

## Agenda

→ feature Importance Code

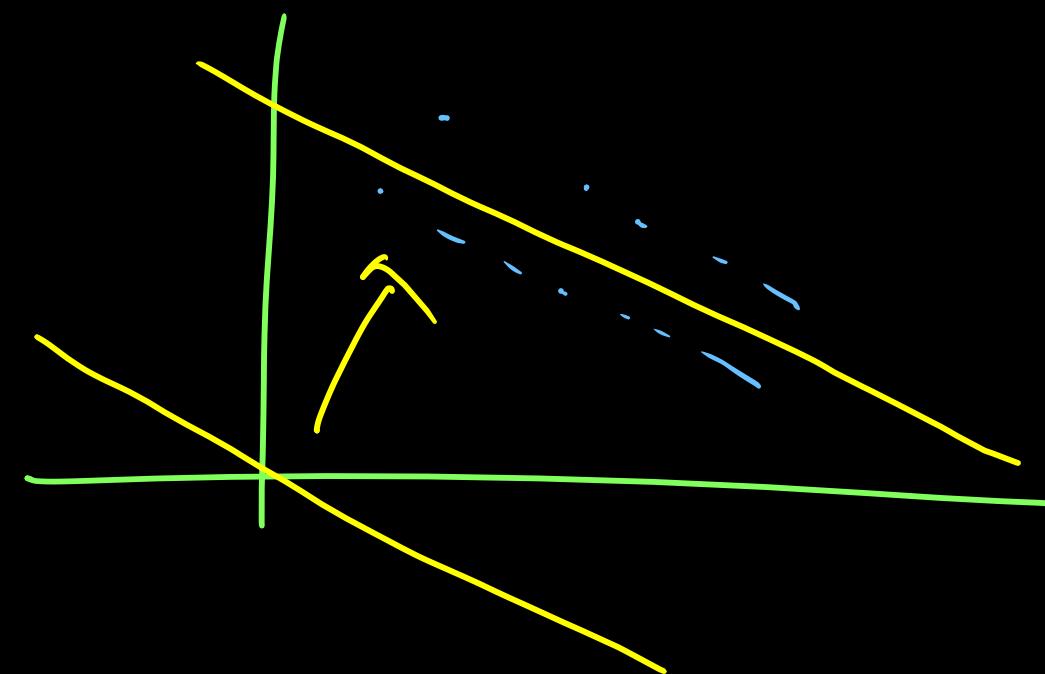
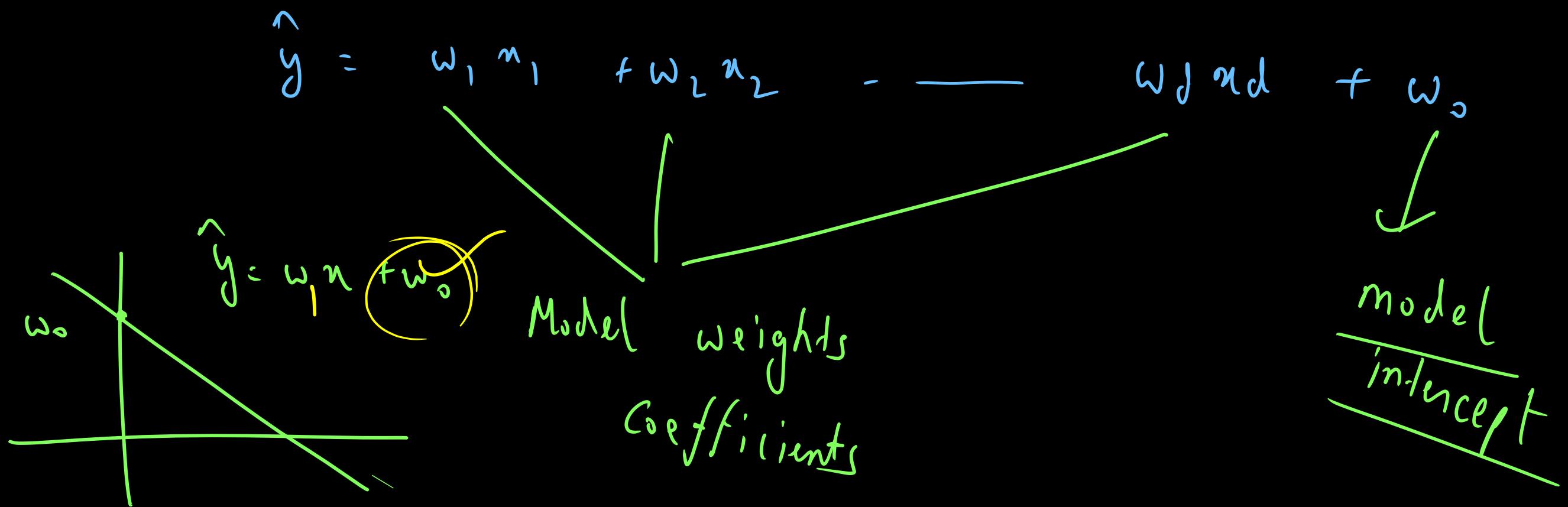
→ Assumption of Linear Regression

