

# İSTATİSTİK-II

## Korelasyon ve Regresyon



1

## Korelasyon ve Regresyon

- Genel Bakış
- Korelasyon
- Regresyon
- Belirleme katsayısı
- Varyans analizi
- Kestirimler için aralık tahminlemesi

2

## **Genel Bakış**

### **İkili veriler**

- ❖ aralarında bir ilişki var mıdır?
- ❖ varsa bu ilişki bir eşitlik ile temsil edilebilir mi?
- ❖ bu eşitliğin kestirimler (öngörüler) için kullanılması

3

## **Korelasyon**

4

# Tanım

## ❖ Korelasyon

**bir değişkenin değeri değişirken diğer bir değişken bununla doğrusal ilişkili olarak değişiyorsa korelasyon vardır denebilir.**

5

# Varsayımlar

- 1.  $(x,y)$  ikili verilerden oluşan örnek bir şans örneğidir.**
- 2.  $x$  ve  $y$ 'lerin dağılışı normaldir.**

6

# Tanım

## ❖Saçılma diyagramı

yatay eksen  $x$ , dikey eksen  $y$  olmak üzere,  $(x,y)$  ikili örnek verilerinin işaretlendiği bir grafiktir. Her bir  $(x,y)$  ikilisi tek bir noktadır.

7

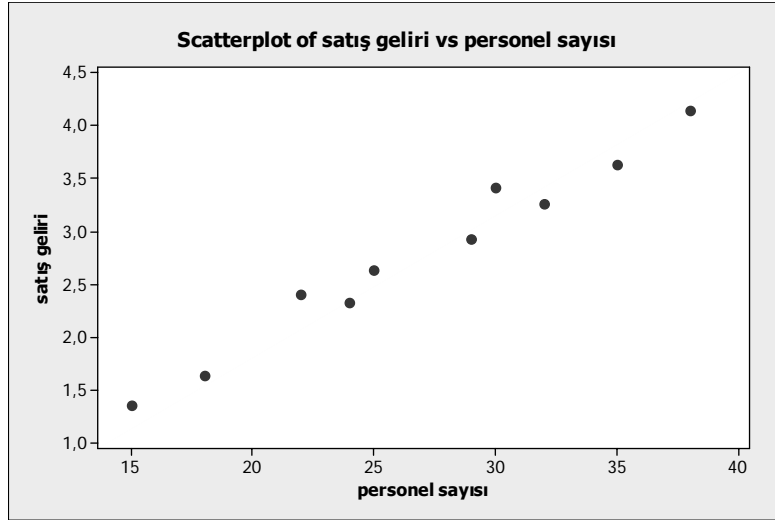
## Örnek

Bir firma bünyesindeki satış personeli sayısı ile satış gelirleri arasındaki ilişkiyi bilmek istemektedir.

<i>Yıllar</i>	<i>Satış Personeli Sayısı (<math>x</math>)</i>	<i>Satış Gelirleri (yüz bin \$) (<math>y</math>)</i>
1999	15	1,35
2000	18	1,63
2001	24	2,33
2002	22	2,41
2003	25	2,63
2004	29	2,93
2005	30	3,41
2006	32	3,26
2007	35	3,63
2008	38	4,15

