

[① Programming languages]

→ Python.

→ R

→ Seaborn, matplotlib...

→ Java

[② Data visualization]

→ Power BI, Tableau:

→ Classification

→ Regression.

→ Reinforcement.

[⑤ Math Statistics]

[Data Science]

[② Machine learning]

→ Deep learning

→ Neural Networks

→ Dimensionality Reduction

→ Clustering.

[④ Web Scraping]

→ BeautifulSoup

→ Scrapy

→ URLLIB

[③ IDE]

→ Jupyter Notebook

→ Spyder.

→ PyCharm.

[⑧ Data Analysis]

→ Feature Engineering

→ Data Wrangling

→ EDA

Classification

fixed category

→ Binary classification. (2 cat)

→ Multiclass classification. (more than 2)

→ Multilabel classification



movie

Action Romance Thriller Comedy

Endgame

1

0

1

1

Dredd

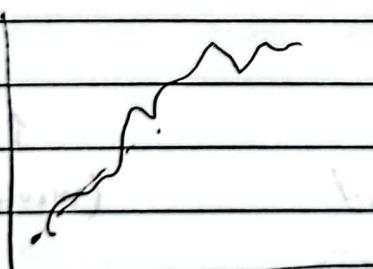
0

0

1

1

forecasting → Regression.



time Series analysis.

Regression → Predict Continuous value.

Unsupervised Machine Learning

→ Clustering

→ Segmentation

→ Reduce Dimension

Semi-Supervised ML

label data → Model → Recommendation (w/o)

→ Reinforcement learning

continuously learning

and new recommendation.

Ex: children learning

way in semi-supervised

(movie recommendation)

or Reinforcement learning.

Step
↓

① Programming Language

→ Python

→ R

→ Julia.

② Exploratory Data Analysis.

③ Feature Engineering

④ Feature Selection.

⑤ Machine Learning Algorithm.

⑥ Regression, Classification, clustering.

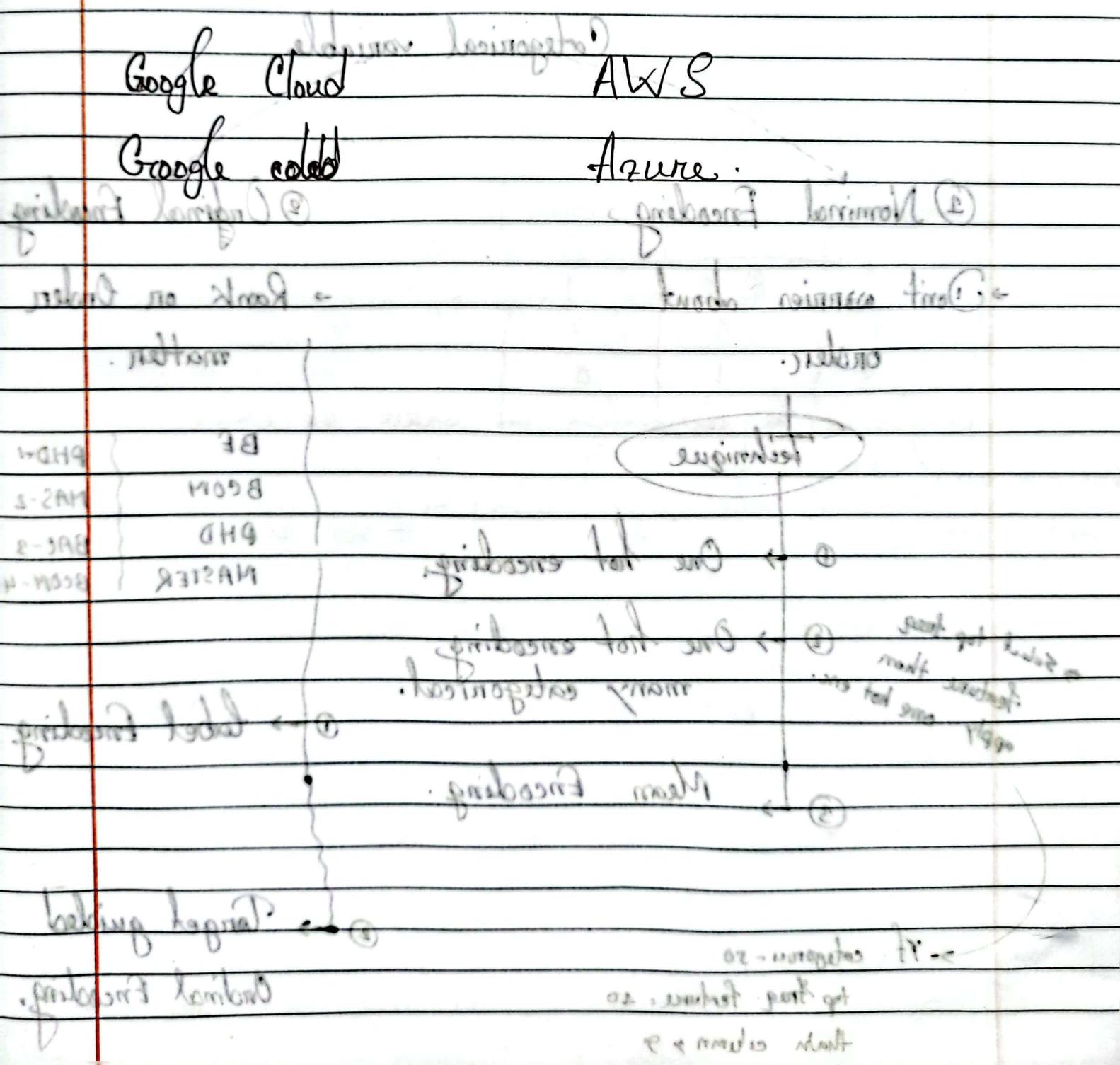
Linear, Logistic Regression, Decision tree