→ OLS using stat model

→ Code

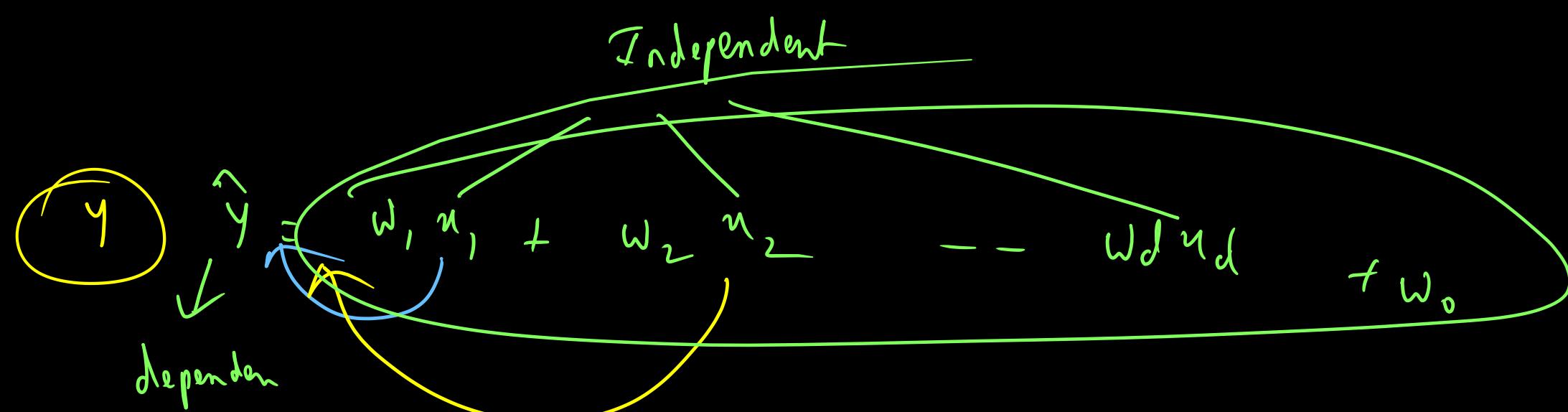
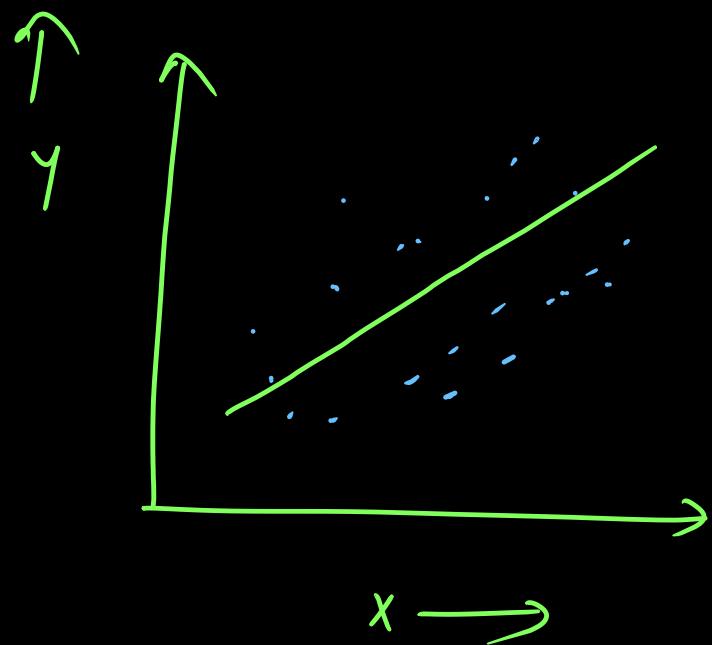
$$\left\{ (y^{(i)} - \hat{y}^{(i)})^2 = S_{SE} \right.$$