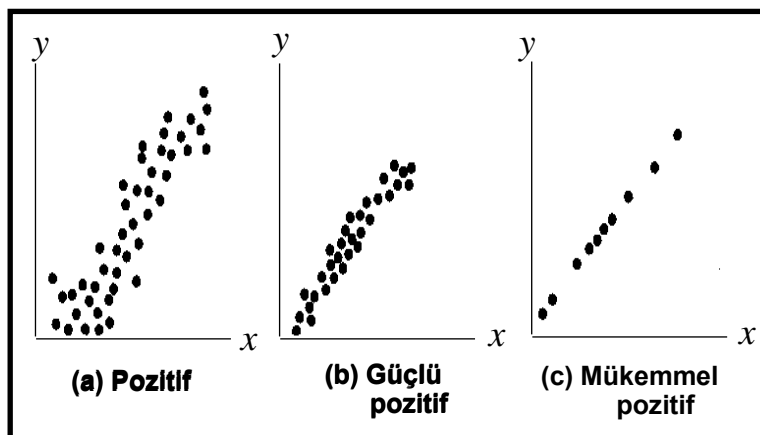
8

## İkili Verilerin Saçılma Diyagramı



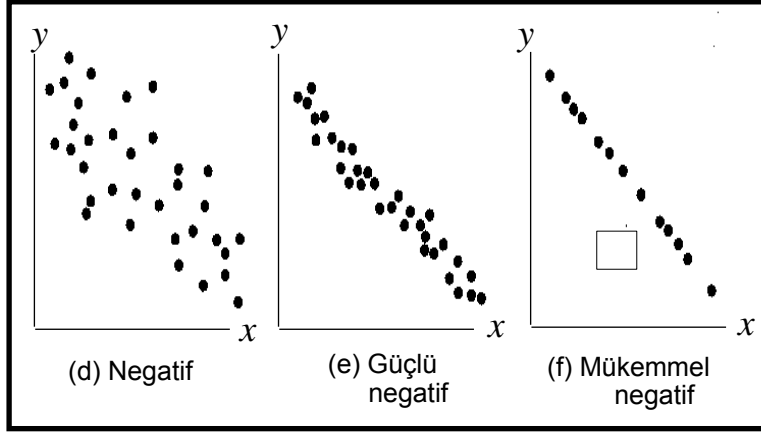
9

## Pozitif Korelasyon

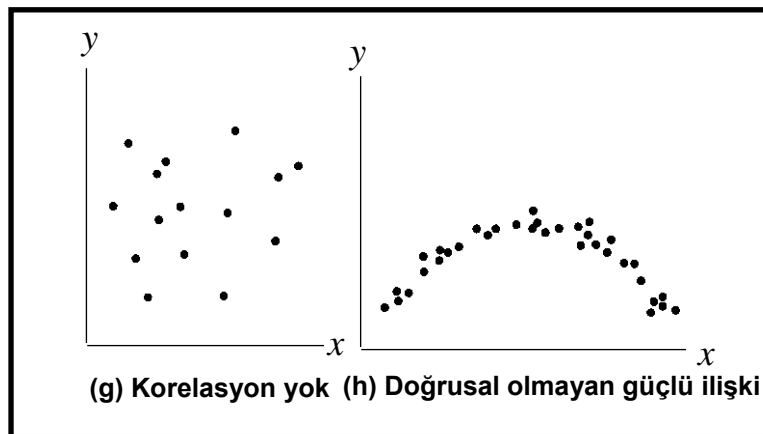


10

## Negatif Korelasyon



11



12

## Tanım

### Korelasyon Katsayısı $r$

Bir örnekteki  $x$  ve  $y$  ikili değerleri arasındaki doğrusal ilişkinin gücünü ölçmektedir.

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

13

## Korelasyon Katsayısı $r$ 'nin Özellikleri

1.  $-1 \leq r \leq 1$
2. Mükemmel pozitif doğrusal ilişki olduğunda  $r = 1$  olur.
3. Mükemmel negatif doğrusal ilişki olduğunda  $r = -1$  olur.
4. Doğrusal ilişki yok ise  $r = 0$  olur.

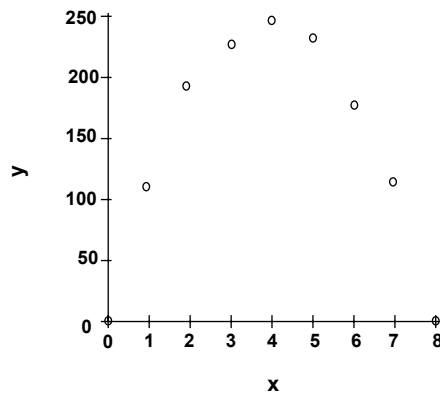
14

## Korelasyon ile ilgili hatalar

1. **Nedensellik:** Korelasyon değişkenler arasındaki sebep sonuç ilişkilerini açıklamaz.
2. **Doğrusallık:** x ile y arasında anlamlı bir korelasyon olmadığı halde, aralarında farklı şekilde bir ilişki olabilir. (Bakınız izleyen slayt)

