

Let's Recall,

- We were predicting the price of sold cars for CARS24.

We already know about Sklearn's linear regression.

from sklearn, — import —
↑



Break: 8:04 AM

Introduction to Statsmodel

What is statsmodel?

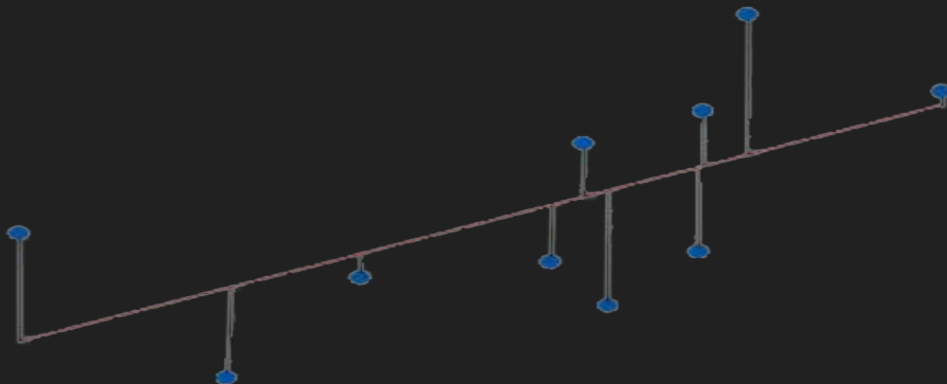
- A python module with statistical functionalities.

$$x_1 | x_2 | x_3 | x_4 | \dots$$

Stats model Library

OLS (Ordinary Least Squares)

- OLS refers to a method of estimating the parameters of a linear regression by minimizing the sum of square residuals.



How is OLS different from sklearn's Linear Regression?

OLS

- Focuses solely on estimating the parameters of a linear regression model

Sklearn

- Offers additional features and functionalities like :

- ➔ Feature scaling
- ➔ Regularization
- ➔ Cross validation
- ➔ Evaluation metrics

$x_1, x_2 \rightarrow R^2 \rightarrow 0.82$
 $\text{Adj } R^2 \rightarrow 0.82$

$x_1, x_2, x_3 \rightarrow R^2 \rightarrow 0.86$
 $\text{Adj } R^2 \rightarrow 0.70$



Assumptions of Linear Regression

Assumption of Linearity

No Multicollinearity

Normality of Residuals

No Heteroskedasticity

No Autocorrelation



Assumptions of Linearity

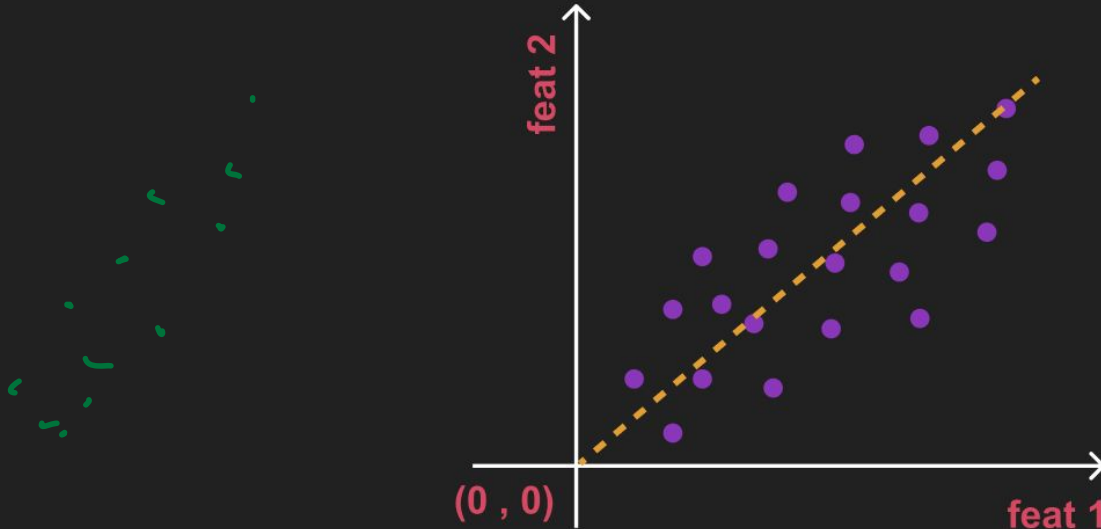
Assumption of linearity

There should be linear relationship between :

- Independent Variables {X}
- Dependent Variables {y}

Straight line fit between variables

chandra salary



No Multicollinearity

What is collinearity?

