

SERİ
SERIE B

CİLT
TOME XX

SAYI
FASCICULE 2

1970

İSTANBUL ÜNİVERSİTESİ
ORMAN FAKÜLTESİ
DERGİSİ

REVUE DE LA FACULTÉ DES SCIENCES FORESTIÈRES
DE L'UNIVERSITÉ D'ISTANBUL



BASİT DOĞRUSAL REGRESYON

Dr. MS., Alptekin GÜNEL

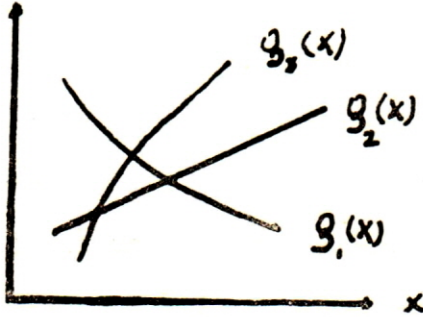
1. GİRİŞ

Bir toplum bir çok karakteristiklere sahiptir ve bu karakteristikler arasındaki ilişkinin mahiyeti tespit edilmek istenir. Tespit edilen bu ilişki-den faydalanarak, daha kolaylıkla elde edilen karakteristikler yardımı ile diğerleri tayin yoluna gidilir.

Burada, en basit hal olan, iki karakteristikli bir toplum ele alınacaktır. Problemin izahı bakımından en basit halin seçilmesi konunun öneminden bir şey kaybettirmemektedir. Zira, daha karışık haller için de uygulanacak teori ve kurallar tamamen aynı olup, fark sadece hesap ayrıntı-larındadır.

Yukarda sözü edilen toplumun karakteristiklerini X ve Y ile gösterelim. Tespiti istenen husus, X ile Y arasındaki ilişkinin mahiyeti ve belirli bir X değerine tekabül eden Y değerinin ne şekilde bulunabileceğidir. Böyle bir problemin tipik örneği, bir ağaçta çapla hacim arasındaki ilişkidir. Aynı şekilde, çapın yaşa bağlı olarak değişimi, asimlasyonun, belirli sınırlar içinde, sıcaklığa bağlı olarak artması iki karakteristikli bir topluma ait örneklerdir.

Çapla hacim arasındaki örnekte olduğu gibi, belirli bir X (çap) değerine birden fazla Y (hacim) değeri tekabül edebilir. Bu durumda, belirli bir X değerine tekabül eden Y değerleri bir «toplum» teşkil edecek-



Şekil — 1. $g(X)$ 'in alabileceği

lerdir; bu toplumun bir ortalaması vardır. Bu ortalamayı $g(X)$ ile gösterelim. Böylece, farklı X değerleri için, Y'lerin ortalamaları da farklı olacaktır. İşte, belirli X değerine tekabül eden Y değerlerinin ortalamasını veren $g(X)$ fonksiyonuna X'in Y'e ait «Regresyon fonksiyonu» adı verilir. Bu fonksiyon çok değişik form-

Bu fonksiyonun formu, genellikle, deneylerle elde edilen verilere dayanılarak tespit edilmek istenir. Söz konusu fonksiyonun çok çeşitli olması, bu amaçla baş vurulan teori ve prensiplerde her hangi bir değişikliği gerektirmemekte, her hal için, gene aynı esaslardan hareket edilmektedir. Bu nedenle, burada

$$g(X) = C + D X \quad (1)$$

formu üzerinde durulacaktır.

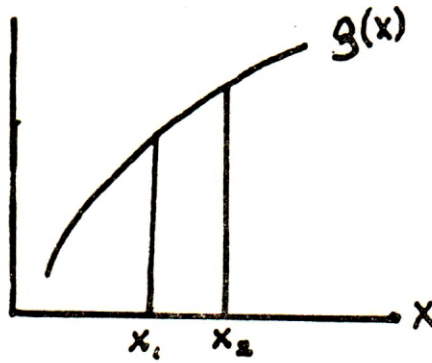
Denklem (1), bilindiği gibi, bir doğru denklemdir. Böyle bir formun seçimini haklı gösteren sebepler vardır: Daha yüksek dereceden eğrilerde bile, belirli sınırlar içinde, X ile Y arasındaki ilişki bir doğru olarak kabul edilebilir (Şekil-2). Ayrıca, birçok fonksiyonlar, logaritmik kâğıt üzerinde doğrusal bir seyir gösterirler. Nihayet, yapılacak değişken değişikliği ile, fonksiyona doğrusal bir karakter verilebilir.

Farz edelim ki, N tane X değerine tekabül eden Y değerleri tespit edilmiş olsun. Yukardaki kabule uygun olarak, X ile Y arasındaki ilişkinin doğrusal olduğunu farz edeceğiz.

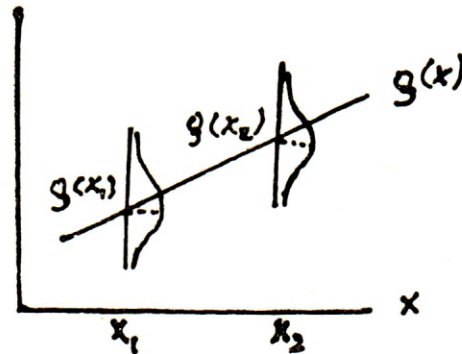
Belirli bir X değeri için bir çok Y değeri tespit edildiğini ve bu Y değerlerinin bir toplum teşkil ettiğini belirtmiştik. Böyle bir toplumun $g(X)$ gibi bir ortalaması olduğu gibi, bir de varyansı olacaktır. Bu varyansın, ortalamaların aksine, bütün Y toplumları için aynı olduğunu ve Y toplumlarının normal bir dağılım gösterdiğini diğer bir kabul olarak ileri süreceğiz (Şekil-3). Söz konusu varyansı $\sigma^2_{y,x}$ ile gösterelim. (1) nolu denklemi

$$g(X) = \mu_{y,x} = A + B(X - \bar{X}) \quad (2)$$

şeklinde yazmak, bazı kolaylık-



Şekil — 2



Şekil — 3. X'e tekabül eden Y toplumları normal dağılımlı ve eşit varyanslıdır. Ortalamaları ise farklıdır.

lar sağlanması nedeni ile, daha çok tercih edilmektedir. (*)Denklemde,

A ve B = Denklem katsayıları

$\bar{X} = \Sigma X/N = X$ 'lerin ortalamasıdır.