15

## Korelasyon ile ilgili hatalar



16



## Örnek Verileri İçin Korelasyon Hesaplamaları

Yıllar	Satış Personeli Sayısı (x)	Satış Gelirleri (yüz bin \$) (y)	x <sup>2</sup>	y <sup>2</sup>	xy
1999	15	1,35	225	1,8225	20,25
2000	18	1,63	324	2,6569	29,34
2001	24	2,33	576	5,4289	55,92
2002	22	2,41	484	5,8081	53,02
2003	25	2,63	625	6,9169	65,75
2004	29	2,93	841	8,5849	84,97
2005	30	3,41	900	11,6281	102,3
2006	32	3,26	1024	10,6276	104,32
2007	35	3,63	1225	13,1769	127,05
2008	38	4,15	1444	17,2225	157,7
<b>Toplamlar</b>	<b>268</b>	<b>27,73</b>	<b>7668</b>	<b>83,8733</b>	<b>800,62</b>

17

## Örnek Verileri İçin Korelasyon Hesaplamaları

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

$$r = \frac{(10)(800,62) - (268)(27,73)}{\sqrt{(10)(7668) - (268)^2} \sqrt{(10)(83,8733) - (27,73)^2}}$$

$$r = 0,987 \quad \text{Güçlü pozitif korelasyon}$$

18

## Anakütle Korelasyon Katsayısının Testi

- ❖  $\rho$  = Anakütle korelasyon katsayısı
- ❖  $H_0: \rho = 0$   
(anlamli bir korelasyon yoktur)
- $H_1: \rho \neq 0$   
(anlamli bir korelasyon vardir)

19

## Test İstatistiği $t$

Test istatistiği:

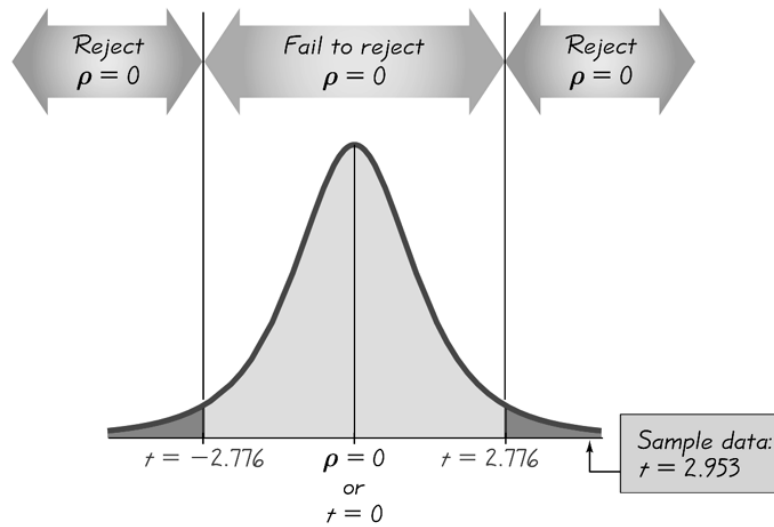
$$t = \frac{r}{\sqrt{\frac{1 - r^2}{n - 2}}}$$

Kritik değerler

serbestlik derecesi =  $n - 2$  olan tablo değerleri dikkate alınarak karar verilir.

20

## Ret Bölgeleri



21

## Anakütle Korelasyon Katsayısının Testi

- ❖  $\rho$  = Anakütle korelasyon katsayısı
- ❖  $H_0: \rho = 0$   
(satış personeli sayısı ile satış gelirleri arasında anlamlı bir korelasyon yoktur)
- $H_1: \rho \neq 0$   
(satış personeli sayısı ile satış gelirleri arasında anlamlı bir korelasyon vardır)

22

## Test İstatistiği $t$

Test istatistiği:

$$t = \frac{r}{\sqrt{\frac{1 - r^2}{n - 2}}} = \frac{0,987}{\sqrt{\frac{1 - 0,987^2}{10 - 2}}} = 17,39$$

### Kritik değer

serbestlik derecesi =  $n - 2 = 10 - 2 = 8$ ,  $\alpha = 0,05$  için  $t_{0,025, 8} = 2,31 < 17,39$

Karar:  $H_0$  ret. Korelasyon anlamlıdır.

