Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

# Methodology Seminar
## The functionalities of the Gretl software

### Hassan OUKHOUYA

Instructor-Researcher: $M^r$.Khalid EL HIMDI
LMSA Lab., FSR, Mohammed V University, Rabat

April 6, 2021

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

## Roadmap

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

## Presentation

$T_{HE}$ aim of this presentation is to propose an initiation, in terms applied to econometrics of time series or panel data, using a Gretl econometrics software with illustrative examples.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Presentation of the interface

# What is GRETL?

### GRETL as per Wikipedia

- Gretl is an open-source statistical package, mainly for econometrics.

- The name is an acronym for Gnu Regression, Econometrics and Time-series Library.

- It is written in C, uses GTK+ as widget toolkit for creating its GUI, and calls gnuplot for generating graphs, and can be integrated with other software such as R, Stat, Python...etc.

- More on http://gretl.sourceforge.net/

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Presentation of the interface

# The main window menus

We find:

- File.
- Tools.
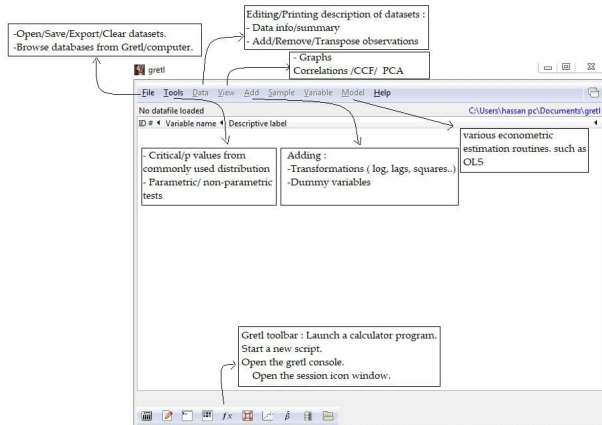- Data.
- View.
- Add.
- Sample.
- Variable.
- Model.
- Help.



Figure: Interface of gretl.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

## Simple linear regression Model

$$y = \beta_0 + \beta_1 x + \epsilon$$

- ⚀ $y$ is the dependent variable.
- ⚁ $x$ is the independent variable.
- ⚂ $\beta_0$ is the constant or intercept.
- ⚃ $\beta_1$ is x's slope or coefficient.
- ⚄ $\epsilon$ is the error term.

⚠ Often we suppose that $\epsilon$ follows the normal distribution $\mathcal{N}(0, \sigma^2)$.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Two Main Objectives

1. Establish if there is a **relationship** between two variables.
   - More specifically, establish if there is a statistically significant relationship between the two.
   - Examples: Income and spending, wage and gender, student height and exam scores.

2. **Forecast** new observations.
   - Can we use what we know about the relationship to forecast unobserved values?
   - Examples: What will be the sales over the next quarter? What will the ROI (Return On Investment) of a new store opening be contingent on store attributes?

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Linear Regression Example

What could explain a family's consumption of a given product?

Table: Data for Linear Regression Example

| ID | Income | Consumption | ID | Income | Consumption |
|----|--------|-------------|----|--------|-------------|
| 1  | 119    | 154         | 11 | 81     | 115         |
| 2  | 85     | 123         | 12 | 81     | 117         |
| 3  | 97     | 125         | 13 | 91     | 123         |
| 4  | 95     | 130         | 14 | 105    | 144         |
| 5  | 120    | 151         | 15 | 100    | 137         |
| 6  | 92     | 131         | 16 | 107    | 140         |
| 7  | 105    | 141         | ⋮  | ⋮      | ⋮           |
| 8  | 110    | 141         | 38 | 96     | 131         |
| 9  | 98     | 130         | 39 | 82     | 127         |
| 10 | 98     | 134         | 40 | 114    | 150         |

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Open and Examine Data in Gretl

In the Gretl program this can be done via the "File, Open Data, User file":

The supported import formats are as follows:

