

Linear Regression (single variable)



$$y = m \cdot x + c$$

slope or gradient.
coefficient

$$\text{price} = m \cdot \text{Area} + b$$

depend

Inde

from sklearn import linear_model

```
df = pd.read_csv('homeprice.csv')
```

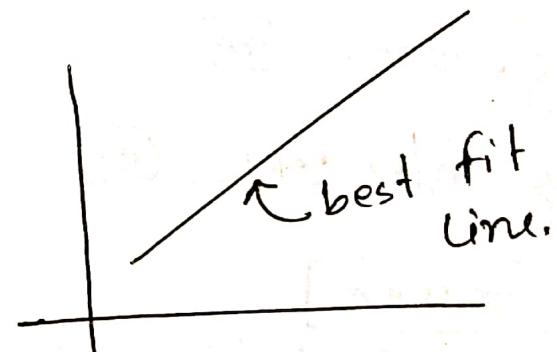
```
plt.xlabel('area')
```

```
plt.ylabel('price')
```

```
plt.scatter(df['area'], df['price'])
```



reg = LinearModel.LinearRegression()
~~reg.fit~~



Fungitac
Sertaconazole Nitrate

```
# neg.fit (df[['area']], df['price'])
```

Linear Regression (copy-xtrue, fit-inception = True, njobs=1, norma=False)

neg.predict ([3000]) →

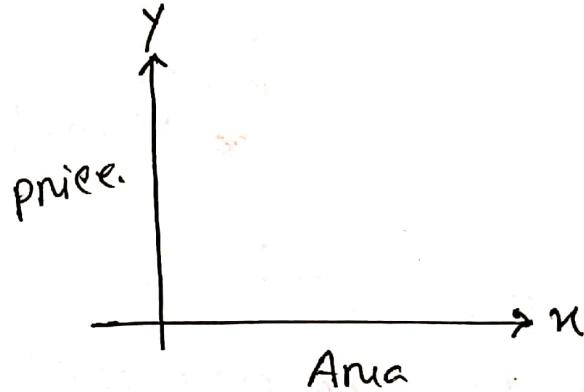
> array (628715.75)

neg. coef - m

neg. intercept - c

$$y = mn + c$$

y = neg.predict(3000)



fill up the column with prediction.

```
d = pd.read_csv('areas.csv')
```

```
p = neg.predict(d)
```

```
d['price'] = p
```

```
df.to_csv('prediction.csv'), index=False)
```

```
plt.scatter(df.area, df.pnuee)
```

```
plt.plot(df.area, neg.predict(df[['area']]), color='blue')
```

excercise:

Given: net national income per capital

Q. predict net income year - 2020.

Fungitac
Sertaconazole Nitrate

Linear Regression (multiple variable) multivariate regression.

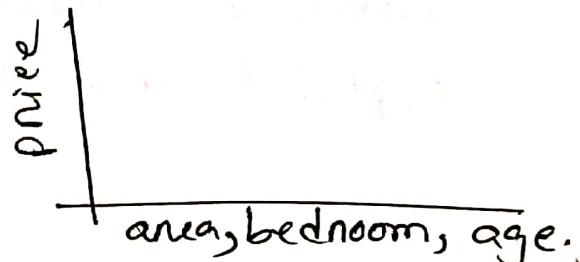
Q. given two home prices. independent variable.

$$\text{price} = m_1 \times \text{area} + m_2 \times \text{bedrooms} + m_3 \times \text{age} + b$$

dependent variable.

$$y = m_1 n_1 + m_2 n_2 + m_3 n_3 + b$$

df = pd.read_csv('homeprice.csv')



① DATA PROCESSING : Handl NA values

new_df = df.fillna ({'bedrooms': 0})

new_df = df.interpolate()

for multiple value.

② LINEAR Regression :

Reg = Linear_model.LinearRegression()

reg.fit (df[['area', 'bedrooms', 'age']], df['price'])

Neg. coef-

$$> m_1 = m_2 = m_3 =$$

neg. intencpt

> c

neg. predict ([[3000, 3, 40]])

$$> \text{price} = 444900.0$$

gradient Descent and cost function

$$y = 2n + 3$$

$$n = [1, 2, 3, 4, 5]$$

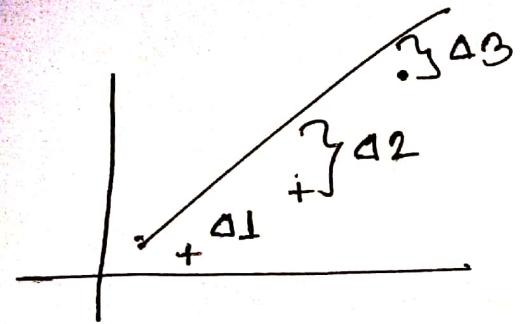
$$y = [5, 7, 9, 11, 13]$$

There given n, y value and find
out the equation.

In plot · scatter() always we can't
get best fit line.

$$\text{price} = m \times \text{area} + c$$

Fungitac
Sertaconazole Nitrate



$$MSE = \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

↓

Sgn. them
they could be
negative.

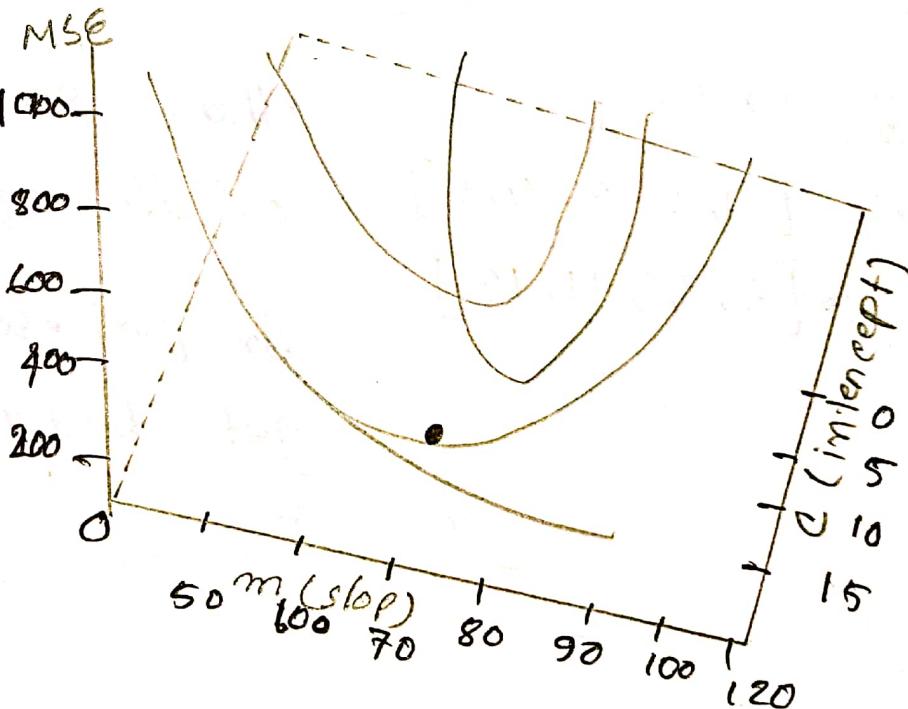
MSE = mean square error.