Buna göre, herhangi bir münferit X_i değerine tekabül eden Y_i gerçek değeri

$$Y_i = A + B(X_i - \bar{X}) + \varepsilon_i ; \quad i = 1, 2, \dots, N \quad (3)$$

denkleminde bulunabilir. (3) nolu denklemde ε_i ler normal dağılımlı, ortalamaları sıfır, birbirlerinden bağımsız ve varyansları ise $\sigma^2_{y.x}$ olan tesadüfi değişkenlerdir. (işaretle $\varepsilon_i : N(0, \sigma^2_{y.x})$) ε_i 'ler hakkındaki bu varsayımlar, özellikle, ε_i 'ler birçok hatların bir sonucu ise gerçeğe daha çok uymaktadır. Bu durum ormancılık problemlerinde belirgin bir niteliktedir.

(2) nolu denklemi yazmakla, X ile Y arasındaki ilişkiyi belirten gerçek bir regresyondogrusunun varlığını kabul etmiş oluyoruz. Amacımız bu doğruyu, deneyle elde ettiğimiz verilere dayanarak, en doğru şekilde tayin etmektir.

Söz konusu doğrunun tayininde en makul yol, bulunacak doğrunun (X_i, Y_i) değer çiftlerine göre elde edilen noktalara mümkün olduğu kadar «yakın» olmasıdır. Böyle bir yakınlık iki şekilde sağlanabilir:

- 1 — Noktalardan doğruya inilen diklerin (noktaların doğruya olan uzaklıklarının) toplamı minimum olsun, veya
- 2 — Bulunacak doğrudan alınacak Y değerleri ile tespit edilen Y değerleri arasındaki farkların kareleri toplamı minimum olsun.

İkinci şekilde yapılacak hesap işlemlerinin, birinciye nazaran daha az karışık olması ve ikinci şekilde elde edilen istatistiklerin, aşağıda da gösterileceği gibi, eğilim hatasından arı olması, ikinci yolun daha çok kullanılmasına sebep olmuştur. İkinci yol, tatbikatta, en küçük kareler metodu olarak tanınır.

2. En Küçük Kareler Metodu:

Yukarda da belirtildiği gibi, en küçük kareler metodunun esası, hesaplanan doğrudan alınan Y değerleri ile, tespit edilen Y değerleri arasındaki farkların kareleri toplamının minimum olması şeklinde ifade edilebilir. Bu şekilde bulunacak doğru, gerçek doğruyu temsil eden birçok doğrudan sadece bir tanesi olup, gerçek regresyon doğrusuna tamamen

* Gerçekte, denklem $g(x) = A + BX$ şeklinde alınır ve (2). bölümde anlatılan esaslar dahilinde A' yı hesaplayıp yerine koyduktan sonra gerekli düzeltmeler yapılırsa denklem 2 elde edilir.

eşit olduğu söylenemez. Bununla beraber, diğer doğrulara üstün tutulması için bir çok sebepler vardır.

En küçük kareler metodu ile hesaplanan regresyon doğrusu denkleminin

$$Y = a + b(X - \bar{X}) \quad (4)$$

şeklinde olduğunu kabul edelim. Bu denklem yardımı ile bulunan ve herhangi bir X_i değerine tekabül eden Y değerini (\hat{Y}_i) ile gösterelim. (4) nolu denklemin belli olması (a) ve (b) değerlerinin bilinmesine bağlıdır. En küçük kareler metodunun prensibine göre,

$$\begin{aligned} \Sigma(Y_i - \hat{Y}_i)^2 &= \Sigma[(Y_i - a - b(w_i - \bar{X}))^2] \\ &= \text{minimum} \end{aligned} \quad (5)$$

olmalıdır, Monoton artan bir fonksiyonun minimumu, türevini sıfır yapan değer olduğundan, (5) nolu eşitliğin bilinmeyenleri olan (a) ve (b) ye göre kısmi türevleri alınır ve $\Sigma(X_i - \bar{X}) = 0$ olduğu göz önünde bulundurulursa,

$$a = \frac{1}{N} \Sigma Y_i = \bar{Y} \quad (6)$$

ve

$$b = \frac{\Sigma Y_i(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2} \quad (7)$$

elde edilir. (b)'nin payı

$$\Sigma X_i Y_i - \frac{1}{N} (\Sigma X_i) (\Sigma Y_i),$$

payadası ise

$$\Sigma X_i^2 - \frac{1}{N} (\Sigma X_i)^2$$

şekillerinde de yazılabileceğinden; ayrıca,

$$S_x^2 = \frac{\Sigma(X_i - \bar{X})^2}{N-1} = \frac{\Sigma X_i^2 - \frac{1}{N} (\Sigma X_i)^2}{N-1} = X\text{'in varyansı}$$

$$S_y^2 = \frac{\Sigma(Y_i - \bar{Y})^2}{N-1} = \frac{\Sigma Y_i^2 - \frac{1}{N} (\Sigma Y_i)^2}{N-1} = Y\text{'in varyansı}$$

$$\begin{aligned} S_{xy} &= \frac{\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{N-1} = \frac{1}{N-1} \left[\Sigma X_i Y_i - \frac{1}{N} (\Sigma X_i) (\Sigma Y_i) \right] \\ &= X \text{ ile } Y\text{'nin kovaryansı} \end{aligned}$$

yazılacak olursa, (b) katsayısı

$$b = \frac{S_{xy}}{S_x^2} \quad (7a)$$

şeklinde de hesaplanabilir.

