

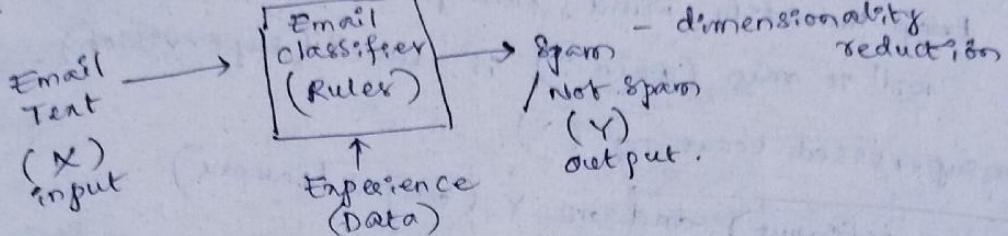
Machine Learning

It has algorithms (so) It is an algorithm

28/01/2025.

Example:

Email classification.



if - domain email id is not repeated.

{ - contains words like discount, free, lottery, win.

then

Spam

else

not spam

Learning: Learn from the experience (data)

Data is of two types Supervised, Unsupervised.

a) Supervised $(X \text{ and } Y)$

email	label
Meeting at 10:00 AM	not spam
Won \$10000	spam

Recognition of different Animals.

Feature	label
	cat
	dog

example for supervised learning

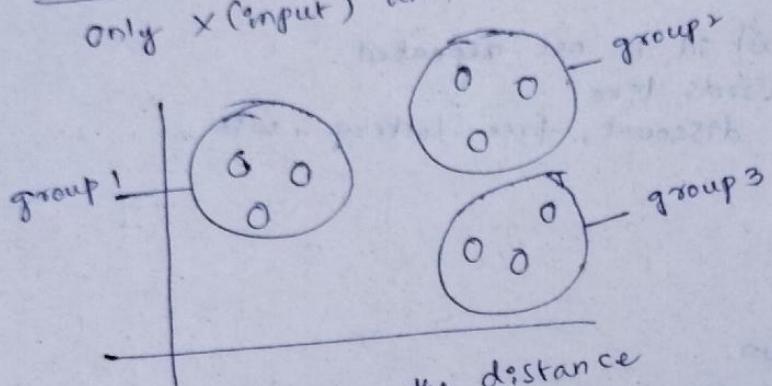
Tomorrow, will it rain? features → label (Y)

temp	humidity	wind	sunlight	Rain?
				Yes
				No

- temp, humidity, wind, sunlight are features
- will it rain (Rain?) is label (supervision)

b) Unsupervised Learning (Y)

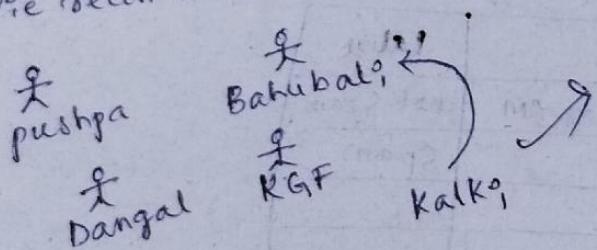
only X (input) and no Y (output / label)



Based on the distance
between the data
we are classified into groups.

Real world Example:

- Movie recommendation system



Features

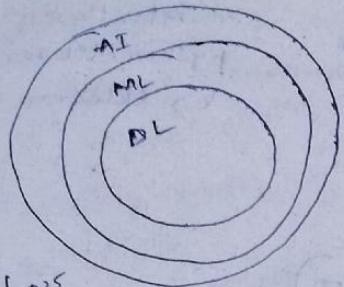
language	Distance	Songs

Budget	Gross

Kalki will be watched by the person
who watcher Babubali.

c) Reinforcement Learning (trial and Error)
learning by trying. (DeepSeekRL)
ex: Bicycling, Swimming

d) Semi-supervised Learning



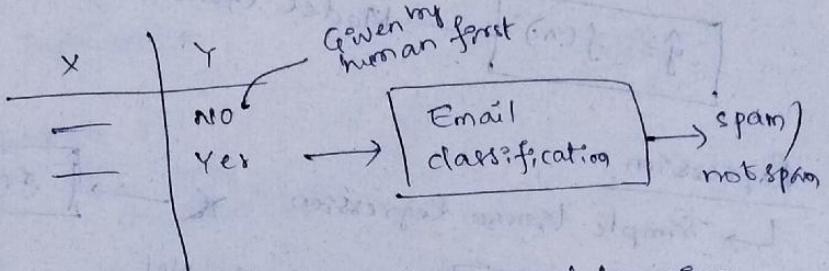
Deep learning - get user
Neural Network
to solve the problem.

What is Neural Network?
collect the signals It is inspired
and send forward. by human brain
(Neuron?)

