

Data Preprocessing.

- ① Handle inappropriate data. (pandas)
→ both numeric / non-numeric supported
- ② Deal with missing values. (pandas, sklearn)
→ colⁿ must be strictly numeric
- ③ Deal with Categorical data. (pandas, sklearn)
→ get dummies
→ LabelEncoder
→ OneHotEncoder
- ④ Deal with numerical data.
- Feature Scaling

* Dealing with Categorical data.

Example →

eid	ecity	esal
1	Mumbai	1000
2	Bangalore	2000
3	Mumbai	1500
4	Chennai	1800

Categorical data doesn't hold any mathematical weightage.

need a strategy that can convert given colⁿ into numerical colⁿ without adding any mathematical weightage.

Algo to deal with Categorical data:

- ① Get the unique value of the categorical column and store the same in the list.

['Mumbai', 'Bangalore', 'Chennai']

- ② Sort the list in ascending order.

['Bangalore', 'Chennai', 'Mumbai']
0 1 2

- ③ Replace the column values with the respective index values from the list. (created in step 2)

Original dataset:

eid	ecity	esal
1	Mumbai	1000
2	Bangalore	2000
3	Mumbai	1500
4	Chennai	1800



eid	ecity	esal
1	2	1000
2	0	2000
3	2	1500
4	1	1800

LABEL ENCODING

(Problem - We introduced mathematical weightage)

- ④ Create different categorical columns and substitute the data accordingly (to remove mathematical weightage)

* DUMMY VARIABLES

eid	ecity	esal
1	2	1000
2	0	2000
3	2	1500
4	1	1800



0	1	2	eid	esal
0	0	1	1	1000
1	0	0	2	2000
0	0	1	3	1500
0	1	0	4	1800

ONE HOT ENCODING

* To remove mathematical impact from labelled columns.

Dealing / Tuning Numerical Data.

Example:

age	salary
27	80000
26	85000
18	10000
12	35000

Given dataset,

range age : $10 \rightarrow 30$

salary : $10000 \rightarrow 85000$

Every ML algo expects your data to follow a standard scale / standard range without impacting the mathematical weightage of the same.

To achieve the above,

feature Scaling: (applicable only for feature columns)

FEATURE ENGINEERING

→ ① Rescale the feature in a limited range which is decided by Data Scientist.

→ ②* **Standardize the data** such that the mean is zero and standard deviation is 1.

→ ③ Normalize the data. (subject to algo performance)
 ④ L1 normalization
 ⑤ L2 normalization

Feature scaling is MANDATORY ONLY FOR

DIMENSIONALITY REDUCTION (PCA, LDA)

Technique Rescale the feature in a limited range
1 → which is decided by Data Scientist.

MinMaxScaler → user can define the range of data.

$$x_i' = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

where,

x → feature column/vector

x_i' → individual element calculated

x_i → individual feature

StandardScaler :

$$x_i' = \frac{x_i - \bar{x}}{\sigma} \quad \dots \quad Z\text{-Score Statistics}$$

where

x_i' → individual element calc.

\bar{x} → sample mean of feature colⁿ

σ → standard deviation of feature colⁿ

Feature Scaling :- Rescaling values such that the mathematical weightage remains the same.

feature Scaling is applied column wise

formula applied @ column axis.

Normalization :- Rescaling values in the range $\beta(0,1)$

Normalization is applied row wise.

formula is applied @ row axis

optimizer
catalyst

Normalization is ideally used to decrease
error in Regression algo.

- Gradient Descent
- Stochastic GD
- Batch GD

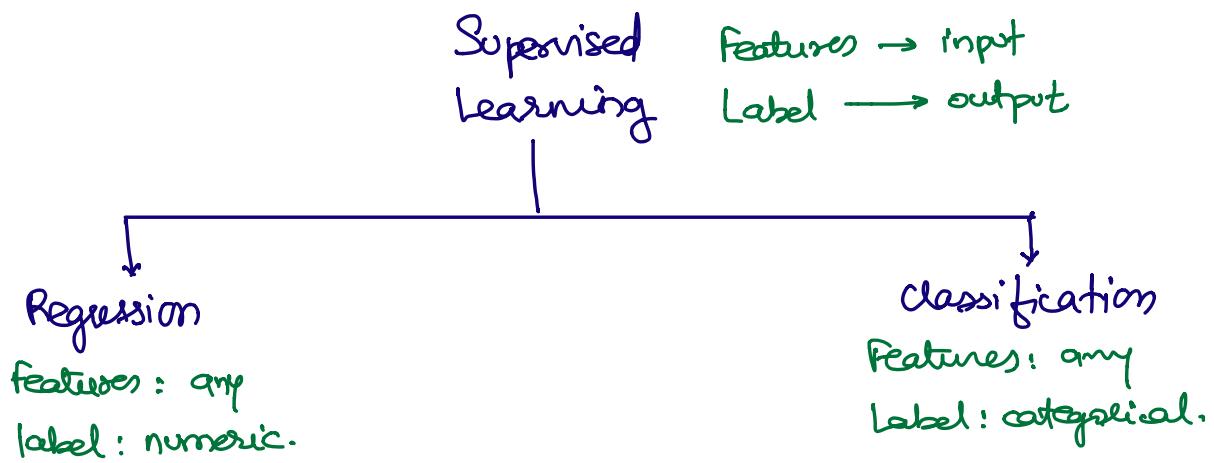
} not used in
ML widely

extensively used in Deep Learning

Artificial Neural Networks

Formula:

- ① Averaging Scale
- ② Range Bound Scale
- ③ L1 normalization \rightarrow Regression (Ridge Regression)
- ④ L2 normalization \rightarrow Regression (Lasso Regression)



① Linear Regression

- Simple linear regression → one feature, one label
- multiple linear regression → 'n' features, one label
- Polynomial linear regression → one feature, one label.
Polyomial features of one feature
 $\frac{x^1}{1}$
one.