Yukardaki hesap işlemlerinden de görüleceği gibi, yapılacak bütün iş, ΣX_i , ΣY_i , ΣX_i^2 , ΣY_i^2 , ve $\Sigma X_i Y_i$ değerlerini hesaplamak ve bunları denklemlerde uygun yerlerine koymaktan ibarettir. Söz konusu beş değer bir defa hesaplandıktan sonra, diğer rakamlara artık ihtiyacımız kalmamıştır.

(a) ve (b) katsayılarından sonra bilinmesi gereken diğer bir ifade $\sigma_{y.x}^2$ dir. Bu amaçla

$$S_{y.x}^2 = \frac{1}{N-2} \Sigma (Y_i - \hat{Y}_i)^2 \quad (8)$$

formülü kullanılabilir. Paydadaki terimin $(N-2)$ olması, (4) nolu denklemde iki parametre (a ve b) bulunması nedeniyledir. Denklemde hesaplanması gereken dört katsayı bulunsa idi, bu ifade $(N-4)$ olacaktı. (8) nolu denklemde $a = \bar{Y}$ ve $\hat{Y}_i = a + b(X_i - \bar{X})$ yazarak

$$\begin{aligned} S_{y.x}^2 &= \frac{1}{N-2} \Sigma (Y_i - a - b(X_i - \bar{X}))^2 \\ &= \frac{1}{N-2} \Sigma ((Y_i - \bar{Y}) - b(X_i - \bar{X}))^2 \\ &= \frac{1}{N-2} (\Sigma (Y_i - \bar{Y})^2 - 2b \Sigma (X_i - \bar{X})(Y_i - \bar{Y}) + b^2 \Sigma (X_i - \bar{X})^2) \\ &= \frac{1}{N-2} [(N-1) S_y^2 - 2b(N-1) S_{xy} + b^2(N-1) S_x^2] \\ &= \frac{N-1}{N-2} (S_y^2 - 2b S_{xy} + b^2 S_x^2) \end{aligned}$$

$b = S_{xy}/S_x^2$ konarak, gerekli kısaltmalardan sonra

$$S_{y.x}^2 = \frac{N-1}{N-2} \left(S_y^2 - \frac{S_{xy}^2}{S_x^2} \right) \quad (9)$$

bulunur.

ÖRNEK:

X_i	2	2	3	3,5	4	4,5	4,5	5	3,6	6
Y_i	1,9	2,1	3,5	2,7	3	3	3,3	3,5	3,6	3,7

$N=10$

$$\begin{aligned}\Sigma X_i &= 40,0 & \Sigma X_i^2 &= 177,0 & \bar{X} &= 4,00 & (\Sigma X_i)^2 &= 1600 \\ \Sigma Y_i &= 29,4 & \Sigma Y_i^2 &= 89,86 & \bar{Y} &= 2,94 & (\Sigma Y_i)^2 &= 864,36\end{aligned}$$

$$\Sigma X_i Y_i = 125,1$$

$$S_x^2 = \frac{1}{9} \left(177,0 - \frac{1}{10} (1600) \right) = 1,889$$

$$S_y^2 = \frac{1}{9} \left(89,86 - \frac{1}{10} (864,36) \right) = 3,804$$

$$S_{xy} = \frac{1}{9} \left(125,1 - \frac{1}{10} (40,0) (29,4) \right) = 0,833$$

$$b = 0,833/1,889 = 0,4409$$

$$a = \bar{Y} = 2,94$$

$$S_{y \cdot x}^2 = \frac{9}{8} \left(3,804 - \frac{(0,833)^2}{1,889} \right) = 4,126$$

Hesaplanan doğru denklemi:

$$\hat{Y}_i = 2,94 + 0,4409 (X_i - 4,00) \quad \text{veya}$$

$$\hat{Y}_i = 1,1764 + 0,4409 X_i \quad \text{dir.}$$

3. a, b, $S_{y \cdot x}^2$ ve \bar{Y} istatistiklerinin özellikleri

A, B, $\sigma_{y \cdot x}^2$ ve $\mu_{y \cdot x}$ parametrelerinin güven sınırlarını bulmak ve söz konusu olabilecek hipotezlerin testi bakımından yukardaki istatistiklerin dağılım fonksiyonlarının bilinmesine ihtiyaç vardır.

3.1 — a'nın dağılım fonksiyonu

a = \bar{Y} olduğundan, (3) nolu denklem yardımı ile,

$$a = \bar{Y} = \frac{1}{N} \Sigma Y_i = \frac{1}{N} \Sigma (A + B(X_i - \bar{X}) + \varepsilon_i)$$

$$= A + B \frac{\Sigma (X_i - \bar{X})}{N} + \frac{\varepsilon_i}{N}$$

$$= A + \frac{\varepsilon_i}{N}$$

$$(\Sigma (X_i - \bar{X}) = 0 \text{ olduğundan})$$

bulunur. Ancak, ε_i lerin ortalaması (beklenen değeri) sıfır kabul edildiğinden, $(\varepsilon_i)/N$ nin beklenen değeri de sıfır olacaktır, dolayısıyla

$$a = \bar{Y} = A$$

bulunacaktır; yani a 'nın beklenen değeri A 'ya eşittir. Böylece, (a) eğilim hatasından aridir.

