

Fazekas István

STATISZTIKA

mobiDIÁK könyvtár

Fazekas István

STATISZTIKA

mobiDIÁK könyvtár

SOROZATSZERKESZTŐ

Fazekas István

Fazekas István

STATISZTIKA

Egyetemi jegyzet
Programtervező és alkalmazott matematikusok részére
Fejlesztés alatt álló változat!

mobiDIÁK könyvtár

Debreceni Egyetem
Informatikai Intézet

Lektor

Nagy Márta
Debreceni Egyetem
Lektorálás alatt!

Copyright © Fazekas István, 2004

Copyright © elektronikus közlés mobiDIÁK könyvtár, 2004

mobiDIÁK könyvtár
Debreceni Egyetem
Informatikai Kar
4010 Debrecen, Pf. 12
<http://mobidiak.inf.unideb.hu>

A mű egyéni tanulmányozás céljára szabadon letölthető. Minden egyéb felhasználás csak a szerző előzetes írásbeli engedélyével történhet.

A mű *A mobiDIÁK önszervező mobil portál* (IKTA, OMF-00373/2003) és a *GNU Iterátor, a legújabb generációs portál szoftver* (ITEM, 50/2003) projektek keretében készült.

Tartalomjegyzék

I. Szórásanalízis	9
1. A Fisher-Cochran-tétel	10
2. Az egyszeres osztályozás	18
3. Kétszeres osztályozás interakció nélkül	23
4. A kétszeres osztályozás az interakció figyelembevételével	27
II. Regresszióanalízis és a lineáris modell	33
1. A lineáris modell	34
2. A lineáris modell normális eloszlás esetén	39
III. A maximum-likelihood módszer	45
1. A Rao-Cramér-féle egyenlőtlenség	46
2. A maximum-likelihood becslés	50
3. A likelihood-hányados próba	59
IV. Nem-paraméteres módszerek	61
1. Hoeffding-féle U -statisztikák	62
2. A Mann-Whitney-féle U -próba	67
Tárgymutató	71

I. fejezet
Szórásanalízis

1. A Fisher-Cochran-tétel

A szórásanalízis alapvető tételei a Fisher-Cochran-tétel segítségével bizonyíthatók. Ebben a fejezetben a Fisher-Cochran-tételnek a nem-centrált χ^2 -eloszlásokat kezelő alakját igazoljuk. Ez alapján pontosan megadható majd a szórásfelbontásban szereplő kvadratikus formák eloszlása. Lineáris algebrai ismeretekkel kezdünk, hiszen minden az ortogonális mátrixokkal történő transzformációkon fog múlni.

1.1. Kvadratikus forma rangja és szabadsági foka

1.1. definíció. Legyen $\mathbf{x} = (x_1, \dots, x_n)^\top$ n -dimenziós vektor, $A = (a_{ij})$ pedig $n \times n$ -es (valós) szimmetrikus mátrix. Azt mondjuk, hogy a

$$Q = Q(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

kvadratikus forma rangja k , ha A rangja k . □

1.2. definíció. Legyen $\mathbf{b}_1, \dots, \mathbf{b}_m$ egy m tagú vektorrendszer \mathbb{R}^n -ben, $y_i = \mathbf{b}_i^\top \mathbf{x}$, $i = 1, \dots, m$, ahol $\mathbf{x} \in \mathbb{R}^n$. Azt mondjuk, hogy a $Q = \sum_{i=1}^m y_i^2 = \sum_{i=1}^m (\mathbf{b}_i^\top \mathbf{x})^2$ kvadratikus forma szabadsági foka k , ha a $\mathbf{b}_1, \dots, \mathbf{b}_m$ vektorrendszer rangja k . □

1.1. tétel. Legyen Q pozitív szemidefinit kvadratikus forma. Ekkor Q rangja megegyezik a szabadsági fokával.