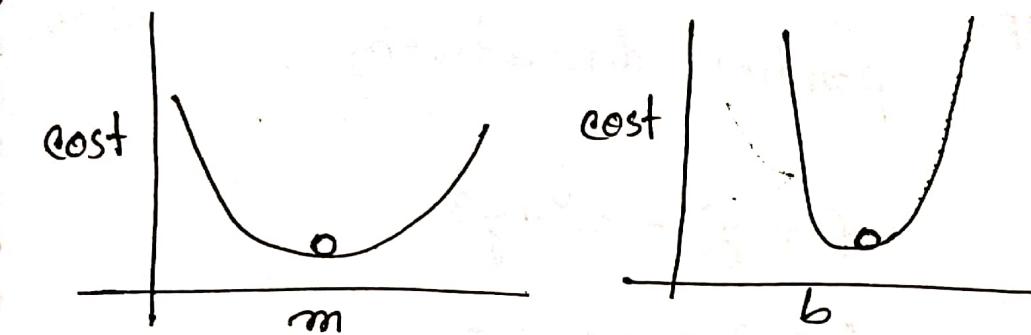
$$MSE = \frac{1}{n} \sum (y_i - (\hat{m}x_i + b))^2$$

↑
cost
function

↑
actual
predicted

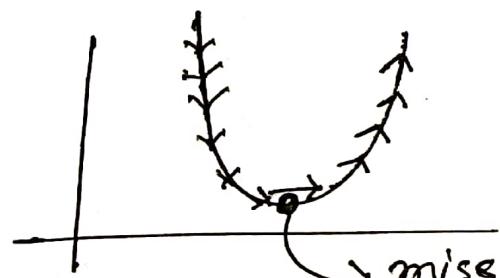
gradient descent is an algorithm
that finds best fit line for given
training dataset





how to do this

if the all steps are
same or equal.
(fixed size)



missed the global minimum.

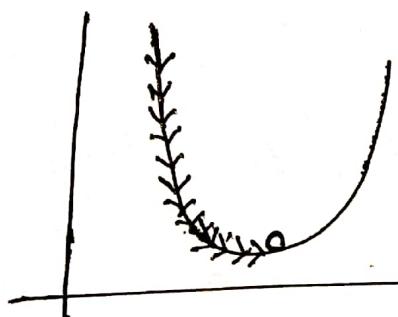
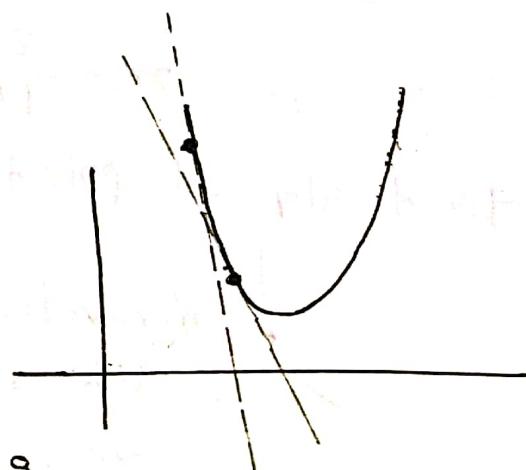
what can work:

here each steps are following
the curvature.

the step size are reducing.

then each point calculate the slop.

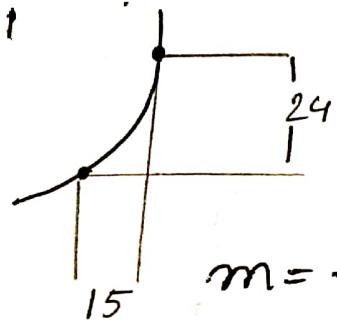
slop also says the direction.



Fungitac
Sertaconazole Nitrate

derivative all about slope.

$$m = \frac{y_2 - y_1}{x_2 - x_1}$$



$$m = \frac{24}{15}$$

- find slope at point:

derivatives are
use in small differences.

$$m = \frac{\Delta y}{\Delta x}$$

$$\text{equation} = f(x) = x^2$$

$$\text{slop} = \frac{d}{dx} x^2$$

$$= 2x$$

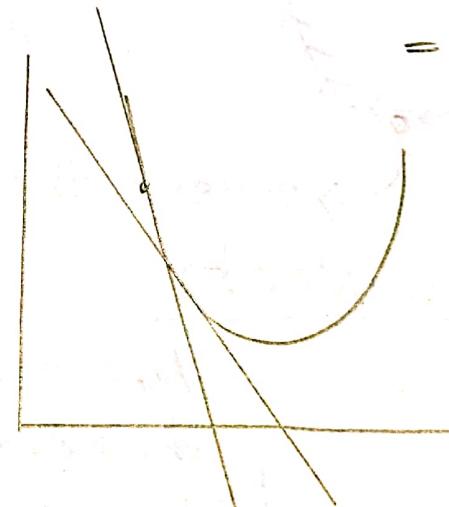
partial derivatives.

$$f(x, y) = x^2 + y^3$$

$$f'(x) = 2x + 0$$

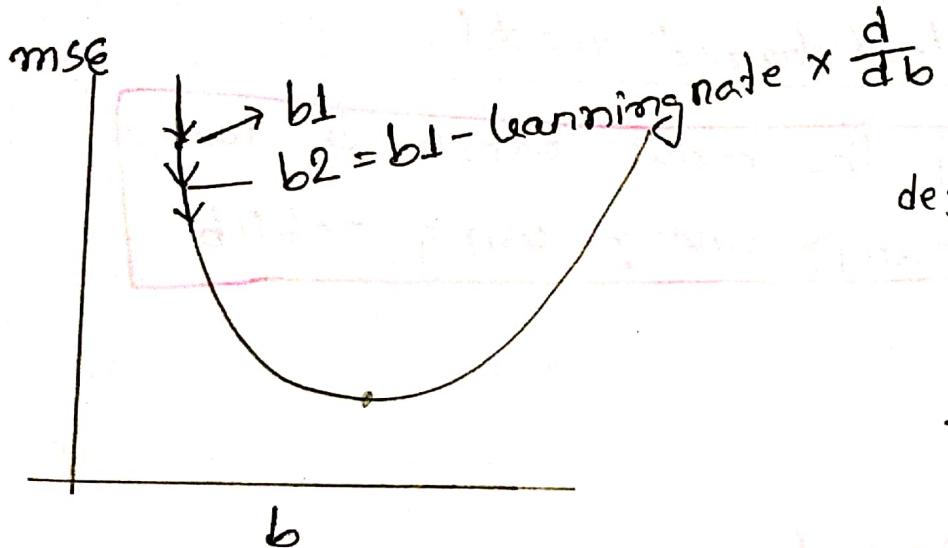
$$f'(y) = 0 + 3y^{3-1}$$

$$\begin{aligned} \frac{d}{dx} x^2 \\ = 2x^{2-1} \end{aligned}$$



partial derivatives.