29/01/2025

Steps

- 1) Train
- 2) Test
- 3) Deploy



30/01/2025

Machine Learning

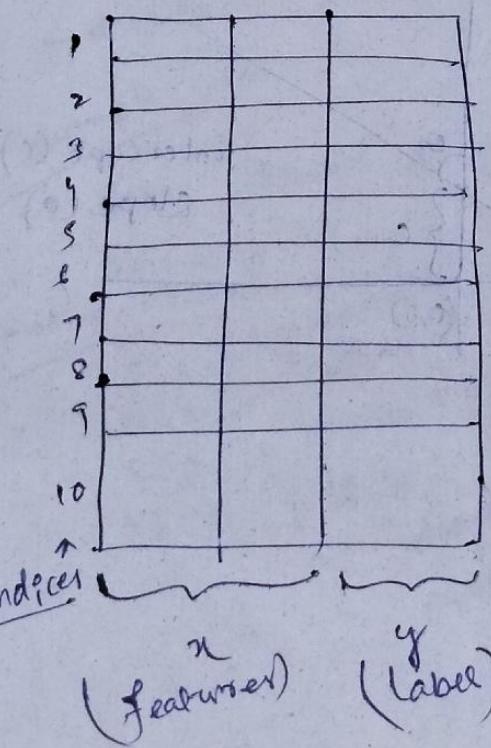
Dataset

labelled

unlabelled

labelled data is
used for
supervised
learning

Unlabelled data
is used for
unsupervised
learning



Supervised learning:

- 1) Train
- 2) Test
- 3) Deployment

Splitting the data

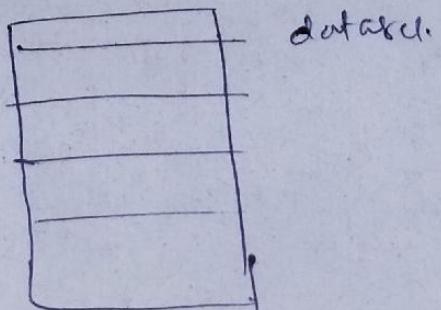
Train	Validation	Test
-------	------------	------

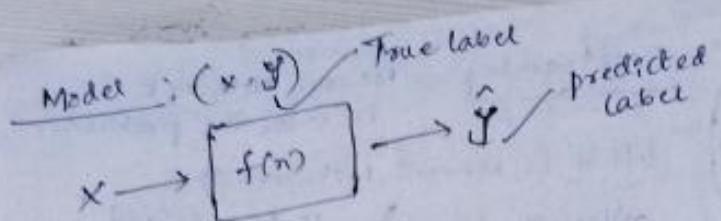
ex: solved examples exercise train
question question question question

- ① Train set $\rightarrow 80\% (4, 7, 2, 3, 5)$
- ② Validation set $\rightarrow 10\% (1, 6)$
development set $\rightarrow 10\%$.
- ③ Test set $\rightarrow 10\% (9, 10)$

Cross-validation

- if dataset size is small





intuitively
python - library
→ sklearn

ex:

$y = \text{true age}$
 $\hat{y} = \text{predicted age}$

$\hat{y} = f(n)$ Model (function)

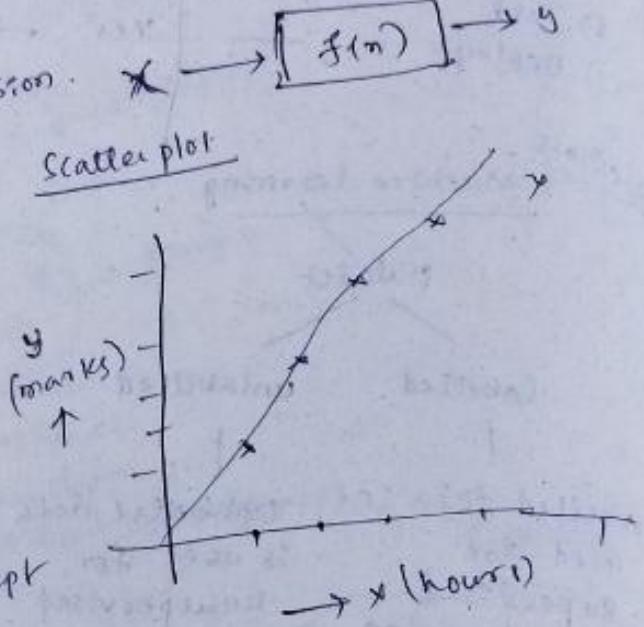
Regression

↳ Simple Linear Regression

Train data:

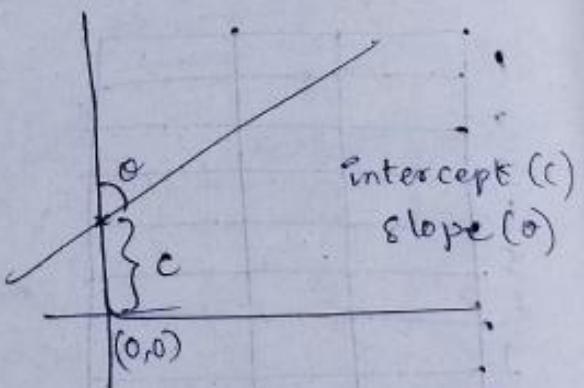
Hours	Marks
1	50
2	55
3	65
4	70
5	80

Scatter plot



output

$y = m x + c$ ~ slope
 $y = p_1 x + p_0$ ~ intercept



Data preprocessing

Categorical - Nominal
- Ordinal

- 1) Missing value
- 2) Encoding
- 3) Feature scaling
- 4) Train val Test split

en:

Colors	T-Shirt size
Red	L
Blue	M
Black	S
White	XL

↳ Nominal value
↳ ordinal value.
 $S < M < L < XL$

iloc

Weight	Weight (a-1)
50	0.5
60	0.6
70	0.7
80	0.8
100	1

we get by dividing by 100

Missing value
pandas

- fillna()

↳ mean
median
mode.