Say we have two features : f_1 and f_2

If, $f_1 = \alpha f_2 + \square$

Then f_1 and f_2 are collinear



What is multicollinearity then?



Collinearity between multiple features

Example : $f_1 = f_1, f_2, f_3$

S.T. $f_1 = \alpha_1 + \alpha_2 f_2 + \alpha_3 f_3$

Then f_1, f_2, f_3 are multicollinear



Why multicollinearity is a problem?

Say we found the optimal weights w^* for a model with 3 features

$w^*=[1,2,3]$ (corresponding to w_1, w_2, w_3) and $w_0 = 5$

$$\text{So, } \hat{Y} = w_1x_1 + w_2x_2 + w_3x_3 + w_0$$

$$= x_1 + 2x_2 + 3x_3 + 5$$

Now, let's say x_1 and x_2 are collinear.

$$\text{Then, } x_2 = 1.5x_1$$

$$\text{Hence, } \hat{Y} = 4x_1 + 3x_3 + 5$$

$$\left\{ \begin{array}{l} \therefore w^* = \langle 4, 0, 3 \rangle \\ \therefore w^* = \langle 1, 2, 3 \rangle \end{array} \right\}$$

Same classifier


This would mess up the weights and we won't be able to do feature importance.



How to deal with multicollinearity?

We will use **Variance Inflation Factor (VIF)**

Say, we have 'd' features $\langle f_1, f_2, f_3, \dots, f_d \rangle$

In, **(VIF)** we treat  one feature as 'y'
remaining features as 'x'



$f_1, f_2, f_3, \dots, f_d$	f_4
$\longleftrightarrow \mathbf{x_i} \longrightarrow$	$\longleftrightarrow \mathbf{y} \longrightarrow$

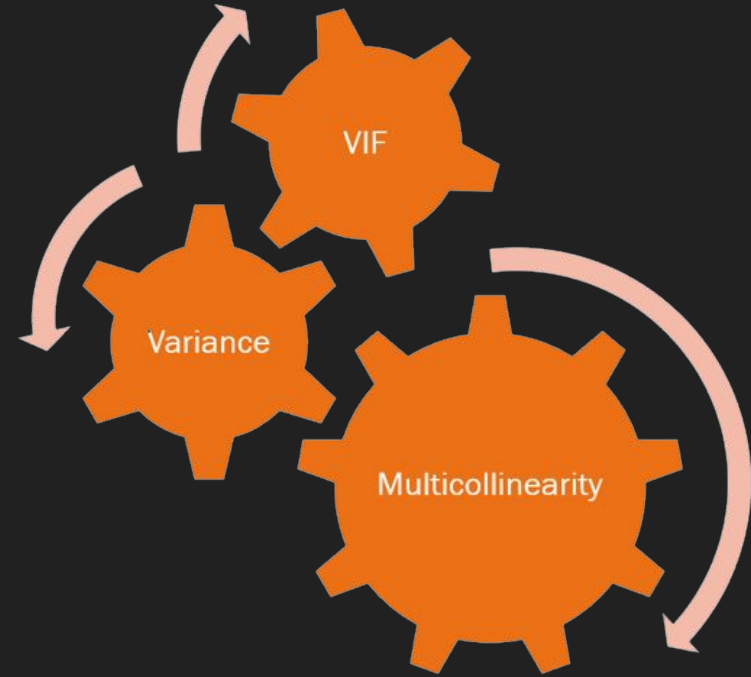
Now,

Train **linear regression** model with (x_i, y)

Find R^2 of the model

To Calculate VIF :

$$VIF = \frac{1}{1 - R_j^2}, R_j^2 : R^2 \text{ for } j^{\text{th}} \text{ feature}$$



Step 1 : Start with a full regression model including all the independent variables.