- Stata files(.dta).
- SPSS files(.sav).
- JMulTi files.
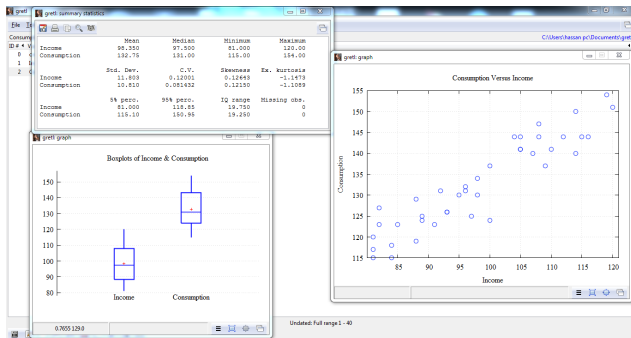- CSV files(.csv).
- Excel fil.(.xlsx).
- ASCII fil.(.txt).

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
**Data Analysis**
Analysis of Residuals
Forecasting

# Summary statistics and Graph X-Y

Shows a full set of descriptive statistics and scatter plot for the variables selected in the main window:

You can click on the graph window for a pop-up menu with the following options.

- Save as PNG.
- Zoom.
- Print.
- Edit.
- Close.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

## Principal Components

We can go further and produces a Principal Components Analysis.



Figure: Principal components; Cross Tabulation; Mahalanobis distances; Correlation matrix.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

## Simple Linear Regression Model

**We can write our model:**

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$consumption \; = \; \beta_0 + \beta_1 \; income + \epsilon$$

**?** Our assumption, which we will test, is that income explains consumption.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Regression Results

We see that $R^2 = 0.86$. We can therefore say that about 86% of the variability of the dependent variable consumption is explained by the income variable; $\quad$ *consumption* $=$ 49.18 + 0.85 *income* $+ \epsilon$
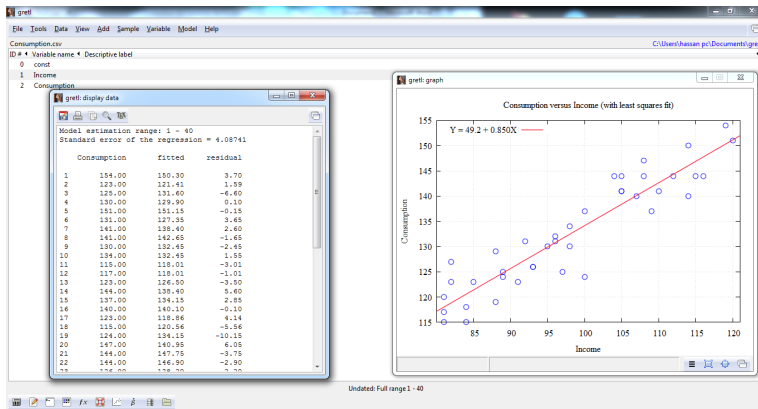
Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Interpreting the coefficients

The estimated model is: $consumption = 49.18 + 0.85\ income + \epsilon$.

49.18 could be interpreted as the consumption level of a family with 0 income.

Most generally, intercept doesn't have intuitive interpretation.

0.85 is the marginal effect of one unit of income on consumption: for every unit more of income a family has, we estimate its consumption grows by 0.85 units.

The slope always has an intuitive interpretation.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
**Analysis of Residuals**
Forecasting

# Estimated vs. Actual Values, Residuals

| ME | RMSE | MAE | MPE | MAPE | Theil's U |
|---|---|---|---|---|---|
| -7.11e-15 | 3.98 | 3.24 | -0.09 | 2.46 | 0.32 |

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
Forecasting

# Distribution of residuals



Figure: Residuals.

Presentation
What is GRETL?
**Simple Linear Regression in Gretl**
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Importing Data
Data Analysis
Analysis of Residuals
**Forecasting**

We want to predict the consumption of 3 families with following incomes: 90, 95 and 100.



Figure: Forecasts for 95% confidence intervals.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# Stock price analysis "Dow Jones"

Stock price analysis is very popular and important in financial study and time series is widely used to trait this topic. The data we use in this presentaion is the Dow Jones Companies from 1st January 2000 to 6th December 2017.