$$\frac{1}{m} \left\{ (y^{(i)} - \hat{y}^{(i)})^2 \right\} = M_{SE}$$



→ Assumption for linear Regression

1  
Linearity



① Pearson test / Spearman Correlation

$r = 1 \rightarrow$  Perfect +ve linear

$r = -1 \rightarrow$  Perfect -ve linear relation

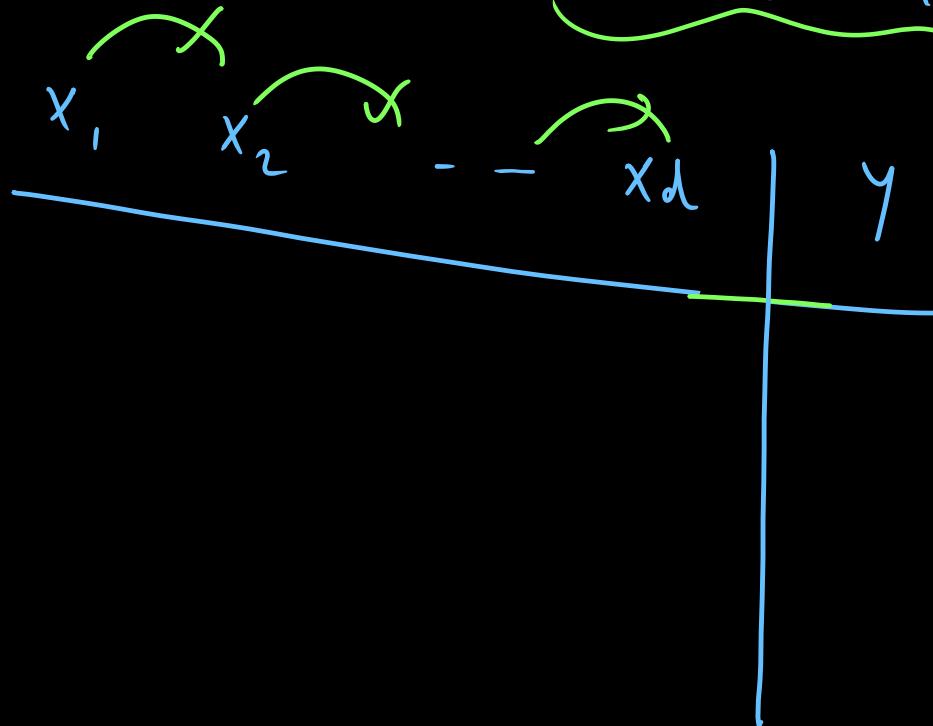
$r = 0 \rightarrow$  No linear relation

② Scatter plot / Line plot

2

## No Multicollinearity

features should be independent



$$x_1 = \alpha x_2 + \beta$$

km driven      Miles      Constant

{  
Gllinear}

$$1 \text{ mile} = 1.6 \text{ km}$$

$$\alpha = \frac{1}{1.6}$$

→ Multicollinearity

$$\left\{ \begin{array}{l} x_1 = \alpha x_2 + \beta x_3 + \gamma x_4 \end{array} \right.$$

$$\hat{y} = x_1 + 2x_2 + 3x_3 + 5$$

$x_3 > x_2 > x_1$

$$= 1 + 2 \cdot 1.5 + 3 \cdot 3 + 5$$

$$= 1 + 3 + 3 + 5$$

$$= 12$$

feature importance

get messed

↑ p : of

Multicollinearity

$$\hat{y} = x_1 + 2 \cdot 1.5 x_1 + 3x_3 + 5$$

$$= 1x_1 + 0x_2 + 3x_3 + 5$$

$x_1 > x_3 > x_2$

$$= 4 + 3 + 5 = 12$$

→ How to remove Multicollinearity ?

Repeat converge

Regression

$x_1, x_2, x_3, x_4$

$y$

↳  $R^2$

Logistic  
Cross

→ VIF for each features  
(Variance Inflation factor)

Linear Regres  
↓  
MSE

↳ Check for feature with highest VIF  $x_1$

Remove the feature with highest VIF

$x_2, x_3, x_4, y$

$VIF$   
  
 $y \approx \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \omega_4 x_4$   
 $x_1 = \alpha x_2 + \beta x_3 + \gamma x_4$