Step 2 : Calculate VIF for each independent variable by regressing it against all other independent variables.

Step 3 : Check for variable with highest VIF, [Thumb rule on next slide]



Step 4 : Remove variable with highest VIF.

Step 5 : Re-fit the model without the removed variable.



Repeat steps 2 - 5 : Continue this until no variable has VIF above threshold.



Case 1

- If $R^2 \approx 1$
 $VIF = 1/1-1 = \infty$
High R^2 means

Feature is highly collinear

Can be removed

Case 2

- If $R^2 \approx 0$
 $VIF = 1-1/0 = 1$
Low R^2 means

Feature is not collinear

Don't remove

Thumb Rule

- $VIF > 10$: Very high multicollinearity, drop
- $5 \leq VIF \leq 10$: High multicollinearity
- $VIF < 5$: Low multicollinearity

****NOTE :** We do this process for each feature.
Calculate the VIF and based on that we keep / remove feature

Normality of Residuals

Residuals/ Errors follow multivariate normal distribution.

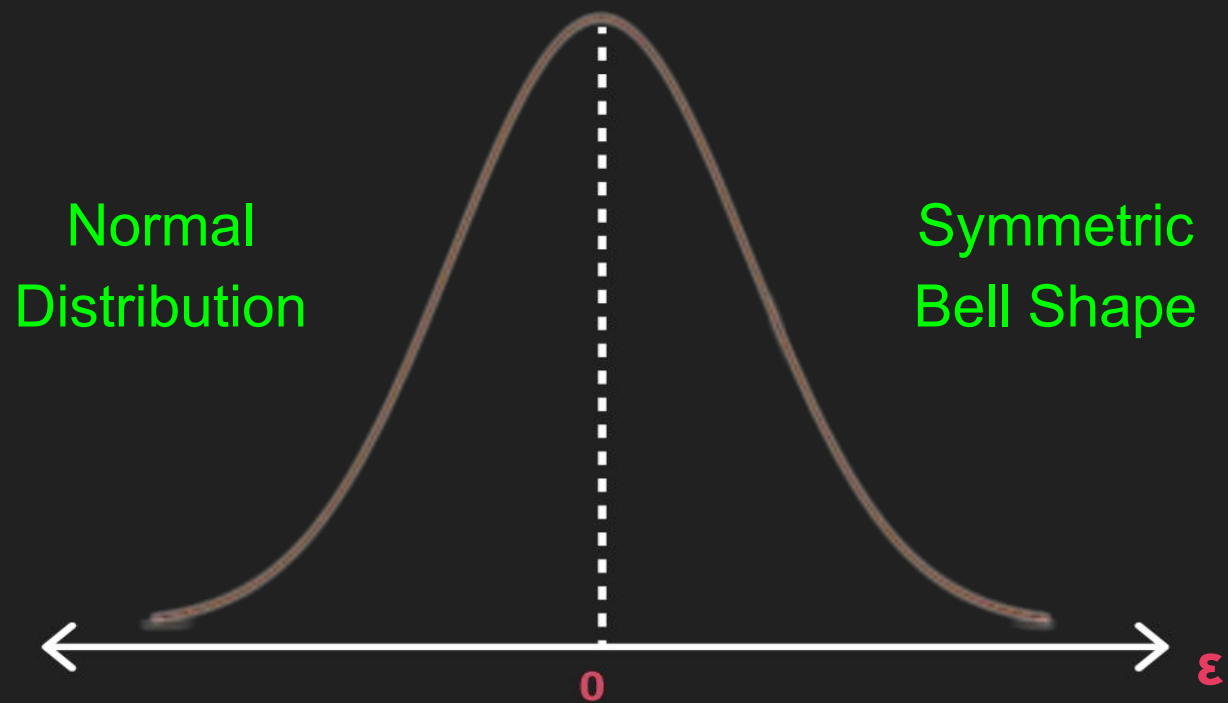
Every linear model has some error.

