01-01"))
> outsample <- (datev >= as.Date("2020-01-01"))
> # Calculate inverse of predictor in-sample
> invreg <- MASS::ginv(predm[insample, ])
> # Calculate coefficients in-sample
> coeff <- drop(invreg %*% respv[insample, ])
> # Calculate forecasts and PnLs out-of-sample
> fcasts <- (predm[outsample, ] %*% coeff)
> pnls <- sign(fcasts)*respv[outsample, ]
```



```
> # Plot dygraph of out-of-sample IR PCA strategy
> wealthv <- cbind(retp[outsample, ], pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv),
+   main="Yield Curve Strategy Out-of-Sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Rolling Yearly Yield Curve Strategy

In the rolling yearly yield curve strategy, the model is recalibrated at the end of every year using a training set of data from the past year. The coefficients are applied to calculate out-of-sample forecasts in the following year.

The rolling yearly strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define yearly dates
> format(datev[1], "%Y")
> yearv <- paste0(seq(2001, 2022, 1), "-01-01")
> yearv <- as.Date(yearv)
> # Perform loop over yearly dates
> pnls <- lapply(3:(NROW(yearv)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > yearv[tday-1]) & (datev < yearv[tday])
+   outsample <- (datev >= yearv[tday]) & (datev < yearv[tday+1])
+   # Calculate coefficients in-sample
+   invreg <- MASS::ginv(predm[insample, ])
+   coeff <- drop(invreg %*% respv[insample, ])
+   # Calculate forecasts and PnLs out-of-sample
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fccasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling yearly IR strategy
> vti <- rutils::diffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Yearly Yield Curve Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Rolling Monthly Yield Curve Strategy

In the rolling monthly yield curve strategy, the model is recalibrated at the end of every month using a training set of the past 11 months. The coefficients are applied to perform out-of-sample forecasts in the following month.

Research shows that looking back roughly a year provides the best out-of-sample forecasts.

The rolling monthly strategy performs better than the yearly strategy, but mostly in periods of high volatility, and otherwise it's flat.

```
> # Define monthly dates
> format(datev[1], "%m-%Y")
> format(datev[NROW(datev)], "%m-%Y")
> months <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-30"), by="month")
> # Perform loop over monthly dates
> pnls <- lapply(12:(NROW(months)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > months[tday-1]) & (datev < months[tday])
+   outsample <- (datev > months[tday]) & (datev < months[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invreg <- MASS::ginv(predm[insample, ])
+   coeff <- drop(invreg %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fccasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling monthly IR strategy
> vti <- rutils::diffit(clossep[zoob::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Monthly Yield Curve Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Rolling Weekly Yield Curve Strategy

In the rolling weekly yield curve strategy, the model is recalibrated at the end of every week using a training set of the past 10 weeks. The coefficients are applied to perform out-of-sample forecasts in the following week.

```
> # Define weekly dates
> weeks <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-
> # Perform loop over weekly dates
> pnls <- lapply(51:(NROW(weeks)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > weeks[tday-10]) & (datev < weeks[tday])
+   outsample <- (datev > weeks[tday]) & (datev < weeks[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invreg <- MASS::ginv(predm[insample, ])
+   coeff <- drop(invreg %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fcasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling weekly IR strategy
> vti <- rutils::difffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Weekly Yield Curve Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Yield Curve Strategy With Shrinkage In-Sample

The technique of *regularization* is designed to reduce the number of parameters in a model, for example in portfolio optimization.

Regularization of the inverse predictor matrix improves the in-sample performance of the yield curve strategy.