② Decision Tree Regression

③ Random Forest Regression

④ Support Vector Regression

⑤ XG boost (extreme gradient boosting)

Linear Regression :-

Algebra \rightarrow Slope - intercept form!

$$y = mx + c$$

where

$m \rightarrow$ slope

$c \rightarrow$ intercept

$x \rightarrow$ feature

$y \rightarrow$ label.

Rewrite the above eqn as;

$$y = b_0 + b_1 x_1$$

where, $b_0 \rightarrow$ intercept

$b_1 \rightarrow$ co-efficient of x_1 (slope of x_1)

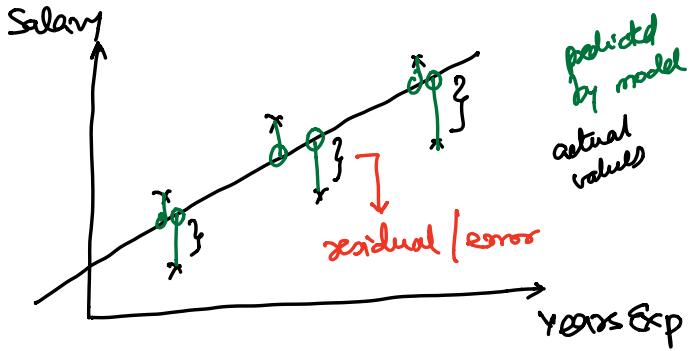
equation of Simple linear regression

$$\text{Salary} = b_0 + b_1 (\text{YearsExperience})$$

equation of multiple linear regression

$$\text{Profit} = b_0 + b_1 (\text{R&D Spend}) + b_2 (\text{Adm}) + b_3 (\text{Mark})$$

— APPLIED MACHINE LEARNING.—

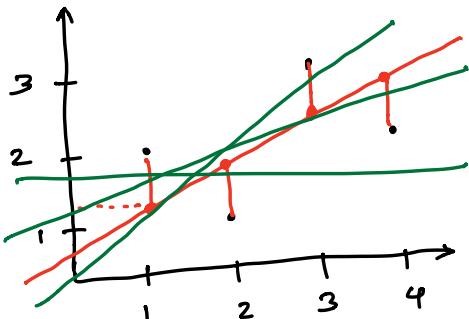


model that can predict the salary of the employee based on his/her experience

Best fit line → the line (model) that has less error!
 $\sum \text{residual}$ is minimum!

Calc. errors:

④ Mean Absolute Error.

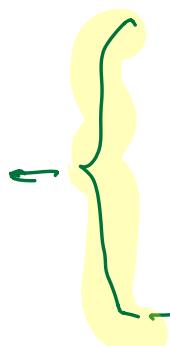


		<u>predicted</u>	
x	y	\bar{y}	error
1	2	1.2	0.8
2	1	1.8	0.8
3	3	2.3	0.7
4	2	2.8	0.8
			<u>3.1</u>

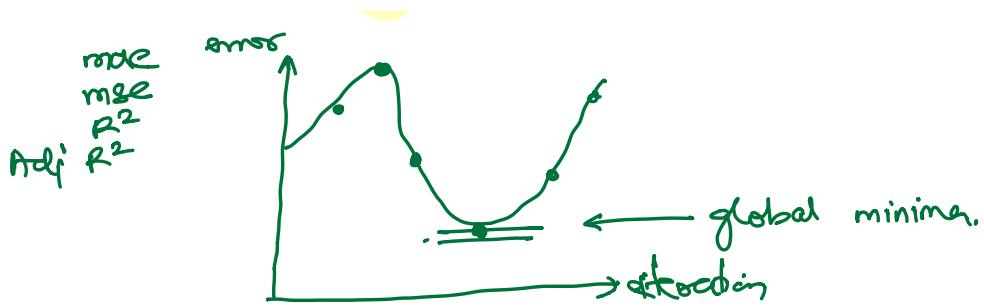
$$\frac{\sum_{i=1}^n |y - \bar{y}|}{n} = \text{MAE} = \frac{3.1}{4} = \underline{\underline{0.775}} \text{ error / loss}$$

Goal → to create and select the line that is minimal loss.

Gradient Descent



iterations	error
1	0.775
2	0.88
3	0.62
4	<u>0.33</u>



Random sampling:

- ① Sequential sampling ← small data but not recommended
- ② Random sampling ← recommended,

train test split ratio → 80 - 20

sequential Sampling

x	y
1	8
2	7
3	6
4	5
5	4
6	3
7	2
8	1
9	0
10	10

training dataset →

x	y
1	8
2	7
3	6
4	5
5	4
6	3
7	2
8	1

testing dataset

x	y
9	0
10	10

seed = 1
 $\rightarrow \text{mod function} = x \% 1 = 1$
 (many algo)

change of seed impact sample.

training dataset

x	y
1	8
3	6
7	2
10	10
:	:

x	y
2	4
5	4
2	7