$x_2$	$x_3$	$x_4$	$x_1$
-------	-------	-------	-------

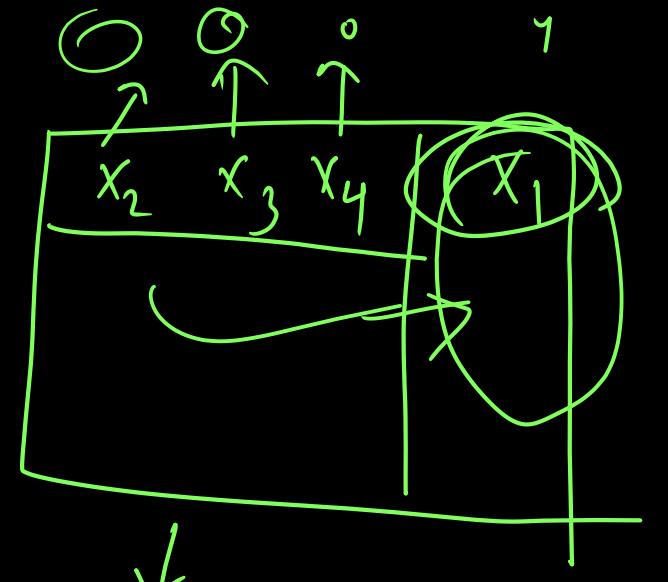
$R^2 \approx 1$   
 $VIF_1 = \frac{1}{1 - R^2}$

$R^2 \approx 0$   
 $VIF = 1$   
 $\downarrow$   
 $VIF = 0$

$VIF > 10 \rightarrow$  <sup>Very</sup> High Multicollinearity

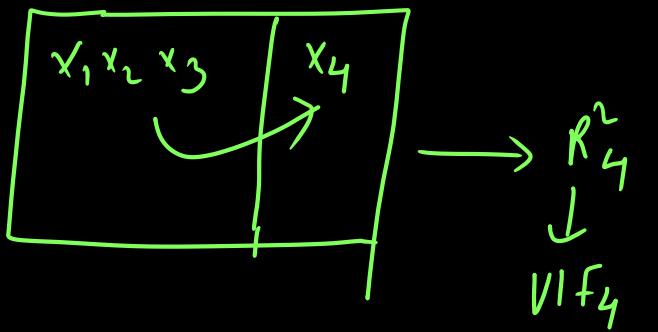
$8 < VIF < 10 : \rightarrow$  Multicoll. exist

