

one hot encoding \Rightarrow converting categorical variable to numeric

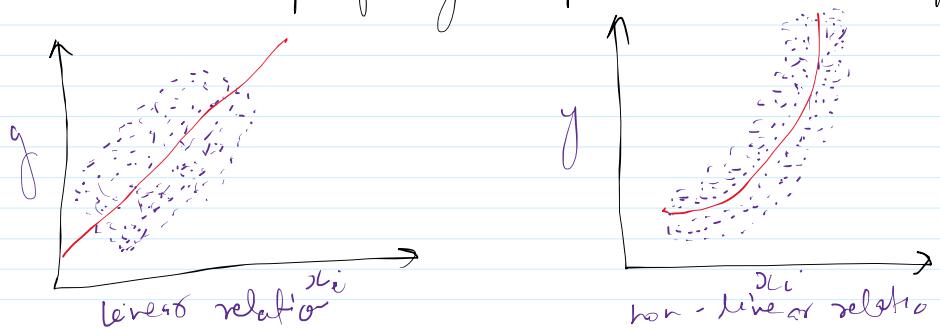
Gr.	origin	Car	origin-USA	origin-Asian	origin-Euro
Tesla	USA	Tesla	1	0	0
Ford	USA	Ford	1	0	0
Hyundai	Asia	Hyundai	0	1	0
BMW	Europe	BMW	0	0	1

only two col. id enough
the now where we have 0 and 0
can be interpreted as Europe

Assumption:

For any data linear regression model works when these assumption hold true

(1) Assumption of linearity \rightarrow i.e. Relationship of any independent variable and dependent variable shall be linear



(2) Features are not multi-collinear:

i) independent variables do not depend on each other

How to check

Treat one independent variable as y and others as x
calculate R^2 and VIF

$$VIF = \frac{1}{1 - R^2}$$

if $VIF > 10$ very high multi-collinearity

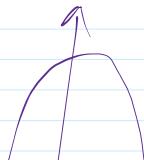
$5 < VIF < 10$ high multi-collinearity

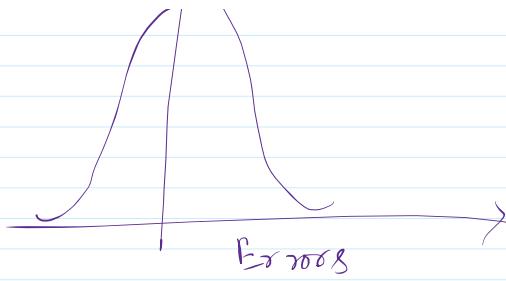
$VIF < 5$ feature can be used

Do the same for each independent variable

VIF \Rightarrow variance inflation factor

(3) Errors are normally distributed:





④ ~~Error~~ Target vs. Error does not show any pattern

⑤ Predicted vs. Error shows linear relationship

Polynomial Regression :-



Data is not linear



In equation polynomial regression tries x^2, x^3, x^4 etc. instead of x .

example $x_1^2 + x_2^2 + x_3^2 + x_4^2 + \dots$

But how do we know what is good x^2 or x^3 or such?

① Start with simple linear equation and calculate R^2 if it is very small e.g. 0.5

from sklearn.preprocessing import PolynomialFeatures

Poly = PolynomialFeatures(degree=1)

X_poly = poly.fit_transform(X)

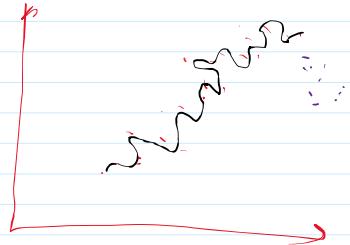
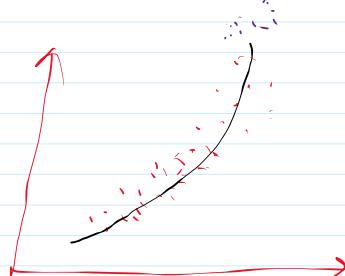
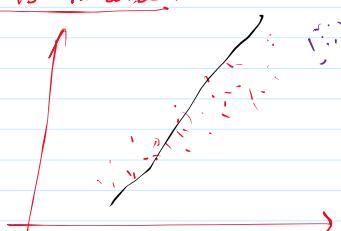
X	degree=1		degree=2			degree=3		
	1	2	1	2	4	1	2	4
1	1	1	1	1	1	1	1	1
2	2	4	4	8	16	2	4	8
3	3	9	9	27	81	3	9	27
4	4	16	16	64	256	4	16	64
5	5	25	25	125	625	5	25	125