BIZONYÍTÁS. 1. Legyen $Q = \mathbf{x}^\top A \mathbf{x}$, ahol A szimmetrikus, pozitív szemidefinit, $n \times n$ -es mátrix. Az, hogy Q rangja k , definíció szerint azt jelenti, hogy A rangja k . De ekkor A -nak k db pozitív sajátértéke van: $\lambda_1, \dots, \lambda_k$; a többi sajátértéke 0 . Az A mátrix előáll $A = U^\top \Lambda U$ alakban, ahol az U ortogonális mátrix sorai az A ortonormált sajátvektorai, a Λ diagonális mátrix főátlójában pedig az A sajátértékei állnak:

$$\Lambda = \begin{pmatrix} \lambda_1 & & & & & \\ & \ddots & & & & \\ & & \lambda_k & & & \\ & & & 0 & & \\ & & & & \ddots & \\ & & & & & 0 \end{pmatrix}.$$

Ekkor $\mathbf{b}_i^\top = \sqrt{\lambda_i} \mathbf{u}_i^\top$, $i = 1, \dots, k$ választással (ahol \mathbf{u}_i^\top az U mátrix i -edik sorvektora)

$$Q = \mathbf{x}^\top A \mathbf{x} = \mathbf{x}^\top U^\top \Lambda U \mathbf{x} = \sum_{i=1}^k \lambda_i (\mathbf{u}_i^\top \mathbf{x})^2 = \sum_{i=1}^k (\mathbf{b}_i^\top \mathbf{x})^2 = \sum_{i=1}^k y_i^2.$$

Itt a $\mathbf{b}_1 = \sqrt{\lambda_1} \mathbf{u}_1, \dots, \mathbf{b}_k = \sqrt{\lambda_k} \mathbf{u}_k$ vektorrendszer rangja k , hiszen az $\mathbf{u}_1, \dots, \mathbf{u}_k$ vektorok ortonormált rendszert alkotnak, és $\sqrt{\lambda_1} \neq 0, \dots, \sqrt{\lambda_k} \neq 0$.

2. Megfordítva, tegyük fel, hogy a $\mathbf{b}_1, \dots, \mathbf{b}_m$ vektorrendszer rangja k , és

$$Q = \sum_{i=1}^m (\mathbf{b}_i^\top \mathbf{x})^2 = \sum_{i=1}^m \mathbf{x}^\top \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x} = \mathbf{x}^\top \left(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top \right) \mathbf{x}.$$

Be kell látnunk, hogy a $\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top$ $n \times n$ -es mátrix rangja k .

Legyen a $\mathbf{b}_1, \dots, \mathbf{b}_m$ vektorrendszer által kifeszített altér L . Ezen altér k -dimenziós. Belátjuk, hogy a $(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top)$ (ún. diádösszeg) mátrix képtere éppen L . Mivel tetszőleges \mathbf{x} vektorra

$$\left(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top \right) \mathbf{x} = \sum_{i=1}^m \mathbf{b}_i (\mathbf{b}_i^\top \mathbf{x}) \in L,$$

így $(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top)$ képtere benne van L -ben. Legyen most \mathbf{v} a $(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top)$ nullteréből való. Ekkor

$$0 = \mathbf{v}^\top \left(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top \right) \mathbf{v} = \sum_{i=1}^m (\mathbf{b}_i^\top \mathbf{v})^2$$

miatt $\mathbf{v} \perp \mathbf{b}_i$, minden i -re. Tehát \mathbf{v} eleme L ortogonális komplementerének. Viszont mind a képtér és a nulltér dimenziójának összege, mind az altér és ortogonális komplementere dimenziójának összege n . Azaz $(\sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^\top)$ képtere éppen L . \square

1.1. megjegyzés. Mivel egy kvadratikus forma egyértelműen meghatározza a szimmetrikus mátrixát, így a kvadratikus forma rangja egyértelműen definiált. Viszont egy kvadratikus forma többféleképpen előállítható lineáris formák négyzetösszegeként. Tehát a szabadsági fok egyértelműségét nem közvetlenül a definícióból, hanem a szabadsági fok és a rang egyenlőségét kimondó tételből olvashatjuk ki. \square

1.1. példa. A

$$Q = Q(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x} = (x_1, x_2) \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 2x_1^2 + 2x_1x_2 + x_2^2$$

kvadratikus forma rangja 2, mivel a mátrixának rangja 2.

