

13/10/22
Thursday

8 Linear Regression — It is a machine learning algorithm based on Supervised Learning.

It is a Statistical method that is used for predictive analysis. Linear regression makes predictions for Continuous / real or numeric variable such as Cost, age, Sales, temperature, product price, etc.

Types —

- i) Linear Regression with Single values.
- ii) Linear Regression with Multiple values.

(Ex)

```
>>> import pandas as pd } for Single variables
       >>> import seaborn as sns
       ①   >>> from sklearn import linear_model. (+ for line)
```

demo (2) -

```
>>> data1 = pd.read_csv('PUBH')
```

	age	Premium	
0	25	18,000	Whole data set
1	30	32,000	
2	35	40,000	
3	40	47,000	
4	45	55,000	

Linear Regression

Only the main code -

(1) `>>> import statsmodels.formula.api as smf
>>> lm = smf.ols(formula = "dependent variable ~ feature1 + feature2 + ... + featureN", data = dataset1).fit()
>>> lm.summary()`

(2) # Using Standard Scales -

`>>> X = df.iloc[:, :-1]
>>> y = df.iloc[:, -1]`

`>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler()`

`>>> scaler.fit_transform(X, y)` or

`>>> X = scaler.fit_transform(X, y)`

Using Lasso, Ridge & elasticnet in Linear Reg.

① \ggg from sklearn.linear_model import Ridge, Lasso, ElasticNet
 \ggg from sklearn.linear_model import RidgeCV, LassoCV, ElasticNetCV

② \ggg lassocv = LassoCV(cv=5, max_iter=20000, normalize=True)
 \ggg lassocv.fit(x_train, y_train)

③ \ggg lassocv.alpha_

\ggg lasso = Lasso(alpha=lassocv.alpha_)
 \ggg lasso.fit(x_train, y_train)

④ \ggg lasso.score(x_test, y_test)

∴ as we use the standardization so before predict we have to transform it.

\ggg t = ~~scaler~~.transform([[1, 2, 3, ..., 577]])

\ggg lasso.predict(t)

Ridge -

⑤ \ggg ridgecv = RidgeCV(alphas=(0.1, 1.0, 10.0), cv=10, normalize=True)

\ggg ridgecv.fit(x_train, y_train)

If alpha & lambda are same then

⑥ $\gg \text{ridge}_{\text{cv}}.\alpha$ -

$\gg \text{ridge}_{\text{cv}} = \text{Ridge}(\alpha=\text{ridge}_{\text{cv}}.\alpha)$

$\gg \text{ridge}_{\text{cv}}.\text{fit}(X_{\text{train}}, y_{\text{train}})$

⑦ $\gg \text{ridge}_{\text{cv}}.\text{score}(X_{\text{test}}, y_{\text{test}})$

(Predict will be same)

Elastic Net -

⑧ $\gg \text{elastic} = \text{ElasticNetCV}(\alpha_{\text{cv}}=\text{None}, \text{cv}=10)$

$\gg \text{elastic}.\text{fit}(X_{\text{train}}, y_{\text{train}})$

$\gg \text{elastic}.\alpha$ (Fit in lambda)

~~$\gg \text{elastic}$~~

$\gg \text{elastic}.\text{l1_ratio}$ (# it is actual alpha on

$\gg \text{elastic}.\text{l1_ratio}$ elastic net)

⑨ $\gg \text{elastic}_{\text{cv}} = \text{ElasticNet}(\alpha=\text{elastic}.\alpha, \text{l1_ratio}=\text{elastic}.\text{l1_ratio})$

$\gg \text{elastic}_{\text{cv}}.\text{fit}(X_{\text{train}}, y_{\text{train}})$

⑩ $\gg \text{elastic}_{\text{cv}}.\text{score}(X_{\text{test}}, y_{\text{test}})$

Note - To predict the value of the train dataset for the user input remember that that if we use Standardization (Standard-Scaler) we have to Transpose the input.

