# Big Data Predictions of Stock Prices in R

## 0. What we want to do

This project aims to introduce a variety of methods for predicting stock prices.
The main tool used is R with packages *forecast* and *fpp,* downloadable on http://cran.r-project.org/web/packages/.

Predicting techniques used and briefly explained are:

- Holt Winters Filtering

- ARIMA model

- Neural Networks

- Pynomial & Linear generalization trends

All of the above are used in a practical demo and finally evaluated on 30 stocks using the
**Mean Absolute Error**.

As a prerequisite, an understanding of basic Machine Learning techniques and methods is
recommended – if this is not the case, I attached sources for more research if you are
interested.

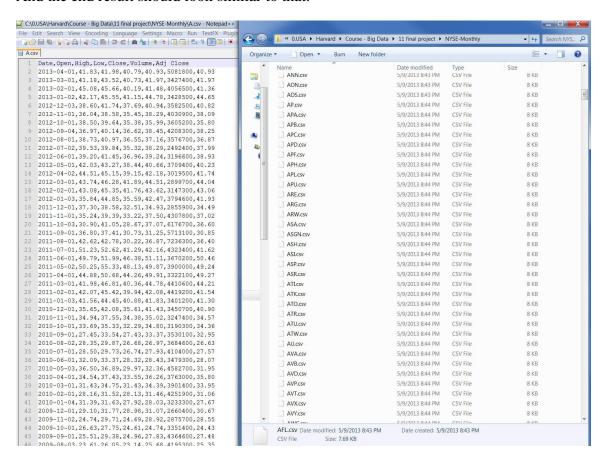The workflow of this project is the following:

- Acquiring all or some stock prices from http://finance.yahoo.com using a *wget*
  script

- Creating three functions of the stock price, on which later predictions are based:
  those are the actual stock graph, as well as two generalized (trend) graphs of it

- Getting familiar with the prediction methods while plotting from attached
  *SCRIPT.R*

- Computing in total 24 different predictions of the stock prices with attached
  *findBestPrediction.R*

- Conducting a test evaluation by using 30 stocks and attached *runTestEvaluation.R*

- Looking at the result and setting a future goal (which you can try for yourself if
  you are interested)

## 1. Getting the data

Before we can start, we need to get the stock prices in CSV file.
I attached a script (getStock.sh) which will download (wget) all stocks (>3000) on
NYSE, monthly prices starting from January 2000.



Here is a snippet of how the download process should look like:

And the end result should look similar to that:



## 2. Trying out different approaches

In order to find a suitable technique for predicting the stock prices accurately, I comprised a pool of 25 prediction techniques which will later compete with each other.

Five main techniques (and variants of those) are each used on the **actual train function** itself and also on two different generalized functions (**trends**).

*In the following I will use the attached script SCRIPT.r with the stock price ABC (AmerisourceBergen Corporation) for demonstration*
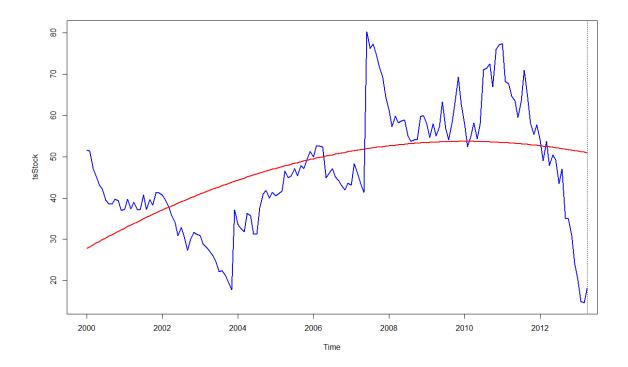
Because the actual functions of stock prices can sometimes be chaotic and misleading for the predictions, I wanted to create two generalizations (trends) of the actual function and also run the prediction techniques on those.

