

Doing Data Science

Unit 1

Faizan Javed
DataScience @ SMU

Grading:

Live Session Assignments (50%) (8 homeworks)

HW1 is due next Tuesday, 1 hour before live session

Case Studies (30%),

Videos and Questions during asynchronous material (BLTs) (15%),

Live Session Attendance (5%).

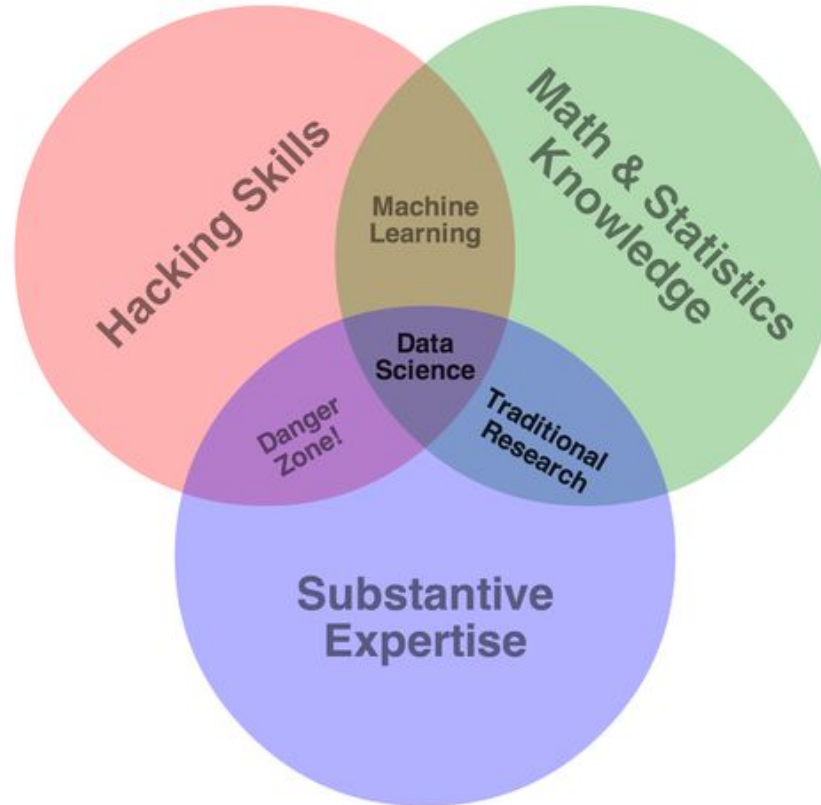
What is Data Science?

