

Feb 1st 2025

①

Data Science as a field. "Jane Wall"

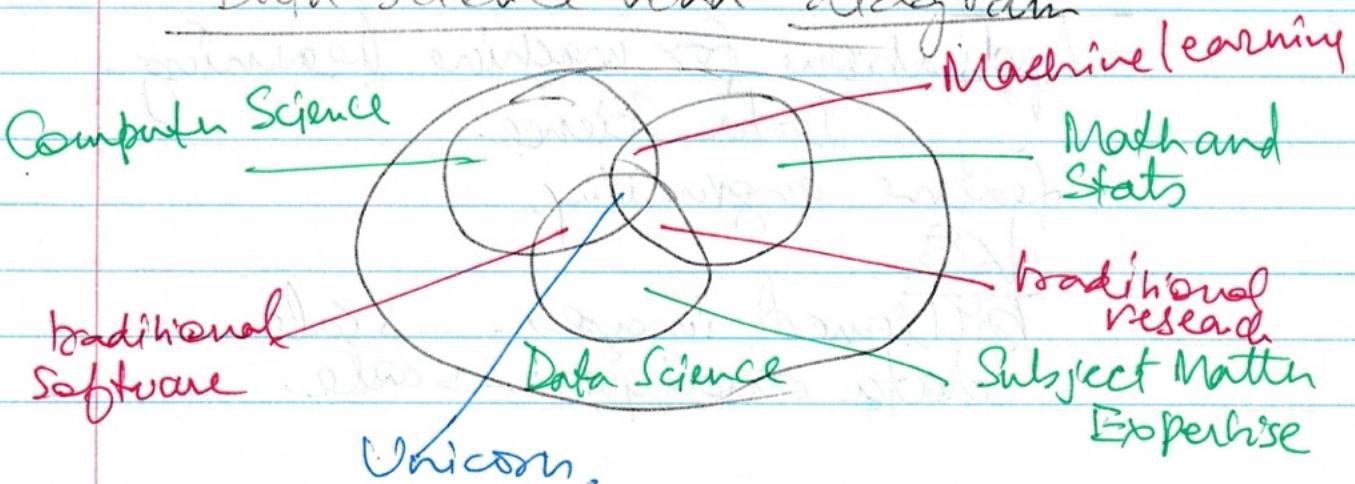
- Data Science applications skills
- → electives to be needed.
- Data Science - reproducible analysis.
- Visualizing.
- Analyzing.
- Pitfalls
 - bias., sources, mitigates
 - ethical.
- communicating results.

Data Science - Origin.

- our ability to count and categorize things
- capacity of data storage has grown
 - capacity to process data has grown.
 - DS as a field has grown.

Fields - Text, Music, click, Purchases
Surveys - Analyse all of it.

Data science Venn diagram



Variety of jobs in Data Science

- getting and cleaning data
- Visualization and Analysis
- Shiny apps / Dashboards
- Modeling & Prediction
- Machine Learning, AI
- Natural language processing
- Ethics & privacy policy.

How does Data Scientists spend their time

- Data Loading - 19%
- Data cleaning - 26%
- Data Visualization - 21%
- Model selection - 11%
- Model training, and saving - 12%
- Deploying Models. - 11%

- Truth from data.
- No perfect answer.

Applications for machine learning.

Data Science.

feature engineering.

LS.

Billions of images. - scale.

data analysis - scale.

③

- What does Google look for when they hire DS
- strong in one area.
 - or
 - balanced.
 - able to explain the concepts in layman terms.
- or product development lifecycle.

Data Management

Data Wrangling

Survey - 2000 responses.

R Package

Recommender System -

Spacial Data Scientist

Data Science on a Map.

RStudio → IDE & GUI

* Rmd (R markdown)

> getwd()

→ R script → *.R

Introduction to R-Basics

Basic Math

- much like calculator.
- Remainder of division.

$$5 \% \% 2 = 1$$
- # This is comments.
- Assignment operator \leftarrow

$$\text{Variable} \leftarrow 100.$$
- Variable starts with small letters
- Functions starts with Capital letters

R Data type

- Numeric
 - floating point
 - Integer.
- Logical
 - Boolean TRUE/FALSE or T/F.
- Character.
- Class - Same as type in python.
 - class(TRUE) = logical.

