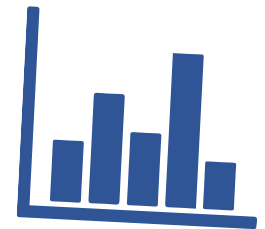# Applied Statistics for Data Scientists with R

# About Me

Md. Ahsanul Islam

Analysis Executive at Kantar Market Research

M.Sc. & B.Sc. in Statistics at University of Chittagong

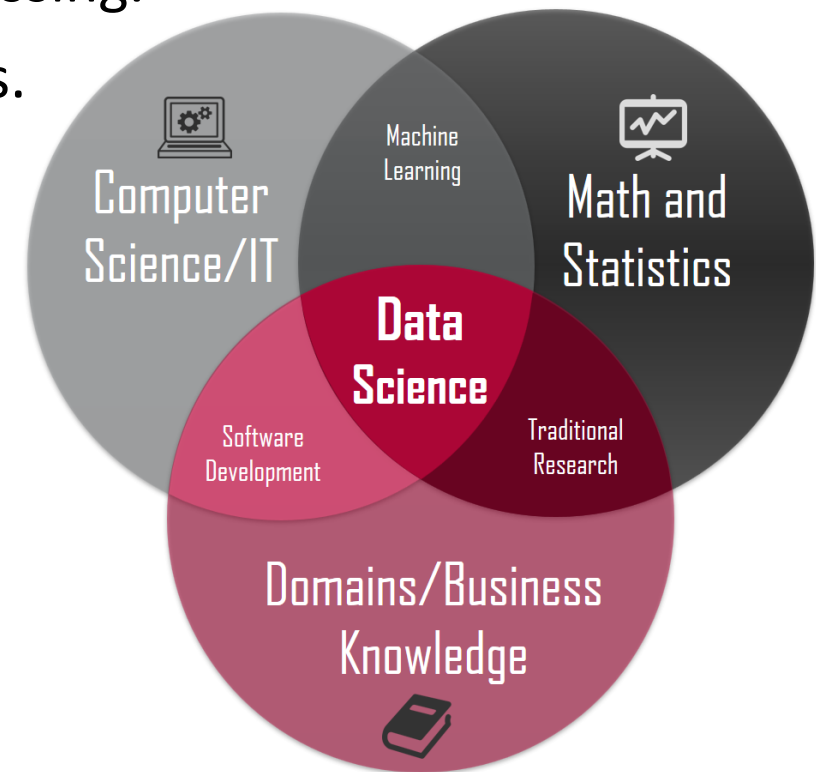Ex-executive (Analysis) at Luminaries Research Ltd.

LinkedIn    GitHub

The capacity to learn is a gift

The ability to learn is a skill

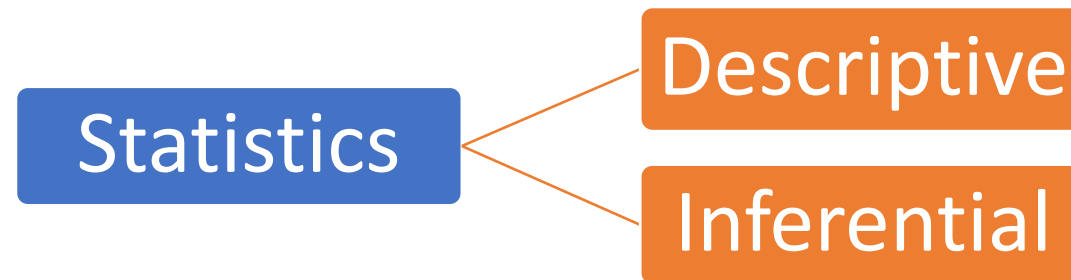The willingness to learn is a choice.


- Brian Patrick Herbert

# Course Prerequisite

- Common sense

- Basic computer skills

- Programming skill is not required

- Ability to use search engines

- Willingness to learn

# Role of Statistics in Data Science

"Statistics is the foundation of predictive modeling, machine learning and AI."

- Exploratory data analysis (EDA) and data preprocessing.

- Statistical Estimation and Optimization Techniques.

- Sampling Techniques.

- Experimental Design.

- Hypothesis Testing.

- Statistical Modeling.

# What is Statistics?

- Statistics is a branch of mathematics that deals with the collection, organization, analysis, interpretation, and presentation of data.

- The focus is in quantitative data to extract meaningful information.

Statistics → Descriptive / Inferential

Note: "Statistics" and "Statistic" are not the same thing.

# A scenario where statistical knowledge helps

- Suppose you are a senior data analyst at a manufacturing company.

- Your boss notices that Machine A's products have fewer defects than Machine B's, but it's not clear if this difference is real or just random variation.

- Your job is to find out if Machine B truly has a problem.

- You may require to use an appropriate sampling technique and statistical test to test the hypothesis.

# Another scenario where statistical knowledge helps

- You are a data analyst for an e-commerce company.

- Your boss wants to predict how much revenue the company will generate during the upcoming Ramadan and what factors might influence the trend.

- This will help in planning marketing campaigns, and inventory management.

- You may require to use an appropriate explainable predictive model.

# Statistics vs Applied Statistics

- **Statistics** is about **creating** the "tools,"
whereas **Applied Statistics** is about **using** those "tools" to address real-world challenges.

# Major Tools Used in Applied Statistics

Applied Statistics for Data Scientists with R

# Why Use R instead of Other Programming Languages?

- Specifically designed for statistics by statisticians.

- Thousands of packages for statistical analysis.

- Excellent for data visualization.

- Easy to learn.

- Mostly academic and research focused.

# What Will You Learn?

- R programming fundamentals

- Data manipulation and visualization

- Applied statistical techniques

- Building predictive models

- Developing Shiny apps


- Extra: SPSS basics for analysis.


*Refer to course outline for details.*

# Usual Analysis Workflow in R

# Usual Analysis Workflow in R



Data file types include:
- CSV
- Excel
- SPSS
- STATA
- JSON … …

# Usual Analysis Workflow in R



Two different ways:
- Base R
- Tidyverse

# Usual Analysis Workflow in R



- Base R
- ggplot2
- plotly
- leaflet
- And many more

# Usual Analysis Workflow in R



A couldn't dare to list anything in this slide.

# Usual Analysis Workflow in R



- PDF from rmarkdown

- Word file from rmarkdown

- HTML file from rmarkdown

- Publish to rpubs.com

- R Shiny Webapp

- Flexdashboard

- Slide using xarigan

# Some Things to Remember During Live Classes

- Take short notes if required

- Stay muted and turn off camera if not asked to turn on

- Ask questions in the chat.
  Answers will be given at the end of each part of the lessons.

- Stay focused on the classes.
  DO NOT BROWSE FACEBOOK or do others tasks during class time.

# To Learn Effectively

- Please attend the live classes regularly

- Reserve some time of your day to practice the codes

- Discussing problems with your peers (classmates) is encouraged

- Share what you have learned with others after each module

- Do projects, solve real problems in your field

- Take good care of your health

We learn more by looking for the answer to a question and not finding it than we do from learning the answer itself.

- Lloyd Alexander

Applied Statistics for Data Scientists with R