- dropna()

example code:

imputer class is imported from sklearn.preprocessing

imputer (missing_value, strategy)

it takes two parameters

from sklearn.preprocessing import Imputer
imputer = Imputer (missing_value="NaN", strategy="mean")

encoding
S:1 Dictionary/map
M:2 - level encoder

L:3 - OneHotEncoder

(if it is categorical & nominal)
XL:4
data is present
equal importance

Scaling

- 1) Standardization
- 2) Min-Max Scaling

Bringing large number to a range

X	$\frac{X - \mu}{\sigma}$
2	$2 - 6 = -4$
4	$4 - 6 = -2$
6	$6 - 6 = 0$
8	$8 - 6 = 2$
10	$10 - 6 = 4$
$\mu = 6$	

$$z = \frac{x - \mu}{\sigma} \quad \text{normal distribution}$$

Stand normal distribution.

input		output
Salary	own house	Loan?
100,000	1	Y

Here the salary is very high & no. of own houses are very small to make them equal we required scaling

Standardisation

from SKL

Min-Max Scaling

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Conclusion

Data preprocessing

- 1) Missing value & handling it

- fillna()
- imputer()

- 2) Encoding

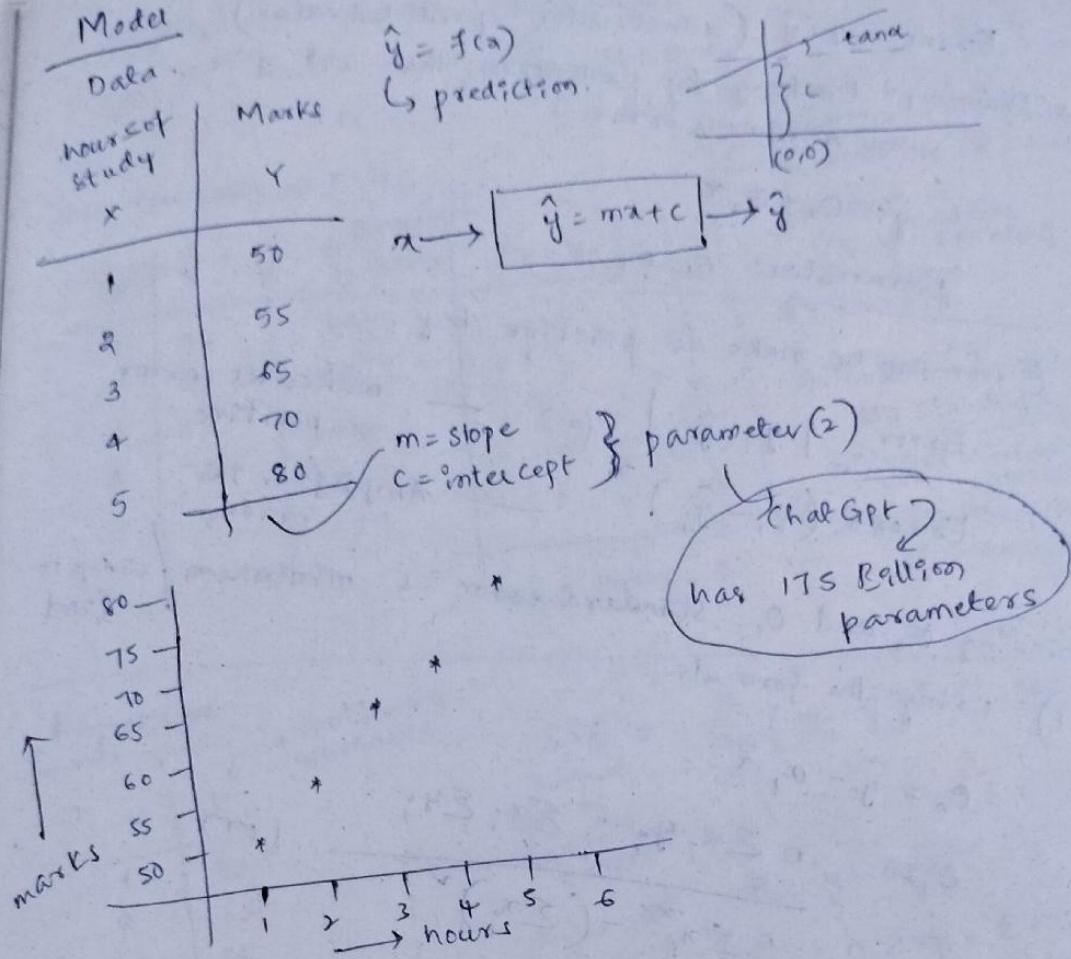
- map()
- LabelEncoder()
- OneHotEncoder()

- 3) Split

You have to learn
and find out the library

- 3) Scaling

- Standard Scaler
- Min-Max Scale



that Gpt
has 175 Billion parameters

learning is always denoted with error.