$$\begin{aligned} \frac{d}{dx} f \\ \frac{d}{dy} f \end{aligned}$$



$$x = df['x']$$

$$y = dt['y']$$

1st step

$$m = 0$$

$$c = 0$$

$$\text{learn-rate} = 0.001$$

$$n = \text{int}(\text{len}(x))$$

$$\text{iteration} = 100$$

2nd step

for i in range

$$\text{prediction} = mx + c$$

def gradient_descent (x, y)

$$m_curr = b_curr = 0$$

$$\text{iteration} n = 1000$$

for i in range (iteration):

$$y_pre = m_curr * x + b_curr$$

$$md = -(2/n) * \sum (x(y - y_pre))$$

Fungitac
Sertaconazole Nitrate

Save and load trained model

1) pickle.

different: if the model contains large numpy array. using joblib

2) sklearn joblib

Sol.

import pickle.

```
with open('model-pickle', 'wb') as f  
    pickle.dump(model, f)
```

```
with open('model-pickle', 'rb') as f  
    mpdlt = pickle.load(f)
```

~~from sklearn.externals import joblib~~

~~joblib.dump(model, 'model-joblib')~~

~~mJ = joblib.load('model-joblib')~~

~~mJ.predict(500)~~

model are saved in
this file.

categorical, Dummy variable,
one hot encoding

build a prediction function to predict price of home.

① 3400 sqft area in west windson

② 2800 sqft home in Robbinsville

Town	area	price
westwindson	2600	58500

↑ how should handle

text data?

in numeric data.

dummy
variable.

Categorical variable.

Nominal → no numeric order.

city

gender

color

red
green
blue

ordinal → having numeric order

satisfied
neutral
dissatisfied

graduate
masters
phd

high
low

Fungitac
Sertaconazole Nitrate

`df = pd.read_csv('house-prices.csv')`

`du = pd.get_dummies(df['town'])` \leftarrow dummy variable.

	town1	town2	town3
1		0	
1	0		0
0		1	0
0		1	0
0	0		0
0	0		1
0	0	0	1

`merged = pd.concat([df, dummien])` axis=1 column

now drop 'the town column, which not

use in neg. model.

dummy variable trap:

when even one variable can be derived from a rest of the vars.

these variables are said to be multi-collinear.

```
final_mange = manged.drop(['town', 'west windson'], axis='column')
```

> from sklearn.linear_model import LinearRegression.

```
model = LinearRegression()
```

Here price is the dependent variable.

```
x = final.drop('price', axis='column')
```

```
y = final.price
```

```
model.fit(x, y)
```

```
model.predict([E2800, 0, 1])
```

```
model.score(x, y)
```

→ compare the predicted value with actual value.

import LabelEncoder.

```
Le = LabelEncoder()
```

df_le = df → taken table column.

```
le.fit_transform(df_le.town) = df_le.town
```

Fungitac
Sertaconazole Nitrate

x is training data set.

x is dependent variable.

$x = df[[\text{'town}', \text{'aged}]] . values$.

$y = df . price$.

import OneHotEncoder. categorical_features=[0])

ohe = OneHotEncoder()

$x = ohe . fit_transform(x) . toarray()$

model . fit (x, y)

model . predict ([[1, 0, 2800]])

ex:

predict price of a mercedes benz that is 4 years old with 45000 mileage.

Training and Testing

Data

A	P	Q
a	1	10
b	2	20
c	3	30
d	4	40
e	5	50
f	6	60
g	7	70

Training
80%

test.

plt.scatter(df['mileage'], df['sell price'])

> A clear relationship in linear reg.

plt.scatter(df['Age'], df['sell price'])

> short of apply of linear reg.

$x = df[['\text{mileage}', '\text{Age}']]$

$y = df['\text{sell price}']$

Fungitac
Sertaconazole Nitrate

```
from sklearn.model_selection import train-test-split
```

~~train-test-split~~

✓ ~~x-train, x-test, y-train, y-test = train-test-split(n, y, test_size=0.2)~~

len(x-train)

clf = linearRegression()

clf.fit(x-train, y-train)

clf.predict(x-test)

clf.score(x-test, y-test)

classification

Logistic Regression. (Binary Classification)

Linear Reg.

Home price

weather

Stock price

predicted value in
continuous

Classification:

Email in spam or not

will customer buy

life insurance?

which party a person is
going to vote for?