The Dow Jones is an index that shows how 30 large publicly owned companies (IBM, Apple, Pfizer...etc.) based in the United States have traded during a standard trading session in the stock market.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Data and Methodology

The data we use in this presentation is the daily stock price of
IBM Holdings:

- Daily (5 days) IBM from 1999-12-31 to 2017-12-05.

- 4512 Days of data.

- In sample training: 1999-12-31 - 2016-02-23 (90%).
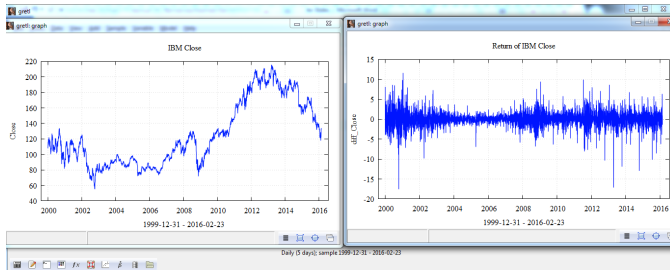
- Run Model ARIMA/GARCH.

- Out-of-sample forecast: 2016-02-24 - 2017-12-05 (10%).

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Load time series dataset in Gretl.

The dataset contains Open, High, Low, Close and Adjusted Close prices of IBM stock each day of this period.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Plot Time Series of IBM Close and Return of IBM

We can find that close price have a general increasing trend over time, and there is a significant drop at the beginning of 2003, 2009 and the beginning of 2016. We decide to analyse returns of IBM which is the differenced close price so that we get a more stationary time serie. We can see that There are periods of persistently high volatility, which shows significantly fluctuations.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Stationarity tests (Augmented Dickey-Fuller test, KPSS test)

After the differenced of the series, we see that the series is stationary
in differences. Using test for unit roots (ADF test, KPSS test). I
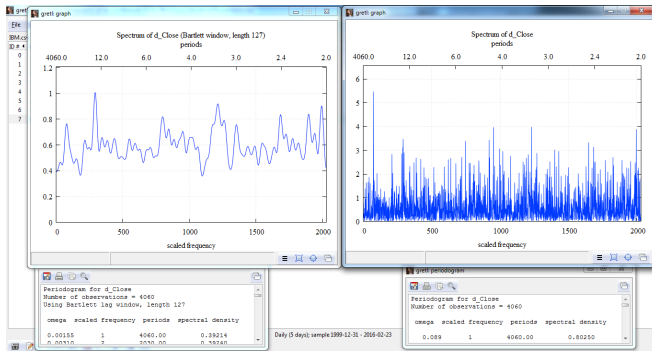chose the constant and trend version, setting the max lag to 30 and
testing down.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# ACF and PACF

The ACF don't die fast which indacates that the serie are not stationary (below left). After differencing data, things seem pretty nice.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# Check seasonal effect

I plot the frequency domain of original data and smoothed data. we could see that there's no significant cycle in a year, which means there's no seasonal effect in this problem. Therefore, we choose to use the ARMA model.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# ARMA Model

**I will start by fitting an ARMA(p,q) model of the form:**

$$Y_t = \mu + \phi_1(Y_{t-1} - \mu) + \ldots + \phi_p(Y_{t-p} - \mu) + \epsilon_t + \theta_1\epsilon_{t-1} + \ldots + \theta_q\epsilon_{t-q}$$

- $\{\epsilon_t\}$ is a white noise process with distribution $\mathcal{N}(0, \sigma^2)$.
- $\phi_1, \ldots, \phi_p$ are the coefficients for the autoregressive part of the model.
- $\theta_1, \ldots, \theta_q$ are the coefficients for moving average part of the model.
- $\mu$ is the population mean.
- $\sigma^2$ is error variance.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# Determine 'p' and 'q'

The orders $p$ and $q$ of the AR and MA models are then obtained by looking at the autocorrelations and partial autocorrelations. Very often, these orders are not obvious. In this case, we can obtain upper bounds for $p$ and $q$ and then select a model by minimizing a penalized criterion of type AIC or BIC.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

