

BATCH: 223

DT/NT : DATA SCIENCE
LESSON : STATISTICS -2
SUBJECT: LINEAR REGRESSION
 R^2



TECHPRO
EDUCATION



techproeducation.com



+1 (585) 304 29 59



STATISTICS - 2

Data Science Program

Session -4

RECAP

**Herkes önceki dersten hatırladığı
1 cümle yazabilir mi?**



Session - 4 Content

► Content

- Linear Regression
- Regression Equation
- Coefficient of Determination
- R^2





LMS Pre-Class'ta bu dersle ilgili kısma çalıştım

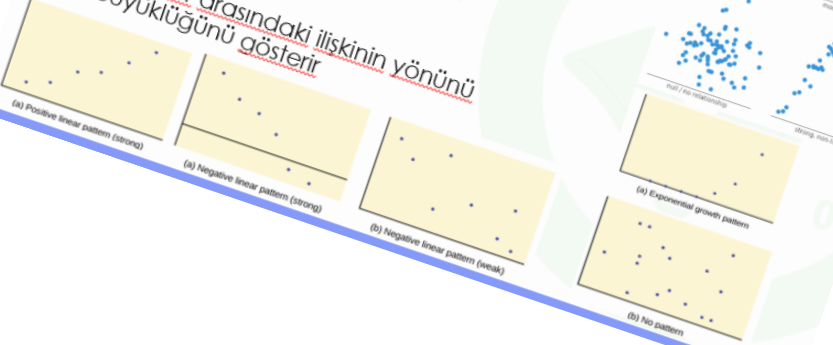
Recap – Previous Lesson



Scatter Plot

Sacılım Grafiği – Serpilme Diyagramı

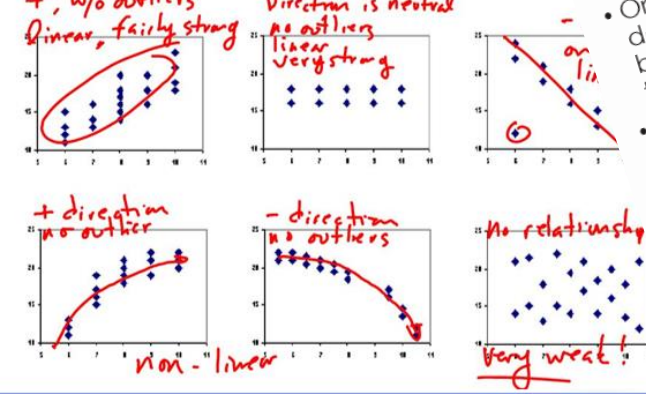
- İki değişkenli bir scatter plot, Y eksenindeki bir değişkenin değerlerini ve X eksenindeki diğer değişkenin değerlerini gösterir.
- Değişkenler arasındaki ilişkinin yönünü ve büyüklüğünü gösterir



Pattern of Data in Scatterplot

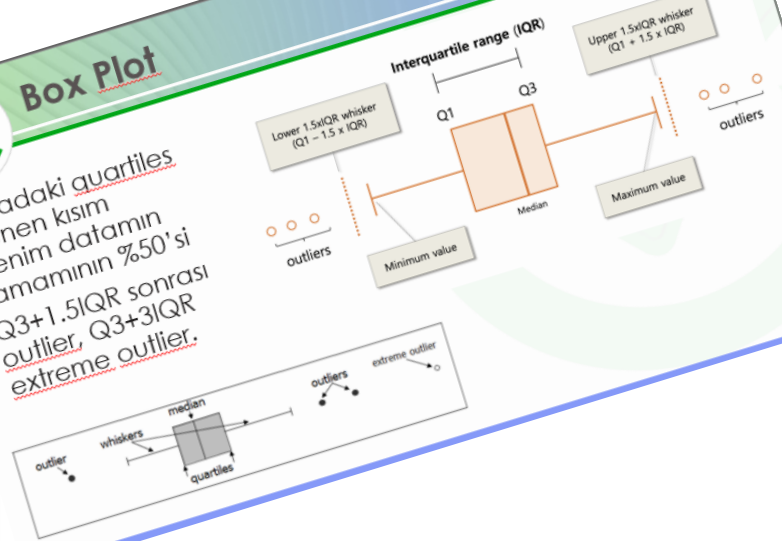
- Metiri

described, in terms of, direction, outliers, linearity, and strength. (DOLS)



Box Plot

- Ortadaki quartilenin kısmı benim datamın tamamının %50'si
- $Q3 + 1.5IQR$ sonrası outlier, $Q3 + 3IQR$ extreme outlier.



Covariance

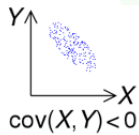
$$\text{Cov}(x, y) > 0$$

- İlişki pozitifdir.
- X artarken Y de artar



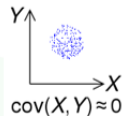
$$\text{Cov}(x, y) < 0$$

- İlişki negatiftir.
- X artarken Y azalır

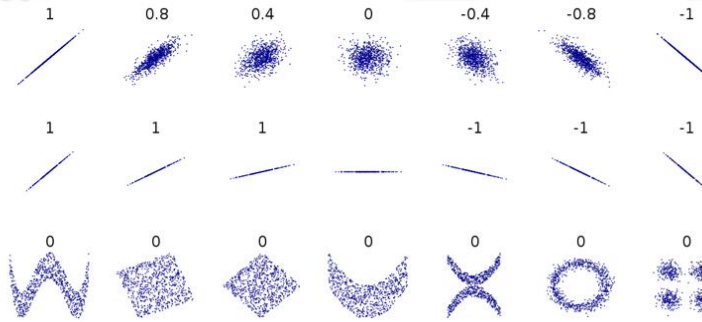


$$\text{Cov}(x, y) = 0$$

- İki değişkenin arasında ilişki yoktur, birbirinden bağımsızdır.



