

Multiple Regression

Date - Feb 13

Multiple linear Regression | model building

One of main importance of ML is that to find best accordance to someone on top of the dataset. So in multiple feature so some feature will get more weight in contrast to another on basis of specific needs.

At such case we are going to use multiple linear regression

Eqn

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

↓ ^{m₁₂ I.i.d₂}
 dependent variable m₁ T
 I.M.I. Slope Independent
 n n variable n

I.M. \Rightarrow Independent variable

So In multiple linear regression each independent variable that is feature has own weight (that is, slope)

Assumptions of linear Regression

1) linearity

there should be linear Relation b/w y and x

2) Homoscedasticity

there should be equal variance in dataset

3) Multivariate normality

There should be some normal relationship in distribution

4) Independence

There should be no pattern ~~to~~ be inside the data

5) Lack of multicollinearity

Predictors should be correlated $X_1 \nparallel X_2$



6) The outlier check (not assumption)

a single outlier can make significant change in line from other points (It is Imp to check for outliers).

dummy variables

So as we know dataset has categorical variable in it so in order to deal with it we have to use dummy variable

⇒ for numerical data we can directly assign some independent variable but for categorical we can't.

Ex:-

State
NYC
California
NYC
California
California

⇒

using dummy variable

NYC	California
1	0
0	1
1	0
0	1
0	1

Dummy Variable works like switches

In above manner we can use categorical variable and as they are in numerical data we can able assign some independent variable to it.

for this chart we can use

~~b~~ b \oplus D \Rightarrow I \oplus d \oplus 0

so if $D_{NYC} = 0 \Rightarrow$ NYC $D_{CA} = 1$ California so it is not case that we can have to use multiple independent variable for a some set of value

→ There exist a case called dummy variable trap which we will learn further

Dummy variable Trap



we learn that we never include all dummy variables in eqn why? (There exist 0 whichable to handle this) # while working if you have n dummy variable only include

(n-1) dummy variable in eqn

so If you have 2 set (different) so apply n-1 dummy variable only to both of the sets.

(Imp) Always omit one dummy variable.

Hypothesis \Rightarrow

P-values

○ statistical significance

Say Toss a coin

p-value probability of
an event happening

H_0 : This is a fair coin

H_1 : This is not a fair coin

\rightarrow we have to find which hypothesis true
 \rightarrow after that consider a hypothesis among
two & try to prove it wrong
 \Rightarrow so if it is wrong another is true
use original is true

α = threshold of p-value & after it we will reject the
hypothesis

		Heads	Tails	Pvalue	assume we are in H_0 environment
Ex.					
Toss 1	①			0.5	
Toss 2	②			0.25	
Toss 3	③			0.125	
Toss 4	④			0.06	
Toss 5	⑤			0.03	
Toss 6	⑥			0.015	

pvalue decreasing
 \downarrow
Say $\alpha = 0.05$
and after it we
will
 $\alpha = 0.5$ reject hypoth

so as H_0 is false according to α value so H_1 is true.

that

so as $\alpha = 0.05$ so $1 - \alpha$ = confidence %

so rejecting hypothesis with 95% Confidence

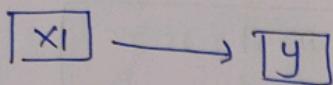
~~fixed~~ # we can set confidence level on our own

so it depends on the experiment

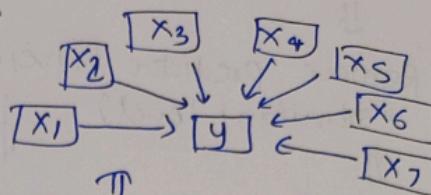
\Rightarrow α is value which ~~start~~ make us suspicious about
the hypothesis

Multiple linear Regression Intuition

Here we will learn to build a model step by step.
for simple linear regression



single independent variable and on basis of it predict dependent variable



↑
In this is multiple linear regression case

⇒ so now we have to select few among all and ignore all on basis of our need.