23

## Regresyon

$x$  bağımsız değişken (açıklayıcı değişken)

$y$  bağımlı değişken (cevap = yanıt değişkeni)

$y = b_0 + b_1x + e$  Basit doğrusal regresyon modeli

$b_1$  = eğim     $b_0$  = kesen

24

# Regresyon

## ❖ Regresyon Eşitliği

Verilen bir ikili veriler topluluğu için regresyon eşitliği,

$$\hat{y} = b_0 + b_1x$$

iki değişken arasındaki ilişkiyi tanımlamaktadır.

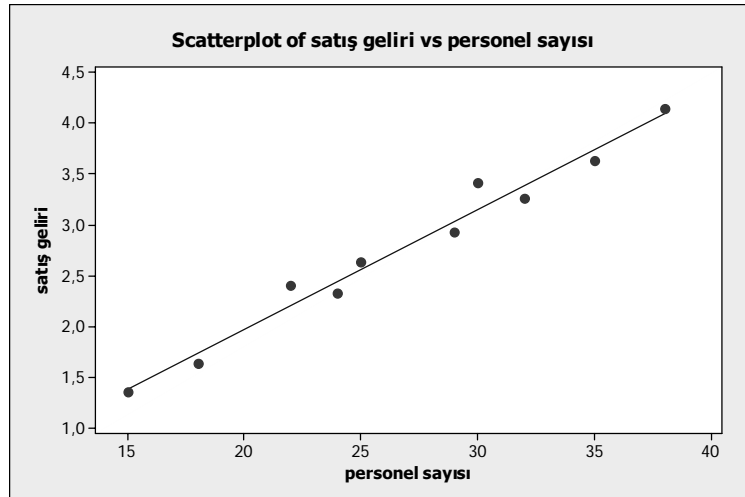
$$b_1 = \text{eğim} \quad b_0 = \text{kesen}$$

## ❖ Regresyon Doğrusu

Regresyon eşitliğinin grafiğidir.

25

## Regresyon Doğrusu



26

## Notasyon

	<u>Anakütle Parametresi</u>	<u>Örnek İstatistiği</u>
Regresyon eşitliğinde kesen	$\beta_0$	$b_0$
Regresyon eşitliğinin eğimi	$\beta_1$	$b_1$
Regresyon modeli ve eşitliği	$y = \beta_0 + \beta_1 x + \varepsilon$	$\hat{y} = b_0 + b_1 x$

27

## Artıklar ve En Küçük Kareler Yöntemi

### ❖ Artıklar

$$e = (y - \hat{y})$$

### ❖ En Küçük Kareler Yöntemi

$\Sigma e^2$ 'yi minimum yapan  $b_0$  ve  $b_1$  değerlerinin bulunmasıdır.

28

## $\beta_0$ and $\beta_1$ için En Küçük Kareler Tahminleyicileri

$$b_0 = \frac{(\sum y) (\sum x^2) - (\sum x) (\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$b_1 = \frac{n(\sum xy) - (\sum x) (\sum y)}{n(\sum x^2) - (\sum x)^2}$$

29

Önce  $b_1$  bulunursa, ardından

$$b_0 = \bar{y} - b_1 \bar{x}$$

30

## Satış geliri için regresyon eşitliğinin tahminlenmesi

$$b_1 = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$b_1 = \frac{10(800,62) - (268)(27,73)}{10(7668) - (268)^2}$$

$$b_1 = 0,118$$

$$b_0 = \bar{y} - b_1\bar{x} = 2,773 - (0,118)(26,8) = - 0,398$$

31

## Kestirimler (Öngörüler)

**Verilen bir x değeri için y'nin değeri ne olur?..**

**Eğer anlamlı bir korelasyon varsa, en iyi öngörülen y değeri, x değerinin regresyon eşitliğinde yerine konulmasıyla bulunur.**