* Error *

By the linear regression, we can solve

Let equation of straight line be $\hat{y} = 10 + 5x$ then data will be

x	y	\hat{y}	E_y	$\hat{y}_2 = 40 + 2x$	E_2	$\hat{y} + E$
1	50	15	35	42	12	57
2	55	20	35	44	14	58
3	65	25	40	46	16	61
4	70	30	40	48	18	62
5	80	35	45	50	20	70
			195			fill the last (e)

Error = $y - \hat{y}$ (actual value - predicted value)

→ Training - Finding the parameter that will give the minimum error

$$\hat{y} = \theta_0 + \theta_1 x$$

parameters θ_0, θ_1

Error can be made as positive by:

$$\text{Error} = |y_i - \hat{y}_i| \quad (\text{or}) \quad \rightarrow \text{makes all error positive}$$

$$\text{Error} = (y_i - \hat{y}_i)^n \quad \rightarrow \text{Amplifies the error.}$$

finding θ_0 and θ_1 , standard error is minimum / ways to find

finding the formula.

1) Using the formula.

$$\theta_0 = \bar{y} - \theta_1 \bar{x} \quad (\text{or})$$

$$\theta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \rightarrow \text{It is little bit more calculative.}$$

Solve

$$y = 10 + 5x$$

x	y	$x_i y_i$	$(x_i)^2$
1	50	50	1
2	55	110	4
3	65	195	9
4	70	280	16
5	80	400	25
(15)	(320)	1035	55

$$\theta_1 = \frac{-5(1035) - 15 \cdot 320}{5(225) - 225}$$

$$= \frac{375}{50} = 7.5$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x} = \frac{320}{5} - 7.5(3)$$

$$= 64 - 22.5 = 41.5$$

$$\theta_0 = 41.5$$

$$\theta_1 = 7.5$$

Prediction

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\hat{y} = 41.5 + 7.5 x_1$$

Conclusion

In Machine Learning

Then the table will be

X	Y	Prediction		$E_i = y_i - \hat{y}_i$	Total Error = 6
		$\hat{y} = \theta_0 + \theta_1 x_1$	$E_i = y_i - \hat{y}_i$		
1	50	49	1		
2	55	56.5	-1.5		
3	65	64.5	1		
4	70	71.5	-1.5		
5	80	79	1		

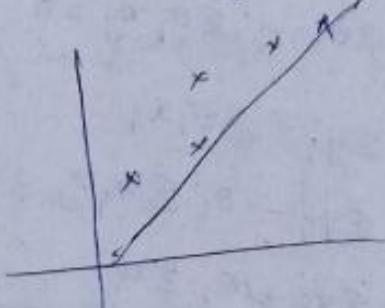
- Model contains parameters
- Train - to reduce the standard error or to get minimum error by finding the parameters.

- Train - to reduce the standard error or to get minimum error by finding the parameters.

Regression

continuous variable

X	Y
1	50
2	55
3	65
4	70
5	80



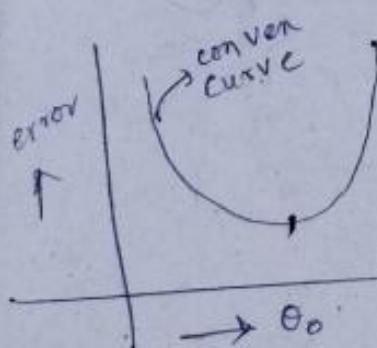
$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Sum of squares of errors

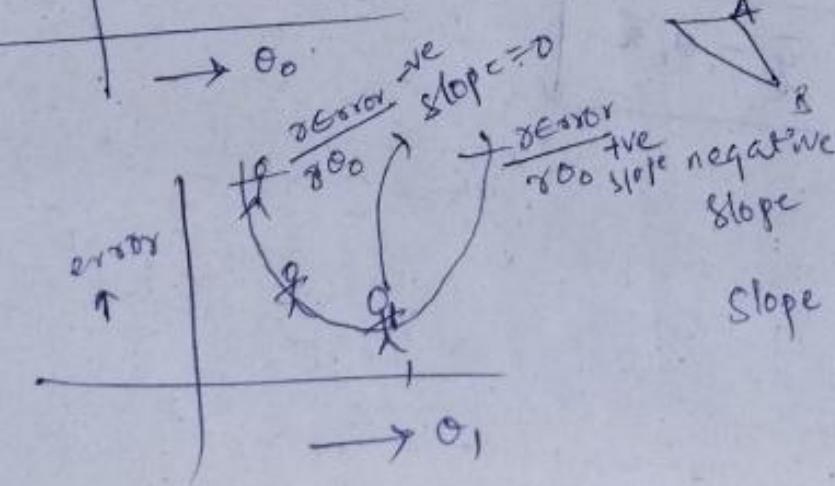
$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_1 = \frac{n \sum xy - \bar{x} \bar{y}}{n \sum x^2 - (\bar{x})^2}$$

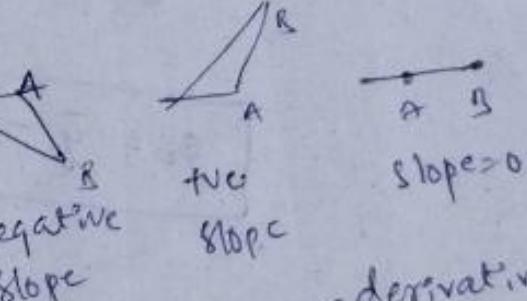
why these attributes will give minimum error? and how