Estimate and interpretation of coefficients ARIMA(p,d,q)

From the result, the autoregressive and moving average term
has a p-value that is less than the significance level of **5%**. You
can conclude that the coefficient for the AR and MA is
statistically significant, and you can choose model ARMA(1,1)
with a minimum of AIC.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Testing ARMA(1,1) Model Root

From the result, we can see that ARMA(1,1) model is causal (or stationary) and invertible because the AR and MA polynomial has all its roots outside of the unit circle. So my final model is ARMA(1,1): $Y_{1:t}$ is defined by $(1 - \phi_1 B)(Y_t - \mu) = (1 + \theta_1 B)\epsilon_t$ i.e.: $(1 + 0.523B)(Y_t - 0.007) = (1 + 0.509B)\epsilon_t$ where $\epsilon_t \sim \mathcal{N}(0, 1.932^2)$

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# GARCH model

## A GARCH(p,q) model can be defined as follows:

$$\epsilon_t = \sigma_t z_t$$

- $z_t \sim$ i.i.d $\mathcal{N}(0,1)$ independent of the past $\epsilon_{t-1}, \epsilon_{t-2}, \ldots \forall t$
- $\sigma_t^2 = \omega + \alpha_1 \epsilon_{t-1}^2 + \ldots + \alpha_p \epsilon_{t-p}^2 + \beta_1 \epsilon_{t-1}^2 + \ldots + \beta_q \epsilon_{t-q}^2$.
- Where $\omega > 0, \alpha_i, \beta_j \geq 0, \ i = 1, \ldots, p, \ j = 1, \ldots, q$ and $\sum_i^p \alpha_i + \sum_j^q \beta_j < 1$.

We can model the trajectory of return by a stationary GARCH(1,1) process.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
**Model selection**
Diagnostics
Forecasting

# GARCH in Gretl

All coefficients of Volatility (alpha(0), alpha(1), beta(1)) are significant. But the omega (const) isn't significant.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

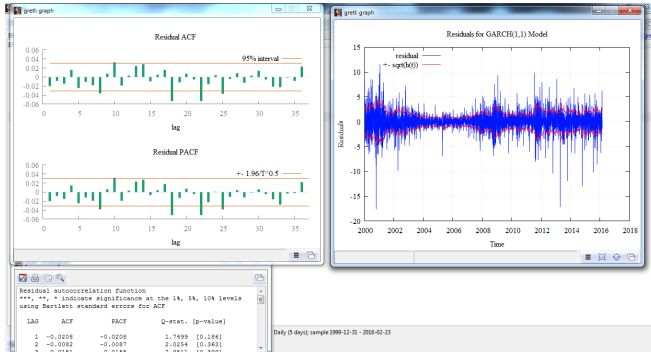Time series analysis
Model selection
Diagnostics
Forecasting

# Residuals for ARMA(1,1) Model

- We can see the mean of residuals is almost zero. However, we can see the not a constant of variance.
- The ACF of the residuals behave like the ACF of a white noise process. The values of ACF almost fall inside the dashed lines except lag 8, 10, 18, 22, 25. So we conclude that the autocorelation of residuals are zero.
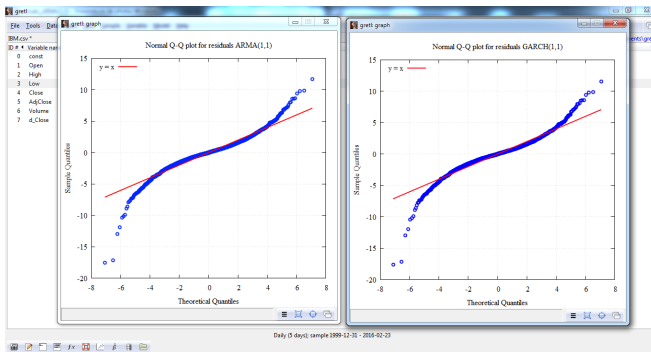
Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# Residuals for GARCH(1,1) Model

Same thing for GARCH model...!

