

STA5092Z – Exploratory Data Analysis

Masters in Science – Data Science Specialisation Program

Mr Stefan Britz

Dr. Şebnem Er

University of Cape Town

Department of Statistical Sciences

Stefan.Britz@uct.ac.za

Sebnem.Er@uct.ac.za – <https://sebnemer.github.io/>



STA5092Z - Exploratory Data Analysis

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Introduction

Introduction - Lecturers

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3



Mr Stefan Britz
Stefan.Britz@uct.ac.za
[Github page](#)
[UCT Research Profile](#)

Dr. Sebnem Er
Sebnem.Er@uct.ac.za
[Github page](#)
[Google Scholar page](#)

Introduction - Structure of the Course

Introduction

Lecture 1 - 2

Software - R

Basic R Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

- First 4 weeks of the first semester
- 24 hours of lectures - in classroom environment
Mon/Wed/Fri 4-6pm.
- 8 lectures by Dr Şebnem Er, and 16 lectures by Mr Stefan Britz

Introduction - Grading

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

There are 3 assignments during the 4 weeks of lectures.
Your final mark will be a weighted average of these assignments.

There will be a late hand-in penalty.

Introduction - Lecture Philosophy

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

- The idea is for the course to be very applied and hands-on, and based around worked examples and case studies, with a focus on the code used to generate an analysis.
- The course is entirely in R, and the goal of the course is as much to introduce students to R and develop R skills as to cover the theory of EDA – which is after all relatively simple.

Introduction - Resources

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

- All material available on VULA
- Exploratory Data Analysis, 1977 Addison-Wesley Publishing Company (John Wilder Tukey)
- Exploratory Data Analysis: Past, Present and Future, 1993 Technical Report (John W Tukey)
<https://apps.dtic.mil/dtic/tr/fulltext/u2/a266775.pdf>
- How to look at data: A review of John W. Tukey's Exploratory Data Analysis, 1979 (Russell M. Church)
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1332871/pdf/jeab0140.pdf>

Introduction - Resources cont.

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

- RP: **Exploratory Data Analysis with R** (Roger Peng):
<https://bookdown.org/rdpeng/exdata/>
- JB: **STA545: Data wrangling, exploration, and analysis with R** (Jenny Bryan): <https://stat545.com/index.html>
- HW: **R for data science** (Hadley Wickham):
<https://r4ds.had.co.nz/>

Introduction - Lecture Outline

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Lecture	Material	Possible resources
1	R installation, basics, workflows	JB2
2	visualizing raw data with ggplot	HW3, RP6
3	managing data frames with dplyr	JB5-7, HW5, RP3
4	(filter, select, arrange, mutate, summarize)	Some from DSFI
5	EDA checklist	RP4, HW7
6	(right questions, correlation, missing values, outliers)	
7	Reshaping, tidying, joining dataframes	JB8, JB14-16, HW12-13
8	(tidyr, more dplyr)	
9-10	Principles of good graphics	JB24-27, RP5, RP14-15
11-12	Exploring time series data	JB10-13, HW14-16
13	Exploring spatial data	
14	(mapview, leaflet, sf)	
15	Functional programming with R	
16	(writing functions, purrr)	
17	Visualising data using animations	RP11-12
18	(ggridanimate)	
19-20	Version control with Git and GitHub	JB3, some from DSFI
21	Dashboards	JB42
22	(flexdashboard, shiny)	
23-24	Dimensionality reduction and Clustering	RP13

**STA5092Z -
Exploratory
Data
Analysis**

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Lecture 1 - 2

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3



- Download R, RStudio, later on Git, account on Github
- Version check

Software - Why R?

- Both R and Python are free.
- R already has all of the statistics support because it was developed by statisticians for statisticians. A lot of statistical modelling research is conducted in R.
- Python was originally developed as a programming language for software development, DS tools (`scikit-learn`, `pandas`, `numpy`) were added on. Though the majority of DL research is done in Python, such as `keras`, `PyTorch`
- R has Tidyverse, a set of packages that makes it easy to import, manipulate, visualise and report data.
- Very easy to generate dashboards using R Shiny.

Ref: Python vs. R for Data Science: What's the Difference (by Richie Cotton)

Ref: R vs Python for Data Science: The Winner is... (by Martijn Theuwissen)

Software - Why R? cont.