Correlation



Correlation - r Calculation

Cigarette (X)	Lung Capacity (Y)
0	45
5	42
10	33
15	31
20	29

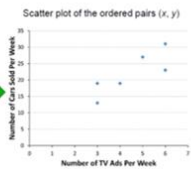
$$r = \frac{n\sum(xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

$$r = \frac{(5)(1585) - (50)(180)}{\sqrt{[(5)(750) - 50^2][(5)(6680) - 180^2]}}$$

$$= \frac{7925 - 9000}{\sqrt{(3750 - 2500)(33400 - 32400)}}$$

$$= \frac{-1075}{\sqrt{(1250)(1000)}} = -0.9615$$

Week	Number of TV Ads	Number of Cars Sold
1	3	13
2	6	31
3	4	19
4	5	27
5	6	23
6	3	19



Week	Number of TV Ads	Number of Cars Sold	xy	x ²	y ²
1	3	13	39	9	169
2	6	31	186	36	961
3	4	19	76	16	361
4	5	27	135	25	729
5	6	23	138	36	529
6	3	19	57	9	361
Σ	27	132	631	131	3110

$$r = \frac{n\sum xy - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

$$= \frac{(6)(631) - (27)(132)}{\sqrt{[(6)(131) - 27^2][(6)(3110) - 132^2]}}$$

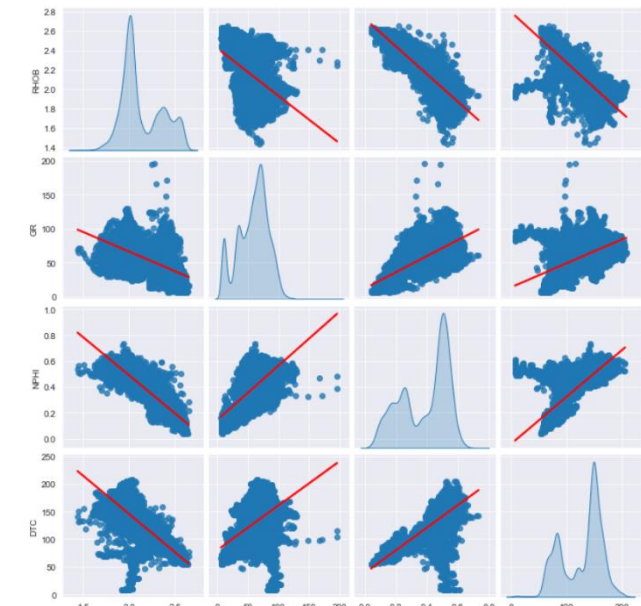
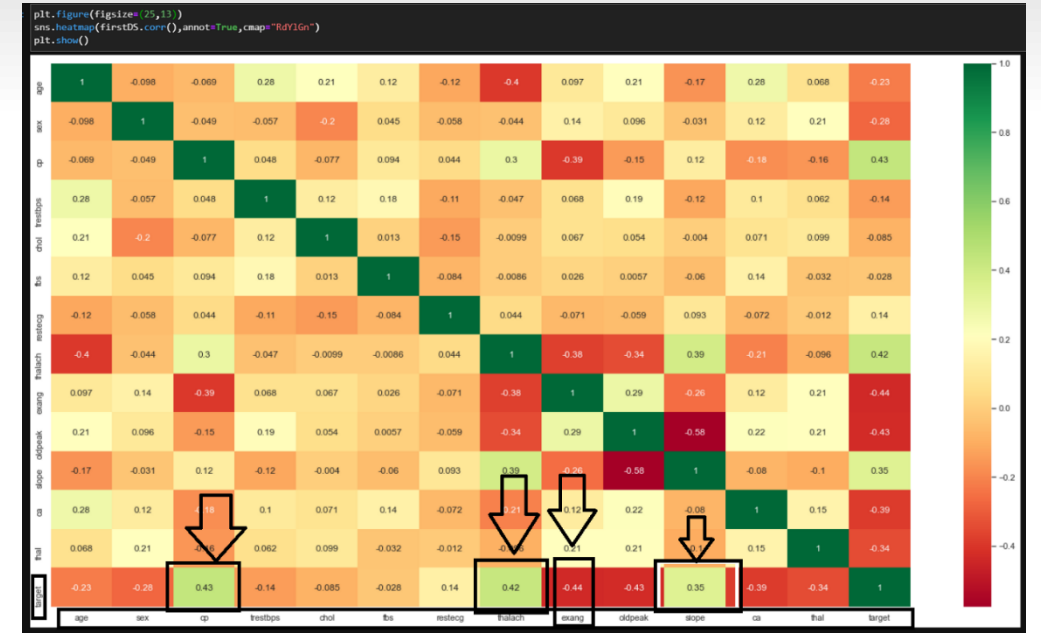
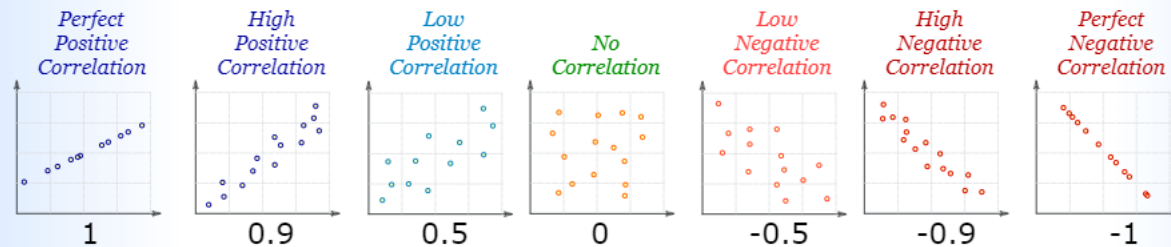
$$= \frac{222}{\sqrt{(37)(1238)}} = \frac{222}{265.43} = 0.836$$