Although it's not realistic to achieve the in-sample performance, it's useful because it provides insights into how the model can be improved.

```
> # Calculate in-sample pnls for different dimax values
> dimv <- 2:7
> pnls <- lapply(dimv, function(dimax) {
+   invred <- HighFreq::calc_invsvd(predm, dimax=dimax)
+   coeff <- drop(invred %*% respv)
+   fcasts <- (predm %*% coeff)
+   sign(fcasts)*respv
+ })
> pnls <- do.call(cbind, pnls)
> colnames(pnls) <- paste0("eigen", dimv)
```



```
> # Plot dygraph of in-sample pnls
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(pnls))
> dygraphs::dygraph(cumsum(pnls), main="In-Sample Returns of Shrinkage YC Strategies")
+ dyOptions(colors=colorv, strokeWidth=2) %>%
+ dyLegend(show="always", width=300)
```

# Yield Curve Strategy With Shrinkage Out-of-Sample

For out-of-sample forecasts, the training set and the test set are separate. The model is calibrated on the training data, and forecasts are calculated using the test data.

The out-of-sample strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define in-sample and out-of-sample intervals
> insample <- (datev < as.Date("2020-01-01"))
> outsample <- (datev >= as.Date("2020-01-01"))
> # Calculate in-sample pnls for different dimax values
> dimv <- 2:7
> pnls <- lapply(dimv, function(dimax) {
+   invred <- HighFreq::calc_invsvd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fcasts)*respv[outsample, ]
+ })
> pnls <- do.call(cbind, pnls)
> colnames(pnls) <- paste0("eigen", dimv)
```

Out-of-Sample Returns of Shrinkage YC Strategies



```
> # Plot dygraph of out-of-sample pnls
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(pnls))
> dygraphs::dygraph(cumsum(pnls), main="Out-of-Sample Returns of Shrinkage YC Strategies") %>%
+   dyOptions(colors=colorv, strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Rolling Monthly Yield Curve Strategy With Dimension Reduction

The shrinkage rolling monthly strategy performs better than the standard strategy because regularization allows using shorter lookback intervals since it suppresses the response noise.

In the rolling monthly yield curve strategy, the model is recalibrated at the end of every month using a training set of the past 6 months. The coefficients are applied to perform out-of-sample forecasts in the following month.

```
> # Define monthly dates
> format(datev[1], "%m-%Y")
> format(datev[NROW(datev)], "%m-%Y")
> months <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-30"))
> # Perform loop over monthly dates
> lookb <- 6
> dimax <- 3
> pnls <- lapply((lookb+1):(NROW(months)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > months[tday-lookb]) & (datev < months[tday])
+   outsample <- (datev > months[tday]) & (datev < months[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invred <- HighFreq::calc_insvsd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fcasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling monthly IR strategy
> vti <- rutils::diffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Monthly Shrinkage YC Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# Rolling Weekly Yield Curve Strategy With Shrinkage

In the rolling weekly yield curve strategy, the model is recalibrated at the end of every week using a training set of the past 4 weeks. The coefficients are applied to perform out-of-sample forecasts in the following week.

```
> # Define weekly dates
> weeks <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-
> # Perform loop over weekly dates
> lookb <- 4
> dimax <- 4
> pnls <- lapply((lookb+1):(NROW(weeks)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > weeks[tday-lookb]) & (datev < weeks[tday])
+   outsample <- (datev > weeks[tday]) & (datev < weeks[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invred <- HighFreq::calc_insvsd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fcasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling weekly IR strategy
> vti <- rutils::diffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Weekly Shrinkage YC Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

## draft: Combined Predictor Matrix

A "kitchen sink" strategy combines many different predictors into a large predictor matrix with many columns.

For example by combining the yield curve predictors with the lagged returns.

```
> # Load the yield curve data
> load(file="/Users/jerzy/Develop/lecture_slides/data/rates_data.RData")
> ratev <- do.call(cbind, as.list(ratesenv))
> namev <- colnames(ratev)
> namev <- substr(namev, start=4, stop=10)
> namev <- as.numeric(namev)
> indeks <- order(namev)
> ratev <- ratev[, indeks]
> closep <- log(quantmod::Cl(rutils::etfenv$VTI))
> colnames(closep) <- "VTI"
> nrow <- NROW(closep)
> datev <- zoo::index(closep)
> ratev <- na.omit(ratev[datev])
> closep <- closep[zoo::index(ratev)]
> datev <- zoo::index(closep)
> retr <- rutils::diffit(log(closep))
> retr <- rutils::diffit(log(ratev))
> # Create a combined predictor matrix
> dimax <- 5
> predm <- sapply(1:dimax, rutils::lagit, input=as.numeric(retr))
> colnames(predm) <- paste0("retslag", 1:NCOL(predm))
> predm <- cbind(predm, rutils::lagit(retr))
> predm <- cbind(rep(1, NROW(predm)), predm)
> colnames(predm)[1] <- "intercept"
> respv <- retr
```

# draft: Combined Strategy With Shrinkage In-Sample

The technique of *regularization* is designed to reduce the number of parameters in a model, for example in portfolio optimization.