$$y^i = w_o + w_x^t{}^{(i)} + \varepsilon$$

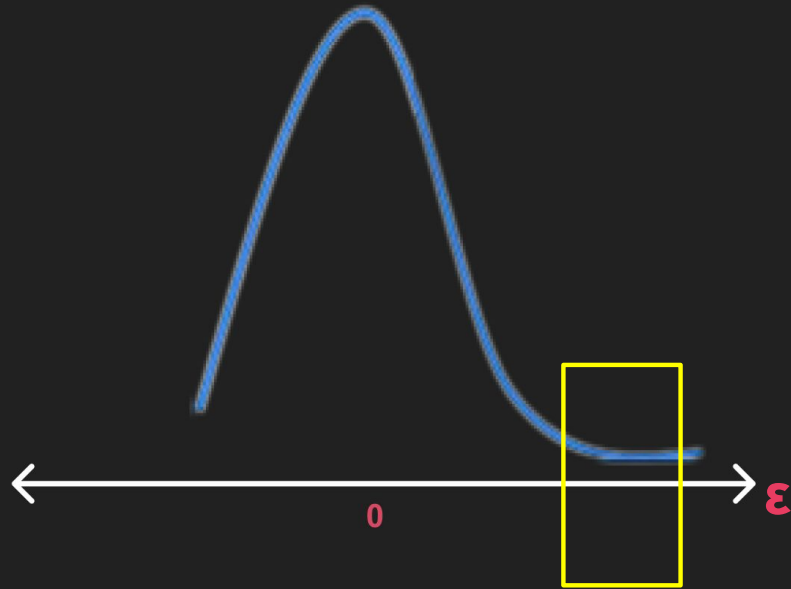
$$\therefore \varepsilon^{(i)} = y^{(i)} - \hat{y}^{(i)}$$



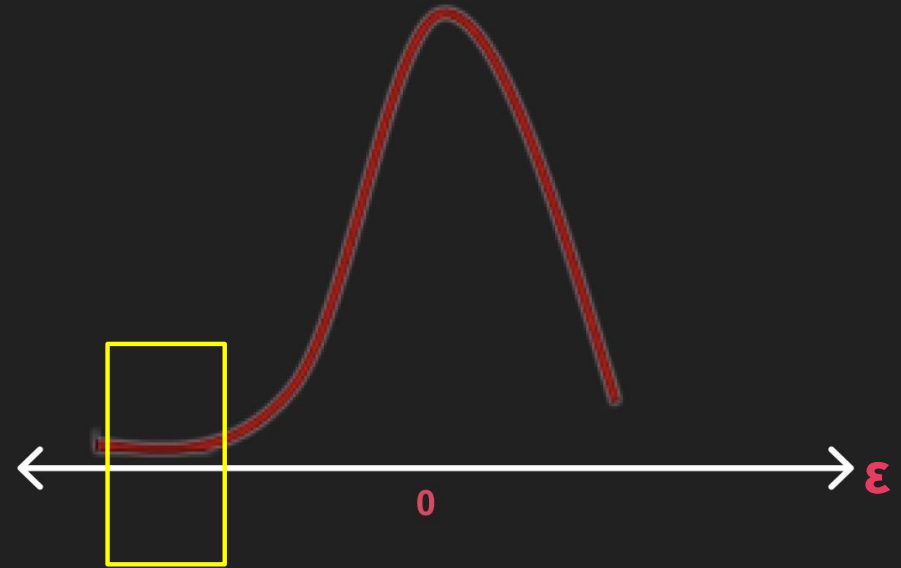
Plotting ' ϵ '