Review

Konular

- Correlation
- Pearson katsayısı
- Sample ve Population corr.
- R hesaplanması

A correlation is assumed to be **linear** (following a line).



Question14: What is the difference between Covariance and Correlation?

Covariance

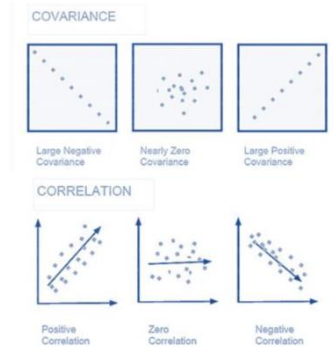
- Signifies the direction of the linear relationship between two variables
- In simple terms, It is a measure of variance between two variables
- It can take any value from positive infinity to negative infinity

Correlation

- It measures the relationship between two variables, as well as the strength between these two variables.
- It can take any value from -1 to 1

Q7. What is correlation and covariance in statistics?

Covariance and Correlation are two mathematical concepts; these two approaches are widely used in statistics. Both Correlation and Covariance establish the relationship and also measure the dependency between two random variables. Though the work is similar between these two in mathematical terms, they are different from each other.



Correlation:

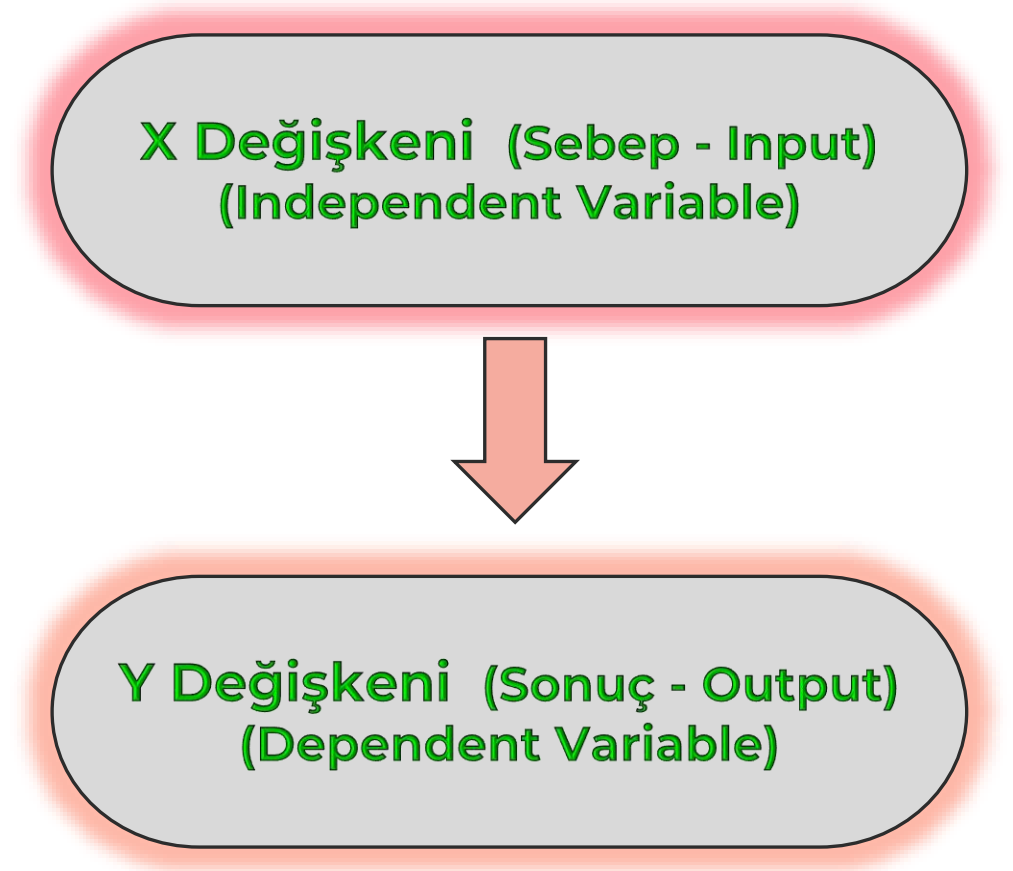
Correlation is considered or described as the best technique for measuring and also for estimating the quantitative relationship between two variables. Correlation measures how strongly two variables are related.