a. Democratic

b. Republic

c. independent

predicted value in
categorical

Logistic neg. is in
one of the technique
used for classification

binary classification

multivalue classification

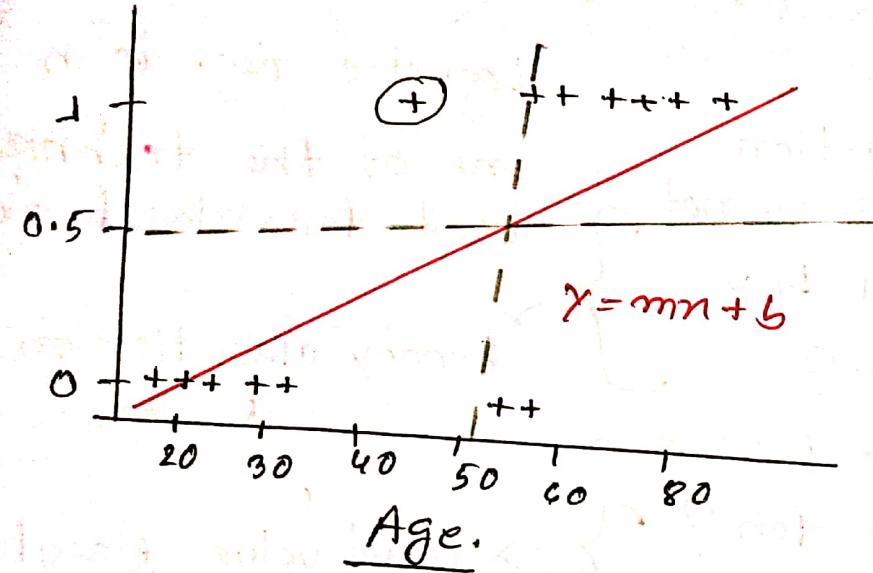
Binary classification

means

Yes $\rightarrow 1$

No $\rightarrow 0$

Fungitac
Sertaconazole Nitrate

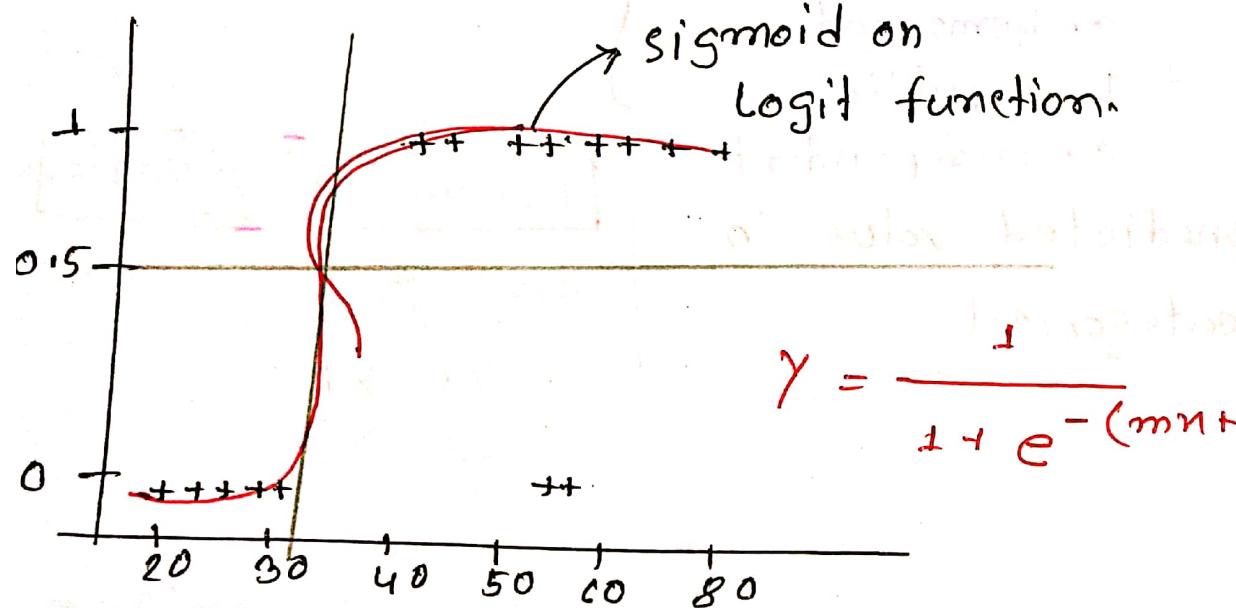


$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

$$e \approx 2.71828$$

↳ Euler's

sigmoid function converts
input into range 0 to 1



$$Y = \frac{1}{1 + e^{-(mn+b)}}$$

plt. scatter (df['age'], df['bought-insurance'])

from sklearn.model_selection import train-test-split

x-train, x-test, y-train, y-test = train-test-split (df[['age']], df['bought-insurance'], test_size=0.1)

	age
x-test	18
x-train	23
	40

(more to more score)

from sklearn.linear_model import LogisticRegression

model = LogisticRegression()

model.fit(x-train, y-train)

model.predict(x-test)?

> 0,0,1

not insurance.

model.score(x-test, y-test)

> 1 ← perfect.

model.predict_proba(x-test)

> 0.6653 0.3346

0.6061 0.3938

0.3918 0.6081 ✓

probability of buying insurance
on mot.

Fungitac
Sertaconazole Nitrate

8
4

Logistic Regression (Multiclass Classification)

```
from sklearn.datasets import load_digits
```

there are 1797 hand written obj
 8×8

```
digiten = load_digits()
```

dir(digiten) → data has numeric data

DESCR, data, images, target, target_names.

digiten.data[0] → img have actual img.

> return array length 8×8

```
plt.gray()
```

```
plt.matshow(digiten.images[0])
```

```
for i in range(5):
```

```
    plt.matshow(digiten.images[i])
```

digiten.target[0:5]

> 0, 1, 2, 3, 4, 5

```
from sklearn.model_selection import train_test_split  
train_test_split(n, y, test_size = 0.2)
```

```
model = LogisticRegression()
```

```
model.fit(x_train, y_train)
```

```
plt.matshow(digit[67]) → 6
```

```
digit.target[67]
```

```
> 6
```

```
model.predict([digit.data[67]])
```

question is how I will know where the error.

Confusing matrix.

```
y_predicted = model.predict(x-test)
```

```
from sklearn.metrics import confusion_matrix
```

Fungitac
Sertaconazole Nitrate

$\text{cm} = \text{confusion-matrix}(y_{\text{test}}, y_{\text{predicted}})$

\downarrow
truth. what model predicted.

import seaborn as sn

plt.figure(figsize

sn. heatmap(cm, annot=True)

plt.xlabel('predict')

plt.ylabel('truth')

Decision tree.
classification
problem

independent

```
input = df.drop('Salary > 100k')
```

Axi'ss = edy

`fenget = df[['galaxy']]`

converting text data into numbers.
import LabelEncoder.

Le-company = cable encoder

$L_e = \pi \rho b =$
 $L_e - \text{degree}$

data we need

```

graph TD
    Company[Company] --> facebook[facebook]
    Company --> ABCpho[ABC pho]
    facebook --> CO[C/O]
    ABCpho --> OneThird[1/3]
  
```

high info gain

mason ement of nandem ness

Degree

Bachelor Master

$\frac{4}{4}$ $\frac{6}{2}$

Low info. gain.

Job - n

degree- n

Fungitac

Sertaconazole Nitrate

inputs_n = inputs.drop(['company', 'job', 'degree'], axis='column')

ready to train classifier.

from sklearn import tree

model = tree.DecisionTreeClassifier()

model.fit(inputs_n, target)

model.score(inputs_n, target)

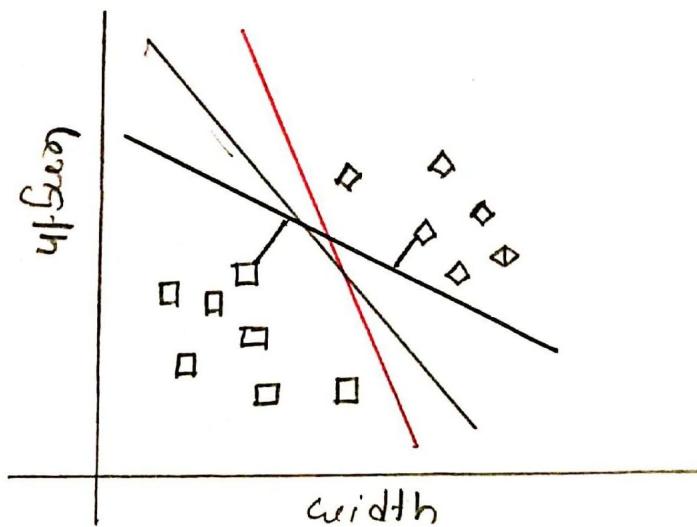
>+

cause of same
dataset.

model.predict([2, 2, 1])

> 0

SVM (Support Vector Classification Algorithm.)



