

## Specialization - 02

### # Week 02 → "Feature Engineering"

- Mapping Raw data into feature
- Mapping numeric value & categorical value
- Empirical Knowledge of data.

→ processing step

Data cleansing / Feature tuning / Representation transform / Feature extraction /

→ Feature Engineering techniques

Feature construction

\* Numerical Range {  
• Scaling  
• Normalizing (min-max)  $x \sim [0, 1]$   
• Standardizing (z-score)  $x_{std} = \frac{x - \mu}{\sigma}$   
 $x_{std} \sim (0, \sigma)$

\* Group {  
• Bucketizing  
• Bag of words

\* Dimensionality Reduction {  
• PCA  
• t-SNE  
• UMAP

Tools to use  
→ Tensorflow  
Embedding  
projection

① Apache Beam run on top of other frameworks like (spark or Flink) on google cloud data flow.

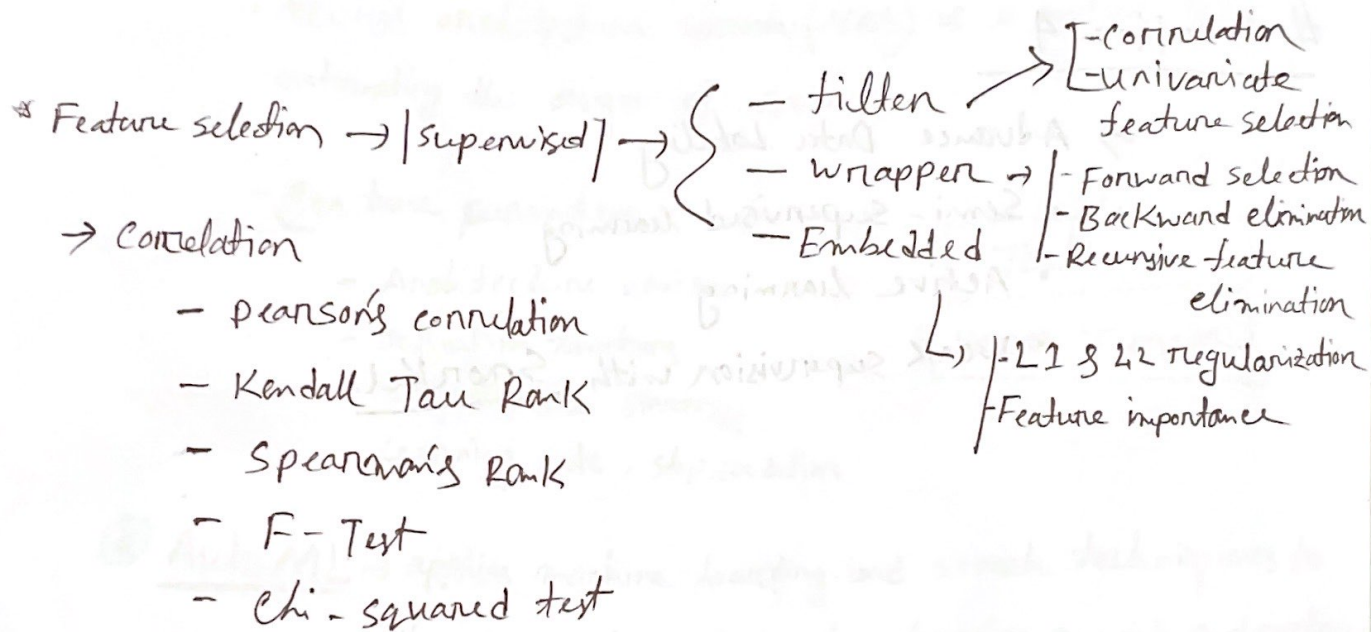
Tools

↳ Tenzent low

Transformation

# Apache Beam →

② Data Affected by → Seasonality, Trend, drift



# week-03 (Data journey)

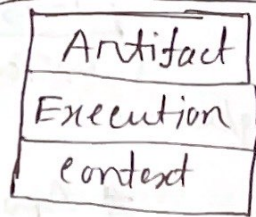
→ Data versioning

- Version control of dataset (Keep track of dataset change over the period)
- Tools = DVC, Git-LFS.

→ Environment version → Docker, Terraform

→ M2 Metadata → Helps to track the data in order to change anything.

ML metadata is a library for recording and retrieving metadata associated with ML production pipelines among other applications.



# Data storage.

# week-04

→ Advance Data Labeling

- Semi-supervised learning
- Active learning

• weak supervision with Snorkel.