ε_i lerin varyansları eşit ve $\sigma^2_{y,x}$ kabul edildiğinden, (a)'nın varyansı

$$S_a^2 = \sigma^2_{x,y}/N$$

olur. ε_i ler normal dağılımlı olduklarından, (a) da normal dağılımlı olacaktır. İşaretle gösterilirse

$$a: N(A; \sigma^2_{y,x}/N)$$

dir.

3.2 — (b)'nin dağılım fonksiyonu

(3) ve (7) nolu formüllerden

$$b = \frac{\sum Y_i (X_i - \bar{X})}{\sum (X_i - \bar{X})^2} = \frac{\sum (A + B(X_i - \bar{X}) + \varepsilon_i) (X_i - \bar{X})}{\sum (X_i - \bar{X})^2}$$

$\sum (X_i - \bar{X}) = 0$ olduğu göz önünde tutulursa,

$$b = B + \frac{\sum \varepsilon_i (X_i - \bar{X})}{\sum (X_i - \bar{X})^2}$$

bulunur. ε_i lerin birbirlerinden bağımsız, ortamları sıfır ve X_i lerin sabit kabul edildikleri hatırlanarak,

$$b = B$$

elde edilir.

(b)'nin varyansının

$$S_b^2 = \frac{\sigma^2_{y,x}}{\sum (X_i - \bar{X})^2} = \frac{\sigma^2_{y,x}}{(N-1)S_x^2}$$

olduğu yukardaki eşitlikten kolayca gösterilebilir. Aynı şekilde, ε_i lerin normal dağılımlı olması nedeni ile, (b) de normal dağılımlıdır. Yani

$$b: N(B; \sigma^2_{y,x}/(N-1)S_x^2)$$

dir.

(b)'nin varyans formülünden de görüldüğü gibi, X 'lerin varyansı ne kadar fazla ise, yani X_i 'ler ne kadar çok birbirlerinden farklı iseler, (b) nin varyansı o kadar azalacaktır. Aynı şekilde, (b)'nin varyansı gözlem sayısı ile ters orantılıdır; gözlem sayısını arttırmak suretiyle, (b)'nin varyansını küçültmek mümkündür. Bu sonuç, arzu edilen bir özelliktir.

3.3 — $S^2_{y..x}$ 'in dağılım fonksiyonu

İspat edilebilir ki

$$S^2_{y..x}/\sigma^2_{y..x}$$

istatistiği ($X^2/d.f$ (khi karesi/serbestiyet derecesi) adı verilen bir dağılım göstermektedir. Bu dağılımın ortalaması (1) olduğundan $S^2_{y..x}$ 'in ortalaması $\sigma^2_{y..x}$ olacaktır.

3.4 — a, b ve $S^2_{y..x}$ arasındaki ilişki

Yukarda elde edilen sonuçlar göz önünde tutulursa, söz konusu üç istatistiğin birbirlerinden bağımsız oldukları görülür. Böylece, meselâ, (a)'nın hesabında yapılacak bir hata, (b)'yi etkilemeyecektir.

3.5 — \bar{Y} 'nin dağılım fonksiyonu

(a) ve (b) birbirlerinden bağımsız ve normal dağılımlı olduklarından, bunların doğrusal bir fonksiyonu olan \bar{Y} de normal dağılımlı olacaktır. (a) ve (b)'ye ait sonuçlardan \bar{Y} nin ortalamasının

$$\mu_{y..x} = A + B(X_i - \bar{X})$$

olduğu kolaylıkla çıkartılabilir. Aynı şekilde, \bar{Y} 'nin varyansı

$$S^2_{\bar{y}} = \sigma^2_{y..x} \left(\frac{1}{N} + \frac{(X_i - \bar{X})^2}{(N-1)S^2_x} \right)$$

yazılabilir.

Parantezin içindeki ikinci terimin payındaki ifadenin de ortaya koyduğu gibi, X_i değeri ortalamadan ne kadar fazla fark ederse, \bar{Y} 'nin sıklığı o kadar daha az olacaktır.

4 — Doğrusal Regresyonla İlgili Hipotezler

4.1 — A ile ilgili hipotezler

Üç çeşit hipotez söz konusudur:

- A belirli bir A_0 değerine eşittir,
- A belirli bir A_0 değerinden küçüktür,
- A belirli bir A_0 değerinden büyüktür.

Bu ifadeler formüle edilecek olursa, aşağıdaki şekilde yazılabilirler:

$$\begin{array}{lll} H_0 : A = A_0 & \text{veya} & H_0 : A \leq A_0 & \text{veya} & H_0 : A \geq A_0 \\ H_A : A \neq A_0 & & H_A : A > A_0 & & H_A : A < A_0 \end{array}$$

Her üç hipotezin de kontrolunda t- testi kullanılabilir.

$$\frac{a - A_0}{S_a} = \frac{a - A_0}{S_{y,x}/\sqrt{N}}$$

ifadesi t- dağılımı gösterir. Dağılımın serbestiyet derecesi (N-2) dir. Birinci hipotez «iki yanlı» hipotez olup diğerleri «tek yanlı» dırlar.

4.2 — B ile ilgili hipotezler

A ile ilgili hipotezlerin tamamen aynı olup hipotezlerin kontrolu gene aynı şekilde yapılır:

$$\begin{array}{llll} H_0 : B = B_0 & \text{veya} & H_0 : B \leq B_0 & \text{veya} & H_0 : B \geq B_0 \\ H_A : B \neq B_0 & & H_A : B > B_0 & & H_A : B < B_0 \end{array}$$

t — değeri

$$\frac{b - B_0}{S_b} = \frac{b - B_0}{S_{y,x}/S_x\sqrt{N-1}}$$

şeklinde hesaplanmaktadır. Dağılımın serbestiyet derecesi (N-2) dir.