Covariance: In covariance two items vary together and it's a measure that indicates the extent to which two random variables change in cycle. It is a statistical term; it explains the systematic relation between a pair of random variables, wherein changes in one variable reciprocal by a corresponding change in another variable.

Linear Regression

► Lineer Regresyon

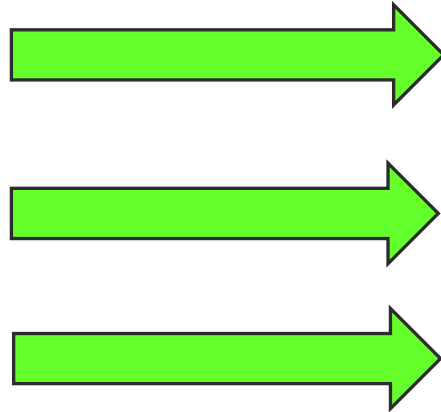
- Amaç: İki değişken arasındaki ilişkiye dayanarak ileri dönük tahmin yapmak
- Sebep-sonuç ilişkisi içinde, Independent variable (bağımsız değişken) sebep, bağımlı değişken ise sonuçtur.



Linear Regression

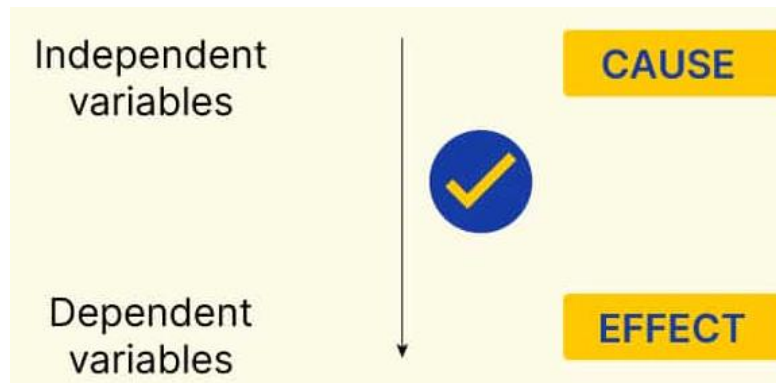
► Independent Variable

- Gelir
- Araç sahipliliği
- IQ değeri
- ???



► Dependent Variable

- Yaşam konforu
- Trafik hacmi
- İş performansı
- ???



RECAP

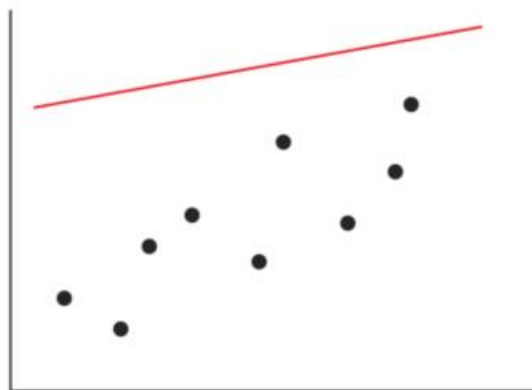
Aklınıza gelen Linear Regresyon örneklerini yazar mısınız?



Q52. What is Linear Regression?

Linear regression is a statistical technique where the score of a variable Y is predicted from the score of a second variable X . X is referred to as the predictor variable and Y as the criterion variable.