Fungitac
Sertaconazole Nitrate

```
from sklearn.datasets import load_iris  
iris = load_iris()  
dir(iris)  
> data, feature_names, target, target_name.  
iris.feature_names  
sepal. length (cm)  
width  
petal length (cm)  
width
```

```
df = pd.read_csv  
df = pd.DataFrame(iris.data, columns=iris.feature_names)  
>  
df['target'] = iris.target  
>  
iris.target_names  
> setosa, versicolor, virginica  
0 1 2  
df[df['target'] == 2].head()
```

Fungitac
Sertaconazole Nitrate

```
df['flower-name'] = df.tanget.apply(lambda x: iris.tanget-namen[x])
```

```
df0 = df[df.tanget == 0]
```

```
df1
```

```
df2
```

```
df2  
> plt.scatter(df0['sepal length'], df0['sepal width']) | x = sepal length  
plt.      df1           df2 | y = sepal width.
```

df0 petal length df0 petal width
df1 df1

import train-test-split.

```
x = df.drop(['tanget', 'flower-name'], axis='columns')
```

```
y = df['tanget']
```

```
xtrain, xtest, ytrain, ytest = train-test-split(x, y, test-size=0.2)
```

from sklearn.svm import SVC

model = SVC() \Rightarrow gamma, kernel

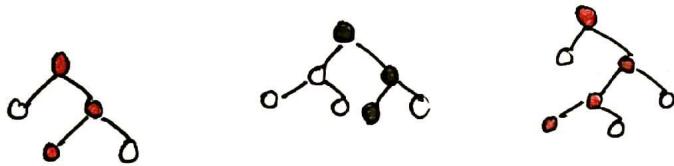
model.fit(x-train, y-train)

model.score(x-test, y-test)

Random forest Algorithm

classification and regression

when there are so many decision trees then there are a random forest.



Decision •

result in •

digiten = load-digiten()

dim(digiten)

plt.gray()

~~digiten.images[1]~~

plt.matshow(digiten.images[:])

df = pd.read_csv('digit.csv')

df.head()

~~df['target'] = digiten.target~~

Fungitac
Sertaconazole Nitrate

AJ Auto regressive language model.

```
x = df.drop(['target'])
```

```
y = digits.target.
```

```
import train-test-split
```

```
x-train, x-test, ytrain, y-test = train-test-split(x, test-size=0.2)
```

```
model(x-test)
```

from sklearn.ensemble import RandomForestClassifier

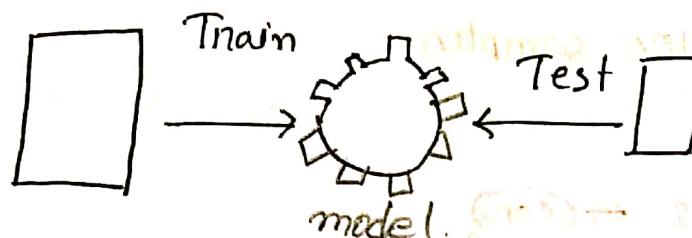
```
model = RandomForestClassifier()
```

```
model.fit(x-train, y-train)
```

kFold cross

validation.

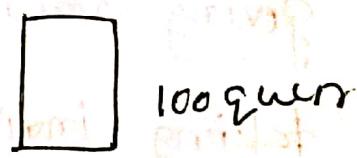
Cross validation is a technique, which allows you to measure the exact



There many several ways you perform this training step.

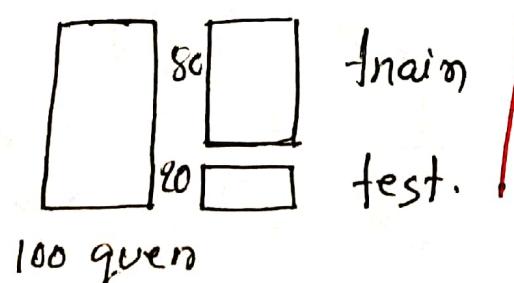
option-1 :
use all available data for training and test on same data set.

100 quer. and make the paper on 100 quer.
this test is not suitable.



option-2 :

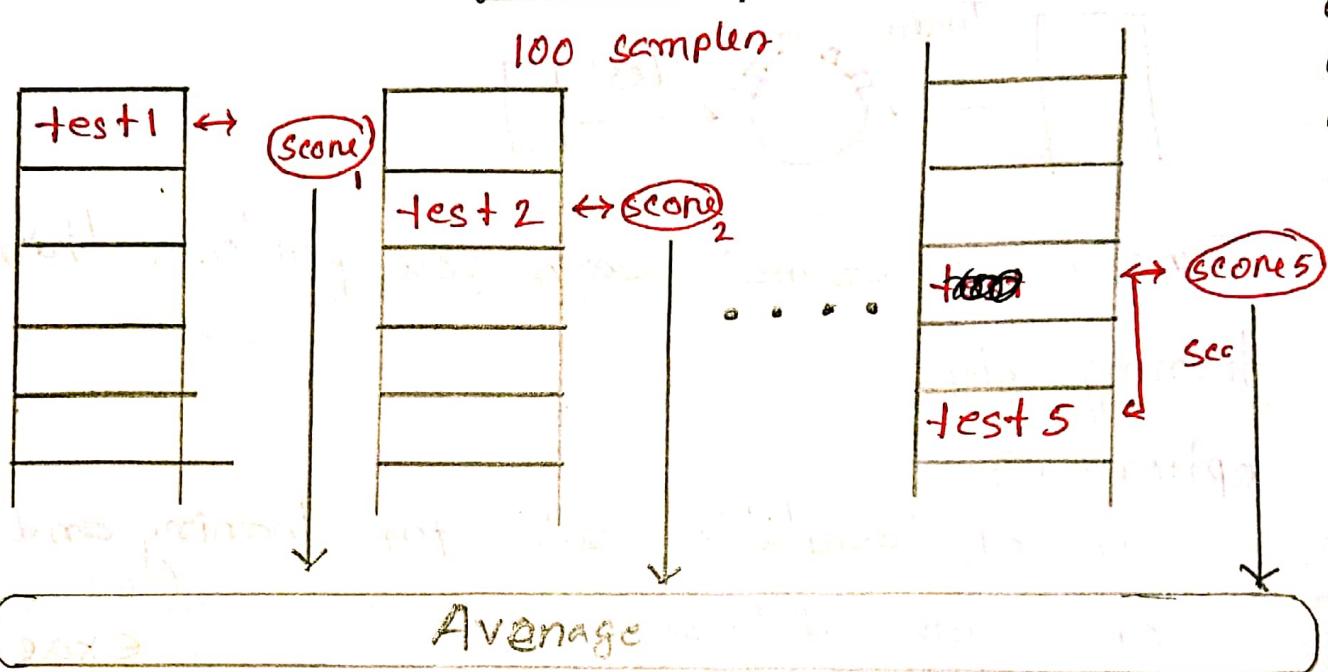
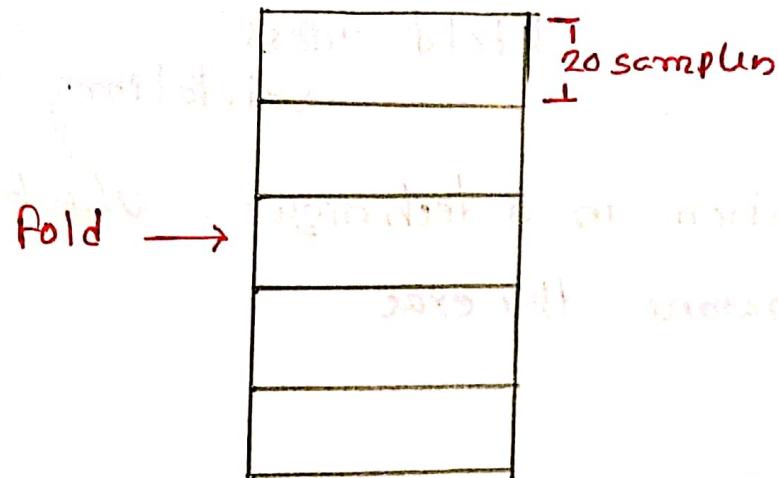
split data into train and test.