- It is the language I know the best, I know very little Python.
- Python and R programmers get inspired from each other, ie. Python's `plotnine` inspired by R's `ggplot2`, and R's `rvest` by Python's `BeautifulSoup`.
- You can also use functions written in Python with `source_python()` function in R.
- You can run R code from Python with `rpy2` package, and you can run Python code from R using `reticulate`. R version of DL package Keras calls Python.
- Though, I do encourage you to learn Python as well. No harm in two languages.

Take away: There is no winner, you are here to learn the skills, your focus should be on skills. If you can program in R, you can do it in any other language.

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Function	Description
<code>version\$version.string</code>	which version of R
<code>RStudio.Version()\$version</code>	which version of RStudio
<code>objects()</code>	list of objects in your environment
<code>ls()</code>	list of objects in your environment
<code>rm(list = ls())</code>	remove everything from your environment
<code>getwd()</code>	what is your working directory
<code>setwd(C:/EDA)</code>	

Dataset resources

Introduction

Lecture 1 - 2

Software - R
Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Before you can work with data you have to get some data:

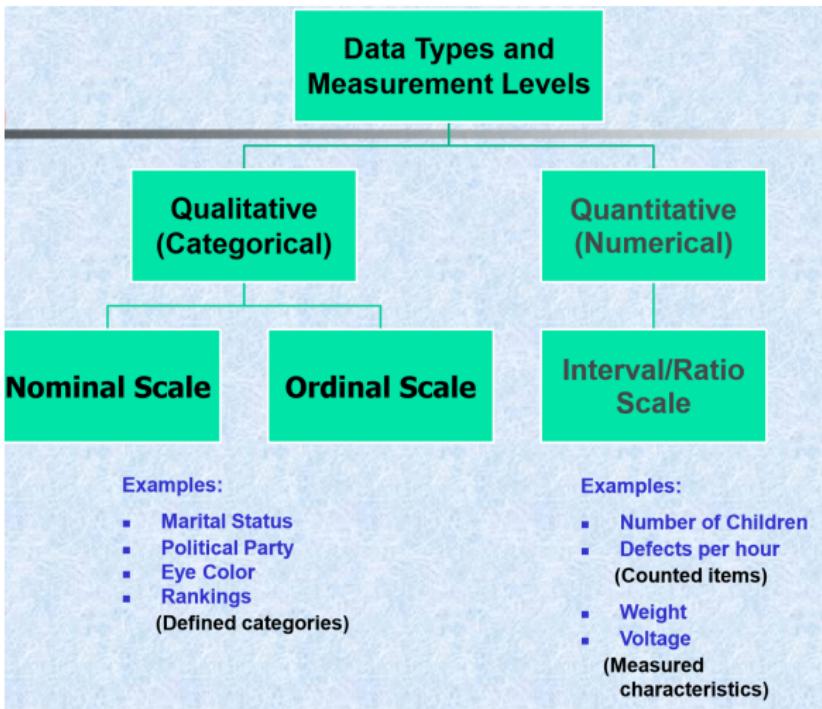
Zindi Competitions: <https://zindi.africa/competitions>

Zindi is the first data science competition platform in Africa.

Zindi hosts an entire data science ecosystem of scientists, engineers, academics, companies, NGOs, governments and institutions focused on solving Africa's most pressing problems.

Kaggle competitions: <https://www.kaggle.com/competitions>

Dataset types - Measurement levels



Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Dataset file formats

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

- Raw files (.csv,.txt, .xlsx, .sav, etc.)
- Databases (mySQL, MongoDB)
- APIs (Twitter)

How do we clean data, how do we tidy data, how do we share data? This course will cover all these as well as exploring the data.

Tukey's EDA book and terminology

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

About the book:

- The book provides techniques and advice about how to explore data.
- The approach of EDA is detective in character, it is a search for clues. Some of the clues may be misleading, but some will lead to discoveries.

Tukey's EDA book and terminology

About the author, Tukey:

- Tukey favors **simplicity** because simple statements are clear.
- Tukey favors **clear visual displays of quantitative facts**.
- Tukey likes **precision**, it is far better to be able to say some response measure is a linear function of a particular stimulus variable than to say it increases with the stimulus variable.
- Tukey favors **depth of analysis**. It is always good to look at the residuals.
- Tukey values **accuracy**. A misplaced decimal vs a misplaced digit.
- Tukey values **replicability** of summary observations in situations containing aberrant observations.