Machine Learning Algorithm
Regression - Alternating Least Squares



$D =$

The regression line is the one with the least value of D

Matching on Peardeck



Daily temperature

Annual Salary

Number of exams passed

Life Expectancy

Amount of time spent studying

Electricity Consumption

GDP per Capita

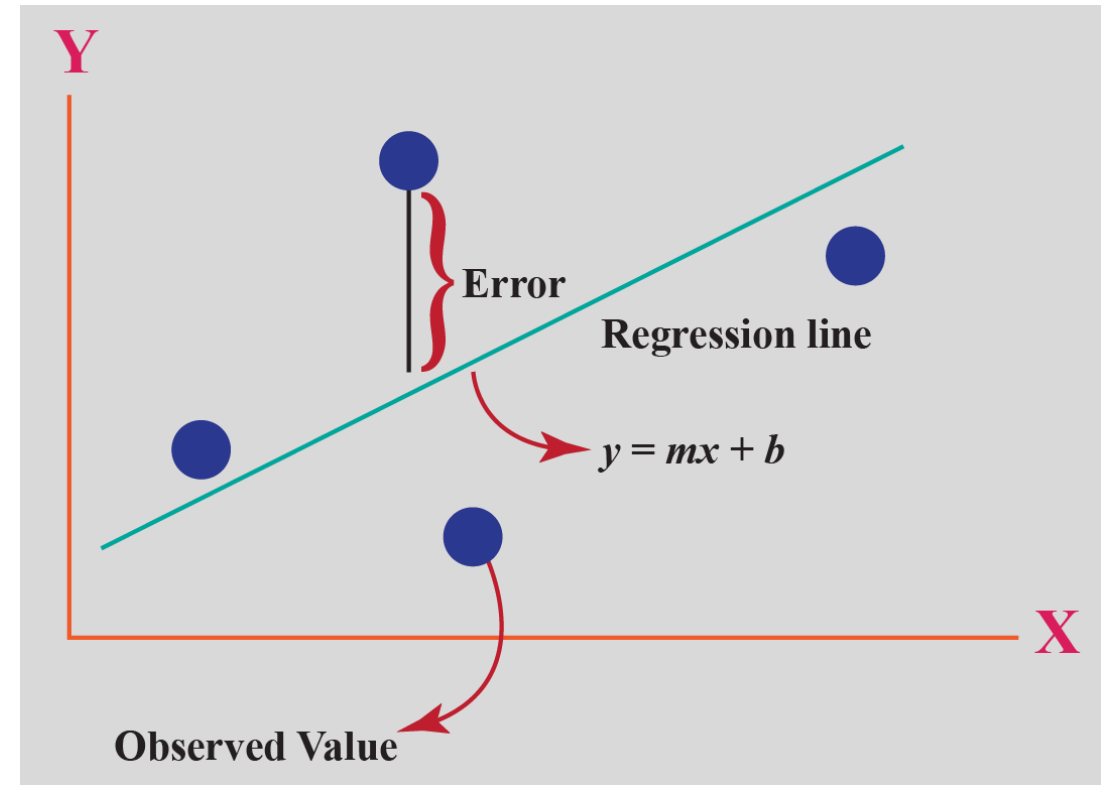
Number of vacations taken

Linear Regression and Equation

► En küçük Kareler Yöntemi

- The least squares (en küçük kareler) yöntemi
- X bağımsız değişkenin değerine bağlı olarak, Y bağımlı değişkenin değerini tahmin etmek için kullanılan bir yöntem

$$Y = aX + b$$



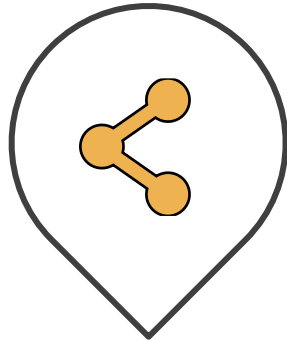
Linear Regression Requirement

Değişken Sayısı



1 Bağımlı değişken
1 Bağımsız değişken

Lineerlik

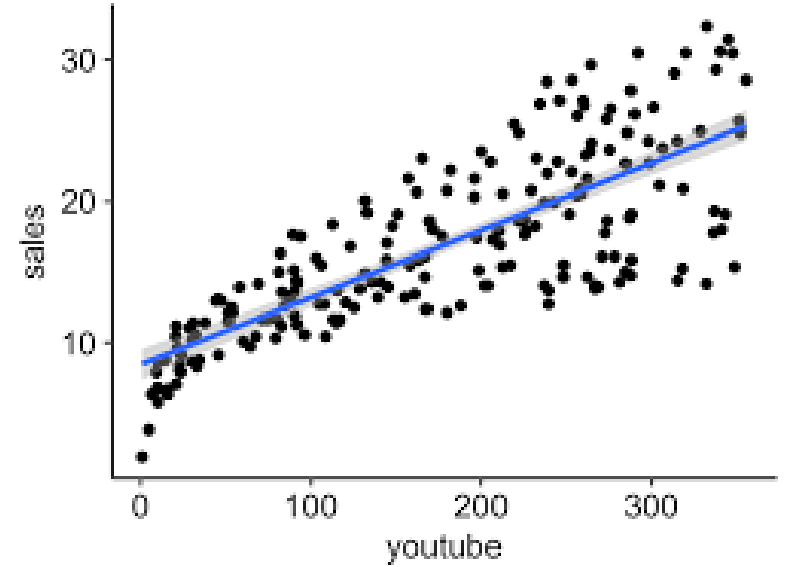


Lineer ilişki olmalı,
nonlinear vb. değil

Ölçülebilirlik

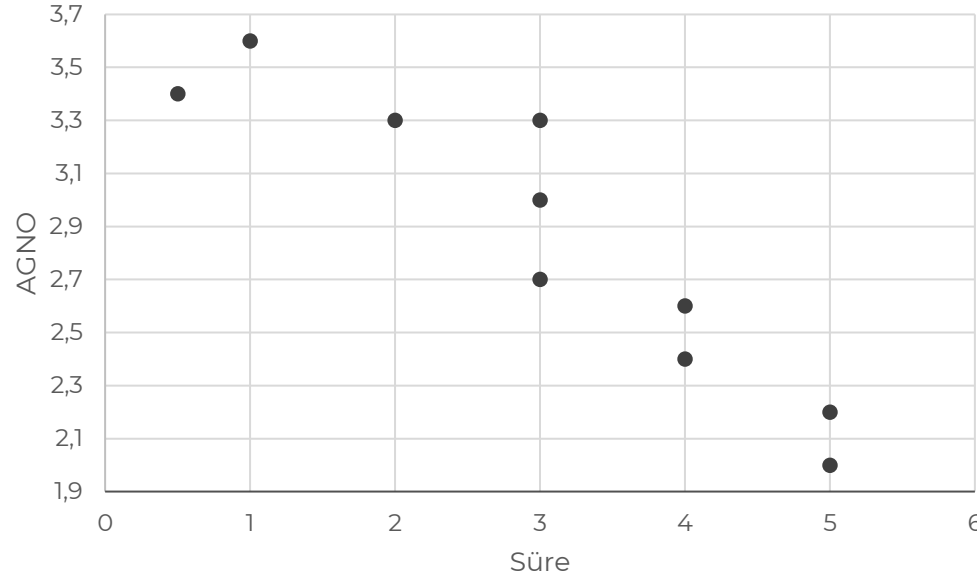


Interval veya
ratio scale



Linear Regression Sample

Herhangi Ekranda geçirilen süre -X	Ağırlıklı Genel Not Ortalaması - Y
3	2,7
5	2,2
2	3,3
0,5	3,4
5	2
3	3
1	3,6
4	2,4
3	3,3
4	2,6
3,05 (ort)	2,85 (ort)