$VIF < 5 : \text{low Multicollinearity}$



$$R^2_1$$

$$VIF_1$$

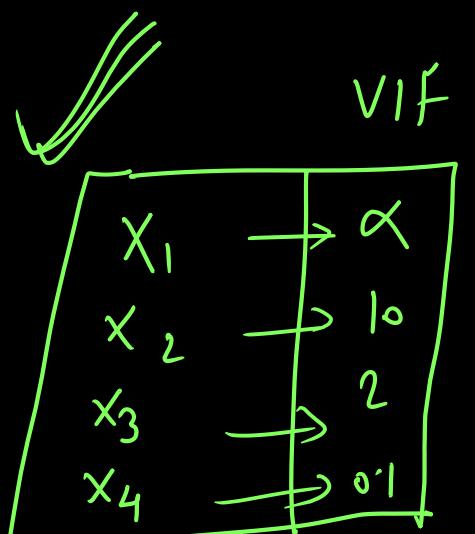


$$R^2_2$$

$$VIF_2$$

$$R^2_3$$

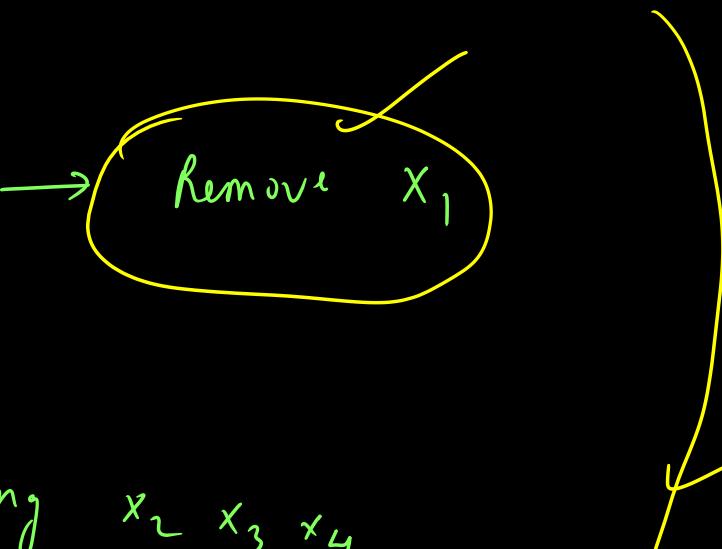
$$VIF_3$$



$\rightarrow x_1, x_2, x_3, x_4, y \rightarrow R^2 = 0.85$

$\rightarrow$  Calculate VIF

VIF
$x_1$
$x_2$
$x_3$
$x_4$