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
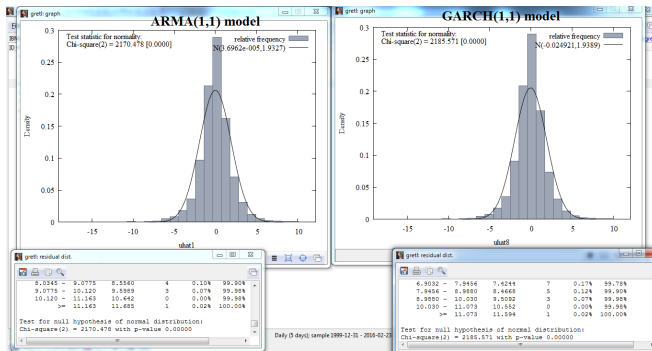Model selection
**Diagnostics**
Forecasting

# Testing for Normality

From QQ-Plot, we can see that the distribution of residuals are
not normally distributed. There are too many extreme positive
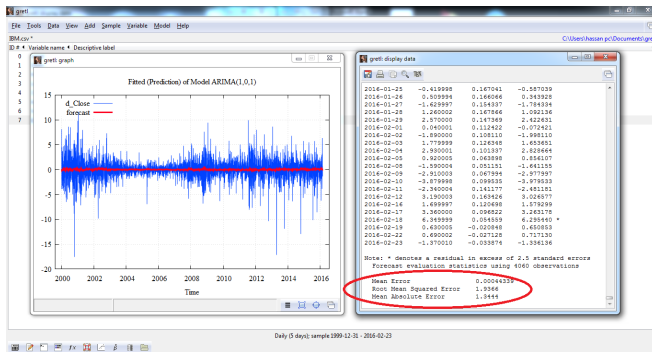and negative residuals. We say the distribution is "**heavy
tailed**".

Presentation
What is GRETL?
Simple Linear Regression in Gretl
**Modeling and Forecasting Time Series in Gretl**
Case Study
Interesting Resources

Time series analysis
Model selection
**Diagnostics**
Forecasting

# Histogram of residuals

Graphically, histogram of residuals suggests that the residuals (and hence the error terms) are normally distributed:

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
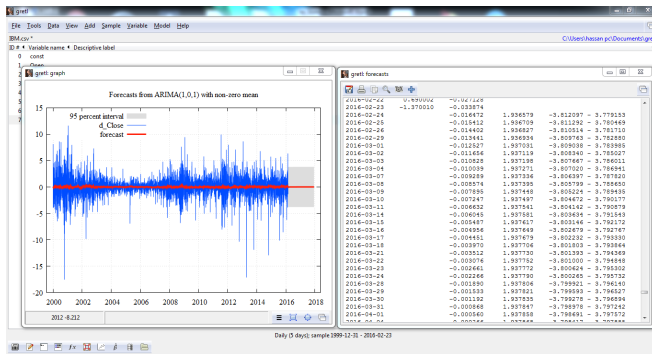Diagnostics
Forecasting

# Forecasting

The value of RMSE is 1.9366, which is pretty small. And this may indicate the ARMA(1,1) model is reasonable although we need to further study this for better analysis.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

# Forecasting

The heavy light gray bar seperately represent the **95%** confidence interval for the forecast. We can see that there will not be a large volatility by forecasting.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Time series analysis
Model selection
Diagnostics
Forecasting

## Conclusion

We find that ARMA(1,1) model with white noise fits the IBM stock price data. During the analysis, we use difference function to Returns. And then, the AIC suggests that ARMA(1,1) is the proper choice. After diagnosing the result using residual plot and ACF, we further made a forecast using the ARMA(1,1) model, which needs more effort for better analysis.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

# Gross domestic product (GDP) of U.S.

## GDP data

Gross domestic product (GDP), the featured measure of U.S. output, is the market value of the goods and services produced by labor and property located in the United States.
For more information, see the Guide to the National Income and Product Accounts of the United States (NIPA) and the Bureau of Economic Analysis.