I began by creating a polynomial model of the function by calling:

```
tl = seq(2000,2013,length=length(tsStock))
tl2 = tl^3
polyStock = lm(tsStock ~ tl + tl2)
tsStocktrend1=ts(polyStock$fit,start=c(2000, 1),frequency=12)
```

*tsStock* is here the time series object of the actual graph.
The function *lm* creates the polynomial model (line 3),
which result is then again converted into a time series object (line4)

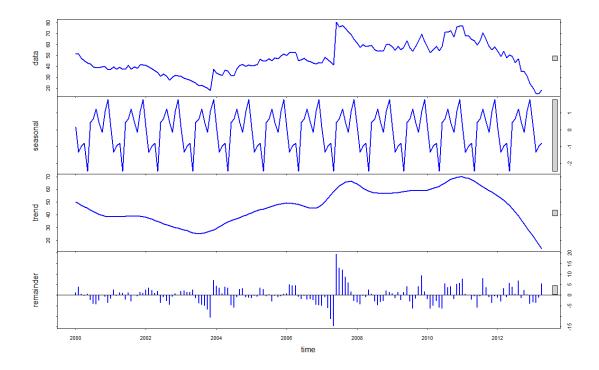The resulting generalized functions can be seen here (red):



Now that we have one generalized graph, let's create another one based on the *loess window for seasonal extraction* called *STL (Seasonal,Trend, Loess)*:

```
stlStock = stl(train,s.window="periodic")
tsStocktrend2 = stlStock$time.series[,2]
```

For more information for the stl function, visit:
http://stat.ethz.ch/R-manual/R-devel/library/stats/html/stl.html

Here you can see the analysis of our stock price graph – we use the **trend component** as our second trend base function (tsStocktrend2).

Let's have an overview of our base functions, before we begin applying the prediction methods on them:



First we begin by using the most general graph (Polynomial), which is marked purple above.

We compute the predictions using several approaches.

- **Holt Winter Filtering**
  (http://stat.ethz.ch/R-manual/R-patched/library/stats/html/HoltWinters.html)
  In this method the predictions with the minimized squared error are selected
  Here two variants are used, the gamma parameter = false means exponential smoothing is used.

```
HWStock1_ng = HoltWinters(tsStocktrend1,gamma=FALSE)
HWStock1 = HoltWinters(tsStocktrend1)
```

- **Neural Networks**
  (http://www.inside-r.org/packages/cran/forecast/docs/nnetar)
  Here a single hidden layer neural network with type feed-forward is used to predict the future prices

```
NETfit1 <- nnetar(tsStocktrend1)
```

The neural network used in this demonstration averaged 20 networks, each of which consisting of a 2-2-1 network with 9 weights - sigma^2 was estimated as 0.0007721

- **ARIMA model**
   (http://stat.ethz.ch/R-manual/R-patched/library/stats/html/arima.html)
   This approach uses an **Auto**regressive **I**ntegrated **M**oving **A**verage model as a generalization of the ARMA model which are used mainly for non-stationarity function predictions.
   Two variants are used, one determining the parameter automatically and the other with a fixed parameter approach.

```
autofit1 = auto.arima(tsStocktrend1)
fit12 <- arima(tsStocktrend1, order=c(1,0,0), list(order=c(2,1,0), period=12))
```

- **Linear Model**
   (http://www.inside-r.org/packages/cran/forecast/docs/tslm)
   This TSLM uses a **l**inear **m**odel as a base and adds **t**rend and **s**eason components to the computation

```
fitl1 <- tslm(tsStocktrend1 ~ trend + season, lambda=0)
```
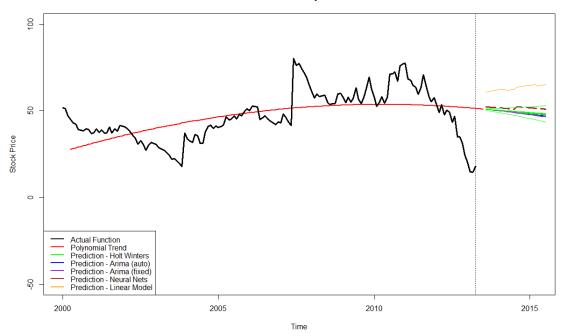
- **STL Model**
   (http://stat.ethz.ch/R-manual/R-devel/library/stats/html/stl.html)
   The Seasonal Trend Loess model, which was introduced above.
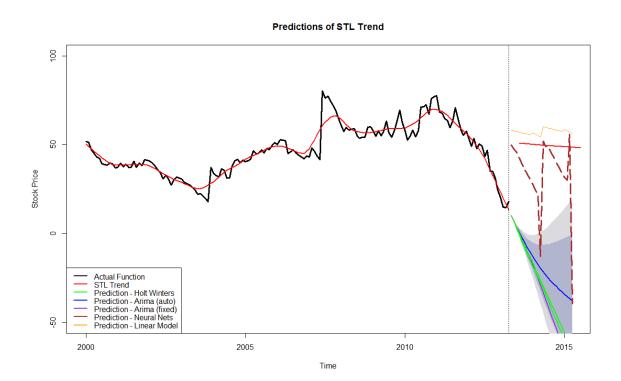
```
stlStock1 = stl(tsStocktrend1,s.window="periodic")
```