$\rightarrow$  feature remaining  $x_2, x_3, x_4$

$\rightarrow x_2, x_3, x_4, y \rightarrow R^2 = 0.89$

$\rightarrow$  VIF  $< 5$

$\rightarrow$  Calculate VIF

VIF
$x_2$
$x_3$
$x_4$

$\rightarrow$  Remove  $x_2$

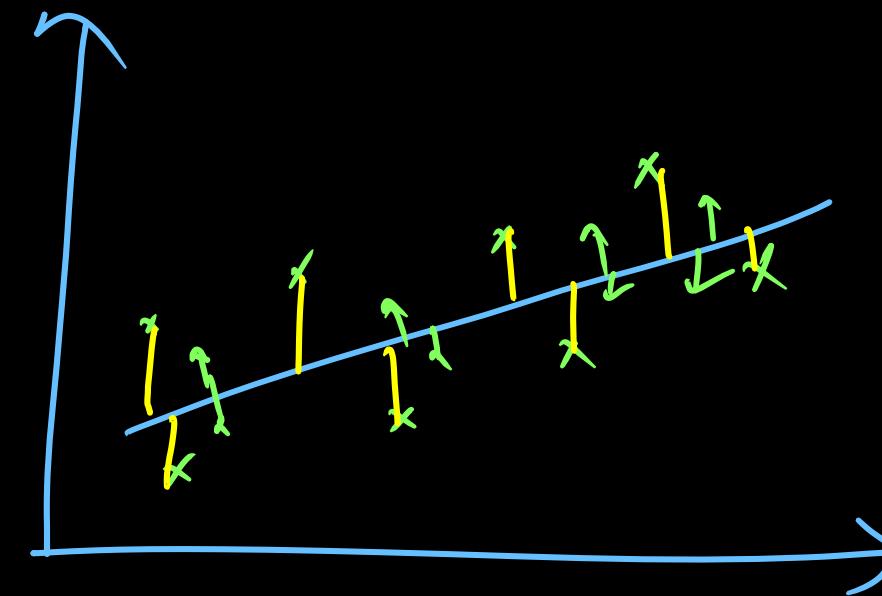
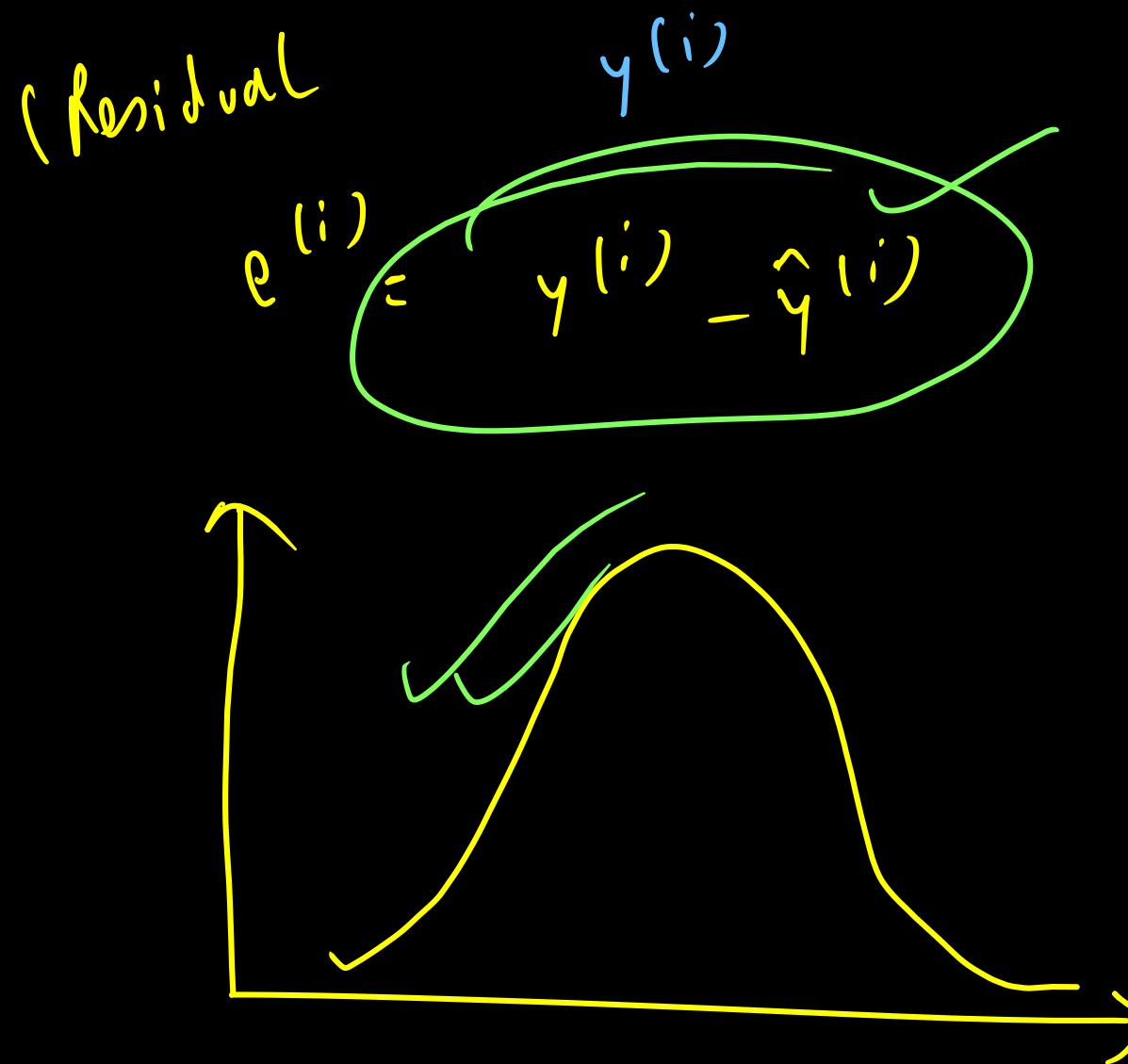
$R^2 \leq 0.75$