Állítsuk elő Q -t lineáris formák négyzetösszegeként kétféle módon. Legyen először $\mathbf{b}_1 = (1, 1)^\top$, $\mathbf{b}_2 = (1, 0)^\top$. Nyilván

$$(\mathbf{b}_1^\top \mathbf{x})^2 + (\mathbf{b}_2^\top \mathbf{x})^2 = (x_1 + x_2)^2 + x_1^2 = Q(\mathbf{x}).$$

Másrészt, $\mathbf{c}_1 = (\sqrt{2}, 1/\sqrt{2})^\top$ és $\mathbf{c}_2 = (0, 1/\sqrt{2})^\top$ választással

$$(\mathbf{c}_1^\top \mathbf{x})^2 + (\mathbf{c}_2^\top \mathbf{x})^2 = (\sqrt{2}x_1 + x_2/\sqrt{2})^2 + (x_2/\sqrt{2})^2 = Q(\mathbf{x}).$$

Mind a $\mathbf{b}_1, \mathbf{b}_2$, mind a $\mathbf{c}_1, \mathbf{c}_2$ vektorrendszer rangja 2. \square

1.2. A kvadratikus forma négyzetösszeg alakja

1.2. megjegyzés. Tegyük fel, hogy a Q kvadratikus forma előáll lineáris formák négyzetösszegeként a

$$Q = (\mathbf{b}_1^\top \mathbf{y})^2 + \dots + (\mathbf{b}_s^\top \mathbf{y})^2$$

és a

$$Q = (\mathbf{c}_1^\top \mathbf{y})^2 + \dots + (\mathbf{c}_r^\top \mathbf{y})^2$$

alakban is. Ekkor a $\mathbf{b}_1, \dots, \mathbf{b}_s$ vektorrendszer által kifeszített altér megegyezik a $\mathbf{c}_1, \dots, \mathbf{c}_r$ vektorrendszer által kifeszítetttel.

Ennek bebizonyítására tekintsük a $Q = \mathbf{y}^\top (\sum_{i=1}^s \mathbf{b}_i \mathbf{b}_i^\top) \mathbf{y}$ és a $Q = \mathbf{y}^\top (\sum_{i=1}^r \mathbf{c}_i \mathbf{c}_i^\top) \mathbf{y}$ előállításokat. Ezeknek a mátrixa megegyezik. Tehát $\sum_{i=1}^s \mathbf{b}_i \mathbf{b}_i^\top = \sum_{i=1}^r \mathbf{c}_i \mathbf{c}_i^\top$. Viszont az első mátrix képtere a $\mathbf{b}_1, \dots, \mathbf{b}_s$ vektorrendszer által kifeszített altér, míg a másodiké a $\mathbf{c}_1, \dots, \mathbf{c}_r$ vektorrendszer által kifeszített. \square

1.2. példa. Állítsuk elő az $y_1^2 + y_2^2$ négyzetösszeget két kvadratikus forma összegeként:

$$y_1^2 + y_2^2 = Q_1 + Q_2 = \frac{1}{2}(y_1 + y_2)^2 + \frac{1}{2}(y_1 - y_2)^2.$$

Itt $Q_1 = (y_1/\sqrt{2} + y_2/\sqrt{2})^2$ és $Q_2 = (y_1/\sqrt{2} - y_2/\sqrt{2})^2$. A Q_1 és a Q_2 kvadratikus formák az $\mathbf{y} = (y_1, y_2)^\top$ vektor lineáris formáinak négyzetösszegei:

$$Q_1 = z_1^2, \quad \text{ahol } z_1 = (1/\sqrt{2}, 1/\sqrt{2}) (y_1, y_2)^\top = \mathbf{b}_1^\top \mathbf{y},$$

$$Q_2 = z_2^2, \quad \text{ahol } z_2 = (1/\sqrt{2}, -1/\sqrt{2}) (y_1, y_2)^\top = \mathbf{b}_2^\top \mathbf{y}.$$

A $\mathbf{z} = (z_1, z_2)^\top$ vektorra

$$\mathbf{z} = A\mathbf{y} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix},$$

tehát az új z_1, z_2 változókat a régi y_1, y_2 változókból az A ortogonális mátrixszal való transzformációval kapjuk meg. \square

A fenti tény sokkal általánosabban is érvényes.

1.1. lemma. Legyen $\mathbf{y} = (y_1, \dots, y_n)^\top$. Legyenek Q_1, \dots, Q_s az y_1, \dots, y_n változók kvadratikus formái. Tegyük fel, hogy $\text{rang}(Q_i) = n_i, i = 1, \dots, s$ és

$$\sum_{i=1}^n y_i^2 = Q_1 + \dots + Q_s.$$

Akkor és csak akkor létezik olyan A ortogonális mátrix, melyre $\mathbf{z} = A\mathbf{y}, \mathbf{z} = (z_1, \dots, z_n)^\top$ jelöléssel

$$Q_1 = z_1^2 + \dots + z_{n_1}^2, \quad Q_2 = z_{n_1+1}^2 + \dots + z_{n_1+n_2}^2, \dots,$$

$$Q_s = z_{n_1+\dots+n_{s-1}+1}^2 + \dots + z_{n_1+\dots+n_s}^2$$

teljesül, ha $n_1 + n_2 + \dots + n_s = n$.

