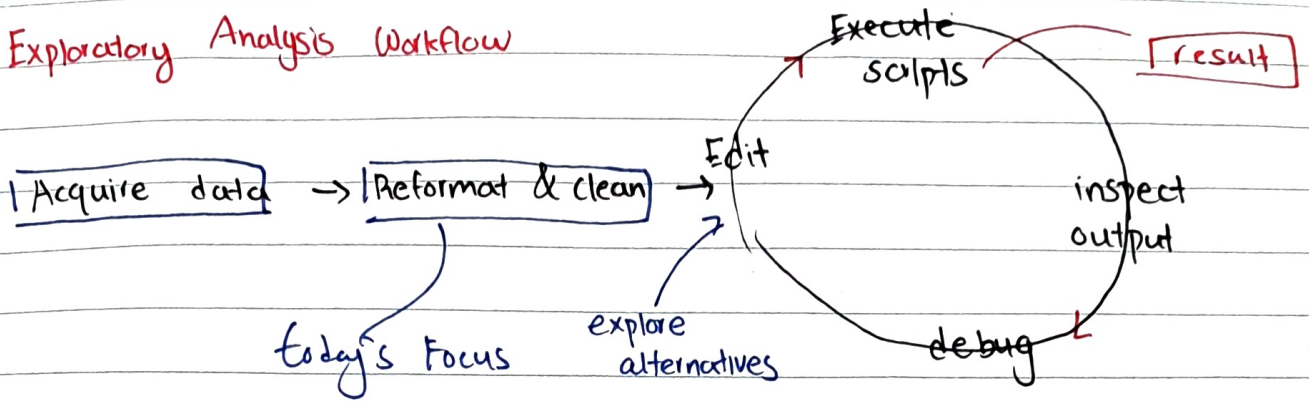


Exploratory Analysis Workflow



Level of Measurements & Types of data

- Categorical (Qualitative)
 - Nominal ex) gender, state
 - ordinal ex) [low, medium, high]
 - Nominal vs ordinal: distance cannot be implied (ordered)
- Quantitative
 - interval: values encode differences, equal interval between variables. No True zero eg. celsius
 - ratio: zero is defined eg. kelvin, experience, weight, pressure

Other

- images, text, videos, etc

~~mode~~
~~median~~
~~mean~~

Countable, equally defined

Order defined

Difference defined (+, -)

Zero defined (x, ÷)

	Nominal	Ordinal	Interval	Ratio
Countable, equally defined	✓	○	○	○
Order defined	✗	○	○	○
Difference defined (+, -)	✗	✗	○	○
Zero defined (x, ÷)	✗	✗	✗	○

Nominal

Ordinal

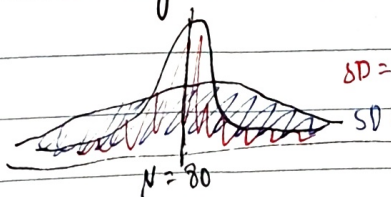
Interval

Ratio

mode	median	mean	count / distribution	min/max	range	percentile	variance, sd
------	--------	------	----------------------	---------	-------	------------	--------------

Q variance / SD ?

- Looking at mean (only) can't tell the performance of a particular student



SD=10 (most are on the same line, few extremes)

SD=30 (more dispersed)

Acquiring Data & cleaning

Acquire from

- File Access (csv, json, excel, XML)
- Programmatically (web scraping, APIs of web services like twitter)
- Database Access
- Collected data (survey)

Clean by (real data often dirty)

- type and name conversion
- filtering of missing/inconsistent data
- unifying semantic data representation
- matching of entites from different sources
- handling of date & time
- scaling & normalization, (advanced)

data collected from different source
may have different name, format, etc

centralise around mean (mean=0, sd=1)

optional dimensionality reduction

for visualisation
often

(2D or 3D)

choose specific vars.
(can't display all vars.)

cleaning with Python

- * look at the data first! before any operations.
- * keep the original data! before modifications.

CSV module

- csv.reader: rows \rightarrow arrays (column num = idx)
- csv.DictReader: rows \rightarrow dict (column name: rows)

pprint (pretty print)

alternative: pandas

- \hookrightarrow open source
- \hookrightarrow data import and analysis functionality
 - \hookrightarrow tabular (df) 2D
 - \hookrightarrow time-series (series) 1D
 - \hookrightarrow matrix ~~1D~~ 2D, 3D... ND

functions

display first/last n rows: head, tail

missing values?: dropna, fillna, replace (null) \rightarrow 0

Cleaning data

- both csv/pandas will likely read-in data as string.
- we need type conversion or else error (mean of strings?)

ex) `data[column] = data[column].astype(int)`

(
or
)

fails if any entry raise Exception
(NaN \rightarrow !int)

define our own function

to loop thru row & ignore NaN (numpy.nan)
(unrecognised values \rightarrow np.NaN (Default))

Descriptive statistics with Pandas

count, sum, abs

min, max

mean, median, mode, std

`data = pd.read_csv('_____'.csv)`

ex) `data['power'].min()
.std()`

prod(): product of values, cumsum(): cumulative sum, cumprod(): cumulative product

groupby(), apply(), applymap()

Filtering

`data.loc[data[col] == '_____', col]`

`data.loc[data['type'] == 'wind-turbine' & data['power'] > 100]`

select

= wind generator with power > 100 mw

`data.groupby('class').size()`
group-by column

Numpy: calculate central tendency & dispersion

`np.nanmax(data) - np.nanmin(data)`

v = row

nan values ignored automatically

`iqr = np.nanpercentile(v, 75) - np.nanpercentile(v, 25)`

Visualising with python

(matplotlib)

```
distribution = data.groupby('fueltype').count().reset_index(name='numstations')
```

create frequency
distribution

↳ \approx select

AS

in sql

```
plt.bar(distribution['fueltype'], distribution['numstations'], alpha=0.5, align='center')
```

• title()

• xlabel()

• ylabel()

• grid()

• xticks(rotation=40)

```
data['power'].hist(bins=10, rwidth=0.9)
```

```
sub = plt.subplot()
```

```
data.plot.scatter(x='numGen', y='power', c='Darkblue')
```

```
sub.set_xlim(0, 50)
```

look for correlation with scatterplot

Correlation statistics to measure dependence

• parametric

pearson's r value

• non-parametric

spearman's rho

kendall's tau

scipy