B ile ilgili en önemli hipotez B=0 hipotezi, yani X ile Y arasında herhangi bir ilişki bulunmadığıdır. Böyle bir hipotezin kontrolunda üzerinde durulması gereken iki hal mevcuttur: Az sayıda ölçme yapılması halinde, B'nin sıfırdan çok farklı olmasına rağmen, B=0 hipotezini red etme ihtimali çok düşük olacaktır. Diğer taraftan, B sıfıra çok yakın olabilir; fakat, gözlem sayısı oldukça yüksek ise, B=0 hipotezini red etme ihtimali atar. Böyle bir durumda, güvenilirlik derecesini küçültmek, yani daha büyük (α) değeri almak yerinde olur. Unutmamak lâzımdır ki, günlük hayattaki güvenilirlik ile istatistikteki güvenilirlik aynı şey değildir.

B=0 hipotezinin kontrolu varyans analizi yolu ile de yapılabilir. Bu amaçla kullanılacak istatistik

$$F_{1,(N-2)} = \frac{b^2 \sum (X_i - \bar{X})^2}{S_{y,x}^2} = t^2_{(N-2)}$$

dir.

Bu testle ilgili varyans analizi tablosu şöyledir:

Varyasyon Kaynağı	Kareler Toplamı	Serbestiyet Derecesi	Kareler Ortalaması	Beklenen Değer
Regresyon	$b^2 \sum (X_i - \bar{X})^2$	1		$\sigma^2_{y,x} + B^2 \sum (X_i - \bar{X})^2$
Hata	$(N-2) S_{y,x}^2$	N-2		$\sigma^2_{y,x}$

$B=0$ hipotezi doğru ise, tablonun son sütunundan da görüldüğü gibi, bu takdirde her iki varyasyon kaynağının beklenen değeri $\sigma^2_{y \cdot x}$ olacağından, $F=1$ bulunacaktır. B ne kadar sıfırdan fark ederse, regresyonun kareler ortalaması hatanın kareler ortalamasından o kadar fark edecek, neticede F değeri 1 den büyük olacaktır. Bu nedenle, hesaplanan F değeri, seçilen (α) değeri için, tablodan alınacak F değerlerinden büyükse, $B=0$ hipotezi red edilmelidir.

5 — Güven Aralıkları

A ve B parametleri normal dağılımlı ve birbirlerinden bağımsız olduklarından, bu parametlere ait güven aralıklarının tayininde, t -dağılımı kullanılabilir. (α) kadar güvenirlilik derecesi için, A 'nın güven aralığı:

$$a - t_{1-\alpha/2} \frac{S_{y \cdot x}}{\sqrt{N}} \leq A \leq a + t_{1-\alpha/2} \frac{S_{y \cdot x}}{\sqrt{N}}$$

B 'nin güven aralığı:

$$b - t_{1-\alpha/2} \frac{S_{y \cdot x}}{S_x \sqrt{N-1}} \leq B \leq b + t_{1-\alpha/2} \frac{S_{y \cdot x}}{S_x \sqrt{N-1}}$$

$\sigma^2_{y \cdot x}$ 'in güven aralığı:

$$\frac{S^2_{y \cdot x}}{K^2_{1-\alpha/2}/(N-2)} \leq \sigma^2_{y \cdot x} \leq \frac{S^2_{y \cdot x}}{K^2_{\alpha/2}/(N-2)}$$

$\mu_{x \cdot y}$ 'in güven aralığı:

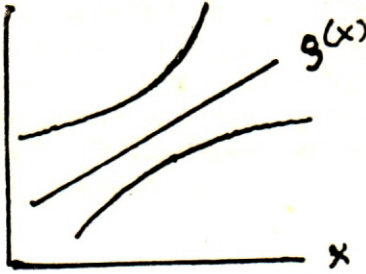
$$\begin{aligned} \bar{Y} - t_{1-\alpha/2} S_{y \cdot x} \sqrt{\frac{1}{N} + \frac{(X_i - \bar{X})^2}{(N-1)S_x^2}} \leq \mu_{y \cdot x} \leq \bar{Y} + \\ + t_{1-\alpha/2} S_{y \cdot x} \sqrt{\frac{1}{N} + \frac{(X_i - \bar{X})^2}{(N-1)S_x^2}} \end{aligned}$$

Münferit bir \bar{Y}_i değerine ait güven aralığı:

$$\bar{Y}_i \pm t_{1-\alpha/2} S_{y \cdot x} \sqrt{1 + \frac{1}{N} + \frac{(X_i - \bar{X})^2}{(N-1)S_x^2}}$$

şeklinde hesaplanırlar.

Son iki denklemden de görüldüğü gibi, X_i \bar{X} 'den uzaklaştıkça güven aralığı genişlemektedir (Şekil-4).



6 — Korelasyon Katsayısı

X ile Y arasındaki doğrusal ilişkinin derecesini tayin amacı ile kullanılabilen bir kriter, Pearson'ın korelasyon katsayısıdır. (r) ile gösterilen bu katsayı

$$r = \frac{S_{xy}}{S_x S_y} \quad (10)$$

Şekil — 4. Güven aralıkları formülünden hesaplanmaktadır.