BIZONYÍTÁS. 1. Tegyük fel, hogy létezik a fenti A mátrix. Ekkor $A^\top A = I$ miatt

$$\sum_{i=1}^n z_i^2 = \mathbf{z}^\top \mathbf{z} = \mathbf{y}^\top A^\top A \mathbf{y} = \mathbf{y}^\top \mathbf{y} = \sum_{i=1}^n y_i^2 = \sum_{j=1}^s Q_j = \sum_{i=1}^{n_1+\dots+n_s} z_i^2$$

minden $\mathbf{z} \in \mathbb{R}^n$ -re. Ezért $n = n_1 + \dots + n_s$.

2. Tegyük fel, hogy $n_1 + \dots + n_s = n$. Mivel Q_1 n_1 rangú kvadratikus forma, így „teljes négyzetté alakításában” n_1 db lineáris forma szerepel:

$$Q_1 = \sum_{i=1}^{n_1} \delta_i z_i^2,$$

ahol $\delta_i = \pm 1$, $z_i = \mathbf{b}_i^\top \mathbf{y}$, $i = 1, \dots, n_1$. Hasonló igaz Q_2, \dots, Q_s -re is. Mivel $n_1 + \dots + n_s = n$, így a Q_i kvadratikus formák előállításában szereplő \mathbf{b}_j^\top n -dimenziós sorvektorokat mind egymás alá írva, egy $n \times n$ -es mátrixot kapunk:

$$A = \begin{pmatrix} \mathbf{b}_1^\top \\ \vdots \\ \mathbf{b}_{n_1}^\top \\ \vdots \\ \mathbf{b}_n^\top \end{pmatrix}.$$

Jelölje most D a δ_i diagonális elemekből álló diagonális mátrixot (δ_i -ket ugyanolyan sorrendben írjuk le, mint a megfelelő \mathbf{b}_i^\top vektorokat az A mátrixban). Ekkor

$$\mathbf{y}^\top A^\top D A \mathbf{y} = \mathbf{z}^\top D \mathbf{z} = \sum_{i=1}^n \delta_i z_i^2 = \sum_{j=1}^s Q_j = \sum_{i=1}^n y_i^2 = \mathbf{y}^\top \mathbf{y}$$

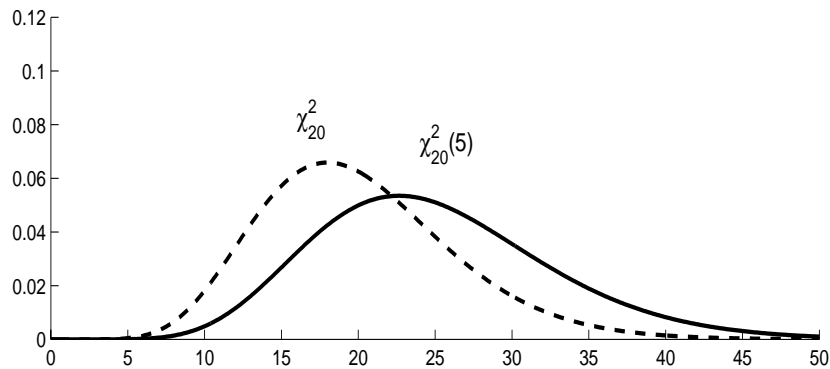
minden $\mathbf{y} \in \mathbb{R}^n$ -re.

