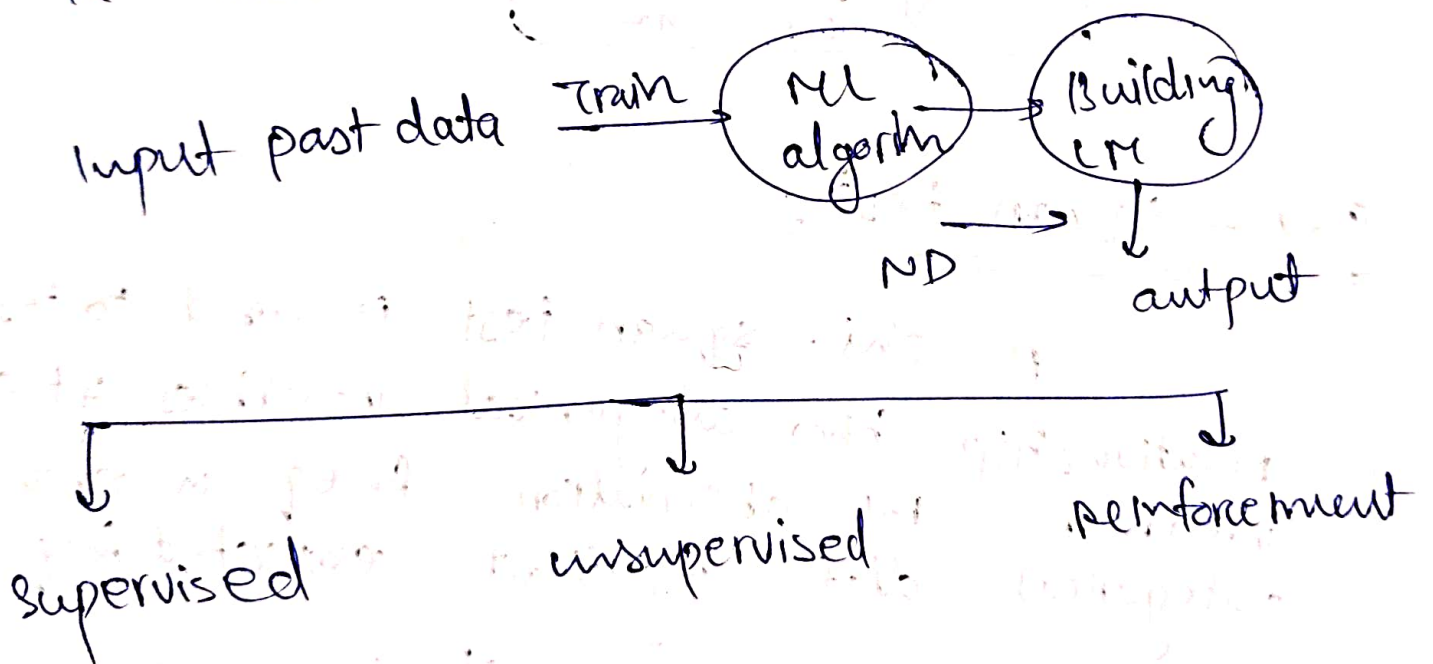
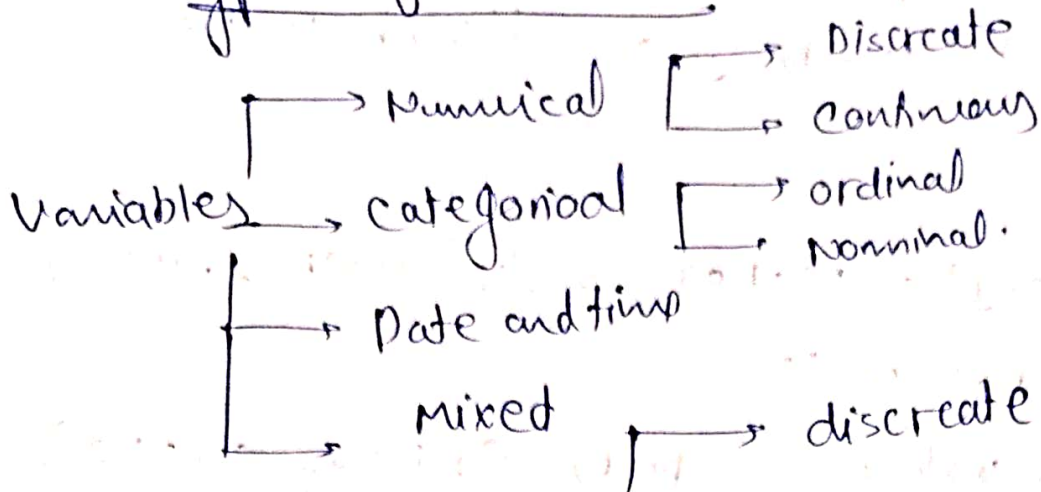