Korelasyon katsayısı aşağıdaki özelliklere sahiptir:

1 — $|S_{xy}| = S_x S_y$ olduğundan, r- değeri -1 ile $+1$ değerleri arasındadır veya bunlardan birine eşittir, yani

$$-1 \leq r \leq +1$$

2 — $r = b \frac{S_x}{S_y}$ yapılabileceğinden, $r=0$, ancak ve ancak, $b=0$ için olur. Ayrıca, r ile b'nin işaretleri aynıdır

3 — Kabaca denebilir ki, r, hesaplanan regresyon doğrusunun noktalara çok yakın olması halinde, -1 veya $+1$ değerlerini alır. Doğrunun noktalardan uzak seyretmesi oranında r de sifıra yaklaşır.

Korelasyon katsayısı r'nin regresyonla ilgi derecesi, Y'lere ait kareler toplamı üzerinde yapılacak ufak bir işlemle daha iyi belirtilebilir. Y'lerin kareler toplamı

$$\begin{aligned} \Sigma(Y_i - \bar{Y})^2 &= \Sigma(\bar{Y}_i - \check{Y}_i + \check{Y}_i - Y)^2 \\ &= \Sigma(Y_i - \check{Y}_i)^2 + \Sigma(\check{Y}_i - \bar{Y})^2 + 2 \Sigma(Y_i - \check{Y}_i)(\check{Y}_i - \bar{Y}) \\ &= \Sigma(Y_i - \check{Y}_i)^2 + \Sigma(\check{Y}_i - \bar{Y})^2; \quad (\Sigma(Y_i - \check{Y}_i)(\check{Y}_i - \bar{Y}) = 0 \text{ dir}). \end{aligned}$$

şeklinde iki kısma ayrılır. Eşitliğin sağ tarafındaki ilk terim hatadan doğan kareler toplamı olup «Hata Kareler toplamı», ikinci terim ise regresyon nedeni ile meydana gelen kareler toplamı olup «Regresyon kareler toplamı» adları ile tanılırlar. (9) nolu formülden, hatalar kareler toplamı

$$(N-1) \left(S_y^2 - \frac{S_{xy}^2}{S_x^2} \right)$$

şeklinde yazılabilir. Diğer taraftan, $\Sigma(Y_i - \bar{Y})^2 = (N-1) S_y^2$ olduğundan, regresyon kareler toplamı

$$\Sigma(Y_i - \bar{Y})^2 = (N-1) \frac{S_{xy}^2}{S_x^2} = (N-1) S_y^2 r^2$$

yardımları ile bulunabilir.

Y'nin kareler toplamına ait eşitliğin her iki yanını $\Sigma(Y_i - \bar{Y})^2$ ile bölmürse,

$$1 = \frac{(N-2) \left(S_y^2 - \frac{S_{xy}^2}{S_x^2} \right)}{\Sigma(Y_i - \bar{Y})^2} + \frac{(N-1) \frac{S_{xy}^2}{S_x^2}}{\Sigma(Y_i - \bar{Y})^2}$$

$$\Sigma(Y_i - \bar{Y})^2 = (N-1) S_y^2 \text{ yazılarak}$$

$$1 = (1 - (S_{xy}^2 / (S_x^2 S_y^2))) + (S_{xy}^2 / (S_x^2 S_y^2))$$

elde edilir.

$$S_{xy}^2 / (S_x^2 S_y^2) = r^2$$

dir. O halde r^2 , Y'lerin kareler toplamının «sebebi bilinen varyasyon» (explained variation) kısmını teşkil etmektedir. Bu durumda, $(1-r^2)$ «Sebebi bilinmeyen varyasyon» (unexplained variation) oranını ifade eder. (r) değeri ne kadar 1'e yakınsa, yani hesaplanan doğru, ölçme ve gözlem sonucunda elde edilen noktalara ne kadar yakınsa, $(1-r^2)$ ifadesi o kadar sıfıra yaklaşacaktır. Bu sonuç, doğrusallık hipotezinin kontrolünde kullanılır.

7 — Regresyonun Doğrusallığının Kontrolü

Yukarıda X ile Y arasındaki ilişkinin doğrusal olduğunu kabul etmiş ve regresyon denklemini

$$g(X) = \mu_{y.x} = A + B(X - \bar{X})$$

şeklinde yazmıştık. En az bir X_i için birden fazla Y değerleri tespit edilmişse (meselâ, aynı çap için birden fazla hacim bulunmuşsa), sözü edilen varsayımın doğruluğunu kontrol etmek mümkündür.

Problemin izahı bakımından aşağıdaki işaretlemeyi kullanacağız:

$$Y_{ij} = X_i \text{ ye tekabül eden } j \text{ inci ölçme,}$$

$$n_i = X_i \text{ ye tekabül eden } Y_{ij} \text{ sayısı}$$

$$T_{y_i.} = \sum_{j=1}^{n_i} Y_{ij}, \quad Y_{y..} = \sum_i \sum_j Y_{ij};$$

$$\bar{Y}_{i.} = (T_{y_i.} / n_i) \quad ; \quad \bar{Y}_{..} = T_{y..} / \Sigma n_i$$

ÖRNEK :

$$\begin{array}{llllll} X_1=5 & Y_{11}=3,0 & X_2=5,5 & Y_{21}=3,8 & X_3=6 & Y_{31}=4,2 \\ & Y_{12}=4,0 & & Y_{22}=4,2 & & Y_{32}=5,0 \\ & & & Y_{23}=4,4 & & \end{array}$$

$$n_1=2; \quad n_2=3; \quad n_3=2$$

$$Ty_{1.}=7,0; \quad Ty_{2.}=12,4; \quad Ty_{3.}=9,2; \quad Ty_{..}=28,6$$

$$Y_{1.}=3,5; \quad Y_{2.}=4,13; \quad Y_{3.}=4,6; \quad \bar{Y}_{..}=4,08$$

Farz edelim ki, k tane farklı X değeri mevcuttur; dolayısıyla, k tane Y değerler topluluğu tespit edilmiş olacaktır. X_i değerine tekabül eden Y değerlerinin gerçek ortalamasını U_i ile gösterelim. Aynı zamanda, yukarıda yaptığımız gibi, Y'lerin normal dağılımlı ve varyanslarının eşit olduğunu da kabul edeceğiz.