Mivel egy kvadratikus forma egyértelműen meghatározza a szimmetrikus mátrixát, ezért $A^\top D A = I$. Ebből adódik, hogy A invertálható, és $D = (A^\top)^{-1} A^{-1} = (A^{-1})^\top A^{-1}$. Tehát D egy pozitív szemidefinit mátrixszal egyenlő, ami alapján a főátlójában szereplő δ_i -k egyike sem lehet -1 , azaz $D = I$. De ekkor $A^\top A = I$, azaz A ortogonális. Továbbá

$$Q_1 = \sum_{i=1}^{n_1} z_i^2,$$

és hasonló áll fenn a többi Q_j -re is. □

1.3. A nem-centrált khi-négyzet eloszlás



ÁBRA 1.1. χ_{20}^2 és $\chi_{20}^2(5)$ sűrűségfüggvénye

1.3. definíció. Legyenek ξ_1, \dots, ξ_n független, normális eloszlású valószínűségi változók: $\xi_i \sim \mathcal{N}(a_i, 1)$, $i = 1, \dots, n$. Ekkor a

$$(1.1) \quad \zeta_n = \xi_1^2 + \dots + \xi_n^2$$

valószínűségi változót n szabadsági fokú, $\lambda = \sum_{i=1}^n a_i^2$ nem-centralitási paraméterű **nem-centráltnégyzet eloszlásának** nevezzük. Jelölése $\chi_n^2(\lambda)$. \square

A nem-centráltnégyzet eloszlás abban a speciális esetben, amikor a kiinduló ξ_i valószínűségi változók 0 várható értékűek, éppen a korábban megismert (centráltnégyzet eloszlás. Azaz $\chi_n^2(0) \equiv \chi_n^2$.

A χ_{20}^2 és $\chi_{20}^2(5)$ sűrűségfüggvénye a 1.1. ábrán látható.

A következő állítás azt mutatja, hogy nem kell a ξ_i -k a_i várható értékeit külön-külön ismerni, az eloszlás csak a $\lambda = \sum_{i=1}^n a_i^2$ nem-centralitási paramétertől függ.

1.1. állítás. Legyenek τ_1, \dots, τ_n független, normális eloszlású valószínűségi változók: $\tau_1 \sim \mathcal{N}(\sqrt{\lambda}, 1)$, $\tau_2 \sim \mathcal{N}(0, 1), \dots, \tau_n \sim \mathcal{N}(0, 1)$, ahol $\lambda \geq 0$. Ekkor

$$(1.2) \quad \tau_1^2 + \dots + \tau_n^2$$

eloszlása megegyezik az (1.1)-ben adott ζ_n eloszlásával, bármilyenek is a $\lambda = \sum_{i=1}^n a_i^2$ feltételt kielégítő a_1, \dots, a_n számok.

BIZONYÍTÁS. Legyen a $B = (b_{ij})$ $n \times n$ -es ortogonális mátrix első sora:

$$b_{11} = a_1/\sqrt{\lambda}, \dots, b_{1n} = a_n/\sqrt{\lambda}.$$

Ilyen B mátrix létezik, hisz az első sorra előírt vektor egységnyi hosszúságú, ez kiegészíthető ortonormált bázissá, és B sorai a bázisvektorok lesznek.

Legyen $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top$, ahol ξ_1, \dots, ξ_n a 1.3. definícióbeliek. Legyen $\boldsymbol{\tau} = B\boldsymbol{\xi}$. Mivel $\boldsymbol{\xi}$ eloszlása $\mathcal{N}_n(\mathbf{m}, I)$, ahol $\mathbf{m} = (a_1, \dots, a_n)^\top$, így $\boldsymbol{\tau} \sim \mathcal{N}_n(B\mathbf{m}, I)$. Ezért $\boldsymbol{\tau}$ koordinátái, τ_1, \dots, τ_n , együttesen normális eloszlású és korrelálatlan valószínűségi változók. Így függetlenek is. A várható értékeik pedig a $B\mathbf{m}$ vektor koordinátái. Azaz $\mathbb{E}\tau_i = \mathbf{b}_i^\top \mathbf{m}$, $i = 1, \dots, n$, ahol \mathbf{b}_i^\top a B mátrix i -edik sora. De $\mathbf{b}_1 = \lambda^{-\frac{1}{2}}\mathbf{m}$, továbbá $\mathbf{b}_2, \dots, \mathbf{b}_n$ erre merőleges. Ezért

$$\mathbb{E}\tau_1 = \lambda^{-\frac{1}{2}}\mathbf{m}^\top \mathbf{m} = \lambda^{-\frac{1}{2}} \sum_{i=1}^n a_i^2 = \lambda^{-\frac{1}{2}}\lambda = \sqrt{\lambda}.$$

$$\mathbb{E}\tau_i = \mathbf{b}_i^\top \mathbf{b}_1 \sqrt{\lambda} = 0, \quad i = 2, \dots, n.$$