Reason can be

- (A) we have to reduce garbage from model (unnecessary features) which make model unreliable
- (B) keep only important among a large set

5 methods for building models

- 1] All in (throw all features in model building)
 - 2] Backward Elimination
 - 3] Forward Elimination
 - 4] Bidirectional Comparison
 - 5] Score Comparison
- Stepwise Regression
among all

1] All in : you have a good data where all things are important or prior knowledge or compulsion to use all

or preparing for backward elimination

2] Backward elimination: to stay model
step 1) select a significance level to fit in model
↓
stay in

step 2) Fit the full model with all predicting

Step 3 : Consider the predictor with highest p-value,
If $p > SL$ go to step 4 else to Find stop → Finish

↓
Step 4 : Remove predictor having value less than SL
(significant level) { All we do in it is removal of feature }

↓
Step 5 : Fit the model without this variable

that is which go out will leave model &
we are going to rebuild the model.

and after rebuilding goto step 3 or until

so do it until our highest p value will be less than significance level

(just we do inside the loop problem as when $p < SL$
we come out and that will be final model).

Forward Selection

Step 1 ⇒ Select Significance level (SL)
↓ to enter the model

→ Create simple regr
for each feature with y

Step 2 ⇒ Fill all simple regression models $y \sim X_n$
select one with lowest p value

↓

Step 3 ⇒ keep this variable and fit all possible models
with one extra predictor added in one(s)

↓

Step 4 ⇒ Consider predictor with lowest p value If $p < SL$
go to Step 3 else go to FIN

In forward all we

do is adding the features

Bidirectional model

Step 1 \Rightarrow Select Significant level to stay ~~and~~ and significant level to enter in model

$$SLENTER = 0.05 \quad \uparrow \quad SLSTAY = 0.05$$

forward selection

backward selection

↓

Step 2 : perform the next step of forward selection
(new variables must have $p < SLENTER$ to enter)

↓ add new variable on basis forward selection

Step 3 : Perform all backward Elim method

(old variables must ~~not~~ have $p < SLSTAY$ to stay)

↓

Step 4 : No new variables can enter or leave so we can exit.

↓ FIN: model is ready

\Rightarrow with forward selection adding variables on basis of SLENTER and removing variable on basis of SLSTAY (i.e backward selection)

All possible

most resource consuming approach

Step 1 : Select criteria of goodness of fit

↓

Step 2 : Construct all possible Regression model : $2^n - 1$ Combinations

↓

Step 3 : Select one with best criteria

↓

FIN model is ready

Tip:

So If you have 10 columns you have 1023 models

\Rightarrow this is research consuming

\Rightarrow we are going to use backward elimination as it is fastest one to build a model.

Coding Part multiple linear Regression

We want to find out best startup to invest on basis of dataset

Step A => Preprocessing

① Encoding

in one hot encoder & column transformer we pass X as a whole

② Feature ~~set~~ scaling

(Ppto)

Performance will give info about model

In multiple linear Regress we don't have to apply for feature Scaling

If your data is not linear so it will perform poorly for multiple linear Regression model

\Rightarrow You don't have to check for assumptions of MLR for dataset, as if assumptions are not true the model will perform badly

\Rightarrow It will be waste of time to check for MLR assumptions

③ do we have to worry about dummy Variable trap?

\Rightarrow no model take care about it

do we have to apply backward elimination?

no scikit learn will take of everything for you.

Code for model building for MLRM

class for MLRM is same for that of SLR

So code

```
from sklearn.linear_model import LinearRegression
```

```
Lr = LinearRegression()
```

```
Lr.fit(xtrain, ytrain)
```

& that will build your model

Predicting Test results

plot real profit vs predicted

$y_{pred} = \text{regressor. predict}(X_{test})$

print through np precision upto 2 points

np.set_printoptions(precision=2)

now concatenation of two vector

print(np.concatenate((yPredicted, reshape(len(yPred), 1)))

↑
Reshape
 y_{pred}

tuple

insert()

↑
tmp for ea

→ for sideways view

m = pd.DataFrame()

m.insert

↑
take 3 arguments

column = "name of column")

value = np array

loc ⇒ column no

* To ~~set~~ round up values

in np array you can use

np.set_printoptions(precision=2)

You can also use round function for this