$\theta_0 \leftarrow \theta_0$
Chaining
constant changing



$\theta_1 \leftarrow \theta_1$
constant changing



Slope = $\frac{dy}{d\theta}$ ✓ derivative

↳ import concept for neural network

Find θ_0 and θ_1 such that [Intuition]
 ∂ -Partial derivative.

$$\frac{\partial \text{Error}}{\partial \theta_0} = 0$$

$$\text{and } \frac{\partial \text{Error}}{\partial \theta_1} = 0$$

we know, $SSE = \sum (y_i - \hat{y})^2$
 $SSE = \sum (y_i - \theta_0 - \theta_1 x_i)^2$ $\left[\because \hat{y} = \theta_0 + \theta_1 x \right]$

$$\begin{aligned} \therefore \frac{\partial \text{Error}}{\partial \theta_0} &= \frac{\partial \sum (y_i - \theta_0 - \theta_1 x_i)^2}{\partial \theta_0} = 0 \\ &= 2 \sum (y_i - \theta_0 - \theta_1 x_i) (-1) \quad \text{error} \\ &= -2 \sum (y_i - \theta_0 - \theta_1 x_i) \quad \text{chain rule} \\ \Rightarrow \sum (y_i - \theta_0 - \theta_1 x_i) &= 0 \end{aligned}$$

$$\sum y_i - \sum \theta_0 - \sum \theta_1 x_i = 0$$

$$\Rightarrow \sum \theta_0 = \sum y_i - \sum \theta_1 x_i$$

$$n \theta_0 = \sum y_i - \theta_1 \sum x_i$$

$$\theta_0 = \frac{\sum y_i - \theta_1 \sum x_i}{n} \rightarrow \text{eq ②}$$

$$\theta_0 = \frac{\sum y_i}{n} - \theta_1 \frac{\sum x_i}{n}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

for θ_1

$$\frac{\partial \text{Error}}{\partial \theta_1} = 0$$

$$\Rightarrow \frac{\partial}{\partial \theta_1} \sum (y_i - \theta_0 - \theta_1 x_i)^2 = 0 \quad \begin{matrix} \text{this } x_i \text{ is within} \\ \text{the sigma}(\Sigma) \end{matrix}$$

$$2 \sum (y_i - \theta_0 - \theta_1 x_i) (-x_i) = 0$$

$$-2 \sum x_i (y_i - \theta_0 - \theta_1 x_i) = 0$$

$$\Rightarrow \sum x_i (y_i - \theta_0 - \theta_1 x_i) = 0$$

$$\Rightarrow \sum x_i y_i - \theta_0 \sum x_i - \theta_1 \sum x_i^2 = 0$$

$$\Rightarrow \sum x_i y_i - \left(\frac{\sum y_i - \theta_0 \sum x_i}{n} \right) \sum x_i - \theta_1 \sum x_i^2 = 0 \quad (\text{from eqn 2})$$

$$\Rightarrow \frac{n \sum x_i y_i - \sum x_i \sum y_i + \theta_0 (\sum x_i)^2 - n \theta_1 \sum x_i^2}{n} = 0$$

$$\Rightarrow n \sum x_i y_i - \sum x_i \sum y_i + \theta_0 (\sum x_i)^2 - n \theta_1 \sum x_i^2 = 0$$

$$\Rightarrow n \sum x_i y_i - \sum x_i \sum y_i = n \theta_0 \sum x_i^2 - \theta_1 (\sum x_i)^2$$

$$\Rightarrow \theta_1 \left(n \sum x_i^2 - (\sum x_i)^2 \right) = n \sum x_i y_i - \sum x_i \sum y_i$$

$$\Rightarrow \boxed{\theta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}}$$

Regression

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

Independent Dependent

x	y
1	50
2	55
3	65
4	70
5	80

Normal Form Equation:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$y = \begin{bmatrix} 50 \\ 55 \\ 65 \\ 70 \\ 80 \end{bmatrix}$$

normal
equation / normal form

$$\theta = (x^T x)^{-1} x^T y$$

numpy:

import numpy as np

x = np.array([[1, 1], [1, 2], [1, 3], [1, 4], [1, 5]]);

y = np.array([50, 55, 65, 70, 80])

theta = np.linalg.inv(x.T @ x) @ x.T @ y

print(theta)

$$\rightarrow \begin{bmatrix} 41.5 \\ 7.5 \end{bmatrix}$$

Vibratory
SKLearn
above can also be learnt
from the SKLearn library.

$$\theta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

GPU.