On the other hand, if



Right Skewed



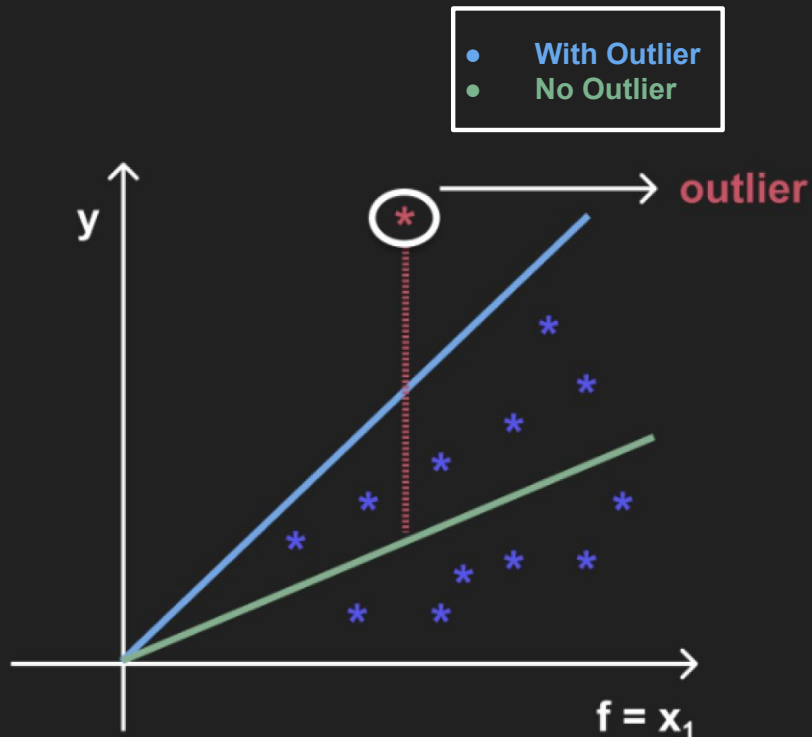
Left Skewed

ϵ is large, outliers are present.

What is the impact of outliers?

If we have outliers,

- The regression line gets pulled towards the outlier to minimize the squared loss.



Q : How to identify outliers ?

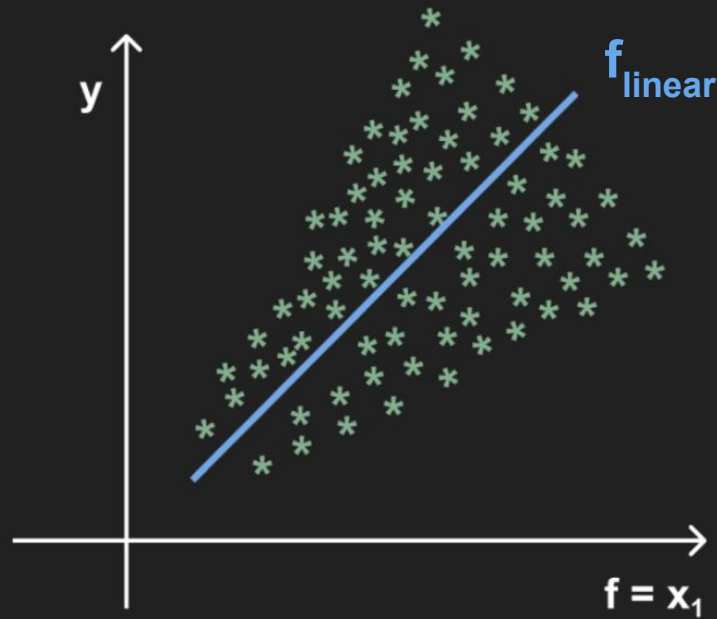
A: Outliers will have high error (ϵ).

Q : How to deal with outliers ?

A: Remove the points with high error as many as you want and fit the model again.

No Heteroskedasticity

When we plot the two features along with the regression line, notice :



As we go from left to right errors are increasing.

How to check Heteroskedasticity ?

By plotting,

$$y^{(i)} \text{ vs. } \varepsilon^{(i)}$$

In maths/stats proof of linear regression, we assume the errors are normally distributed