Tukey's EDA book and terminology

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

The methods described in the book:

- Frequency distributions
- Measures of central tendency and variability
- Scale transformations
- Graphical displays of single variable and of relationships
- Smoothing techniques
- Analysis of tables

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Tukey's approach to data analysis is highly visual and he has numerous suggestions for graphical displays. Tukey emphasizes the value of graphs for the following:

- Graphs can be used to **store** quantitative data,
- Graphs can be used to **communicate** conclusions,
- Graphs can be used to **discover** new information. Tukey particularly emphasizes the value of graphs for discovery.

Some types of plots are better for one purpose, others are better for another.

Distribution of a Single Quantitative Variable

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Tukey's novel distribution tools:

- Stem and Leaf: The measures Tukey proposes involve no arithmetic, only counting. - not widely used.
- Histogram: One should note,
 - its height,
 - where it is centered,
 - how spread out it is,
 - whether it is asymmetric,
 - whether there are any discontinuities.
- Box-Whisker Plots: These plots show medians, quartiles, and two extreme values in a format that is easy to grasp quickly. Very powerful when comparing several frequency distributions.
- Q-Q plots

Visual Display of a Single Qualitative Variable

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Relative Frequency Distribution -
compare percentages and probability

- Frequency distribution
- Pie chart
- Bar chart

Summarizing data

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

In order to visually display data, sometimes we need to summarize our data with summary measures:

- Compute a few "key" numbers: 5 number summaries.
- Our aim is to reduce a large data set to a few numbers which will help us understand the important features of the data.
- Let's begin with the concept of ranked data.
- In a sample of size n , the smallest number has a rank of 1; the second smallest number has a rank of 2; ; the largest number has a rank of n .

$x_1, x_2, x_3, \dots, x_n$ and $x_{(r)}$ is the number with rank r .

Summarizing data - cont.

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

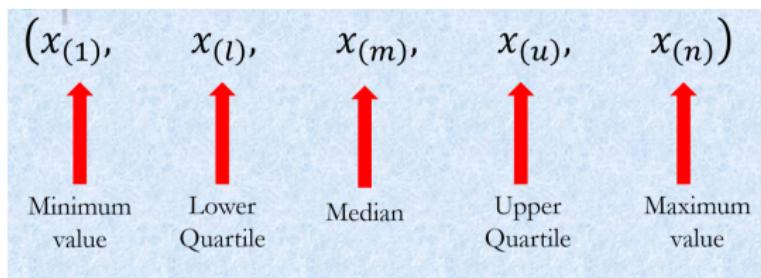
Graphs

Transform

Tools

Example 1

Lecture 2 - 3



Which transformations and Why?

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Very often it is more convenient to look at some transform of the original variable. If the distribution is far from symmetrical, one end of the distribution will be too crowded to permit careful inspection.

Tukey deals extensively with scale transformations. He gives three main reasons for transformations:

- A transformation may be selected to produce a symmetrical distribution,
- A transformation may increase the similarity of the spread of the different sets of numbers,
- A transformation may straighten out a line.

Types of transformations

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

The transformations discussed range on:

$$x^n$$

$$x^{n-1}$$

$$x^2$$

$$\log x$$

$$-\frac{1}{x}$$

$$-\frac{1}{x^2}$$

$$-\frac{1}{x^n}$$

Other dependent variables, e.g. counts and latencies, are occasionally transformed by taking the square root, the logarithm, or the reciprocal.

Data Science Toolbox

Visualising plays an important role in exploring your data, and you would know that Tukey favours analysis of data with four-color pen, graph paper, few tables etc.:



Though we will use R and its functions for this purpose:

- ggplot2
- tidyr
- dplyr

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Example 1 - Female headed households in SA

Let us look at the “Womxn in Big Data South Africa: Female-Headed Households in South Africa” competition at <https://zindi.africa/competitions/womxn-in-big-data-south-africa-female-headed-households-in-south-africa>

The datasets are provided in a .csv file format, test.csv, train.csv, variable_descriptions.csv.

The target variable of interest is the percentage of households per ward that are both female-headed and earn an annual income that is below R19,600 (approximately \$2,300 USD in 2011).

```
> Train <- read.csv("C:/Users/01438475/Google Drive/  
UCTcourses/DataScience/LectureNotes/EDA/Datasets/Woman/  
Train.csv", header = TRUE)  
> str(Train)  
> summary(Train)
```

Now we will explore this dataset.

tidyverse ecosystem

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

The majority of the packages that you will use are part of the so-called tidyverse package:

```
> install.packages("tidyverse")
> library(tidyverse)

> Attaching packages tidyverse 1.2.1
> ggplot2 3.2.1      purrr   0.3.3
> tibble  2.1.3      dplyr    0.8.3
> tidyr   1.0.0      stringr 1.4.0
> readr   1.3.1     forcats  0.4.0
> Conflicts tidyverse_conflicts()
> dplyr::filter() masks stats::filter()
> dplyr::lag()   masks stats::lag()
```