Tehát τ_1, \dots, τ_n kielégíti a tétel feltételeit. Másrészt

$$\sum_{i=1}^n \xi_i^2 = \boldsymbol{\xi}^\top \boldsymbol{\xi} = \boldsymbol{\xi}^\top B^\top B \boldsymbol{\xi} = \boldsymbol{\tau}^\top \boldsymbol{\tau} = \sum_{i=1}^n \tau_i^2.$$

Tehát az (1.2) előállítás teljesül. \square

1.2. állítás. A nem-centráltnégyzet $\chi_n^2(\lambda)$ khi-négyzet eloszlás karakterisztikus függvénye:

$$\varphi(t) = (1 - 2it)^{-n/2} \exp[it\lambda(1 - 2it)^{-1}];$$

várható értéke:

$$\mathbb{E}\zeta_n = n + \lambda;$$

szórásnégyzete:

$$\mathbb{D}^2\zeta_n = 2(n + 2\lambda);$$

ferdesége:

$$\beta_1(\zeta_n) = \frac{\sqrt{8}(n + 3\lambda)}{(n + 2\lambda)^{3/2}};$$

lapultsága:

$$\beta_2(\zeta_n) = \frac{12(n+4\lambda)}{(n+2\lambda)^2}.$$

BIZONYÍTÁS. Legyen $\xi \sim \mathcal{N}(m, 1)$. Ekkor ξ^2 karakterisztikus függvénye:

$$\begin{aligned} \varphi_{\xi^2}(t) &= \int_{-\infty}^{+\infty} e^{itx^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2}} dx = \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left(x\sqrt{1-2it} - \frac{m}{\sqrt{1-2it}} \right)^2 \right\} dx \exp \left\{ \frac{itm^2}{1-2it} \right\} = \\ &= \frac{1}{\sqrt{1-2it}} \exp \left\{ \frac{itm^2}{1-2it} \right\}, \end{aligned}$$

ahol az utolsó lépésben az

$$y = x\sqrt{1-2it} - \frac{m}{\sqrt{1-2it}}$$

helyettesítést végeztük el. Az így kapott karakterisztikus függvényekből szorzással adódik $\chi_n^2(\lambda)$ karakterisztikus függvénye.

Ha $\xi \sim \mathcal{N}(m, \sigma^2)$, akkor

$$\mathbb{E}\xi^2 = \sigma^2 + m^2, \quad \mathbb{D}^2\xi^2 = 2\sigma^4 + 4\sigma^2m^2$$

(ez utóbbit a standard normálisra visszavezetve kaphatjuk). Ezután alkalmazzuk $\chi_n^2(\lambda)$ definícióját! \square

1.2. tétel. (χ^2 addíciós tétel.) Legyenek η_n és η_m független χ^2 -eloszlású valószínűségi változók n , ill. m szabadsági fokkal, valamint λ_1 , ill. λ_2 nem-centralitási paraméterrel. Ekkor az $\eta_n + \eta_m$ változó $n+m$ szabadsági fokú, $\lambda_1 + \lambda_2$ nem-centralitási paraméterű χ^2 -eloszlású.

A tétel szavakban kifejezve: független χ^2 -ek összege χ^2 , a szabadsági fokok és a nem-centralitási paraméterek pedig összeadódnak.

BIZONYÍTÁS. Legyenek ξ_1, \dots, ξ_{n+m} független normális eloszlásúak: $\xi_i \sim \mathcal{N}(a_i, 1)$, $i = 1, \dots, n+m$. η_n -et $\xi_1^2 + \dots + \xi_n^2$ -ként, η_m -et pedig $\xi_{n+1}^2 + \dots + \xi_{n+m}^2$ -ként reprezentálva,

$$\eta_n + \eta_m = \xi_1^2 + \dots + \xi_{n+m}^2$$

adódik, ez pedig $n+m$ szabadsági fokú, $\lambda_1 + \lambda_2$ nem-centralitási paraméterű khi-négyzet eloszlású. \square

1.4. definíció. Legyenek ξ és η független χ^2 -eloszlású valószínűségi változók: $\xi_i \sim \chi_n^2(\lambda)$, $\eta_i \sim \chi_m^2$. Ekkor a

$$(1.3) \quad \frac{\xi}{n} / \frac{\eta}{m}$$

valószínűségi változót n és m szabadsági fokú, λ nem-centralitási paraméterű **nem-centrál** F -eloszlásúnak nevezzük. \square

A nem-centrál F -eloszlás abban a speciális esetben, amikor a kiinduló ξ valószínűségi változó is centrál, éppen a korábban megismert (centrál) F -eloszlás.