**Önemli Not: Regresyon doğrusu yalnızca tahminlemede kullanılan x uzayı içinde geçerlidir. Mevcut x'lerden uzak bir noktada öngörümler yapılmamalıdır.**

32



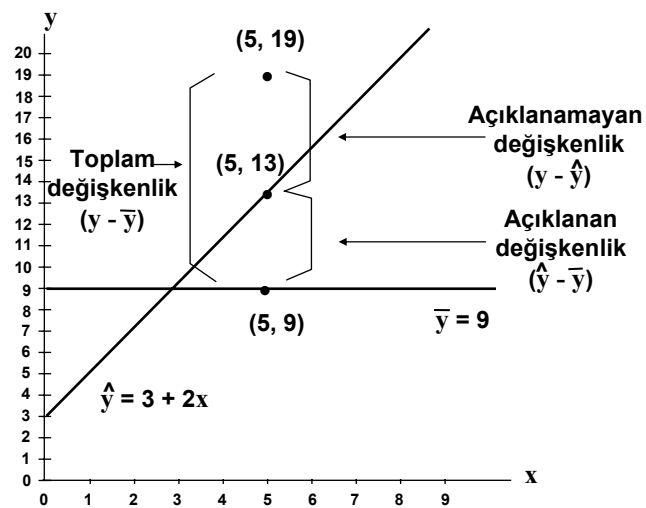
**30 satış personeli çalıştığında satış gelirinin kestirilmiş değeri nedir?**

$$\hat{y} = - 0.398 + 0.118 (30)$$

$$\hat{y} = 3.1516, \quad 315160 \$$$

33

### **Toplam Değişkenlik, Açıklanan Değişkenlik ve Açıklanamayan Değişkenlik**



34

(toplam değişkenlik) = (açıklanan değişkenlik) + (açıklanamayan değişkenlik)

$$(y - \bar{y}) = (\hat{y} - \bar{y}) + (y - \hat{y})$$

(toplam değişkenlik) = (açıklanan değişkenlik) + (açıklanamayan değişkenlik)

$$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2$$

(Genel kareler toplamı) = (regresyon kareler toplamı) + (artık kareler toplamı)

35

## Tanım

### Belirleme Katsayısı

y'deki değişkenliğin ne kadarının regresyon doğrusu tarafından açıklanabildiğini söyler.

$$r^2 = \frac{\text{Regresyon kareler toplamı}}{\text{Genel kareler toplamı}}$$

$$r^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2} = \frac{\text{RKT}}{\text{GKT}}$$

36

$$r^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2} = \frac{b_1^2(\sum x^2 - (\sum x)^2/n)}{\sum y^2 - (\sum y)^2/n}$$

$$r^2 = \frac{0,118^2(7668 - (268)^2/10)}{83,873 - (27,73)^2/10} = \%97,4$$

y'deki deęişmelerin %97,4'ü regresyon doğrusu ile açıklanabilmektedir.

37

## Varyans Analizi Tablosu (VAT)

Deęişkenlik Kaynağı	Kareler Topamları (KT)	Serbestlik Derecesi	Kareler Ortalaması (KO)	F-Oranı
<b>Regresyon</b>	$RKT = b_1^2(\sum x^2 - (\sum x)^2/n)$	<b>1</b>	Regresyon KO = $RKO = RKT / 1$	$F = \frac{RKO}{AKO}$
<b>Artık</b>	Artık Kareler Toplamı $AKT = GKT - RKT$	<b>n - 2</b>	Artık KO = $AKO = AKT / (n - 2)$ $= S^2$	
<b>Toplam (Genel)</b>	Genel Kareler Toplamı $GKT = \sum y^2 - (\sum y)^2/n$	<b>n - 1</b>		