Random Forest, K-means -----

⑦ Hyperparameter Tuning

GridSearch, RandomizedSearch, Hyperopt
Optuna, Genetics Algorithms.

- (9) Docker And Kubernetes.
- (10) Model Deployment
- (11) End to End ML Projects.



Types of Encoding

categorical variable

as like. Gender

Female }
Male }

Categorical variable

① Nominal Encoding

→ Don't worry about
order.

② Ordinal Encoding

→ Rank or Order
matter.

technique

① → One hot encoding.

② → One hot encoding
many categorical.

③ → Mean Encoding.

BF	PHD-1
BCOM	NAS-2
PHD	BAL-3
MASTER	BCOM-4

→ Select top freq
feature than
apply one hot enc.

① → Label Encoding.

→ If. categories = 50
top freq. feature = 40
than column > 9

② → Target guided
Ordinal Encoding.

One-hot encoding

→ 4 Category

$$4-1 = 3 \text{ column}$$

→ 5 Category

$$5-1 = 4 \text{ Column}$$

One column skipped

(y) Yes (no) (can't)

Suppose State

German

Spain

France

German Spain France

0	0	0
0	0	1
0	1	0

we can skip one column.

order primary alt. fibering at below is kind ⑤

(fibering is at multiple nationally flags etc)

and the flag is not aligned ⑥

primary becomes new flag ⑦

thus all will align in given

④ Why we need feature scaling?

When

and

why

we should perform feature Scaling.

And

when should
not

why
not

→ Handle missing values of "categorical variables".

① Ignore the observation or Delete the row.

if we have large dataset → if not this is not good way.

② Build a model to predict the missing value
(can apply classifier algorithm to predict.)

③ Replace with the most frequent values.

④ Apply unsupervised learning
using an algorithm like KNN.

High cardinality \rightarrow In a Column \rightarrow Categorical variable
 high num of unique values. \downarrow high label category.

Related to the time series data.

Decision Tree.

$$\text{Entropy}(S) = \sum_{i=1}^n -P_i \log_2 P_i$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum (\text{Prob} \times \text{Entropy})$$

flow \rightarrow higher Information Gain \rightarrow more entropy removed

Partition more criteria \rightarrow Target \rightarrow Target variable.

Partition more criteria \rightarrow Target \rightarrow

than \rightarrow

Information Gain \rightarrow root node.