Reproducible vs Replicable

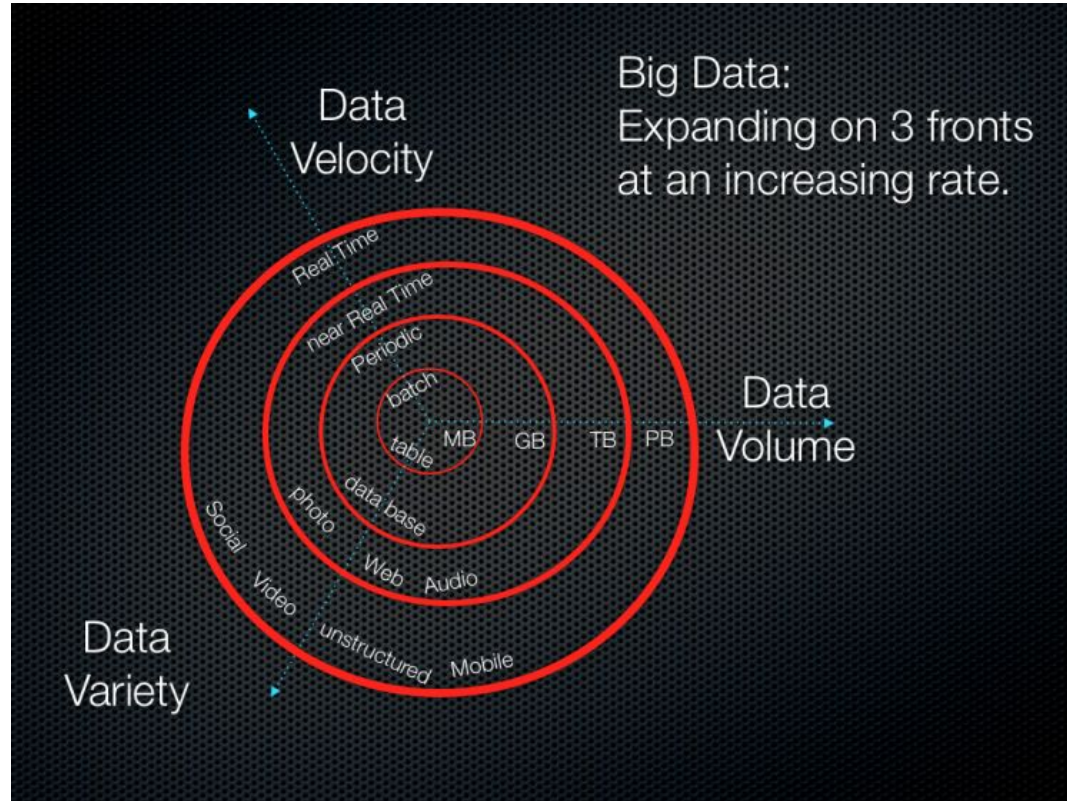
Tools for Data Science

Data Science Venn Diagram

https://s3.amazonaws.com/aws.drewconway.com/viz/venn_diagram/data_science.html



Big Data: the 3 Vs



Data Science profile

Computer Programming

Math

Statistics

Machine Learning

Domain Expertise

Communication, Presentation Skills

Data Visualization

Reproducible research

Use specific set of computational functions/analyses (usually specified in terms of code) and original data to recreate findings

Replicable: recreate findings with new data

A study is only replicable if you perform the exact same experiment (at least) twice, collect data in the same way both times, perform the same data analysis, and arrive at the same conclusions.

Why is replicability important?

Open Source

Public research

Business/Private Sector

Tools for reproducible research

R

Knitr

RMarkdown

RStudio

Section 1.5.1, install R and RStudio, and all packages on page xix

Workflow of Reproducible research

Data Gathering

Data Analysis

Results Presentation

Machine learning

Training data (e.g. implicit and explicit feedback on job ads)

Target function (e.g. probability of user clicking and/or applying for a job)

Metric (e.g. precision vs. recall, or any ranking metric that correlates to AB test metrics)

Practical tips

Document everything

Ensure compatibility of your software/libraries/packages
(troubleshooting assistance)

Comment your code (variable/argument names should not be cryptic)

Source code comment header

R data structures

<https://jamesmccaffrey.wordpress.com/2016/05/02/r-language-vectors-vs-arrays-vs-lists-vs-matrices-vs-data-frames>

Vector: collection of items of fixed size (all of same type)

List: can hold items of different types and list size can be increased on the fly

Matrix: 2-D vector (fixed size, fixed type)

Array: a vector with one or more dimensions

Data frame: each column has items of same types and columns can have headers

Getting started with R

#Print R session info -- why is this useful?

sessionInfo()

c(combine) function : create vectors

NumVec ← c (2, 3, 4)

CharVec ← c (“doing”, “data”, “science”)

cbind()/rbind() : combine vectors side-by-side

```
StringNumObj ← cbind(NumVec, CharVec)
```

data.frame() : create an object with rows and columns

```
StringNumObj ← data.frame(NumVec, CharVec)
```

Reassign row.names

```
row.names(StringNumObject) ← c("First", "Second", "Third")
```

\$: component selection

NewNumeric ← StringNumObj\$NumVec

head()/tail() : select first/last few rows

head(cars)

[rows,columns] subscript operators, : sequence operator

cars[3:7,]

str(..)

compactly display the structure of an R object

E.g.,

3 obs. of 2 variables:

\$ NumVec : num 2 3 4

\$ CharVec: Factor w/ 3 levels "data","doing",...: 2 1 3

summary(..)

display summary statistics for analysis (mean, quantiles, etc)

E.g.,

NumVec CharVec

Min. :2.0 data :1

1st Qu.:2.5 doing :1

dim(..)

retrieve or set the dimensions of an object (array, matrix, dataframe)

E.g.,

[1] 3 2

Missing/extreme values:

NA = not available

NaN = undefined

Inf = extremely small/large (infinity)

What did you learn today?

Questions?