38

## Hata Varyansının Tahmini

$$S^2 = \frac{\sum (y - \hat{y})^2}{n - 2} = \text{Artık Kareler Ortalaması}$$

## Tahminin Standart Hatası

$$S = \sqrt{\frac{\sum (y - \hat{y})^2}{n - 2}}$$

39

## F - Testi

❖  $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$   
(Model anlamsızdır)

$H_1$ : en az bir  $i$  için  $\beta_i \neq 0$   
(Model anlamlıdır)

40

## F – Testi (Basit Doğrusal Regresyon İçin)

❖  $H_0: \beta_1 = 0$   
(Model anlamsızdır)

$H_1: \beta_1 \neq 0$   
(Model anlamlıdır)

Test İstatistiği =  $F$  – oranı

Ret Bölgesi =  $F > F_{\alpha, 1, (n-2)}$  ise  $H_0$  RET.

41

## Varyans Analizi Tablosu (VAT) (Satış Gelirleri Örneği)

Değişkenlik Kaynağı	Kareler Toplamları (KT)	Serbestlik Derecesi	Kareler Ortalaması (KO)	F-Oranı
<b>Regresyon</b>	$RKT = b_1^2(\sum x^2 - (\sum x)^2/n)$ $= 0,118^2(7668 - (268)^2/10)$ $= 6,7982$	<b>1</b>	Regresyon KO = $RKO = RKT / 1$ $= 6,7982 / 1$ $= 6,7982$	$F = \frac{RKO}{AKO}$ $F = \frac{6,7982}{0,0225}$ $= 302,41$
<b>Artık</b>	Artık Kareler Toplamı $AKT = GKT - RKT$ $= 6,9780 - 6,7982$ $= 0,1798$	$n - 2 =$ $10 - 2 = 8$	Artık KO = $AKO = AKT / (n - 2)$ $= 0,1798 / 8$ $= 0,0225$	
<b>Toplam (Genel)</b>	$GKT = \sum y^2 - (\sum y)^2/n$ $= 83,873 - (27,73)^2/10$ $= 6,9780$	$n - 1 =$ $10 - 1 = 9$		

42

## F – Testi (Satış Gelirleri Örneği İçin)

❖  $H_0: \beta_1 = 0$   
(Model anlamsızdır)

$H_1: \beta_1 \neq 0$   
(Model anlamlıdır)

Test İstatistiği =  $F$  – oranı = 302,41

Karar =  $F = 302,41 > F_{0,05, 1, 8} = 5,32$   $H_0$  RET.

43

## Anakütle Regresyon Katsayılarının Testi

❖  $\beta_1$  = Anakütle regresyon katsayısı ( $X_1$  için)

❖  $H_0: \beta_1 = 0$   
( $\beta_1$  anlamsızdır)

$H_1: \beta_1 \neq 0$   
( $\beta_1$  anlamlıdır)

44

## Test İstatistiği $t$

**Test istatistiği:**

$$t = \frac{b_1}{S_{b_1}}$$

$S_{b_1}$  =  $b_1$ 'in standart hatasıdır.

$$S_{b_1} = \frac{S}{\sqrt{(\sum x^2 - (\sum x)^2/n)}}$$

45

## Kritik değerler

**serbestlik derecesi =  $n - 2$  olan  
tablo değerleri dikkate alınarak  
karar verilir.**

$|t| > t_{\alpha/2, n-2}$  ise  $H_0$  RET.

46

## Anakütle Regresyon Katsayılarının Testi (Satış Gelirleri Örneği)

❖  $\beta_1$  = Anakütle regresyon  
katsayısı ( $X_1$  için)

❖  $H_0: \beta_1 = 0$   
( $\beta_1$  anlamsızdır)

$H_1: \beta_1 \neq 0$   
( $\beta_1$  anlamlıdır)

47

## Test İstatistiği $t$

Test istatistiği:

$$t = \frac{b_1}{S_{b_1}} = \frac{0,118}{0,006804} = 17,39$$

$S_{b_1}$  =  $b_1$ 'in standart hatasıdır.

$$S_{b_1} = \frac{S}{\sqrt{(\sum x^2 - (\sum x)^2/n)}} = \frac{0,1499}{\sqrt{(7668 - (268)^2/10)}} = 0,006804$$

48



**Kritik değerler**  
**serbestlik derecesi =  $n - 2$  olan**  
**tablo değerleri dikkate alınarak**  
**karar verilir.  $\alpha = 0,05$  olsun.**

$|17,39| > t_{\alpha/2, n-2} = t_{0,025, 8} = 2,306$   
 $H_0$  RET.  $\beta_1$  anlamlıdır.