1) Graphical processing

Unit

Color
Machine learning

requires GPU

CPU-Z

$$\text{SSE} : \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad \hat{y} = x\theta$$

$$= (y - \hat{y})^T (y - \hat{y})$$

$$= (y - x\theta)^T (y - x\theta)$$

$$= y^T y - y^T x\theta - \theta^T x^T y + \theta^T x^T x\theta$$

$$\hat{y} = \begin{bmatrix} 49.5 \\ 56.5 \\ 62.5 \\ 71.5 \\ 75 \end{bmatrix}$$

$$x^T x\theta = \theta^T x^T y$$

$$1 \times 5, 5 \times 2, 2 \times 1 \quad | \quad 1 \times 2, 2 \times 5, 5 \times 1$$

$$1 \times 1 \qquad \qquad \qquad 1 \times 1$$

$$= y^T y - 2 \cdot \theta^T x^T y + \theta^T x^T x\theta$$

$$\frac{\partial \text{SSE}}{\partial \theta} = 0 - 2x^T y + 2x^T x \cdot \theta = 0$$

$$\Rightarrow 2x^T x \theta = 2x^T y$$

$$\Rightarrow x^T x \theta = x^T y$$

$$\theta = \frac{x^T y}{x^T x}$$

$$\boxed{\theta = (x^T x)^{-1} x^T y}$$

Multiple Linear Regression

$$\theta = (X^T X)^{-1} \cdot X^T Y$$

size, bedrooms, age (Independent) $\rightarrow f(x)$ \rightarrow price (Dependent)

$$\hat{y} = \theta_0 + \theta_1 * \text{size} + \theta_2 * \text{bedrooms} + \theta_3 * \text{age}$$

parameter need to find are $\{\theta_0, \theta_1, \theta_2, \theta_3\}$

i.e. $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$

Example:
 independent variables (size, bedrooms, age) and
 dependent variable (price)

size (sq ft)	Bedrooms	Age (years)	price (Y)
1500	3	15	400,000
1800	4	20	460,000
2400	3	10	56000
3000	4	8	60000
3500	5	5	2000

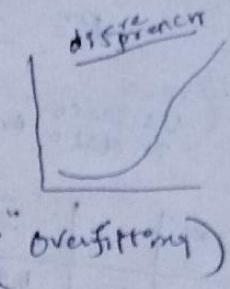
$$X = \begin{bmatrix} 1 & 1500 & 3 & 15 \\ 1 & 1800 & 4 & 20 \\ 1 & 2400 & 3 & 10 \\ 1 & 3000 & 4 & 8 \\ 1 & 3500 & 5 & 5 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \quad Y = \begin{bmatrix} 400000 \\ 460000 \\ 56000 \\ 60000 \\ 2000 \end{bmatrix}$$

5×4

- polynomial ~~linear~~ Regression.
- Evaluating Model performance

PLR: $\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_n x^n$

hours (x)	x	x^2	y (marks)
1	1	1	4
2	2	4	14
3	3	9	24
4	4	16	34
5	5	25	44



Evaluation:

↳ Test

↳ Error

|

→ Mean Absolute Error

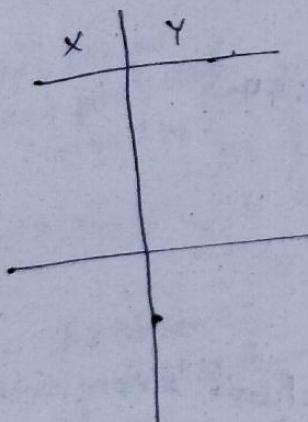
→ MSE

→ RMSE

→ R^2

more error → less learn
less error → more learn.

Evaluation Metrics



Train

Error

$$\rightarrow SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\rightarrow MSE = \frac{SSE}{n}$$

Test

$$\rightarrow RMSE = \sqrt{MSE}$$

Every ML algorithm will have any one of these errors based on problem.

Evaluation

↳ Evaluation Metrics

↳ R^2

→ Adjusted R^2

$$R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \rightarrow$$

Variance in data ↙

variance is explained by model.

Variance: how the data is distributed around the mean.

if $\text{var}(x)$ is less, near to mean
else if $\text{var}(x)$ is more, far from mean

eq:	X	Y	\hat{y}
	1	50	49
	2	55	56.5
	3	65	64
	4	70	71.5
	5	80	79
			$\bar{y} = 64$

$$R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = 0.986$$

$$\begin{aligned}
 R^2 &= \frac{(50-49)^2 + (55-56.5)^2 + (65-64)^2 + (70-71.5)^2 + (80-79)^2}{(50-64)^2 + (55-64)^2 + (65-64)^2 + (70-64)^2 + (80-64)^2} \\
 &= \frac{1 + (-1.5)^2 + 1 + (-1.5)^2 + 1}{(-14)^2 + (-9)^2 + 1 + (6)^2 + (16)^2} = \frac{3+2(2.25)}{196+81+1+36+256} \\
 &= \frac{7.5}{570} = 0.013
 \end{aligned}$$

Note: If $R^2 = 1 \rightarrow$ Good Model / perfect fit
 $R^2 = 0 \rightarrow$ Bad Model / not fitting.

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{n-k-1} = 0.94$$

where $n = \text{no. of data points} = 5$
 $k = \text{no. of features} = 1$

→ Adjusted R^2 will decrease when unnecessary columns are added to the dataset (ID: B200652)

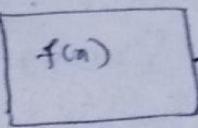
↳ It does not require to predict the marks

R^2 value remains same but adjusted R^2 value decreases

Classification

Evaluation Metrics:

Email



Binary Classification

$f(n)$

$y = \{ \text{spam, not spam} \}$

1 0

discrete!