Regularization of the inverse predictor matrix improves the in-sample performance of the yield curve strategy.

Although it's not realistic to achieve the in-sample performance, it's useful because it provides insights into how the model can be improved.

```
> # Calculate in-sample pnls for different dimax values
> dimv <- 2:11
> pnls <- lapply(dimv, function(dimax) {
+   invred <- HighFreq::calc_invsvd(predm, dimax=dimax)
+   coeff <- drop(invred %*% respv)
+   fcasts <- (predm %*% coeff)
+   sign(fcasts)*respv
+ })
> pnls <- do.call(cbind, pnls)
> colnames(pnls) <- paste0("eigen", dimv)
```

In-Sample Returns of Combined Strategies With Shrinkage



```
> # Plot dygraph of in-sample pnls
> colrv <- colorRampPalette(c("blue", "red"))(NCOL(pnls))
> dygraphs::dygraph(cumsum(pnls), main="In-Sample Returns of Combined Strategies With Shrinkage")
+ dyOptions(colors=colrv, strokeWidth=2) %>%
+ dyLegend(show="always", width=300)
```

# draft: Combined Strategy With Shrinkage Out-of-Sample

For out-of-sample forecasts, the training set and the test set are separate. The model is calibrated on the training data, and forecasts are calculated using the test data.

The out-of-sample strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define in-sample and out-of-sample intervals
> insample <- (datev < as.Date("2020-01-01"))
> outsample <- (datev >= as.Date("2020-01-01"))
> # Calculate in-sample pnls for different dimax values
> dimv <- 2:11
> pnls <- lapply(dimv, function(dimax) {
+   invred <- HighFreq::calc_invsvd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fcasts)*respv[outsample, ]
+ })
> pnls <- do.call(cbind, pnls)
> colnames(pnls) <- paste0("eigen", dimv)
```

Out-of-Sample Returns of Combined Strategies With Shrinkage



```
> # Plot dygraph of out-of-sample pnls
> colorv <- colorRampPalette(c("blue", "red"))(NCOL(pnls))
> dygraphs::dygraph(cumsum(pnls), main="Out-of-Sample Returns of Com")
+ dyOptions(colors=colorv, strokeWidth=2) %>%
+ dyLegend(show="always", width=300)
```

# draft: Rolling Monthly Combined Strategy With Dimension Reduction

The shrinkage rolling monthly strategy performs better than the standard strategy because regularization allows using shorter lookback intervals since it suppresses the response noise.

In the rolling monthly yield curve strategy, the model is recalibrated at the end of every month using a training set of the past 6 months. The coefficients are applied to perform out-of-sample forecasts in the following month.

```
> # Define monthly dates
> format(datev[1], "%m-%Y")
> format(datev[NROW(datev)], "%m-%Y")
> months <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-30"))
> # Perform loop over monthly dates
> lookb <- 6
> dimax <- 3
> pnls <- lapply((lookb+1):(NROW(months)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > months[tday-lookb]) & (datev < months[tday])
+   outsample <- (datev > months[tday]) & (datev < months[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invred <- HighFreq::calc_invsd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fccasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling monthly IR strategy
> vti <- rutils::diffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Monthly Shrinkage YC Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Rolling Weekly Combined Strategy With Shrinkage