Basit doğrusal regresyonda  $t^2 = F$   
olmaktadır.

49

## **Anakütle Regresyon Katsayılarının Testi**

❖  $\beta_0$  = Anakütle regresyon  
modelinde sabit terim

❖  $H_0: \beta_0 = 0$   
( $\beta_0$  anlamsızdır)

$H_1: \beta_0 \neq 0$   
( $\beta_0$  anlamlıdır)

50

## Test İstatistiği $t$

**Test istatistiği:**

$$t = \frac{b_0}{S_{b_0}}$$

$S_{b_0}$  =  $b_0$ 'in standart hatasıdır.

$$S_{b_0} = \frac{S \sqrt{\sum x^2}}{\sqrt{n(\sum x^2 - (\sum x)^2/n)}}$$

51

## Kritik değerler

**serbestlik derecesi =  $n - 2$  olan  
tablo değerleri dikkate alınarak  
karar verilir.**

$|t| > t_{\alpha/2, n-2}$  ise  $H_0$  RET.

52

## Anakütle Regresyon Katsayılarının Testi (Satış Gelirleri Örneği)

❖  $\beta_0$  = Anakütle regresyon  
modelindeki sabit terim

❖  $H_0: \beta_0 = 0$   
( $\beta_0$  anlamsızdır)

$H_1: \beta_0 \neq 0$   
( $\beta_0$  anlamlıdır)

53

## Test İstatistiği $t$

**Test istatistiği:**

$$t = \frac{b_0}{S_{b_0}} = \frac{-0,398}{0,1884} = -2,11$$

$$S_{b_1} = \frac{S \sqrt{\sum x^2}}{\sqrt{n(\sum x^2 - (\sum x)^2/n)}} = \frac{(0,1499)\sqrt{(7668)}}{\sqrt{(10)(7668 - (268)^2/10)}}$$

$$= 0,1884$$

54

**Kritik değerler**  
**serbestlik derecesi =  $n - 2$  olan**  
**tablo değerleri dikkate alınarak**  
**karar verilir.  $\alpha = 0,05$  olsun.**

$$|-2,11| < t_{\alpha/2, n-2} = t_{0,025, 8} = 2,306$$

$H_0$  REDDEDİLEMEZ.  $\beta_0$  anlamsızdır.

55

## **E(y) Değeri İçin Kestirim Aralığı**

$$\hat{y} - E < E(y) < \hat{y} + E$$

**Burada**

$$E = t_{\alpha/2, n-2} S \sqrt{\frac{1}{n} + \frac{n(x_0 - \bar{x})^2}{n(\sum x^2) - (\sum x)^2}}$$

- $x_0$ ,  $x$ 'in verilen bir değeridir.
- Karekök içindeki ifade ile  $S$ 'nin çarpımı ise  $x_0$ 'daki  $\hat{y}$  değeri için standart hatadır.
- Standart hata en düşük değerini  $x_0 = \bar{x}$  olduğunda alır.

56

## **E(y) Değeri İçin Kestirim Aralığı**

$x_0 = 30$  personel için satışların beklenen değeri  
%95 güven ile hangi aralıkta gerçekleşir?

$$3.1516 - E < E(y) < 3.1516 + E$$

$$E = (2,306)(0,1499) \sqrt{\frac{1}{10} + \frac{(10)(30 - 26,8)^2}{(10)(7668) - (268)^2}}$$

$$E = (2,306)(0,01815) = 0,04186$$

$$3,1097 < E(y) < 3,1935$$