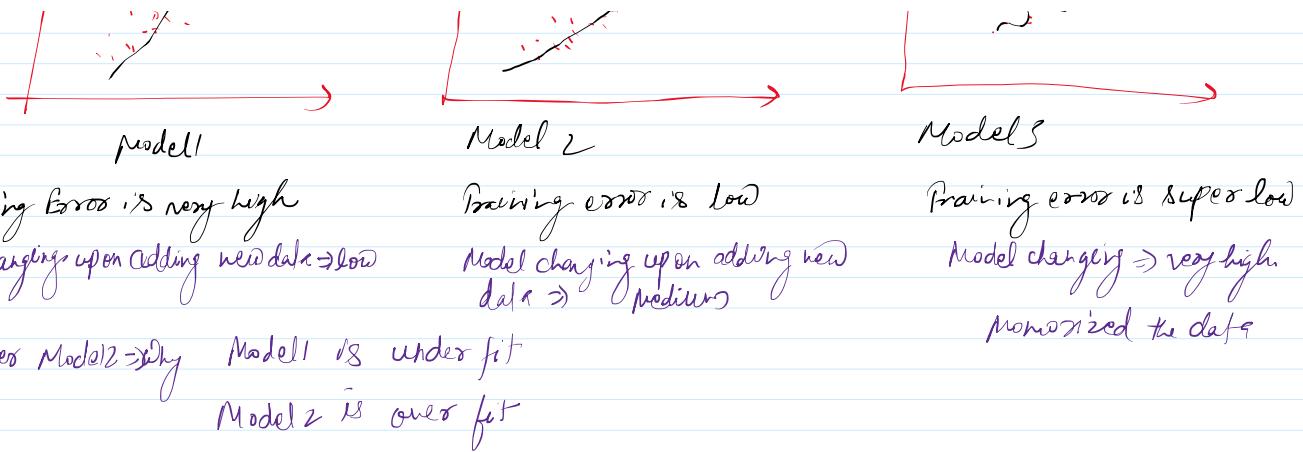
model = LinearRegression()
model.fit(X_poly)

Data \Rightarrow LR $\Rightarrow R^2 \Rightarrow$ Degree 2 $\Rightarrow R^2 \Rightarrow$ Degree 3 $\Rightarrow R^2$

Overset and underset

Error vs. Variance :-





on test data Model 1 and 2 will perform bad.

in overfitted model for 2 data points which are very close prediction will be very far apart while in good model predictions would be quite close as well

overfitted model performs very well in testing but performs very bad in testing data

underfitted model performs very bad on both training and testing data

for good model or best fit model both training and testing errors will be low

Polynomial leads to overfit

$$y = c + m_1 x^1 + m_2 x^2 + m_3 x^3 + m_4 x^4 + m_5 x^5 + m_6 x^6 + m_7 x^7 \dots$$

if we keep m_i very low and keep it decreasing as power goes up it will solve overfitting problem caused by high power value.

This is called regularization

Regularization: -
normally

$$\text{Loss} = \text{MSE}$$

in regularization

$$\text{loss} = \text{MSE} + m^2 \quad \text{for } y = mx + c$$

$$\text{loss} = \text{MSE} + (\beta)^2 \quad \text{for } y = \beta x + \gamma$$

Regularization adds a penalty m^2 while calculating loss which make loss very high

model tries to reduce the loss again while keeping m very low which is desirable to avoid overfitting caused by high order of power in polynomial regression making m low leads towards underfitting

$$\text{MSE} = \frac{\sum (y - (mx + \beta))^2}{n} + \sum (m_i^2)$$

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + m_4 x_4 + m_5 x_5 \dots$$

Algo will try to find a combination where m is small and MSE is also small

$$\text{MSE} = \sum (y - (mx + \beta))^2 \rightarrow \sum m^2$$

$$MSE = \frac{\sum (y - (mx + c))^2}{n} + \lambda \sum (m_i^2)$$

$\lambda \rightarrow$ Super high will lead to underfit as it will lead model to choose m value very low

$\lambda \rightarrow$ Super low will lead to overfit as it will lead model to choose m value very high

if training $R^2 >$ testing $R^2 + 5\%$. then model is overfitting