In the rolling weekly yield curve strategy, the model is recalibrated at the end of every week using a training set of the past 4 weeks. The coefficients are applied to perform out-of-sample forecasts in the following week.

```
> # Define weekly dates
> weeks <- seq.Date(from=as.Date("2001-05-01"), to=as.Date("2021-04-
> # Perform loop over weekly dates
> lookb <- 8
> dimax <- 4
> pnls <- lapply((lookb+1):(NROW(weeks)-1), function(tday) {
+   # Define in-sample and out-of-sample intervals
+   insample <- (datev > weeks[tday-lookb]) & (datev < weeks[tday])
+   outsample <- (datev > weeks[tday]) & (datev < weeks[tday+1])
+   # Calculate forecasts and PnLs out-of-sample
+   invred <- HighFreq::calc_invsd(predm[insample, ], dimax=dimax)
+   coeff <- drop(invred %*% respv[insample, ])
+   fcasts <- (predm[outsample, ] %*% coeff)
+   sign(fccasts)*respv[outsample, ]
+ }) # end lapply
> pnls <- do.call(rbind, pnls)
```



```
> # Plot dygraph of rolling weekly IR strategy
> vti <- rutils::diffit(clossep[zoo::index(pnls),])
> wealthv <- cbind(vti, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Rolling Weekly Shrinkage YC Strategy") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Forecasts Using Aggregated Predictor

Needs more work to improve performance

Aggregating the predictor reduces its noise and increases the significance of correlations.

The optimal aggregation number can be found by maximizing the regression t-values.

```
> # Find optimal nagg for predictor
> nags <- 5:100
> tvalues <- sapply(nags, function(nagg) {
+   predm <- HighFreq::roll_mean(retr, lookb=nagg)
+   predm <- cbind(rep(1, NROW(predm)), predm)
+   predm <- rutils::lagit(predm)
+   regmod <- lm(respv ~ predm - 1)
+   regsum <- summary(regmod)
+   max(abs(regsum$coefficients[, 3][-1]))
+ }) # end sapply
> nagg[which.max(tvalues)]
> plot(nags, tvalues, t="l", col="blue", lwd=2)
> # Calculate aggregated predictor
> nagg <- 53
> predm <- HighFreq::roll_mean(retr, lookb=nagg)
> predm <- rutils::lagit(predm)
> predm <- cbind(rep(1, NROW(predm)), predm)
> regmod <- lm(respv ~ predm - 1)
> summary(regmod)
```



```
> # Calculate forecasts and PnLs in-sample
> invreg <- MASS::ginv(predm)
> coeff <- drop(invreg %*% respv)
> fcasts <- (predm %*% coeff)
> pnls <- sign(fcasts)*respv
> # Plot dygraph of in-sample IR strategy
> wealthv <- cbind(retp, pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Aggregated YC Strategy In-sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```

# draft: Aggregated Forecasts Out-of-Sample

Needs more work to improve performance

For out-of-sample forecasts, the training set and the test set are separate. The model is calibrated on the training data, and forecasts are calculated using the test data.

The out-of-sample strategy performs well in periods of high volatility, but otherwise it's flat.

```
> # Define in-sample and out-of-sample intervals
> insample <- (datev < as.Date("2020-01-01"))
> outsample <- (datev >= as.Date("2020-01-01"))
> # Calculate forecasts and PnLs out-of-sample
> invreg <- MASS::ginv(predm[insample, ])
> coeff <- drop(invreg %*% respv[insample, ])
> fcasts <- (predm[outsample, ] %*% coeff)
> pnls <- sign(fcasts)*respv[outsample, ]
```



```
> # Plot dygraph of out-of-sample YC strategy
> wealthv <- cbind(retp[outsample, ], pnls)
> colnames(wealthv) <- c("VTI", "Strategy")
> colv <- colnames(wealthv)
> dygraphs::dygraph(cumsum(wealthv)[endd],
+   main="Aggregated YC Strategy Out-of-Sample") %>%
+   dyAxis("y", label=colv[1], independentTicks=TRUE) %>%
+   dyAxis("y2", label=colv[2], independentTicks=TRUE) %>%
+   dySeries(name=colv[1], axis="y", col="blue", strokeWidth=2) %>%
+   dySeries(name=colv[2], axis="y2", col="red", strokeWidth=2) %>%
+   dyLegend(show="always", width=300)
```