This tells you that tidyverse is loading the ggplot2, tibble, tidyr, readr, purrr, and dplyr packages.

ggplot2 functions1



ggplot2 part of the tidyverse

3.2.1

	<code>geom_abline()</code> <code>geom_hline()</code> <code>geom_vline()</code>	Reference lines: horizontal, vertical, and diagonal
	<code>geom_bar()</code> <code>geom_col()</code> <code>stat_count</code> (<code>)</code>	Bar charts
	<code>geom_bin2d()</code> <code>stat_bin_2d()</code>	Heatmap of 2d bin counts
	<code>geom_blank()</code>	Draw nothing
	<code>geom_boxplot()</code> <code>stat_boxplot()</code>	A box and whiskers plot (in the style of Tukey)
	<code>geom_contour()</code> <code>stat_contour()</code>	2d contours of a 3d surface
	<code>geom_count()</code> <code>stat_sum()</code>	Count overlapping points
	<code>geom_density()</code> <code>stat_density()</code> (<code>)</code>	Smoothed density estimates
	<code>geom_density_2d()</code> <code>stat_density_2d</code>	Contours of a 2d density estimate

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

ggplot2 functions2



ggplot2

part of the [tidyverse](#)

3.2.1

	<code>geom_dotplot()</code>	Dot plot
	<code>geom_errorbarh()</code>	Horizontal error bars
	<code>geom_hex()</code> <code>stat_bin_hex()</code>	Hexagonal heatmap of 2d bin counts
	<code>geom_freqpoly()</code> <code>geom_histogram()</code> <code>stat_bin()</code>	Histograms and frequency polygons
	<code>geom_jitter()</code>	Jittered points
	<code>geom_crossbar()</code> <code>geom_errorbar()</code> <code>geom_linerange()</code> <code>geom_pointrange</code> (<code>)</code>	Vertical intervals: lines, crossbars & errorbars
	<code>geom_map()</code>	Polygons from a reference map
	<code>geom_path()</code> <code>geom_line()</code> <code>geom_step</code> (<code>)</code>	Connect observations

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

ggplot2 functions3

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3



ggplot2 part of the tidyverse

Refer

	<code>geom_point()</code>	Points
	<code>geom_polygon()</code>	Polygons
	<code>geom_qq_line()</code> <code>stat_qq_line()</code>	A quantile-quantile plot
	<code>geom_qq()</code> <code>stat_qq()</code>	
	<code>geom_quantile()</code> <code>stat_quantile()</code>	Quantile regression
	<code>geom_ribbon()</code> <code>geom_area()</code>	Ribbons and area plots
	<code>geom_rug()</code>	Rug plots in the margins
	<code>geom_segment()</code> <code>geom_curve()</code>	Line segments and curves
	<code>geom_smooth()</code> <code>stat_smooth()</code>	Smoothed conditional means
	<code>geom_spoke()</code>	Line segments parameterised by location, direction and distance
	<code>geom_label()</code> <code>geom_text()</code>	Text

ggplot2 functions 4

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

The screenshot shows the ggplot2 documentation page. At the top, there's a logo of a hexagon containing a line graph, followed by the text "ggplot2 part of the tidyverse 3.2.1". Below this, there are two sections:

- Rectangles**: Contains icons for a grid (geom_raster()), a rectangle (geom_rect()), and a tiled grid (geom_tile()).
- Violin plot**: Contains icons for a violin shape (geom_violin()) and a density plot (stat_ydensity()).
- Visualise sf objects**: Contains an icon of a map (coord_sf(), geom_sf(), geom_sf_label(), geom_sf_text(), stat_sf()).