Machine learning :

getting comp learn and act like human do
and improve their learning over time
in autonomous fashion.



Types of variables



Numerical value

Continuous

Categorical

- Ordinal
- Nominal

Data cleaning :

Data cleaning is a process of preparing data for analysis / ML/DL by removal or modifying data that is incorrect, incomplete, irrelevant, duplicated, or improperly formatted.

Missing value :

↳ Mandatory and non-mandatory

↳ This data is stored as blank content

↳ Removal of missing value is important because it doesn't work on blank value.

↳

```
dataset = pd.read_csv("Address.csv")
```

```
data.head(3)
```

Null value → NaN (Missing content)

```
dataset.isnull().sum()
```

⇒ heatmap

Encoding

one-hot encoding ÷ one-hot encoding is

a technique used to convert categorical variables into a format that can be provided to machine learning algorithms. Each category of a variable is transformed.

↳ one-hot encoding is used when there is no ordinal relationship.

Dummy variables

Dummy variables are a variation of one-hot encoding. When performing regression analysis, one of the categories in a categorical feature is left out to avoid the dummy variable trap.

Similar to one-hot encoding but removing one category to serve as baseline.

Label encoding

Label encoding is used to convert categorical text values to number. This is generally used for ordinal data, where there is a clear order or ranking.

Low = 1

Median = 2

High = 3

The problem arise when there is no ordinal relationship b/w categories, as Machine may interpret the no. as ordinal.

Ordinal encoding:

Ordinal encoding is a technical used when there is a natural order between categories. Each category is assigned an increase integer based on order.

To assign no. to categories, reflecting their order

Outliers:

Definition: An outlier is a data point that is significantly different from the rest of data. It can affect the model's performance and lead to skewed results.

Why to remove outliers: Outliers can distort statistical measures like mean and standard deviation and affect algorithms like regression, which assume normal distribution.

outliers removal using IQR:

IQR (Interquartile Range) is a measure of variability defined as the difference between 75th percentile and 25th percentile (QD).

$$IQR = Q_2 - Q_1$$

$$\text{Lower bound} \rightarrow Q_1 - 1.5 * IQR$$

$$\text{Upper bound} \rightarrow Q_3 + 1.5 * IQR$$

outliers removal using z-score:

A statistical measurement that describes a value's relationship to the mean of a group of values. It measures how many SD a data point is from the mean.

$$Z = \frac{(X - \mu)}{\sigma}$$

→ data point with a z-score less than -3 or greater than +3 are considered outliers.

Feature Scale (Normalization)

Normalization scale the data to a fixed range, typically $[0, 1]$ which is useful when the value need to be constrained.

$$X_{\text{new}} = \frac{(X - X_{\text{min}})}{X_{\text{max}} - X_{\text{min}}}$$

Handling duplicate data:

Definition \div Duplicate data refer to identical or almost identical rows of data in a data set, which can lead to bias in machine learning model.

Replace and data type change:

Replacing refers to the process of handling missing or erroneous data by replacing it.

Feature selection Techniques:

Feature selection techniques are used to select important features that contribute the most to the output, improving model performance and reducing overfitting.

A) Forward Elimination:

Process: Start with no variables in model, and add on features at a time. After adding each feature, evaluate the model's performance. Stop adding variable when performance doesn't improve.

B) Backward Elimination:

Process: Start with all the features and iteratively remove the least significant feature, one at a time, based on some criteria. Continue until no further features can be removed without harming the model's performance.