$$\varepsilon^{(i)} \sim N(0, \sigma)$$

mean = 0

Std. Dev. is
not constant

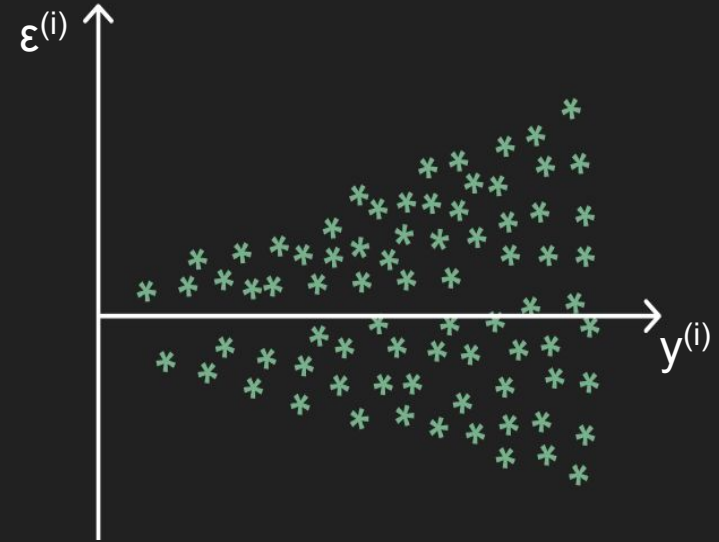
Spread of $\varepsilon^{(i)}$ is not same for all values of $y^{(i)}$, this is known as Heteroskedasticity.



What does it tell us ?

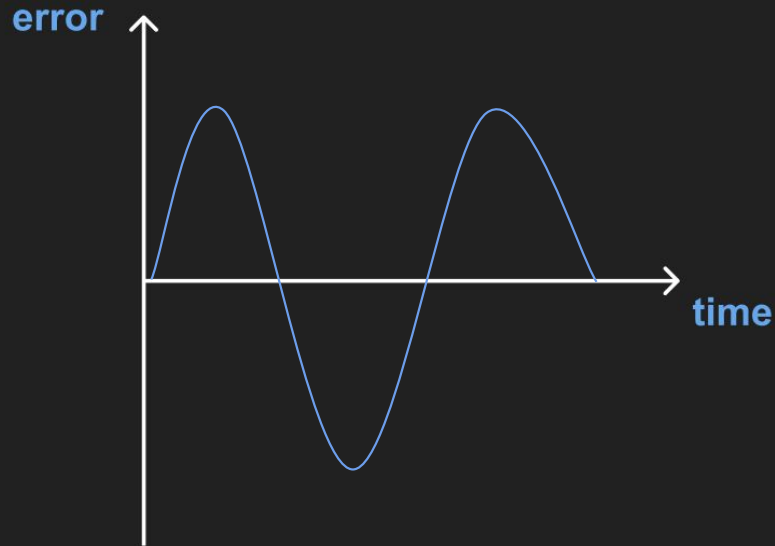
Outliers in data

Linear model isn't right



No Autocorrelation

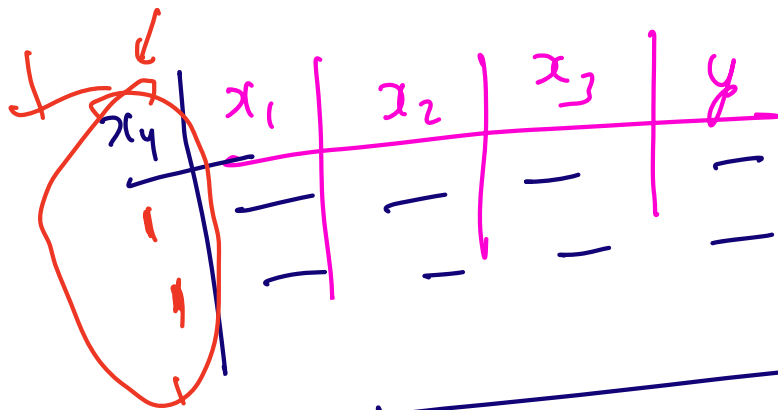
Autocorrelation plays a role only when “**Time Series**” data is involved



Example : Predicting sales on the same day

When we plot the error w.r.t. time, if some pattern is observed then autocorrelation exists.

]



x_0	x_1	x_2	x_3	y
1	—	—	—	—
1	—	—	—	—

! $\rightarrow w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0 \underline{x_0} = y$

$x_0 = 1$

Linear Reg:

Assumption: All columns are irrelevant

$x_1 \quad x_2 \quad x_3 \quad x_4$

reject $\xrightarrow{H_0}$ H_0 : All columns are irrelevant

H_A : All columns are relevant

< 0.05

$< 5\%$

		p-value
Rejecting H_0 →	x_1	0.01
Accept H_0 ←	x_2	0.8
	x_3	0.03

Batch : whole dataset

Mini-Batch : sample

Idc : stochastic GD

→ Regularisation