Example 1 with ggplot2

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1



```
> library(ggplot2)
```

Warning message:

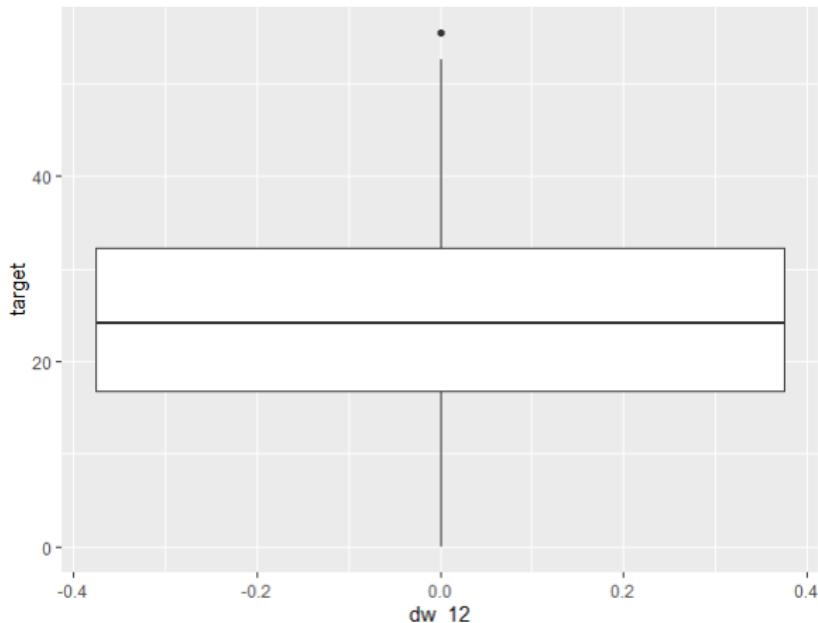
package ‘ggplot2’ was built under R version 3.6.1

```
> R.version.string
```

```
[1] "R version 3.6.0 (2019-04-26)"
```

Example 1 with ggplot2 - boxplot

```
> plot1 = ggplot(data = Train, aes(dw_12,target))  
> plot1 + geom_boxplot()
```



How to interpret this: IntroStat p:21-26

Example 1 with ggplot2 - boxplot options

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

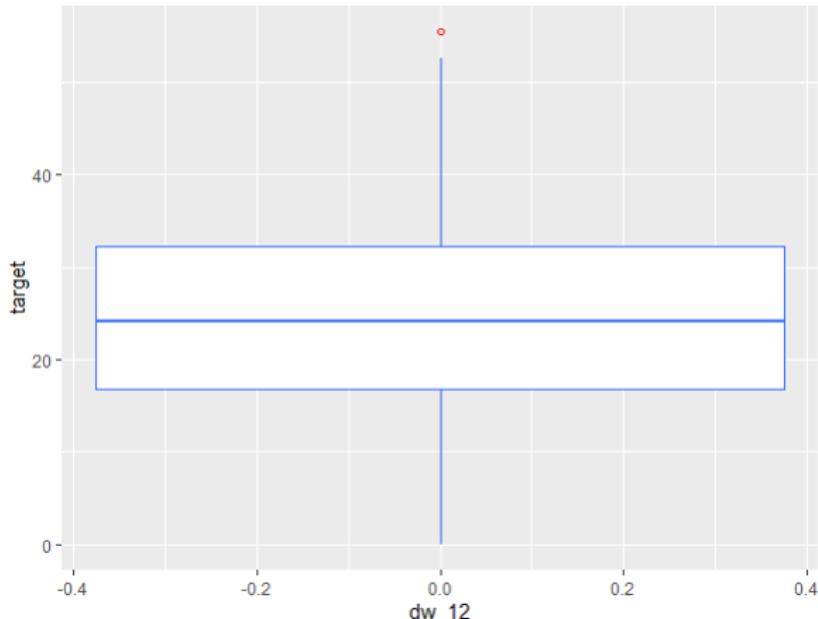
Transform

Tools

Example 1

Lecture 2 - 3

```
> plot1 + geom_boxplot(fill = "white", colour = "#3366FF",  
outlier.colour = "red", outlier.shape = 1)
```



Example 1 with ggplot2 - boxplot options

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

A cheat sheet:

<https://rstudio.com/wp-content/uploads/2016/11/ggplot2-cheatsheet-2.1.pdf>

<https://ggplot2.tidyverse.org/reference/>

ggplot2: Elegant Graphics for Data Analysis (Use R!) by (Hadley Wickham)

Example 2

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Rmd file... Reproducible results, don't save output!

Lecture 2 - 3

STA5092Z - Exploratory Data Analysis

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

Lecture 2 - 3

Getting and Cleaning Data

Introduction

Lecture 1 - 2

Software - R

Basic R
Functions

Datasets

What is EDA?

Graphs

Transform

Tools

Example 1

Lecture 2 - 3

```
> install.packages("swirl")
> packageVersion("swirl") # make sure it is at least 2.2.21
> library(swirl)
> install_from_swirl("Getting and Cleaning Data")
> swirl
```

```
[1] "R version 3.6.0 (2019-04-26)"
```