- Vector - 1d array.
 - created using combine function .(c)
 - $\text{hvec} \leftarrow \text{c}(1, 2, 3, 4, 5)$
 - class(hvec) \Rightarrow numeric.

(5)

> `cvec <- c('U', 'S', 'A')`

> `class(cvec) > character.`

- Cannot mix data types in vectors.

> `v <- c(TRUE, 20, 40)`

> ~~class~~ `v > [1] 20 40`

> `class(v) > numeric.`

- another example.

> `cvec <- c('USA', 20, 40)`

> `cvec > 'USA', '20', '40'`

> `class(cvec) > character.`

- `temp <- (72, 71, 68, 73, 69, 75, 76)`

> `temp > 72, 71, 68, 73, 69, 75, 76`

> `names(temp) <- c('Mon', 'Tue', 'Wed', ...)`

> `temp >`

Mon Tue Wed

72 71 68

> `days <- ('Mon', 'Tue', 'Wed', ...)`

> `names(temp) <- days.`

> `temp >`

Mon Tue Wed

72 71 68

Vector Operations.

> created vector operations. R under MAC

Vector Indexing and Slicing.

> indexing starts from 1

(6)

Stack overflow

Getting help.

> `help('vector')`

> `?vector`

> `help.search('vector')`

Dataframe

> `data()` → all available dataframe in

> `head(state.x77)` → display first few rows

> `tail(state.x77)`

> `str(state.x77)` → structure of dataframe

> `Summary(state.x77)`

Data Manipulation

> Dplyr package - Manipulating

> Tidyr package - cleaning data

> Pipe operator %>%

Data visualisation - ggplot2

— Built on the idea of adding layers

— First 3 layers are

- Data — `ggplot(dataframe)`

- Aesthetics — columns/features

- Geometries. — `geom_point()`

- Facets — multiple plots

- Statistics — `stat_smooth()`

- Coordinates

- theme

Visualisation

- Coordinates and - Resize of our Coordinate
- faceting - have multiple plots.

Bias

- feeling.
- How are you getting Survey.
- Outliers. - Ignore them? Introducing data.
- ML algorithm.
- Recognize bias - mitigate

Ethics → Why?

- It impacts many people or everyone.
- facial recognition
- search on web browser
- ethical issues related to data
- ethical issues related to model results from analysis and modeling that data using ML.
- ethical issues within Data Science.

NYC Shooting data

incident_id ✓
 occurr_date (char)
 Occur_time (S:thms)
 BORO - (char)

28,562 rows. 21 columns

When records are read using `read_csv`, above details were obtained.

Occur_date ✓

Occur_time ✓

Mindingflag ✓

Perpetrator's age

Sex

Race

Victim's age

Sex

Race

→ 3 profiles.

Perpetrator's profile

Victim's profile

Location profile.