$\rightarrow x_3, x_4, y \rightarrow R^2 = 0.75$

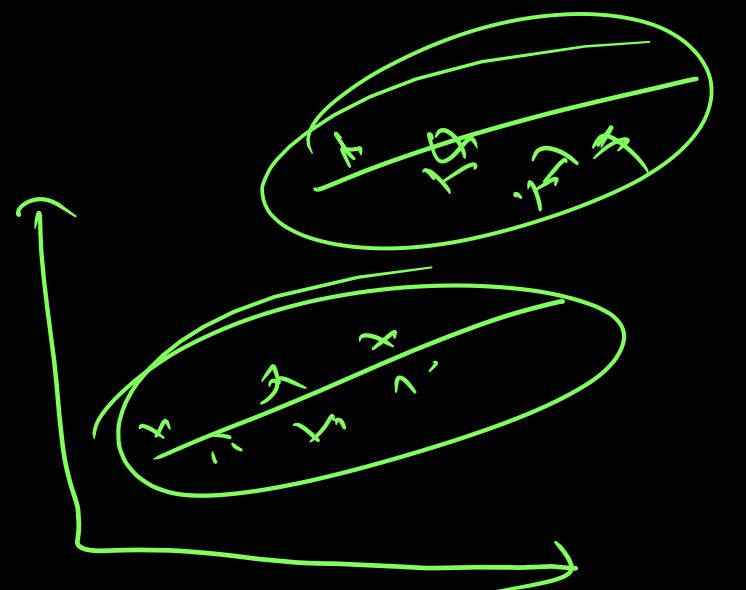
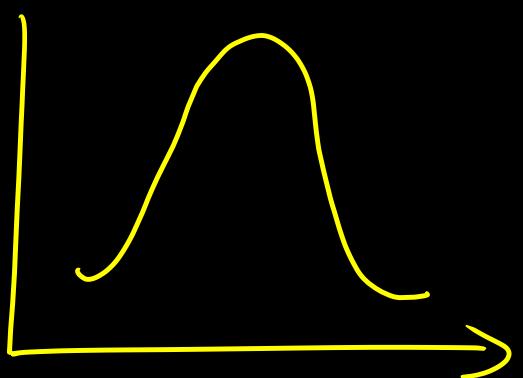
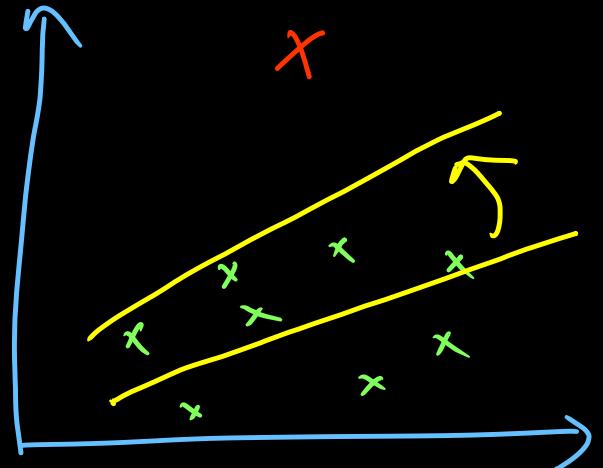
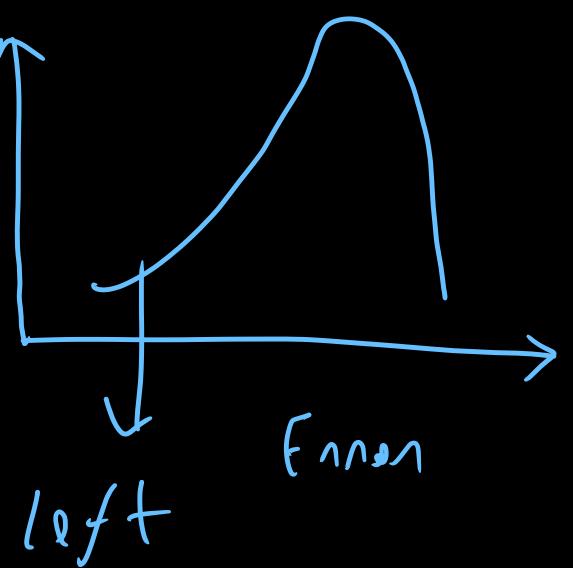
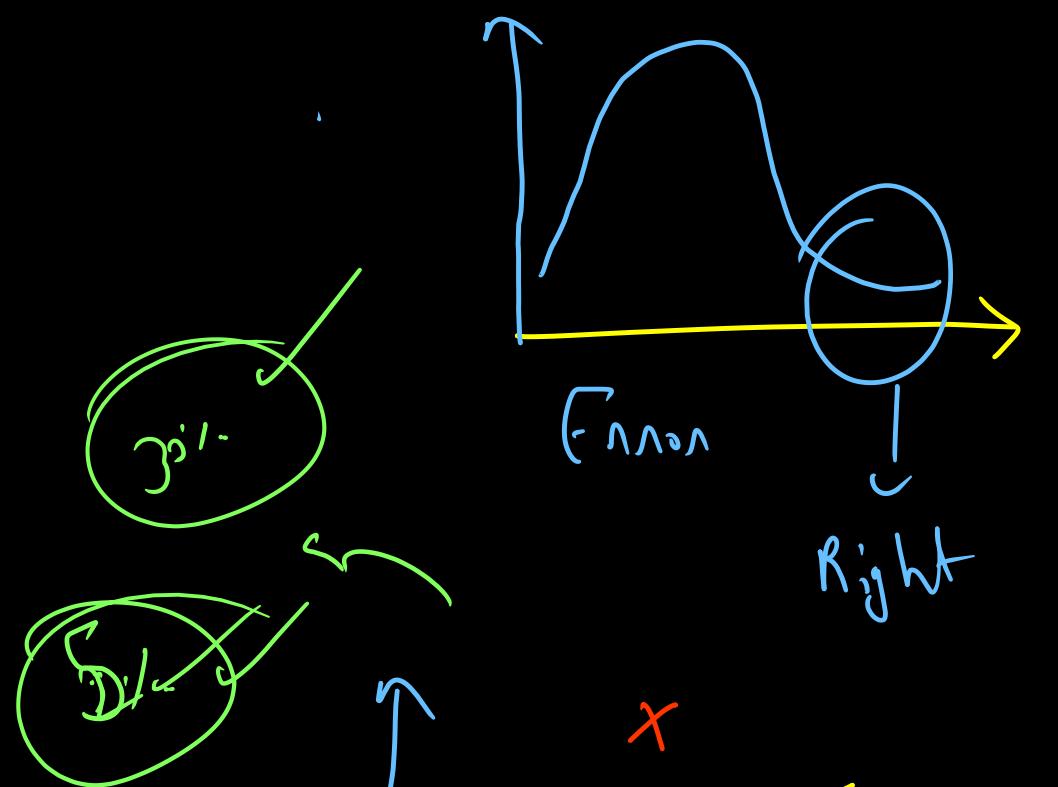
$\rightarrow$  Stop it

$\rightarrow$  non linear and non normally distributed

$$\hat{y}^{(i)} = \omega_1 x_1^{(i)} + \omega_2 x_2^{(i)} - \dots - \omega_d x_d^{(i)}$$