all the type of quer
are not in train

option - 3

k fold



giving vanity of sample. to model, and taking individual score. and averaging them.

After this we can take average of all the scores to apply with the model of form $y = w_0 + w_1x_1 + \dots + w_nx_n$.

Final

```
digitr = load_digits()
```

```
ttrain, test, split (digitr.data, digitr.target, test_size = 0.3)
```

①

In logistic Regression score **95%**.

In SVM score **40%**.

In Random forest Classifier score **97%**.

Thus model score value will change because of the change of test size/on split.

from sklearn.model_selection import kFold

kf = KFold(n_splits=3) → How many folds.

kf → method.

for train_index, test_index in kf.split([1,2,3,4,5..9]):
 print (train_index, test_index)

[3 4 5 6 7 8] [0 1 2]

[0 1 2 6 7 8] [3 4 5]

[0 1 2 3 4 5] [6 7 8]

↓ i.e. it is divided into 3 folds.

$$9/3 = 3$$

Sergel®
Esomeprazole USP

```
def get_score(model, X_train, X_test, Y_train, Y_test):  
    model.fit(X_train, Y_train)  
    return model.score(Y_train, Y_test)
```

```
[from sklearn.modelselection import StratifiedKFold  
foldn = StratifiedKFold(n_splits=3)]
```

score-logi = []

real datasets.

score-SVM = []

score-Random = []

foldn. [1, 2, 3, ..., 9]

for train_index, test_index in foldn.split(digits.data):

X-train, X-test, Y-train, Y-test

ablot = digits.data[train_index], digits.data[test-index]
q. / digits.target[train_index], digits.target[test-index]

score-logi.append(get_score(LogisticRegression(), X-train,
X-test, Y-train, Y-test))

score-SVM.

score-nf.

(m) from sklearn.model_selection import cross_val_score
cross_val_score(logisticRegression(), digits.data, digits
.target)

ex:-

- ① use inner flower datasets.
- ② use different classifiers.
- ③ measure the score.

k Means
clustering Algorithm.
→ step-1) scaling
step-2) find centroid
step-3) SSE + Snapshot

Machine learning algorithms are three main categories:

1. supervised → given data set, class label or target.
2. unsupervised → set of features, don't know about class label / target
3. reinforcement learning.

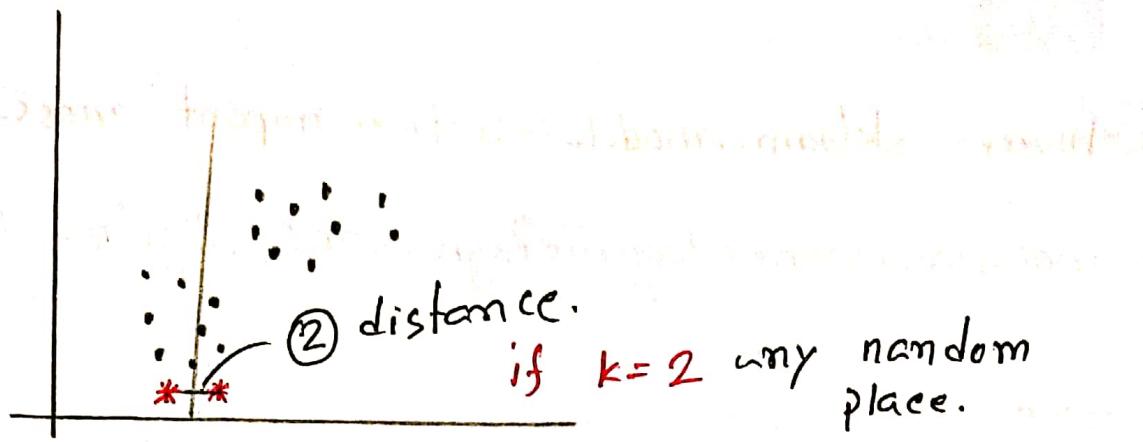
to find the
useful prediction.

} try to identify the underlying structure in that data.

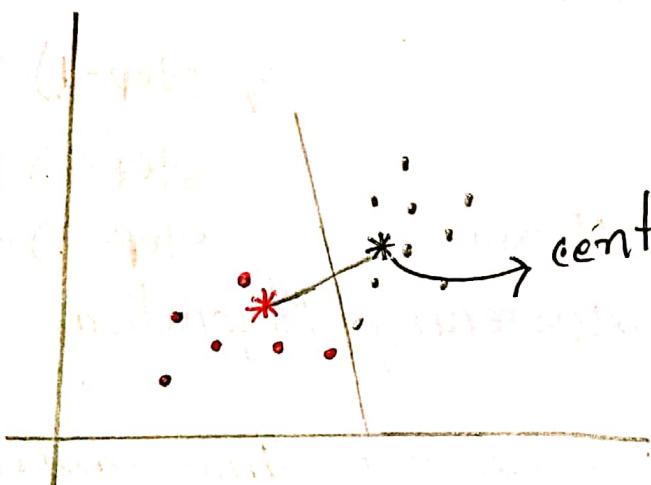
try to find data.clusters.

Sergel®
Esomeprazole USP

①



imperfect clunky clusters and now we try to improve. improve this cluster.



③ if the red point is more near to centroid. set the cluster red.