1.4. A Fisher-Cochran-tétel

1.3. tétel. Fisher-Cochran. Legyenek ξ_1, \dots, ξ_n független valószínűségi változók, $\xi_i \sim \mathcal{N}(a_i, 1)$, $i = 1, \dots, n$. Legyenek Q_1, \dots, Q_s a ξ_1, \dots, ξ_n kvadratikus formái, $\text{rang}(Q_i) = n_i$, $i = 1, \dots, s$. Tegyük fel, hogy

$$\xi_1^2 + \dots + \xi_n^2 = Q_1 + \dots + Q_s.$$

Ekkor Q_1, \dots, Q_s akkor és csak akkor függetlenek és $\chi_{n_1}^2, \dots, \chi_{n_s}^2$ eloszlásúak, ha $n_1 + \dots + n_s = n$.

Ebben az esetben Q_i nem-centralitási paramétere: $Q_i(\mathbf{a})$, ahol $\mathbf{a} = (a_1, \dots, a_n)^\top$.

BIZONYÍTÁS. 1. Legyenek Q_1, \dots, Q_s függetlenek, Q_i eloszlása $\chi_{n_i}^2$, $i = 1, \dots, s$. Ekkor a χ^2 -addíciós tétel alapján $Q_1 + \dots + Q_s$ eloszlása $\chi_{n_1 + \dots + n_s}^2$. De $Q_1 + \dots + Q_s = \xi_1^2 + \dots + \xi_n^2$, aminek az eloszlása definíció szerint χ_n^2 . Tehát $n_1 + \dots + n_s = n$. Megjegyezzük, hogy itt nem volt külön szükség a $\text{rang}(Q_i) = n_i$ feltételre.

2. Legyen most $n_1 + \dots + n_s = n$. Az 1.1. lemma miatt létezik olyan A ortogonális mátrix, hogy az $\boldsymbol{\eta} = A\boldsymbol{\xi}$ -re (ahol $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n)^\top$, $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top$)

$$Q_1 = \eta_1^2 + \dots + \eta_{n_1}^2, \quad Q_2 = \eta_{n_1+1}^2 + \dots + \eta_{n_1+n_2}^2, \dots, \\ Q_s = \eta_{n_1+\dots+n_{s-1}+1}^2 + \dots + \eta_{n_1+\dots+n_s}^2.$$

Viszont $\boldsymbol{\xi} \sim \mathcal{N}_n(\mathbf{a}, I)$, ezért $\boldsymbol{\eta} = A\boldsymbol{\xi} \sim \mathcal{N}_n(A\mathbf{a}, AA^\top) = \mathcal{N}_n(A\mathbf{a}, I)$, hiszen A ortogonális mátrix.

Tehát $\boldsymbol{\eta}$ koordinátái független, 1 szórású normálisak. Q_i pedig n_i db ilyenek a négyzetösszege, tehát $\chi_{n_i}^2$ eloszlású. Mivel különböző Q_i -k előállításában azonos η_j -k nem vesznek részt, így Q_i -k függetlenek.

3. A nem-centralitási paraméter. A lemmában adott előállítás $Q_1 = z_1^2 + \dots + z_{n_1}^2$; részletesen kiírva:

$$Q_1(\mathbf{y}) = (\mathbf{b}_1^\top \mathbf{y})^2 + \dots + (\mathbf{b}_{n_1}^\top \mathbf{y})^2, \quad \mathbf{y} \in \mathbb{R}^n,$$

ahol $\mathbf{b}_1^\top, \dots, \mathbf{b}_{n_1}^\top$ az A mátrix sorvektorai. Ezt a $\boldsymbol{\xi}$ várható érték vektora, azaz $\mathbf{y} = \mathbf{a}$ helyén véve:

$$Q_1(\mathbf{a}) = (\mathbf{b}_1^\top \mathbf{a})^2 + \dots + (\mathbf{b}_{n_1}^\top \mathbf{a})^2 = (\mathbb{E}\eta_1)^2 + \dots + (\mathbb{E}\eta_{n_1})^2.$$

Ez utóbbi pedig, a $Q_1 = \eta_1^2 + \dots + \eta_{n_1}^2$ alapján, éppen Q_1 nem-centralitási paramétere. Hasonló áll Q_2, \dots, Q_s -re. \square