①
 »» t = Scaler.transform([[337, 118, 4, 4.5 ... , 15]])
 »» elastic_lr.predict(t)

- To export the model in pk1 -

»» import pickle
 »» pickle.dump(elastic_lr, open("elastic.model.pk1", "wb"))

$$\begin{bmatrix} \text{wb} = \\ \text{rb} = \end{bmatrix}$$

$$y = m_1 x_1 + m_2 x_2 + \dots + c$$

$$m = \text{coef / Slope} = \theta,$$

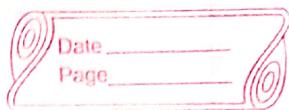
$$c = \text{Intercept} = \theta_0$$

$$y = m x + c \text{ or}$$

$$y = \beta_1 x + \beta_0$$

~~16/10/22~~
~~Sundays~~

Treatise



Logistic Regression (Binary classification)

ii) Logistic Regression:- It is a machine learning algorithm based on Supervised learning.

It is a Statistical method that is used for predicting probability of target variable.

Logistic Regression makes probability for classification problems that are discrete in nature.

Ex - English or Hindi, True or False, 1 or 0, Right or wrong, cat or dog or goat

Types of Logistic Regression -

i) Binary classification - Ex - Win or loss, dead or alive.

ii) Multiclass classification - Ex - Onion or Potato or Sweet Potato, Lily or Sunflower or Rose

iii) Real life example of Logistic Regression
are used in hospital to check whether what is the chance of getting heart attack of patients.

→ Heart Attack Prediction /
→ Cancer Prediction

- Logistic Regression -

- # In Logistic Regression we Solve Classification Problem.
- # OVO → One Versus One → Binary Class Classification
OVR → One Versus Rest → Multi-Class Classification
- # If your dataset is Skewness that it is generally a good idea to replace Non value or (any value or '0') in mean
- >> df['BM'] = df["BMT"].replace(0, df["BMT"]
• mean())
- # If there are skewness in the data if mean the column has a "Outliers".
- # If there is so many Outlier in a single feature that we will try to remove reduce the data from that particular feature.

drop 10% data.

```

>>> q1 = df["Insulin"].quantile(.90)
>>> df_new = df[df["Insulin"] < q1]
    # This is new data set after dropping
  
```

Q How to remove Outlier?

Ans) To experiment with quantile range we can remove outlier as it ~~is~~ one of the reason.

» $q_1 = df["Insulin"] . quantile (.95)$

» $df_new = df[df["Insulin"] < q_1]$
if after this

» ~~age effect~~

» $q_1 = df["Age"] . quantile (.95)$

» $df_new = df[df["Age"] < q_1]$

(H Quantile is to drop the data from data set to remove outliers)

Q Do we do Regularization in Logistic Regression?

Ans Yes we do Regularization in Logistic Regression but by default.

• Solver = "liblinear" \rightarrow that it do L2 regula. by default and

• Solver = "saga" \rightarrow that it is elastic net (L1 + L2) regularization.

• Solver = "liblinear" \rightarrow It support both L1 and L2 reg. with a dual formulation under the L2 penalties

Note → for small dataset, "liblinear" is a good choice, whereas 'sgd' 'Sag' are faster for large ones.

→ liblinear is limited to one-versus-rest

→ We will not use "liblinear-Solver" for multi-class classification

Logistic Regression work on Sigmoid function:-

$$\text{Sigmoid} = \frac{1}{1+e^{-x}}$$

If we feed an output value to the Sigmoid function, it will return the probability of the outcome between 0 & 1.

If the value is below 0.5, then the output returns as No / Fail. If it above 0.5 that mean the output is returned as Yes / Pass.

Parameters

Verbose - It internal think which we want to see.

Confusion Matrix :- used in order to evaluate the performance of a classification Model
It is $N \times N$ matrix

		Actual Value	
		Positive	Negative
Predicted Value	Positive	T.P	F.P \leftarrow Type 1 error
	Negative	F.N	T.N

Type 2 error

- when Predicted Value Matches with Actual \rightarrow T.P
- when Value is False & model is also saying false \rightarrow F.N
- when the Actual Value is False & Predicted Value is saying True \rightarrow F.P (Type 1-error)
- when the Actual Value is True & model predict False \rightarrow F.N

$$\text{Accuracy} = \frac{T.P + T.N}{\text{Total}}$$

$$\text{Precision} = \frac{T.P}{(\text{pred. Yes}) T.P + F.P}$$

$$\text{Error Rate} = 1 - \text{accuracy}$$

$$\text{recall} = \frac{T.P}{(\text{actual yes}) F.N + T.P}$$

$$\frac{(\text{Type 1 error} + \text{Type 2 error})}{T.P + F.N}$$

Prediction Value

		(Pos) 1		(Neg) 0	
		Y (Pos)		F.N	
Actual Value	1	T P			F.N \leftarrow Type B error
	0	F P		T N	

