

Data → Data Preparation

- Data Cleaning
- Data Transformation
- Data Splitting

Feature → Feature Engineering

- Feature selection
- Feature Creation
- Feature Transformation
- Dealing with missing / inconsistent data

Training & Testing set

- 80-20 split b/w training data & testing data.

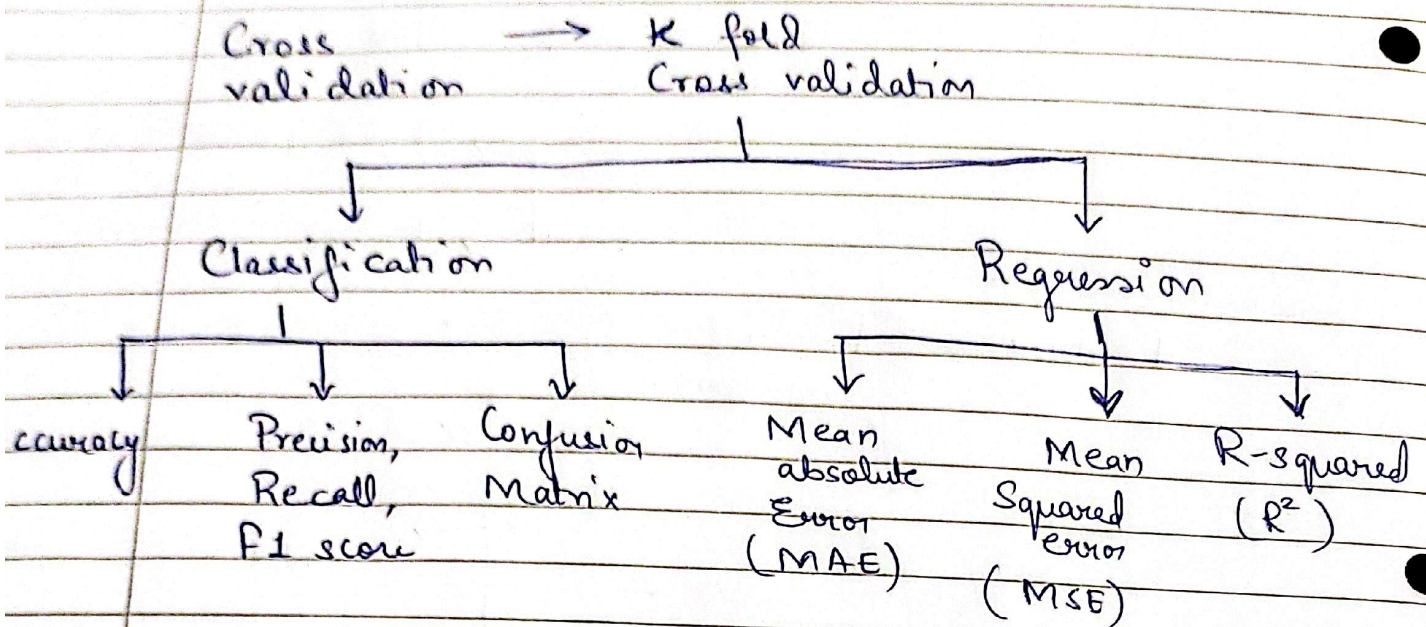
* validation set



helps tune ^{hyper}parameters without affecting dataset

- 70-30 split b/w training data & testing data.

- 60-20-20 split b/w training, test & validation



Model & Algorithm

Algorithms → steps or procedures used to learn patterns from data.

They define how the model will adjust its internal parameters to improve its predictions.

Model → o/p of the algorithm after training on data.

Example Project	Data Type	Feature Type	Label Type	Model Type	Algorithm Type	Training, Testing + Validation
Sentiment Analysis	Text	Word embeddings (eg. TF-IDF, BERT)	Categorical eg. (positive, negative)	Text Classification	<p><u>Naive Bayes</u>: Calculates the probability of each class based on feature occurrence assuming feature independence.</p> <p>Transformer Transformer (BERT): Uses attention mechanisms to capture long range dependencies & word context</p>	<p>Train, test \neq split with k fold cross validation</p> <p>Transformers may use pre training + fine tuning approach on large datasets..</p>

Example Project	Data Type	Feature Type	Label Type	Model Type	Algorithm Type	Training & validation testing
Customer Segmentat ⁿ	Tabular	Categorical Continuous	Categorical	Clustering	<p><u>k-means</u> :</p> <p>Minimizes within cluster variance by assigning data points to nearest cluster center</p> <p><u>Hierarchical Clustering</u> :</p> <p>Merges data point iteratively based on similarity</p>	<p><u>Hold out validation</u> on a separate test set to assess model performance.</p> <p><u>Elbow method</u> used for optimal cluster number selection.</p>

Example Project	Data Type	Feature Type	Label Type	Model Type	Algorithm Type	Training, testing & validation.
Price Prediction	Tabular	Numerical Features	Continuous	Regression	<u>Linear Regression:</u> Finds a linear relationship b/w features & target <u>XGBoost:</u> Uses boosted decision trees to capture non linearities in data.	Train - test - split MAE & RMSE used as metrics.