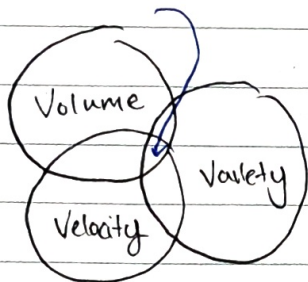


# Wk 1 - lecture

## Data Science Skills

- Math & Science
- Programming & Statistics
- Domain knowledge
  - ex) ↓ costs through route optimisation <sup>greater impact considering more vehicles, time</sup>
  - ↳ get knowledge from experts, domain changes by research topics.
- Communication

## Big Data



Volume: what size

Velocity: how fast data <sup>added</sup> ~~changing~~

Variety: how various the structure is

<sup>biggest challenge</sup>

how many ID vids  
twitter twits  
review  
csv, json, text, image, emails, etc  
structured semi- unstructured  
≈ 80% of data  
from enterprises

never ending  
with manual processing

ex) Panama paper leak  
(TBS of data)

From certain volume, manual processing  
is NOT AN OPTION.

## Data-Driven Journalism

data → filter → visualize → story

Data can be

- too big
- too fast
- needs to be combined from diverse sources
  - ↳ scripting &

## Why Python?

R: for research

- interpreted  $\times$  compiling Python: for industry projects
- dynamically typed  $\leftarrow$  data types may vary (errors occur in Java)
- readable 😊 for devs
- deployable  $\leftarrow$  server-side language
- productivity  $\leftarrow$  facilitates rapid & interactive prototyping

## Most common IDE

- Jupyter Notebook  $\leftarrow$  download in Anaconda & preferred.  
or google colab  
or ~~heroku~~ uni server
- VS Code
- R studio
- etc
- ✓ live code share
- ✓ code by 'cells'

## Importing packages.

```
import csv (as —)
```

```
from csv import DictReader
```

) b easy importing? 📖

## Open-source libraries

- scipy: scientific & technical computing - scipy.org
- numpy: arrays & matrices
- matplotlib: visualization
- scikit-learn: machine learning
- nltk: natural language processing
- pandas: R-like dataframe manipulation