Type A error

```
>>> from sklearn.metrics import confusion_matrix
>>> plot_confusion_matrix
```

```
>>> from sklearn.linear_model import LogisticRegression
```

```
>>> logr = LogisticRegression(verbose=1)
```

```
>>> logr_liblinear = LogisticRegression(verbose=1, solver="liblinear")
```

```
>>> logr.fit(X_train, y_train)
```

```
>>> logr.liblinear.fit(X_train, y_train)
```

```
>>> y_pred = logr.predict(X_test)
```

```
>>> y_pred_liblinear = logr.liblinear.predict(X_test)
```

```
>>> confusion_matrix(y_test, y_pred)
```

```
>>> confusion_matrix(y_test, y_pred_liblinear)
```

Now we will make function -

```
>>> def model_eval(y_test, y_pred):
    [tn, fp, fn, tp = confusion_matrix(y_test,
                                         y_pred).ravel()]
    accuracy = (tp+tn)/(tp+tn+fp+fn)
    precision = tp/(tp+fp)
    recall = tp/(tp+fn)
    result = {"accuracy": accuracy, "precision": precision,
              "recall": recall}
    return result
```

```
>>> model_eval(y_test, y_pred)
```

```
>>> model_eval(y_test, y_pred_liblinear)
```

ROC AUC Score -

The ROC AUC Score tells us how efficient the model is. The higher the AUC, the better the model performance.

```
>>> from sklearn.metrics import roc_curve, roc_auc_score
```

```
>>> roc_auc_score(y_test, y_pred)
```

```
>>> roc_auc_score(y_test, y_pred_liblinear)
```

ROC & AUC (Threshold = 0, 0.2, 0.4, 0.6, 0.8)

① \hat{y} (Actual value)

	\hat{y}	$\hat{y}(0)$
	1	1
0.	0.8	1
0.	0.96	1
1	0.4	1
1	0.3	1
0	0.2	1
1	0.7	1

②

	1	0	Actual
Pred Rely	1	TP	FP
	0	FN	TN

If Threshold = 0

\therefore Sensitivity \rightarrow TPR (True Positive Rate) = $\frac{TP}{TP+FN}$

$$\Rightarrow \frac{4}{4+0} = 1$$

\therefore Specificity \rightarrow FPR (False Positive Rate) = $\frac{FP}{FP+TN}$

$$\Rightarrow \frac{2}{2+0} = 1$$

TPR

1 -

\rightarrow FPR

• ROC & AUC should be only used for classification problem, to describe threshold and how the model performed.

$$\text{ROC} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

• Classification threshold

• Confusion matrix

• Precision

• Recall

• F1 score

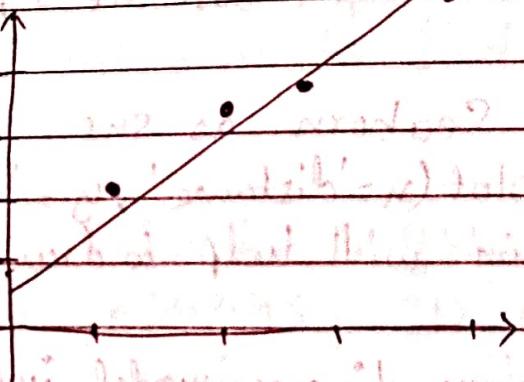
• Area under ROC curve

~~11/10/22~~
Fridays

Date _____
Page _____

10 Polynomial Regression: It is a form of linear regression in which the relationship between the independent variable x and dependent variable y is modeled as an n^{th} degree polynomial.