Bu takdirde, iki ayrı varsayımın kontrolü söz konusudur:

— Bütün U_i 'ler birbirlerine eşittir, yani

$$U_1=U_2= \dots = U_k$$

— U_i 'ler X'lerin doğrusal bir fonksiyonudur.

Bu iki varsayımın kontrolünde varyans analizi kullanılabilir. Aynı X_i değerine tekabül eden Y değerlerine «Grup» diyecek olursak, aşağıdaki eşitlikler yazılabilir:

Gruplar içi kareler toplamı

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 \\ &= \sum_{i=1}^k \sum Y_{ij}^2 - \sum \frac{T_{y_{i.}}^2}{n_i} \end{aligned} \quad (11)$$

Gruplar arası kareler toplamı

$$= \sum n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2 = \sum \frac{T_{y_{i.}}^2}{n_i} - \frac{T_{y_{..}}^2}{\sum n_i} \quad (12)$$

Regresyon Kareler Toplamı

$$\begin{aligned} &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2 = (N-1) \frac{S_{xy}^2}{S_x^2} \\ &= (N-1) S_y^2 r^2 \end{aligned} \quad (13)$$

Regresyon etrafındaki Kareler Toplamı

$$= \text{Gruplararası Kareler toplamı} - \text{Regresyon Kareler toplamı} \quad (14)$$

Doğrusallık kontrolü amacı ile hazırlanacak varyans analizi tablosu şöyledir:

Varyasyon Kaynağı	Kareler Toplamı	Serbestiyet Derecesi	Kareler Ortalaması	Kareler Ortalamasının Beklenen Değeri
1	2	3	4	5
Gruplar içi			$\Sigma n_i - k$	$\sigma^2_{y \cdot x}$
Regresyon			1	$\sigma^2_{y \cdot x} + \Sigma n_i \sigma^2_u R^2_{xu}$
Regresyon etrafı			$k - 2$	$\sigma^2_{y \cdot x} + \frac{\Sigma n_i}{k - 2} \sigma^2_u (1 - R^2_{xu})$
Gruplar arası			$k - 1$	$\sigma^2_{y \cdot x} + \frac{\Sigma n_i}{k - 2} \sigma^2_u$

Tablonun son sütunundaki

$$\sigma_u^2 = \frac{1}{N} \Sigma n_i (U_i - \bar{U})^2 = \text{Ortalamalar varyansı}$$

$$\bar{U} = \frac{1}{N} \Sigma n_i U_i = \text{Genel ortalama}$$

$R_{xu} = X_i$ ile buna tekabül eden U_i 'ler arasındaki korelasyon katsayısıdır.

Birinci varsayımı kontrol için

$$F_1 = \frac{\text{Gruplararası Kareler ortalaması}}{\text{Gruplar içi kareler ortalaması}}$$

istatistiği kullanılır. Varsayımımız doğru ise, $\sigma_u^2 = 0$ olacağından, pay ve paydanın beklenen değerleri, tablonun son sütunundan da görüldüğü gibi, $\sigma^2_{y \cdot x}$ olacak ve $F_1 = 1$ bulunacaktır. Hiç değilse, ortalamalardan biri diğerlerine eşit değilse, $F_1 \neq 1$ olacaktır. F_1 'in serbestiyet derecesi ($k - 1$) ve $(\Sigma n_i - k)$ 'dir. Hesaplanan değer, seçilen (α) güvenirlilik derecesi ve yukarıdaki serbestiyet dereceleri için, tablodan alınan değerden büyükse ortalamaların eşitliği varsayımı red edilecektir. **İkinci varsayım ancak birinci varsayım red edildikten sonra kontrol edilir.** İkinci varsayımı kontrol için kullanılacak istatistik

$$F_2 = \frac{\text{Regresyon etrafında kareler ortalaması}}{\text{Gruplar içi kareler ortalaması}}$$

dir.

İkinci varsayım doğru ise, $R_{xu}=1$ olacaktır. Bu takdirde, pay ve paydanın beklenen değerleri aynı ve $\sigma_{y \cdot x}^2$ olacak, dolayısıyla $F_2=1$ bulunacaktır. $R_{xu} \neq 1$ ise, F_2 'nin değeri 1'den farklıdır. F_2 'nin serbestiyet dereceleri $(k-2)$ ve $(\sum n_i - k)$ 'dir. Tablodan alınan değer hesaplanan değerden küçükse, varsayım red edilir.

Evvence de belirtildiği gibi, her iki tez birbirinden bağımsız değildir. Testin kuvveti k ve n_i sayıları yükseldikçe artar.

