# DS 501 Data scientist express bootcamp

*Week 1 [Ella]*

# Copyright Policy

All content included on the Site or third-party platforms as part of the class, such as text, graphics, logos, button icons, images, audio clips, video clips, live streams, digital downloads, data compilations, and software, is the property of BitTiger or its content suppliers and protected by copyright laws.

Any attempt to redistribute or resell BitTiger content will result in the appropriate legal action being taken.

We thank you in advance for respecting our copyrighted content.

For more info:
see https://www.bittiger.io/termsofuse
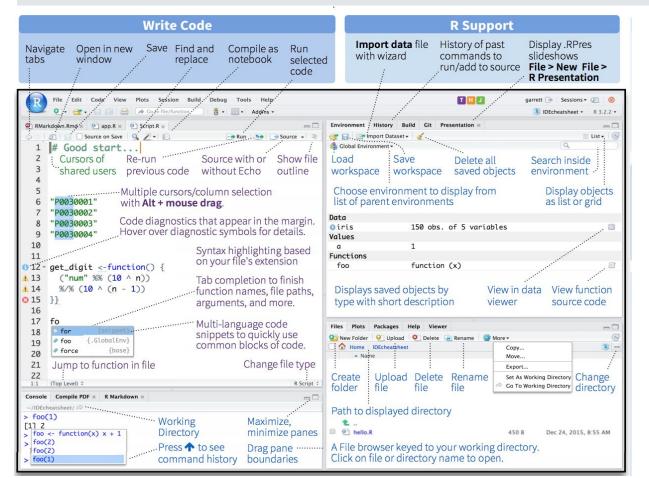and https://www.bittiger.io/termsofservice

## Summary

- Get familiar with R studio

- Common data structures in R

- Data exploration

  - Numeric variables

  - Categorical variables

  - Numeric variable with numerical response

  - Categorical variable with numerical response

  - Numeric variable with categorical response

  - Categorical variable with categorical response

# R Studio



R material list

# Data structure

| | Homogeneous | Heterogeneous |
|---|---|---|
| 1d | Atomic vector | List |
| 2d | Matrix | Data frame |
| nd | Array | |

- Dimensions
- homogeneous: all columns must be of the same type
- heterogeneous: columns can be of different types

## Load data in R

- Download [data](#)
  - [https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data](https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data)

- Prepare workspace
  - Getwd(), setwd()

- [Load](#) in data
  - .txt file: read.table()
  - **.csv file: read.csv()**
  - Excel file: readWorksheetFromFile from **library**
  - Json file, XML file, HTML table,
  - Other stats software files
  - Relational [database](#) and non-relational [database](#)

# Getting to know the data

- Metadata
  - summary(), str(), dim(), head(), colnames()/rownames(), length(), unique()

- Categorical variable
  - table(), barplot(), pie()

- Continuous variable or ordinal categorical
  - **by**(), apply()...
  - mean(), median(), sd(), quantiles(), density(), boxplot()
  - plot(), lines() to visualize results

## How to understand data?

- Data exploration
  - Numeric variables
    - mean(), sd(), quantile(), boxplot(), density()...
  - Categorical variables
    - Sort(), table(), barplot()...
  - Numeric variable with numerical response
    - Cor, library(corrplot)...
  - Categorical variable with numerical response
    - boxplot(), by(), apply(), library(lattice)...
  - Numeric variable with categorical response
    - boxplot(), library(tabplot)...
  - Categorical variable with categorical response

## Summary

- Get familiar with R studio
- Common data structures in R
- Data exploration