X	Y	\hat{y}
TN	0	0 ✓
TP	1	1 ✓
FN	0	0 ✗
FP	0	1 ✗
TN	1	1 ✓
TP	0	1 ✗
FP	1	0 ✗
FN	0	0 ✓

$$\text{Accuracy} = \frac{5}{9} \times 100$$

Confusion Matrix

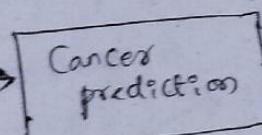
		predicted	
		1	0
actual	1	TP 2	FN 2
	0	FP 2	TN 3

Not Spam & but marked as Spam

TP - True positive
TN - True Negative
FN - False Negative

Example:

Features
(dimensions)



Cancer prediction

$\{ \text{yes, NO} \}$

1 0

has Cancer and predicted correctly.

	1	0
1	TP	FN
0	FP	TN

has Cancer but not predicted (more dangerous)

no Cancer but predicted Cancer

no Cancer and predicted no Cancer

Classification

x	y	\hat{y}
0	0	0
0	0	0
1	0	+ve Spam
0	1	-ve not spam
1	1	Cancer / not cancer.
0	1	
1	0	
0	0	
1	1	

↳ Two labels { 1, 0 }

↳ +ve / -ve

↳ Spam / not spam

↳ Cancer / not cancer.

↳ Accuracy: Near to 1 = Good model
 Near to 0 = Bad model

↳ Precision: Near to 1 = Good model
 Near to 0 = Bad model

↳ Recall: Near to 1 = Good model
 Near to 0 = Bad model

↳ F1 Score: Near to 1 = Good model
 Near to 0 = Bad model

↳ Confusion Matrix

		Prediction	
		1	0
Actual	1	TP 3	FN 2
	0	FP 2	TN 3

→ Type 2 error.

→ Type 1 error

$$\text{Accuracy} = \frac{\text{Total correct prediction}}{\text{Total prediction}} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$= \frac{6}{10} = 0.6 = 60\%.$$

$$P(\text{precision}) = \frac{\text{No. of correct positive predictions}}{\text{No. of positive predictions}} = \frac{TP}{TP + FP}$$

$$= \frac{3}{5} = 0.6$$

$$R(\text{Recall}) = \frac{\text{No. of correct positive prediction}}{\text{No. of actual positive cases}} = \frac{TP}{TP + FN}$$

$$= \frac{3}{5} = 0.6$$

$$F1 = \frac{2 * P * R}{P + R} = \frac{2 * 0.6 * 0.6}{0.6 + 0.6} = 0.6$$

To understand the complete picture, we need precision, recall and F1 along with accuracy.

eg:

	\hat{y}	y
0	0	0
0	0	0
0	0	0
1	0	0
0	0	0
0	0	0
0	0	0
1	0	0
0	0	0
0	0	0

	TP	FN
TP	2	0
FN	0	8

$$\text{Accuracy} = \frac{8}{10} = 80\%$$

$$P = 0$$

$$R = 0$$

F1 = not defined.

eg:

email classification

	\hat{y}	y
1	0	1
0	1	2

	TP	FN
TP	1	2
FN	3	14

$$\text{Accuracy} = \frac{15}{1+2+3+14} = \frac{15}{20} = \frac{3}{4} = 0.75$$

$$\text{precision} = \frac{1}{1+3} = \frac{1}{4} = 0.25$$

$$\text{Recall} = \frac{1}{1+2} = \frac{1}{3} = 0.33$$

$$F1 = \frac{2 \times 0.25 \times 0.33}{0.25 + 0.33} = 0.27$$

Multiclass Classification

eg:

0 - Regular Email
1 - Social Email
2 - Spam Email

	\hat{y}	y
0	0	0
0	0	2
0	2	0
0	0	0
0	2	2
1	1	1
1	1	1
1	1	1
1	0	1
1	1	1
2	1	1
2	1	2
2	2	2
2	2	2
2	2	2

actual

predicted

	0	1	2
0	3	0	2
1	1	4	0
2	0	2	3

$$\text{Accuracy} = \frac{10}{15} = 0.66$$

class 0: TP=3, TN=9, FP=1, FN=2

$$P = 0.75 \quad R = 0.6 \quad F1 =$$

class 1: TP=4, TN=8, FP=2, FN=1

$$P = ? \quad R = ? \quad F1 = ?$$

class 2: TP=3, TN=8, FP=2, FN=2

$$P = ? \quad R = ? \quad F1 = ?$$

$$\text{Macro } P = \frac{10}{3}$$

$$\text{Macro } R = \frac{10}{3}$$

$$\text{Macro } F1 = \frac{10}{3}$$

Q9:

y	5
0	2
0	2
0	0
0	0
0	1
1	1
1	1
1	2
1	0
1	1
2	0
2	0
2	0
2	2
2	2

Evaluation Metrics

- ↳ classification
- ↳ Accuracy
- precision
- Recall
- F1 score

Classification Algorithm

① Decision Tree

$\{X, Y\}$ discrete

{ spam, not spam }

sample code:

```

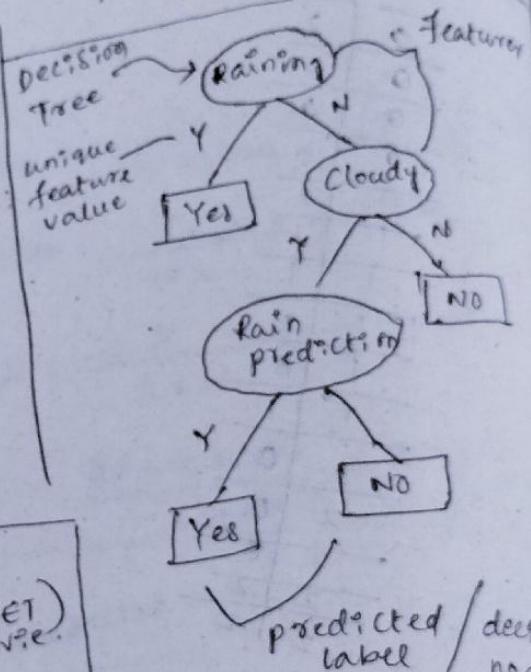
if raining
then yes
else cloudy:
    if rain prediction
    then yes
else
    no
else
    no;

```

Cloudy, weather prediction, raining

	X	Y	
X	Cloudy	weather prediction	raining
Y			Umbrella?
N	N	N	N
Y	N	Y	Y
Y	N	Y	Y
Y	Y	Y	Y

Accuracy = 100%



1. Uncertainty and Entropy (TENET movie)

basically entropy measure
uncertainty

example: will the students wear uniform tomorrow?

C1	C2	C3	C4
Y	Y	Y	Y
Y	Y	N	Y
Y	Y	N	N
Y	N	N	N
Y			

$$\begin{aligned} C_2 & \\ p(U) &= \frac{3}{4} \\ p(NU) &= \frac{1}{4} \end{aligned}$$

$$\begin{aligned} C_3 & \\ p(U) &= \frac{1}{4} \\ p(NU) &= \frac{3}{4} \end{aligned}$$

$C_1 \rightarrow p(U) = \frac{4}{4} = 1$
More certainty $p(U) = \frac{0}{4} = 0$
(Entropy is less)

\uparrow less certainty (Entropy is more) C_4

$p(U) = \frac{2}{4} = \frac{1}{2}$
 $\frac{p(U)}{p(NU)} = \frac{2}{4} = \frac{1}{2}$

My calculate $p(U), p(NU)$
for $C_2 \& C_3$

$$\text{Entropy} = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$\{P_1, P_2\} = \{U, NU\}$$

uniform \swarrow not uniform.

$$\log 1 = 0$$

$$E(c_1) = -1 \log_2 1 - 0 \log_2 0 = 0$$

$$E(c_4) = -\log_2 \frac{1}{2} = 1$$

$$E(c_4) = -y_2 \log_2 y_2 - \bar{y}_2 \log_2 \bar{y}_2 = -0.81$$

$$E(C_2) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.81$$

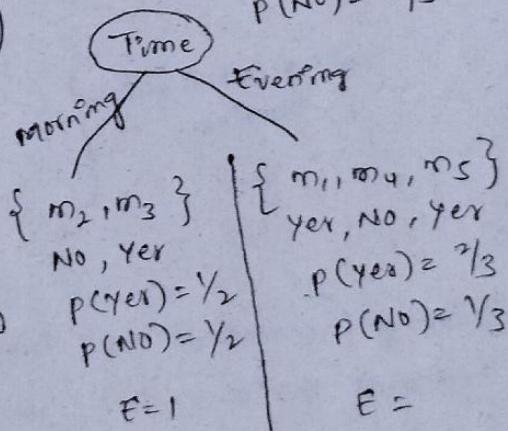
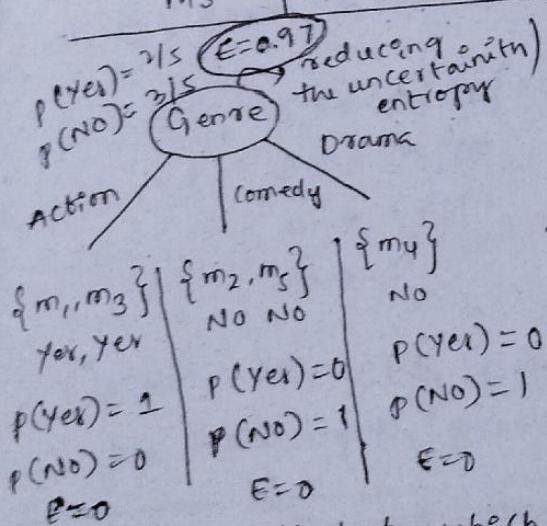
$$E(C_2) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.81$$

Note: Uncertainty is more \Rightarrow Entropy is more
i.e. Uncertainty \propto Entropy.

Example: watch movie?

Movie	Genre	Time	Watch?
M1	Action	Even	Yes
M2	Comedy	Morn	No
M3	Action	Morn	Yes
M4	Drama	Even	No
M5	Comedy	Even	No

Entropy (1) = High uncertainty / entropy
Entropy (0) = low uncertainty / entropy



$E=0$ Select the attribute which help you take the decision.

Select the attribute which help us to classify the data. here, gender is suitable attribute (feature) which decreases the uncertainty (entropy) more than the attribute time.