ÖRNEK

X_i	5,0	5,5	6,0	6,5	7,0	7,5	8,0	8,5	9,0	9,5
Y_{ij}	3,0 4,0	3,8 4,2 4,4	4,2 5,0	5,4 6,0 6,2	5,6 6,1 6,1 6,6	6,3	7,0 8,0 8,4	7,1 8,0 8,5 8,5	8,6 9,0	8,0 9,4 9,8
n_i	2	3	2	3	4	1	3	4	2	3
$T_{y \cdot i}$	7,0	12,4	9,2	17,6	24,4	6,3	23,4	32,1	17,6	27,2
$\bar{Y}_{\cdot i}$	3,5	4,13	4,6	5,86	6,1	6,3	7,8	8,02	8,8	9,06

$$\sum_{i=1}^k n_i = 27; \quad k=10; \quad i=1,2,\dots,10; \quad J=1,2,\dots,n_i$$

$$\sum n_i X_i = 198,0; \quad \sum n_i X_i^2 = 1505,5; \quad (\sum n_i X_i)^2 = 39204$$

$$\sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij} = 177,2; \quad \sum_i \sum_j Y_{ij}^2 = 1257,54; \quad \left(\sum_i \sum_j Y_{ij} \right)^2 = 31399,84$$

$$\sum_i (T_{y \cdot i}^2 / n_i) = 1251,471; \quad \frac{(\sum \sum Y_{ij})^2}{n_i} = \frac{T_{y \cdot \cdot}^2}{n_i} = \frac{31399,84}{27} = 1162,957$$

$$\sum_i \sum_j X_i Y_{ij} = 1367,7; \quad \left(\sum_i n_i X_i \right) \left(\sum_j Y_{ij} \right) = 35085,6$$

$$S_r^2 = \frac{1}{26} \left(1505,50 - \frac{39204}{27} \right) = 2,0576; \quad S_r = 1,4344$$

$$S_y^2 = \frac{1}{26} \left(1257,54 - \frac{31399,84}{27} \right) = 3,6378; \quad S_y = 1,9073$$

$$S_{xy} = \frac{1}{26} \left(1367,7 - \frac{35085,6}{27} \right) = 2,6244$$

$$r = \frac{2,6244}{(1,4344)(1,9073)} = 0,959; \quad r^2 = 0,919; \quad (1-r^2) = 0,081$$

Varyans Analizi Tablosu

Varyans Kaynağı	Kareler Toplamı	Serbestiyet Derecesi	Kareler Ortalaması	
Grupları içi	1257,52 - 1251,471 = 6,049	17	0,357	$F_1 = \frac{9,8349}{0,357}$ = 27,55**
Regresyon	(N - 1) S _y ² r ² = 86,9218	1	86,9218	
Regresyon Etrafı	1,5922	8	0,199	
Gruplar arası	1251,471 - 1162,957 = 88,514	9	9,8349	

Hesaplanan F_1 değeri, (9) ve (17) serbestiyet dereceleri ve $\alpha = \%1$ için tablodan alınan F değerinden büyüktür. Bu durumda, bütün ortalamaların eşitliği hakkındaki varsayım red edilecektir.

Birinci varsayım red edildiği için, ikinci varsayımı kontrol edebiliriz;

$$F_2 = \frac{\text{Regresyon etrafında kareler ortalaması}}{\text{Gruplar içi kareler ortalaması}}$$

Ancak, gruplar içi kareler ortalaması, regresyon etrafında kareler ortalamasından büyüktür. Bu takdirde F_2

$$F_2 = \frac{0,257}{0,15} = 1,70$$

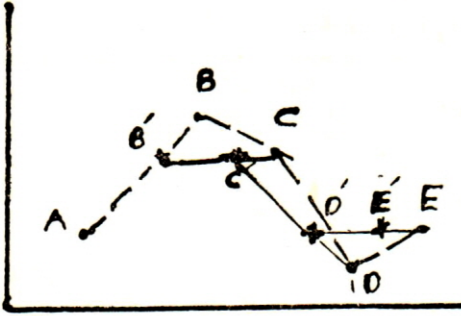
şeklinde hesaplanmalıdır.

Hesaplanan F_2 değeri, (17) ve (8) serbestiyet dereceleri ile $\alpha = \%1$ için bulunan tablo değeri 3,71'den küçüktür. Bu durumda, X ile Y arasındaki ilişkinin doğrusal olmadığı hakkında elimizde yeter delil yoktur. Söz konusu varsayım kabul edilecektir.

8. Regresyon Doğrusunun Geometrik Tayini

Regresyon doğrusunu karakterize eden (a) ve (b) değerlerinin bilinmesine ihtiyaç duyulmadığı hallerde, söz konusu doğruya geometrik olarak tayin mümkündür. Çizim metodu Şekil-5 yardımı ile açıklanmıştır:

Yapılan ölçmelerle A, B, C, D ve E gibi beş nokta elde edildiğini farz edelim. AB doğru parçası üzerinde ve AB uzunluğunun $(2/3)$ 'ündeki



Şekil - 5

nokta B' diyelim. B' ile C arasında ve gene bu uzaklığın $(2/3)$ 'ündeki nokta C'; C' ile D arasında ve benzer uzaklıktaki nokta D'; D' ile E arasında, aranan nokta E' olsun

E' noktası, regresyon doğrusunun geçmesi gereken noktalardan biridir. Doğrunun çizilebilmesi için tayini zorunlu diğer nokta aynı şekilde hareket edilerek, fakat bu sefer E',

D, C, B ve A noktaları için, bulunur. Elde edilecek son nokta B'' ile gösterilirse, regresyon doğrusu B'' E' doğrusudur.

AB'' uzaklığı ile E'E uzaklığı birbirine eşitse, çizimde kaba hata yapılmamış demektir.

FAYDALANILAN ESERLER

- 1 — ANDERSON, R. L., and T. A. BANCROFT, 1952, Statistical Theory in Research, McGraw-Hill, New-York.
- 2 — DRAPER, N. R., and H. SMITH, 1968, Applied Regression Analysis, John Wiley, New-York.
- 3 — JEROME, C. R. Li., 1966, Statistical Inference, Edwards Bothers, Ann Arbor, Mich.
- 4 — PRODAN, M., 1968, Forest Biometrics, Pergamon Press, London.
- 5 — SELBY, S. M., 1965, Standard Mathematical Tables, The Chemical Rubber Co., Ohio.
- 6 — SNEDECOR, G. W., 1953, Statistical Methods, State Colege Press, İowa.