→ Linear model:

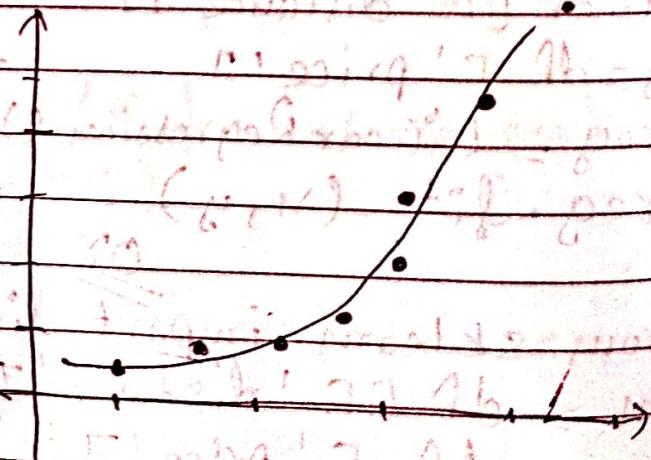


→ Polynomial Model:

0 Degree Polyno.: $y = \text{constant}$

1 Degree Polyno.: $y = mx + c$

2 degree Polyno.: $y = ax^2 + bx + c$



Generalized Polynomial Equation

↓

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots + a_nx^n$$

Column ① $\text{dataset} = \text{df1} = \text{pd.read_csv(r'path')}$

<u>Position</u>	<u>Level</u>	<u>Salary</u>
0	Bussin-Analy	1 45,000
1	Junior Cons.	2 50,000
2	Senior Cons.	3 60,000
3	Manager	4 80,000
4	Country Manager	5 110000
5	Region Manager	6 150000
6	Partnert	7 200000
7	Senior Partnert	8 230000
8	C-Level	9 300000
9	CEO	10 1000000

Column ② $\ggg x = \text{df1.iloc[0:10].values}$ (# independent Variable)

$\ggg x$

Output = array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

Column ③ $\ggg y = \text{df1.iloc[0:10].values}$

$\ggg y$

Output = array([45000, 50000, 60000, 80000, 110000, 150000,

200000, 230000, 300000, 1000000])

Code no ④ \ggg import matplotlib.pyplot as plt

Code no ⑤ \ggg plt. Scatter (x, y)

Output:

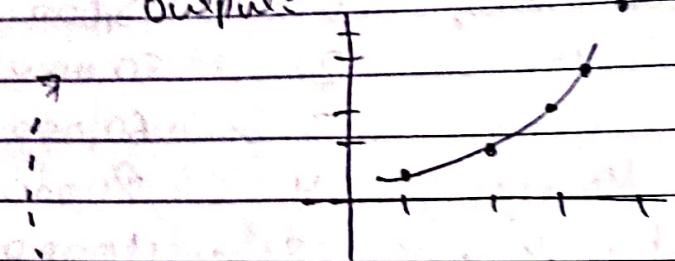


Fig - 1 (Polynomial)

Code no ⑥ Now we will check whether the data suit in 'Polynomial regression' or 'Linear regression' by drawing Scatter plot

\ggg sns. lmplot (x='level', y='Salary', data=df)

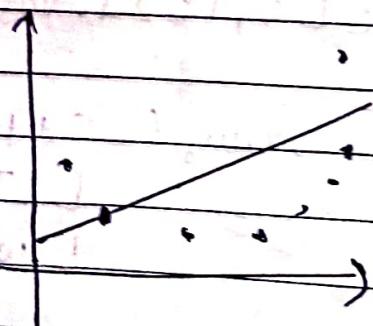


Fig - 2
(Linear Reg.)

Hence Fig-1 is Polynomial Reg. and Fig-2 is Linear Reg. If it is clear that point data is close to best with Polynomial Reg as best fit line touch most value in Polynomial.

import Linear regression

```
>>> from sklearn import linear_model
>>> reg = linear_model.LinearRegression()
>>> reg.fit(X, y)
    } We know already
    } that linear
    } regression will be
    } false, but we just
    } check it.
>>> reg.predict([[6.5]])
    } hence the linear regression gives the
    } wrong prediction, so we will use polynomial
    } regression.
```

import Polynomial regression

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(degree=2)
>>> X_poly = poly.fit_transform(X)
>>> reg2 = linear_model.LinearRegression()
>>> reg2
    } output will be a Linear Regression() but it's a Poly linear()
```

`>> reg2.predict(poly.fit_transform([6.5]))`

Output - array ([189498.107])

∴ here the prediction is almost correct

Q) In mean why do we have to use Mean Square Error (MSE)

Ans The reason is that we have only one Global Minima no local minima