intercept

$$Y = a + bX$$

slope

Linear Regression Sample

$X - X_{ort}$	$(X - X_{ort})^2$	$y - y_{ort}$	$(y - y_{ort})^2$	$(X - X_{ort}) * (y - y_{ort})$	
-0,05	0,0025	-0,15	0,0225	0,0075	
1,95	3,8025	-0,65	0,4225	-1,2675	
-1,05	1,1025	0,45	0,2025	-0,4725	
-2,55	6,5025	0,55	0,3025	-1,4025	
1,95	3,8025	-0,85	0,7225	-1,6575	
-0,05	0,0025	0,15	0,0225	-0,0075	
-2,05	4,2025	0,75	0,5625	-1,5375	
0,95	0,9025	-0,45	0,2025	-0,4275	
-0,05	0,0025	0,45	0,2025	-0,0225	
0,95	0,9025	-0,25	0,0625	-0,2375	
	21,225		2,725	-7,025	Toplam
	SSx		SSy	SP	

SSx: Sum of Square for independent variable

SSy: Sum of Square for dependent variable

SP : Sum of products

$$b = SP / SSx = -7.025/21,225 = -0,3310$$

$$a = y_{ort} - b * x_{ort} = 2,85 - (-0,3310) * 3,05 = 3,85$$

$$Y = a + bX$$

$$Y = 3,85 - 0,331X$$

intercept

slope

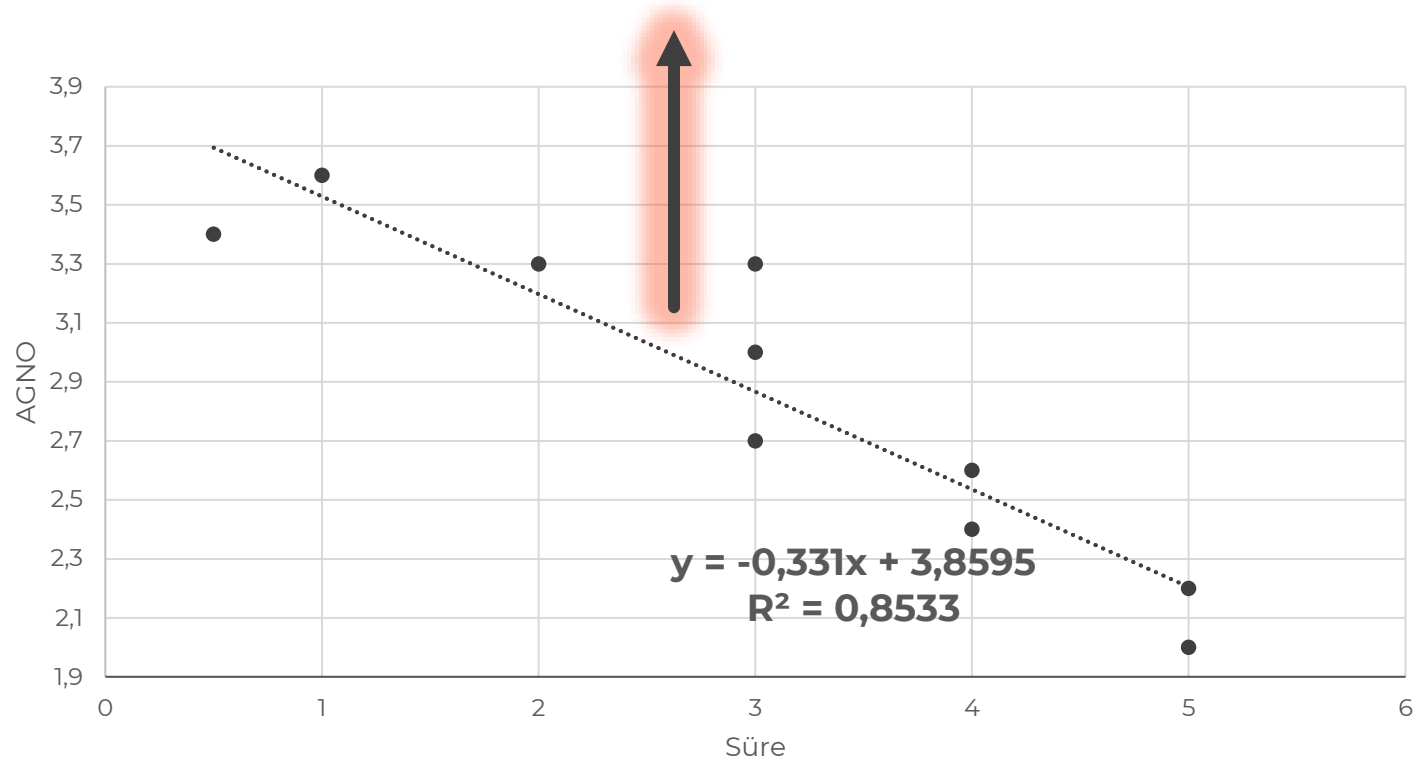
Linear Regression Sample

En iyi eğilim çizgisi (Line of Best Fit)

$$Y = 3,85 - 0,331X$$

intercept

slope



Regresyon çizgisi, gerçek değerler ve tahmin edilen değerler arasındaki 'sum of square farklarını' minimize eder