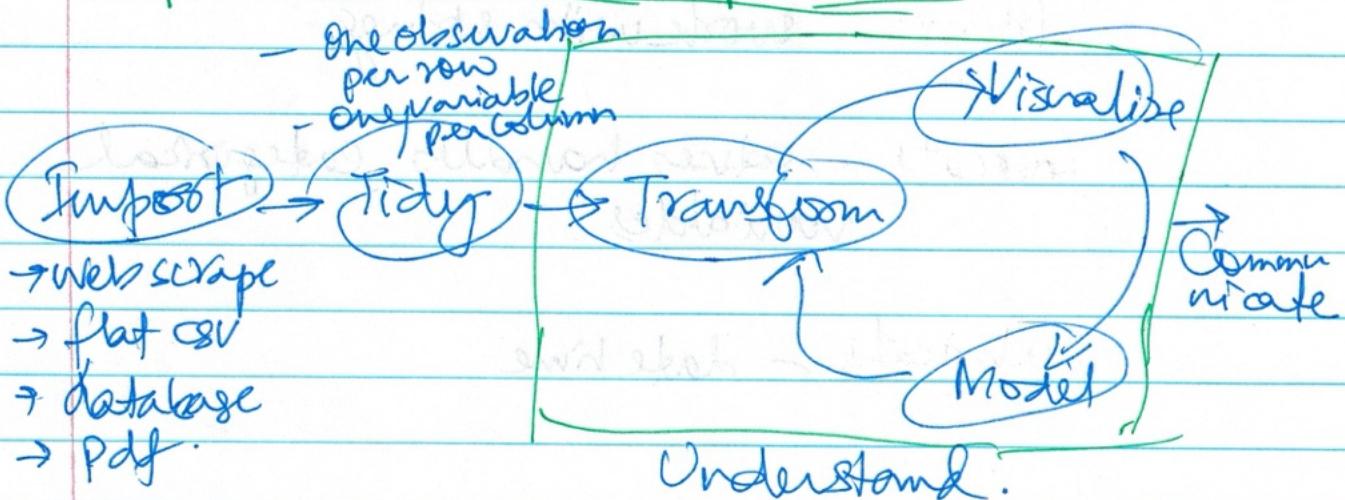
Reproducible Research

- Why?
- ① Makes your conclusion carry more weight
 - ② Others can build off your work
 - ③ Do it for your future self

How?

- ① Document everything
 - start with question and raw data
- ② All files are text files
 - platform independence
 - future proof
- ③ Make code readable
 - comments
 - formatting (indentation)
 - meaningful names.
 - each script does one task.
- ④ Use relative path.

Steps in data science process



Tidyverse packages

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

ggplot2 - graphics

dplyr - data manipulation

tidy - tidy data

readr - read rectangular data

purr - enhances R's functional
programming

bibble - modern re-imaging of the
data frame

stringr - work with strings

forcats - solves handles categorical
variables

lubridate - date time

Covid

- > library(tidyverse)
- > urls ← str_c(url_in, filenames)
 - ↳ vector of urls.
- > create data frame variables


```
global_cases<- read_csv(url[1])
```
- > different ways to explore datasets .
 - just print
 - head()
 - str()
 - summary() - Summary by column/variable. (very useful)
- > use of pivot_longer - to convert columns into single column , and values to another column.


```
name_to =  
values_to.
```

pivot_longer - reshape data from a wide format to a long format.

syntax :

```
pivot_longer(data, cols, name_to, values_to)
```

data - data frame

cols - columns you want to pivot This specifies which columns should be collapsed into key-value pairs.

- names_to - name of new column

(new_name - values_to - values)

→ then join data. (death data and
cases data)

→ global <- global.cases %>%

fill_join(global.deaths) %>%

rename(Country_Region = "Country Region")

Province_State = "Province (state)" %>%

mutate(date = mdy(date))

mutate - is used to create new
column or modify existing coln.

example - new column.

df contains x and y

df <- df %>%

mutate(z = x + y) creates z

example - modify existing column.

df <- df %>%

mutate(y = y^2).

library(tidyverse) - for using

date = mdy(date)

- > global & global %>% filter(cases > 0)
- > Then analyse data and then visualize it.
 - > group by is used for categorical data.
 - > Summarise - used for creating summary statistics.
- > visualisation
 - Covid 19 in "US"
 - Covid 19 in "new york"
 - new cases.
- > Use of lag() function - used to create a lagged version of a vector, shifting elements by a specified number of positions. It is used in time series analysis to compare current values with past values.
- > exercise .
 - take the case of different country.
 - worst and best state.
- > linear model .
 - lm.

Bias

Add bias in the project report.