Now let's see the predictions:
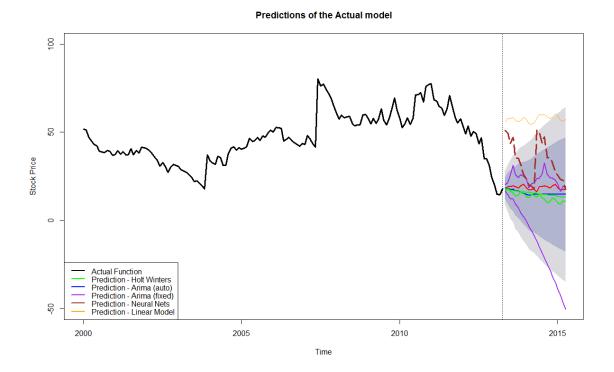
**Predictions of the Polynomial Trend**



As we can see, the predictions do not vary a lot, because the Polynomial Trend model has not a lot of variance in the Y Axes.

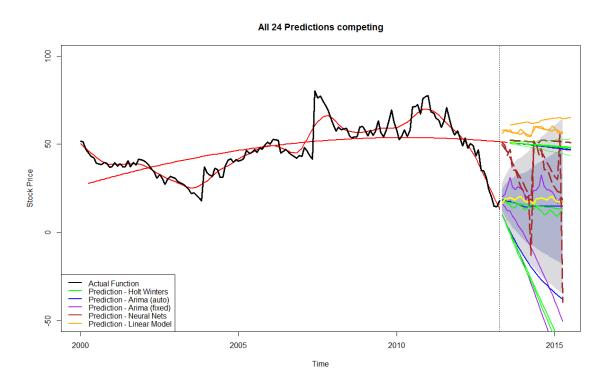This is different, when we look at the **STL Trend predictions**:

**Predictions of STL Trend**

Also the predictions for the Actual Graph vary a lot more:

**Predictions of the Actual model**



## 3. Evaluating the techniques

Before we begin with the evaluation, here are all predictions which are in the race:

**All 24 Predictions competing**

To evaluate this, I wrote a function *findBestPrediction.R*

Please open this attached function, to better understand my explanations.

The function is divided in mainly **4 steps**:

- Reading in the data and splitting it into train and test data

- Computing all 24 predictions of the various techniques using the three base functions of the train data

- Comparing the predictions with the actual test data, which gives as the error

- Choosing the prediction ID which gave the smallest total error

- This *findBestPrediction.R* is called in another script *runTestEvaluation.R*, which reads in 30 sample stocks, stores them in a vector and iteratively queries *findBestPrediction.R* while saving the result of each in another vector (*BestPrediction*).