- Seasonal Adjustment: Seasonally Adjusted Annual Rate.

- Frequency: Quarterly.

- Units: Billions of Dollars.

- Date Range: 1990-01-01 to 2020-10-01.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## The Box-Jenkins Univariate ARIMA Approach:

- Identification.
- Estimation.
- Diagnostic Checking.
- Forecast.

Data collected:

1. Use for model construction (Sample Range): 1990 Q1 - 2014 Q3 (80%).

2. Use for forecast performance assessment (Forecast Range): 2014 Q4 - 2020 Q4 (20%).

Sample size for each observation: 124.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Identification

- 📈 There is a trend (perhaps exponential).

- ☹ The ACF graph for GDP dies out slowly (exponentially decaying), with one spike in PACF that cuts off after lag 1. Data is nonstationary.



🐻 It is possible to remove the trend since differencing should remove any deterministic(linear) trend from the series. Let's plot the first order differencing of log!

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Identification

📊 The difference of logs GDP (below right) may show a slight
upward trend until the bottom dropped out in late 2008.



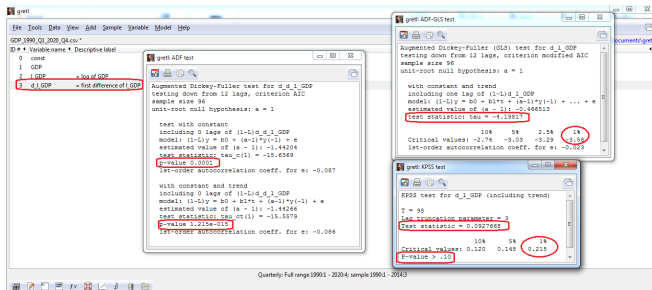From ACF &
PACF we can
choose a models:

- $ARIMA(1, 1, 1)$

- $ARIMA(1, 1, 0)$

- $ARIMA(2, 1, 0)$

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Identification

Now test for unit roots using the ADF-GLS test. I chose the constant and trend version, setting the maximum lag to 12 and testing down.



☞ The test results indicate that d_l_GDP has not unit root with significantly asymptotic p-value < 10%. This tells us that serie is stationary.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Estimation

Below are the estimated results of ARIMA $(1, 1, 0)$/ARIMA $(2, 1, 0)$ of model d_l_GDP ($Y_{1t}$).

Replace the symbol with estimated parameters, we have:

$$\Delta^{1d}Y_{1t} = 4.314 \times 10^{-5} - 0.44\Delta^{1d}Y_{1,t-1} + \epsilon_{1t} \text{ where } \Delta^{1d}Y_{1t} = Y_{1t} - Y_{1,t-1}$$

$$\Delta^{1d}Y_{1t} = 4.482 \times 10^{-5} - 0.52\Delta^{1d}Y_{1,t-1} - 0.176\Delta^{1d}Y_{1,t-2} + \epsilon_{1t}$$

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Diagnostic Checking

🔧 For diagnostic checking, I will perform additional tests on all the models passing the normality test, for instance Ljung-Box Q test for autocorrelation and ARCH test for heteroskedasticity (lag order 4). The selection technique for the best fitted model is to have the best overall white noise, which is normally distributed, independently distributed (NO autocorrelation), and homoskedastic (NO ARCH effect).

Below is the summarized testing results:

|  | Normality ($\chi^2$) Test | Autocorrelation Test | ARCH Effect Test |
|---|---|---|---|
| ARIMA(1,1,0) | 8.398 (0.015) | 4.461 (0.216) | 7.859 (0.097) |
| ARIMA(2,1,0) | 6.821 (0.033) | 3.466 (0.177) | 10.516 (0.033) |

p-value for the testing statistics is in the parenthesis. 0.02 is used as significance level for all the hypothesis testing.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
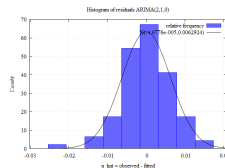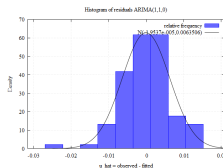Box–Jenkins model diagnostics
Forecasting

## Diagnostic Checking

✎ On a Q-Q plot normally distributed residuls appears as roughly a straight line.