Linear Regression Sample - Python

Linear Regression Sample



```
In [2]: import numpy as np  
from scipy import stats
```

```
In [3]: Ekran_sure = np.array([3,5,2,0.5,5,3,1,4,3,4])
```

```
In [4]: AGNO = np.array([2.7,2.2,3.3,3.4,2,3,3.6,2.4,3.3,2.6])
```

```
In [5]: reg = stats.linregress(Ekran_sure, AGNO)
```

```
In [7]: print("a: ", reg.intercept)  
print("b: ", reg.slope)
```

```
a: 3.859481743227327  
b: -0.330977620730271
```

$$Y = a + bX$$

$$Y = 3,85 - 0,331X$$

intercept

slope



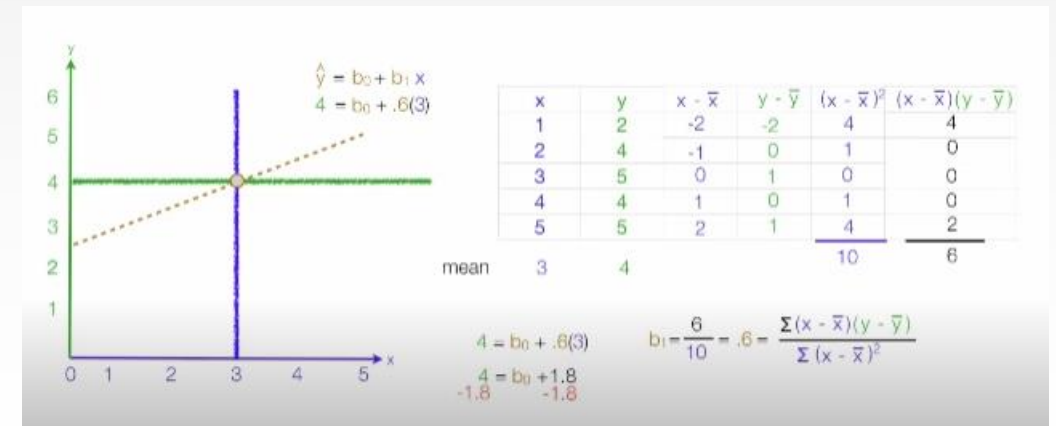
Google
Sheets



YOUTUBE VIDEO ONERI

<https://www.youtube.com/watch?v=JvS2triCgOY>

- How to calculate linear regression using least square method



Pearson's r Calculation

$X - X_{ort}$	$(X - X_{ort})^2$	$y - y_{ort}$	$(y - y_{ort})^2$	$(X - X_{ort}) * (y - y_{ort})$	
-0,05	0,0025	-0,15	0,0225	0,0075	
1,95	3,8025	-0,65	0,4225	-1,2675	
-1,05	1,1025	0,45	0,2025	-0,4725	
-2,55	6,5025	0,55	0,3025	-1,4025	
1,95	3,8025	-0,85	0,7225	-1,6575	
-0,05	0,0025	0,15	0,0225	-0,0075	
-2,05	4,2025	0,75	0,5625	-1,5375	
0,95	0,9025	-0,45	0,2025	-0,4275	
-0,05	0,0025	0,45	0,2025	-0,0225	
0,95	0,9025	-0,25	0,0625	-0,2375	
	21,225		2,725	-7,025	Toplam
	SSx		SSy	SP	

Formula of Pearson's Correlation Coefficient

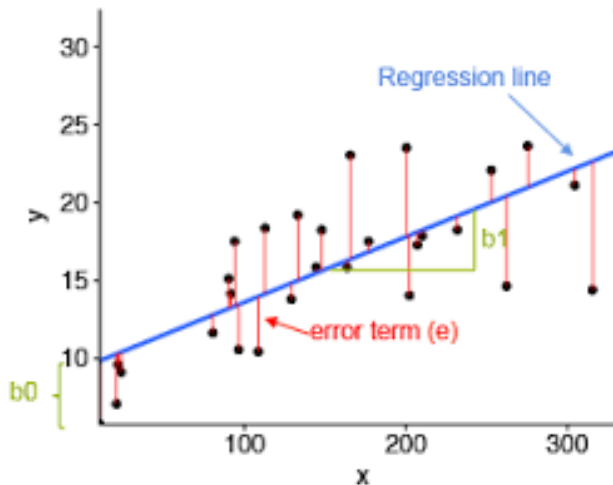
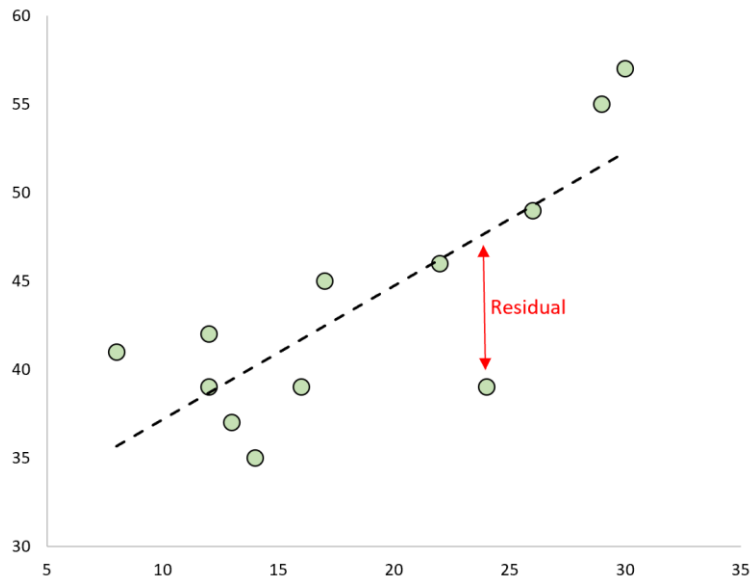
$$r = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$$r = \frac{SP}{\sqrt{SS_x SS_y}}$$

$$r = -0,92$$

Residual term (e)