1. starts with k centroids by putting them at random place $k=2$

2. compute distance of every point from centroid and cluster them accordingly.

3. Adjust centroids so that they become center of gravity for give cluster.

- 4) again re-cluster every point based on their distance.
- 5) Again adjust centroids.
- 6) Recompute clusters

After final clustering, most important point here, need to supply k .

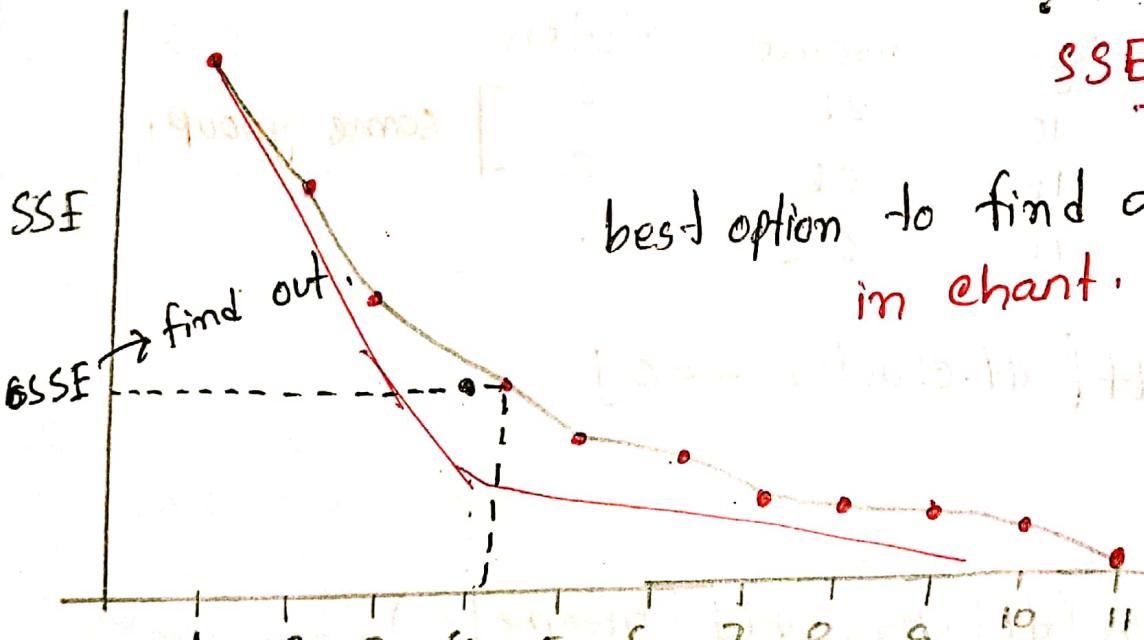
Q. What is the good number of k ?

A. There is a technique called elbow method.

They try to compute sum of square error (SSE)



$$SSE_j = \sum_{i=1}^n \text{dis}(x_i - c_j)^2$$



best option to find an elbow in chart.

```
df = pd.read_csv('income.csv')
```

```
plt.scatter(df['Age'], df['income'])
```

↳ 3 clusters.

```
from sklearn.cluster import KMeans
```

```
km = KMeans(n_clusters=3)
```

```
km.fit_predict(df[['Age', 'income']])
```

y-predicted.

```
[2, 2, 0, 0, 1, 1, ..., 1, 0, 0, 0, 0, 2, 2, 0]
```

$\frac{2}{2}$ $\frac{0}{0}$ $\frac{1}{1}$ ↳ three clusters.

```
df['cluster'] = y-predict
```

Name	Age	income	cluster
x	10	2k	2
y	12	5k	2
z	14	8k	0

↳ same group.

```
df1 = df[df['cluster'] == 0]
```

```
df2
```

```
df3
```

```
plt.scatter(df1['Age'], df1['income'])
```

df2 df2

df3 df3

from sklearn.preprocessing import MinMaxScaler.

scaler = MinMaxScaler()

scaler.fit(df[['income']])

df['income'] = scaler.transform(df[['income']])

df (the income now scaled 0 to 1)

scaler.fit(df['Age'])

then fit the data sets into KMeans

km = KMeans(n_clusters=3)

y_predited = km.fit_predict(df[['Age', 'Income']])

y_predited

df['cluster'] = y_predited.

if df.drop('cluster')

Clustering Aging.

df1 = df[df['cluster'] == 0]

2

plt.scatter(df1['Age'], df1['income'])

df2

df2

How to get the centroid.

kM.clusten - center -

plt.scatter($\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$) plot on centroid.

(kM.clusten - center - [:, 0], kM.clusten - center - [0:, 1])

k_nag = range(1, 10) ← Assuming no. of cat.

SSE = []

for k in k_nag:

km = KMeans(n_clusters=k)

km.fit(df[['Age', 'income']])

SSE.append(km.inertia_)

SSE

plot-ing SSE And K

plt.xlabel('k')

plt.ylabel(SSE)

plt.plot(k_nag, SSE)

Nive Bayer

Application: spam email, hand knitting, weather pre
face detection, Neuron, antiepile.
it contains independent variable. Sex, Age, class

```
df = pd.read_csv('titanic.csv')
```

```
df.drop(['PassengerId', 'Name', 'Class', 'Parch', 'Ticket'], axis=1)
```

```
target = df['Survived']
```

```
inputn = df.drop(['Survived'])
```

```
[dummien = pd.get_dummies(input['Sex'])]
```

```
inputn = pd.concat([inputn, dummien])
```

```
inputn.drop('Sex')
```

```
inputn.columns [inputn.isna().any()]
```

```
[inputn['Age'] = inputn['Age'].fillna(inputn['Age'].mean())]
```

```
x-train, x-test, y-train, y-test =
```

```
train-test-split (inputn, target, test_size = 0.2)
```

Create naive Bayes model.

```
from sklearn.naive_bayes import GaussianNB  
model = GaussianNB()  
model.fit(x-train, y-train)  
model.score(x-test, y-test)
```

part-2

Email spam

ham → good message.
spam

```
df = pd.read_csv('spam.csv')
```

```
df.groupby('category').describe()
```

category and message → numben.

```
df['spam'] = df['category'].apply(lambda x: 1 if x == 'spam'
```

```
else 0 )
```

→ message(text)

x-train, x-test, y-train, y-test =

```
train-test-split(df.message, df.spam)
```

test-size = 0.25

Healthcare Date: _____

countVectorize → the message column convert into
numbenn. → comput.

each of this doc. find the unique word.

This is the
first doc

This doc is the
second doc

And this is the
third one.

and, doc, first, in, one, second
the, third, this
unique

Is this the first
doc.

3 types of classification.

1. Bernoulli → feature (0,1)

2. Multinomial → discrete data (movie rating).

3. Gaussian → normal distribution
(bell curve)

from sklearn.feature_extraction.text import

CountVectorizer.

v = CountVectorizer()

x_train_count = v.fit_transform(x_train.values)

x_train_count.toarray()[:3]

↳ Num of unique words.