$\rightarrow$  Shapiro, qq plot



Ordinary Least Square  
 $\sum (y^{(i)} - \hat{y}^{(i)})^2$  → Break until 22:25

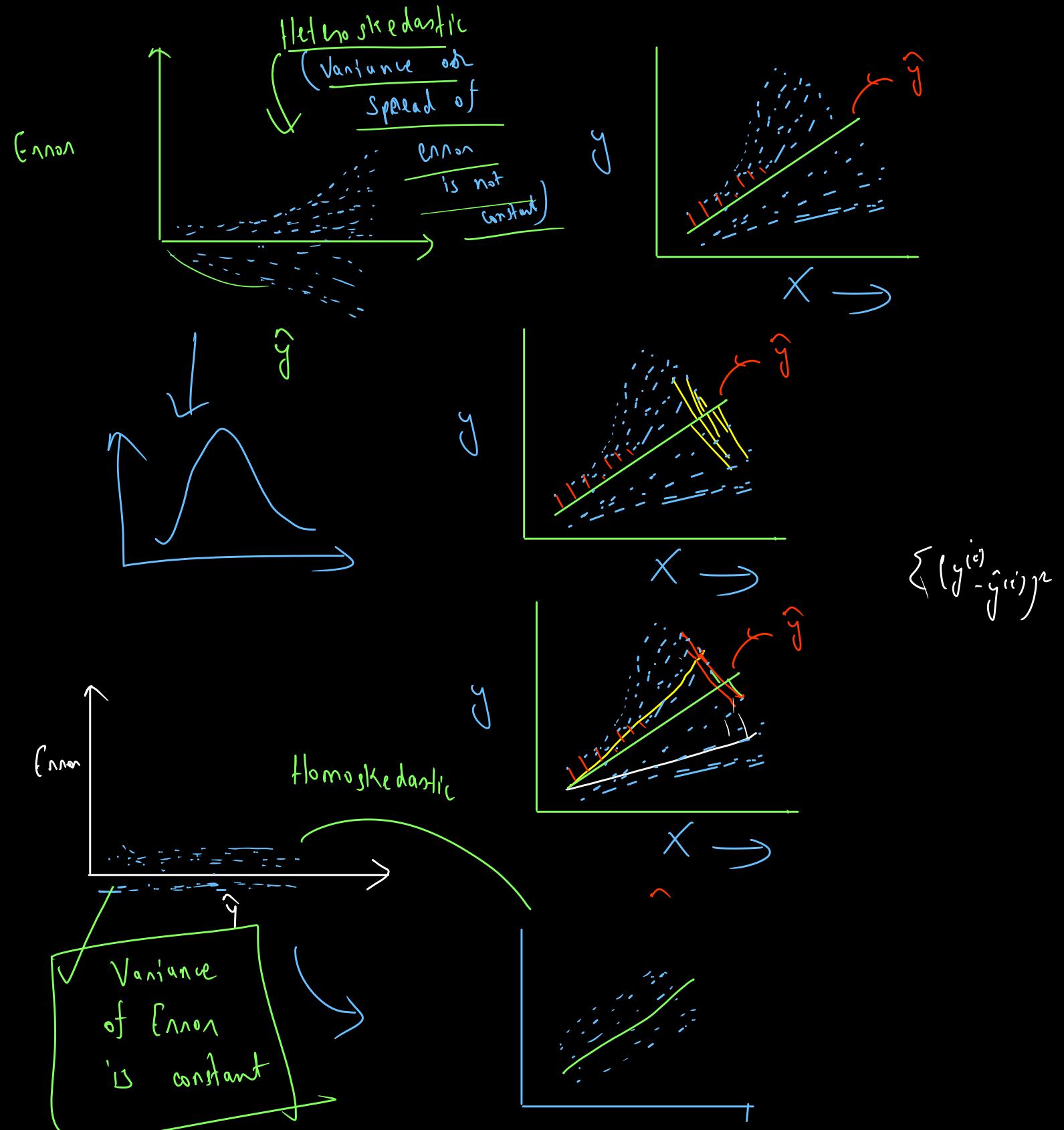
Regression  $n = m - k - 1 = 13874 - 16 - 1$   
 $x_0$   
 $+ w_0 \cdot 1$

SKLearn  $\hat{y} = w_1 x_1 + \dots + w_d x_d + w_0 \cdot 1$   
 $= [w^T x + w_0]$   
 Statsmodel  $y = w^T x + \epsilon$   
 $w^T x$   
 $bias$

$x = [x_1, \dots, x_d]$   
 $w = [w_1, \dots, w_d, w_0]$   
 $w_0 \cdot 1$

$w_1 x_1 + w_2 x_2 + \dots + w_d x_d + w_0 \cdot 1$

4) No Heteroskedasticity



```
# Performing the Goldfeld-Quandt test to check for Homoscedasticity
from statsmodels.compat import lzip
import statsmodels.stats.api as sms

name = ['F statistic', 'p-value']
test = sms.het_goldfeldquandt(y_train, X2_sm)
lzip(name, test)

[('F statistic', np.float64(1.0176613865185007)),
 ('p-value', np.float64(0.21807877686818042))]
```

$H_0$  : HomosKedas

$H_A$  : Not Homos