Residual = Observed value – Predicted value



Regression Model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Dependent Variable \rightarrow Y_i

Population Y intercept \rightarrow β_0

Population Slope Coefficient \rightarrow β_1

Independent Variable \rightarrow X_i

Random Error term \rightarrow ϵ_i

Linear component: $\beta_0 + \beta_1 X_i$

Random Error component: ϵ_i

Coefficient of Determination – R^2

Determinasyon – Belirlilik Katsayısı

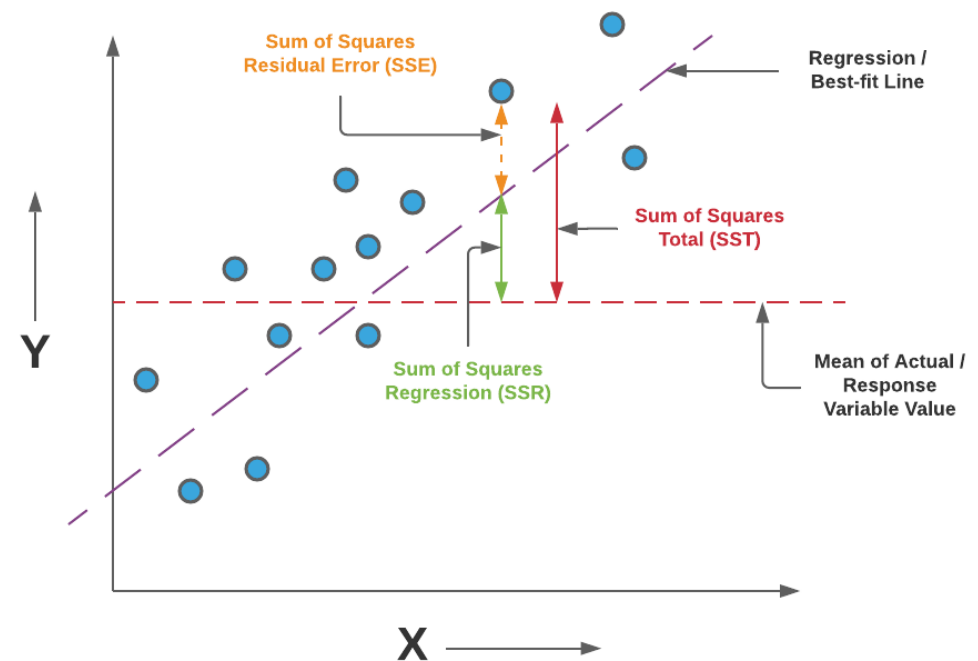
Coefficient of Determination – R²

➤ R² R-square

- Analizimizde iki değişken arasındaki ilişki hakkında fikir sunar
- R² değeri bize bağımlı değişkendeki toplam varyansın yüzde kaçının bağımsız değişken tarafından açıklandığını söyler.
- R² 0 -1 arasında değişir

$$r = -0,92$$

$$R^2 = 0,85$$

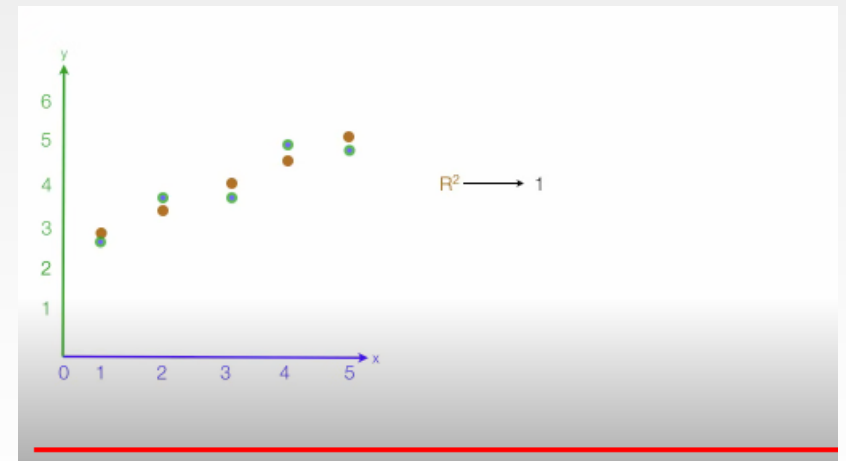


$$R^2 = \frac{SSR}{SST} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

YOUTUBE VIDEO ONERI

<https://www.youtube.com/watch?v=w2FKXOa0HGA>

- How to Calculate R Squared Using Regression Analysis



Kahoot Uygulaması

Python Calculation



- It is time to code by Python...