Simplifying acid control

Sergel®
Esomeprazole USP

```
from sklearn.naive-bayes import MultinomialNB  
model = multinomialNB() ← for classification  
model.fit(x-train-count, Y-train)  
  
Email converted into a num. matrix  
stopwords
```

```
Email = ['first text', 'second text']  
email_count = v.tnann form (email)  
model.predict(email_count)  
[0, 1] ← 2nd email is spam.
```

In pipeline where can define a transformation

pipeline is super useful.

from sklearn.pipeline import Pipeline.

```
clf = pipeline([  
    ('vectorizer', CountVectorizer()),  
    ('nb', MultinomialNB())  
])  
clf.fit(x-train, Y-train)
```

Hypen para

How to choose the best model you given machine learning problem?

* can be classify own model, by using kfold and loop but it is not convenient.

- from `sklearn.model_selection import GridSearchCV`

```
clf = GridSearchCV(SVM.SVC(gamma='auto'))  
for hyperparameter tuning.  
param_grid: [C: [1, 10, 20],  
            kernel: [rbf, linear]]  
of this parameter.  
}, cv=5, n_jobs=-1, verbose=10)
```

```
clf.fit(iris.data, iris.target)
```

```
clf.cv_results_
```

```
df = pd.DataFrame(clf.cv_results_)
```

```
df = df[['param_C', 'param_kernel', 'mean_test_score']]
```

```
df
```

```
clf.best_score_ ← best of model.  
clf.best_params_ ← best of model.
```

```
from sklearn.model_selection import RandomizedSearchCV
RS = RandomizedSearchCV(svm.SVC(gamma='Auto'), {
    'C': [1, 10, 20], not try even single
    'kernel': ['rbf', 'linear'] permutation and combination
}, cross_val. of parameters.
cv=5, try on random params
return_train_score = False,
n_iter = 2 ← num of iteration.
```

)
ns.fit(liris.data, liris.target)
pd.DataFrame(ns.cv_results)[['params', 'mean_fit_time', 'mean_test_score']]

Choosing best model.

```
import SVM
import RandomForestClassifier
import Logistic Regression.
```

Model_Parms = {
 'SVM': {
 'model': svm.SVC(),
 'params': {
 'C': [1, 10, 20]
 'kernel': ['rbf', 'linear']
 }
 }
}



Date: _____ / _____ / _____

'random_forest': {

'model': RandomForestClassifier(),

'Pnames': {

'n_estimators': [1, 5, 10]

}

},

Logistic_Regression: {

'model': LogisticRegression()

'Pnames': {

'C': [1, 5, 10]

}

},

Scans = []

for model name, np, in model_Pnam

x

|

```
scons = []
```

```
for model-name, mp in model-param.items():
```

```
    clf = GridSearchCV(mp[model], mp[pname], cv=5,  
                       return_train_score=False)
```

```
    clf.fit(inis.data, inis.target)
```

```
scons.append({})
```

```
    'model': model-name,
```

```
    'best-score': clf.best_score_,
```

```
    'best-pname': clf.best_param_
```

```
})
```

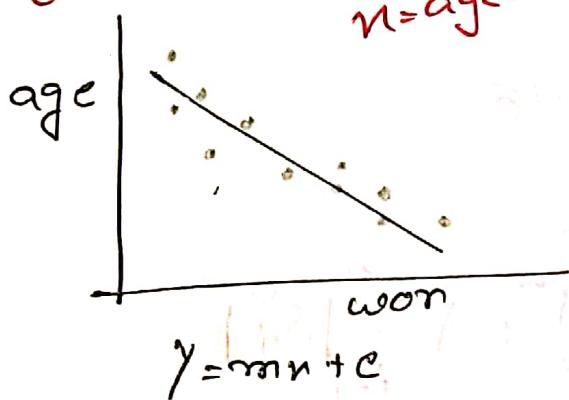
```
df = pd.DataFrame(scons, columns=['model', 'best-score',  
                                   'best-pname'])
```

Regularization

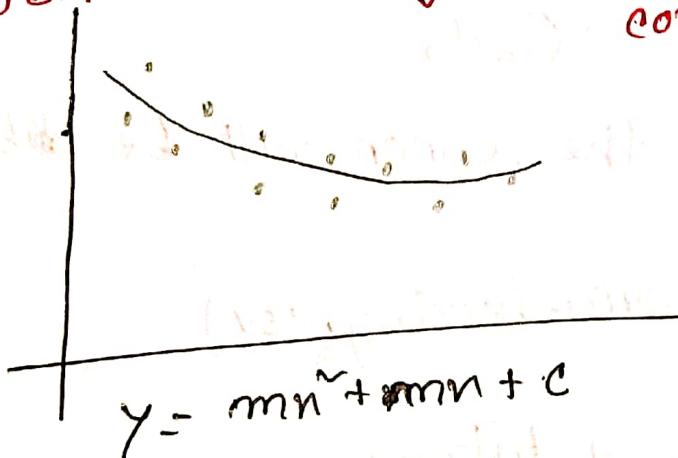
L1 and L2 are used in overfitting issue.

How a model is overfit and can be addressed
on accuracy can be improved.

① Under fitting ② over fitting



③ Balanced fitting.



then equ. is too much complicated.

thin in the balance of 1 and 2

thin in the balance of 1 and 2

How to reduce overfitting.

In overfitting eqn.

$$y = \theta_0 + \theta_1 x_{age} + \theta_2 x_{age^2} + \theta_3 x_{age^3} + \theta_4 x_{age^4}$$

try to make θ_3, θ_4 almost close to zero.

$$y = \theta_0 + \theta_1 x_{age} + \theta_2 x_{age^2} \quad [\text{shrink the parameters}]$$

we did this by using mse

L1

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$$

$$h_\theta(x_i) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3^3$$

if the θ is biggen the error will be biggen.

```
df = pd.read_csv('melbourne-housing.csv')
```

df

Data cleaning + fillna.

```
df.isna().sum()
```

```
cols_to_fill_zero = ['propertycount', 'Distance', 'Bedroom2',  
                     'Bathroom', 'Car']
```

use categorical feature using dummy variable

$x = df.drop(['price', axis=1])$

$y = df['price']$

$xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.3)$

import LinearRegression

reg = LinearRegression()

reg.fit(xtrain, ytrain)

reg.score(xtest, ytest) underfitting

> 0.13

reg.score(xtrain, ytrain) overfitting.

> 0.68

from sklearn import linear_model

Lasso-reg = linear_model.Lasso(alpha=50, max_iter=50, tol=0.1)

Lasso-reg.fit(xtrain, ytrain)

Lasso-reg.score(xtest, ytest)

L1