NYPD Shooting

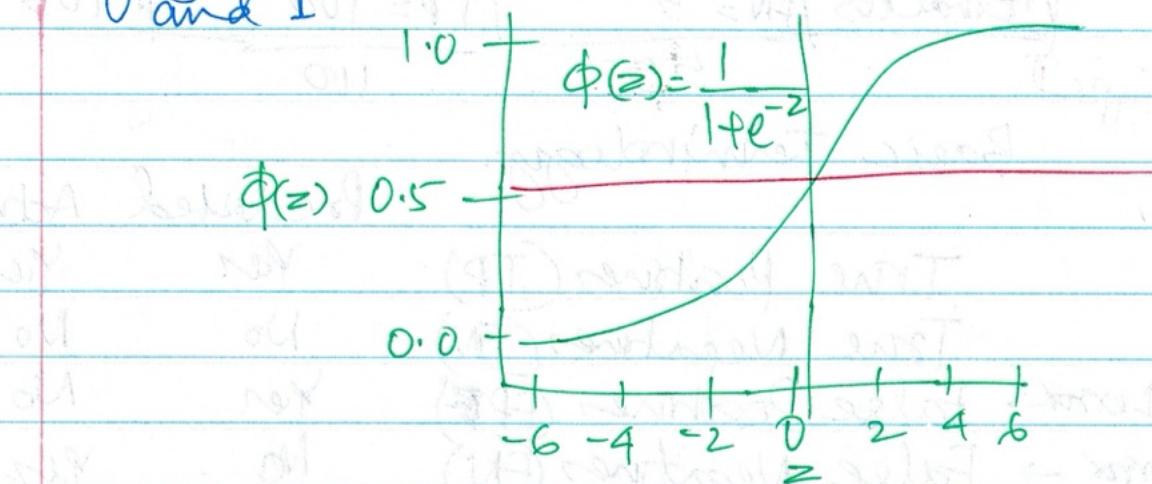
- Complete all the steps in data science process in reproducible manner.
- Took at NYPD shooting data and analyse
- grading
 - Data science process
 - Reproducibility,
 - clear writing
 - validity of analysis
 - identification of bias.
- Rubric

✓ import	Reproducible.
✓ tidy	visualisation
✓ analyse.	analysis .
✓ 2 visualisation	
- bias bias	
- the model.	
- upload html / pdf.	

→ plots are good way to find correlations

Logistic Regression

> The sigmoid (aka Logistic) Function takes in any value and outputs it to be between 0 and 1



$\phi(z)$ is always between 0 and 1.

Linear Model $y = b_0 + b_1 x \Rightarrow$ transforms into

$$\text{Logistic Model } p = \frac{1}{1+e^{-(b_0+b_1 x)}}$$

> We usually set a cutoff point (say 0.5)
 anything below probability is (0) zero
 anything above probability is (1) one.

Evaluation of classification models

(16)

> Confusion Matrix:

		PREDICTED NO	PREDICTED YES	Type I
		TN = 50	FP = 10	60
		FN = 5	TP = 100	105
Type II		55	110	

Basic Terminology:

Predicted Actual

True Positives (TP) Yes Yes.

True Negatives (TN) No No

TYPE I error → False Positives (FP)

Yes No

TYPE II error → False Negatives (FN) No Yes.

Accuracy → how often it is right

$$\text{Correct} = \frac{\text{TP} + \text{TN}}{\text{Total}} = \frac{100 + 50}{165} = 90\%$$

Misclassification Rate - how often it is wrong.
(error rates).

$$\text{Misclassification} = \frac{\text{FP} + \text{FN}}{\text{Total}} = \frac{15}{165} = 9\%$$

Communicating Results

- Reports
- Presentations
- Elevator Pitch

Written report -

- be organised
- Easy to read - tell the story
- key visualizations
- Format, others can open
- good grammar / sentence structure
- Why should I care?
- Where is the data from.
- What does it tell you - analysis and visualization
- What do you conclude.
- How could you be wrong.

Presentations -

- don't be boring
- be organised
- tell the story
- key visualizations
- Make eye contact (or camera)

Presentations (DONT)

- Don't read slides.

- Don't have lots of words on slide
- Don't be monotone
- Don't go over time limit.
- Don't overwhelm with many similar visualizations.

Presentations (DO)

- Include pictures to make your main point
- Target talk to audience
- Be enthusiastic
- Practice
- Maintain eye contact
- Use whiteboard or other option

Elevator pitch

- Explain to non-technical person.
- Target is 13 year old.
- Short approx: 1 minute

Build portfolio

→ GitHub

Meetup

→ Data Science workshop

25th FebSequenceQuestion of Interest

Import Libraries

Import data.

Read global & US data.

Explore Columns

Tidy and Transform Global data.

Tidy and Transform US data.

Visualization of US data-

- US
- state - New York
- State - Arizona
- New cases per day US
- New Deaths - State
- Best state
- Worst state

Bias Sources

Political affiliation

Population density.

Extent of lockdown.

Climate of area.