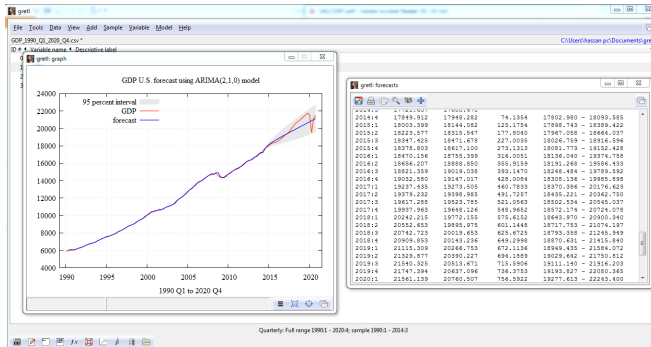The normal distribution is symmetric, so it has no skew (the mean is equal to the median).

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

## Diagnostic Checking

📈 We will utilize as model 2 (ARIMA(2,1,0)) based on their forecasting performance and overall the goodness of fit for forecasting.



| | ME | RMSE | MAE | MPE | MAPE | Theil's U |
|---|---|---|---|---|---|---|
| ARIMA(1,1,0) | 0.188 | 75.809 | 53.845 | -0.041 | 0.465 | 0.467 |
| ARIMA(2,1,0) | 0.453 | 74.135 | 52.789 | -0.032 | 0.454 | 0.453 |

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
Box–Jenkins model diagnostics
Forecasting

# Forecast



Forecast evaluation statistics using 25 observations:

|  | ME | RMSE | MAE | MPE | MAPE | Theil's U |
|---|---|---|---|---|---|---|
| ARIMA(2,1,0) | 212.37 | 590.08 | 447.83 | 0.942 | 2.184 | 1.027 |

Hassan OUKHOUYA          Methodology Seminar

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

Data
The Box–Jenkins method
Box–Jenkins model identification
Box–Jenkins model estimation
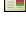Box–Jenkins model diagnostics
Forecasting

## Conclusion

Overall, we can see that ARIMA (2, 1, 0) provides a good fit for GDP data. It gives a fairly accurate forecasting. However, although forecasts from 2014 Q3 to 2020 Q4 are within the 95% percent interval, the graph shows that the red line of actual data has in of the confidence interval starting from 2020 Q1 and climbing back up toward the interval from 2020 Q2. Such trend exactly coincides with the way of how the economy has evolved since great recession of 2020 (COVID-19). But, the weakness of ARIMA model is that it could not predict such trend but rather assume the same pattern from 1990 Q1-2014 Q3, that is, the ARIMA model is not good for volatility analysis, but GARCH is. Finally, GDP does strive back after a deep dip in 2020 Q1.

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

References

## Sources

- Gretl Tutorial 1: Simple Linear Regression

- Time Series Analysis for Stock Data

- Data of Dow Jones

- Time Series Analysis for Log Returns of Nasdaq

- ARIMA(p,d,q) Models (Video 6 of 7 in the gretl Instructional Video Series)

- Les processus ARIMA 2020-2021

- Time Series Analysis for Log Returns of S&P500

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

References

## Sources

📄 Further reading:
  - Interpret the key results for ARIMA

📄 External links:
  - Gretl Tutorial 6: Modeling and Forecasting Time Series Data
  - ARIMA and its Diagnostics in gretl
  - ARCH & GARCH in gretl

📄 Guide d'utilisation de GRETL

📄 Gross Domestic Product Data

Presentation
What is GRETL?
Simple Linear Regression in Gretl
Modeling and Forecasting Time Series in Gretl
Case Study
Interesting Resources

References

## Sources

- Comprehensive Time-Series Regression Models Using Gretl-U.S. GDP and Government Consumption Expenditures & Gross Investment from 1980 to 2013

- Using gretl for Principles of Econometrics, 5th Edition

- Introduction to Gretl