The results of the first 30 stocks beginning with the letter A are shown here:

```
RStudio
File  Edit  Code  View  Plots  Session  Project  Build  Tools  Help

     Go to file/function

Console ~/
> stockvector[[6]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ABM.csv", sep=",", header=TRUE)
> stockvector[[7]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ABT.csv", sep=",", header=TRUE)
> stockvector[[8]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ABV.csv", sep=",", header=TRUE)
> stockvector[[9]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ABX.csv", sep=",", header=TRUE)
> stockvector[[10]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ACE.csv", sep=",", header=TRUE)
> stockvector[[11]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ACI.csv", sep=",", header=TRUE)
> stockvector[[12]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ACO.csv", sep=",", header=TRUE)
> stockvector[[13]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ACT.csv", sep=",", header=TRUE)
> stockvector[[14]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ADC.csv", sep=",", header=TRUE)
> stockvector[[15]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ADM.csv", sep=",", header=TRUE)
> stockvector[[16]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/ADX.csv", sep=",", header=TRUE)
> stockvector[[17]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEC.csv", sep=",", header=TRUE)
> stockvector[[18]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEE.csv", sep=",", header=TRUE)
> stockvector[[19]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEG.csv", sep=",", header=TRUE)
> stockvector[[20]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEM.csv", sep=",", header=TRUE)
> stockvector[[21]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEO.csv", sep=",", header=TRUE)
> stockvector[[22]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AEP.csv", sep=",", header=TRUE)
> stockvector[[23]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AES.csv", sep=",", header=TRUE)
> stockvector[[24]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AET.csv", sep=",", header=TRUE)
> stockvector[[25]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AF.csv", sep=",", header=TRUE)
> stockvector[[26]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AFG.csv", sep=",", header=TRUE)
> stockvector[[27]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AFL.csv", sep=",", header=TRUE)
> stockvector[[28]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AGCO.csv", sep=",", header=TRUE)
> stockvector[[29]] = read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AGM.A.csv", sep=",", header=TRUE)
> stockvector[[30]] =  read.table("C:/0.USA/Harvard/Course - Big Data/11 final
project/NYSE-Monthly/AGN.csv", sep=",", header=TRUE)
>
> for(i in 1:length(stockvector))
+ {
+     BestPrediction[i] = findBestPrediction(stockvector[[i]])
+ }
There were 19 warnings (use warnings() to see them)
> BestPrediction
 [1]  5 15 20 11 20  7  4  4  5 20 11 22 16  7  3 22  5 18 25 14  7 16  6 20  3
[26] 15  7 23 16  7
> |
```

After some conversion (see below), we get the frequencies:

```
freq <- ave(rep(1, times=BestPrediction), BestPrediction, FUN=sum)
results = data.frame(BestPrediction,freq)
```

```
> results
   BestPrediction freq
1               5    3
2              15    2
3              20    4
4              11    2
5              20    4
6               7    5
7               4    2
8               4    2
9               5    3
10             20    4
11             11    2
12             22    2
13             16    3
14              7    5
15              3    2
16             22    2
17              5    3
18             18    1
19             25    1
20             14    1
21              7    5
22             16    3
23              6    1
24             20    4
25              3    2
26             15    2
27              7    5
28             23    1
29             16    3
30              7    5
>|
```

Now we have to look up the IDs of the prediction techniques in the *findBestPrediction.R*
function and get the resulting ranking from 30 samples:

1.    **ARIMA** model with auto parameters (5 times the best)

2.    **Linear** model with trend, season components (4 times the best)

3.    **Holt-Winters** Filtering (3 times the best)

      **Holt-Winters** Filtering with exponential smoothing (3 times the best)

## 4. Discussing the results and Outlook

Because only 30 stocks were used in the evaluation phase, the results are merely hints to which techniques are best suitable for stock predictions.

Using the attached scripts, however, it should not be a problem for you to try it out with 300 or more… ☺

Also it has to be said, that most of the used techniques are heavily dependent on the parameters used – here I also recommend trying out different values and/ or other parameters all together.

A future goal could be to find the top 10 stocks which are worth investing in. Then the evaluation test should be done again with maybe 500 stocks to acquire the most suitable prediction techniques. Afterwards these prediction methods can be used to predict all stocks which have historic data of at least 10 years and a **criterion** for promising trends then should select the most promising stocks. For example the most rapid and constant increase of stock price from the point when the prediction began to the prediction end is feasible.

You are free to use my scripts and methods as a starting point.

Let me know which stocks you found and which predicting method you used. ☺

---

Dietmar Aumann, 05/09/2013, Malmö (Sweden)