1.3. megjegyzés. A fenti bizonyításból kiderült, hogy a Fisher-Cochran-tétel feltételei esetén a $\xi_1^2 + \dots + \xi_n^2$ teljes négyzetösszeg előáll a ξ_1, \dots, ξ_n lineáris formái négyzetösszegeként:

$$\xi_1^2 + \dots + \xi_n^2 = (\mathbf{b}_1^\top \boldsymbol{\xi})^2 + \dots + (\mathbf{b}_n^\top \boldsymbol{\xi})^2,$$

ahol az $\eta_i = \mathbf{b}_i^\top \boldsymbol{\xi}$, $i = 1, \dots, n$, lineáris formák függetlenek, és $\eta_i \sim \mathcal{N}(\mathbf{b}_i^\top \mathbf{a}, 1)$. Sőt, Q_j ezek közül az első n_1 db, Q_2 a következő n_2 db, ... négyzetösszege.

Ha Fisher-Cochran-tétel feltételei mellett létezik a Q_i kvadratikus formáknak valamilyen lineáris formák négyzetösszegeként történő előállítása, akkor a különböző Q_i -k előállításában szereplő formák függetlenek.

Ezt bizonyítjuk arra az esetre, ha Q_1 előáll a $Q_1 = (\mathbf{b}_1^\top \boldsymbol{\xi})^2 + \dots + (\mathbf{b}_{n_1}^\top \boldsymbol{\xi})^2$ mellett $Q_2 = (\mathbf{c}_1^\top \boldsymbol{\xi})^2 + \dots + (\mathbf{c}_m^\top \boldsymbol{\xi})^2$ alakban is. Egy előző megjegyzés szerint a $\mathbf{c}_1, \dots, \mathbf{c}_m$ által kifeszített altér megegyezik a $\mathbf{b}_1, \dots, \mathbf{b}_{n_1}$ által kifeszítettel. De az utóbbi ortogonális a $\mathbf{b}_{n_1+1}, \dots, \mathbf{b}_n$ által kifeszített altérre. Ezért $\mathbf{c}_i^\top \mathbf{b}_j = 0$, amiből

$$\text{cov}(\mathbf{c}_i^\top \boldsymbol{\xi}, \mathbf{b}_j^\top \boldsymbol{\xi}) = \mathbf{c}_i^\top I \mathbf{b}_j = \mathbf{c}_i^\top \mathbf{b}_j = 0,$$

ha $1 \leq i \leq m$, $n_1 + 1 \leq j \leq n$. A $(\mathbf{c}_1^\top \boldsymbol{\xi}, \dots, \mathbf{c}_m^\top \boldsymbol{\xi})^\top$ és a $(\mathbf{b}_{n_1+1}^\top \boldsymbol{\xi}, \dots, \mathbf{b}_n^\top \boldsymbol{\xi})^\top$ vektorok együttesen normális eloszlásúak, korrelálatlanok, tehát függetlenek. Ezért pl. $\mathbf{c}_1^\top \boldsymbol{\xi}$ és Q_2 független. \square

2. Az egyszeres osztályozás

Ebben a szakaszban a Fisher-Cochran-tételt alkalmazzuk az egyszeres osztályozásra.

2.1. Az egyszeres osztályozás feladata

A szórásanalízis (ANOVA=ANalysis Of VARIance) alapkérdése: több minta esetén a várható értékek egyenlők-e. Alapvető feltétel: minták egymástól is függetlenek, normális eloszlásból származnak, és a szórásaik egyenlők! Tehát a minták között csak a várható értékekben lehet eltérés.

A legegyszerűbb szórásanalízisbeli modell az egyszeres osztályozás (one-way classification, one-way layout). Itt egyetlen **tényező szintjeit** kell összehasonlítani. Mivel a megfigyelések eredményeit tényezőnként egy-egy oszlopban szokták elhelyezni, a tényezők szintjeinek hatását oszlophatásnak nevezzük. Példaként tekintsünk egy fogyókúra kísérletet.

diéta tényező		
1. diéta	2. diéta	3. diéta
9.40	22.84	17.35
9.48	15.32	16.36
7.56	11.04	15.88
11.52	17.92	14.28
11.56	19.68	18.60
12.12	26.20	19.32
11.36		14.20
4.60		17.52
14.48		

2.1. példa. Három különböző diéta hatását mérték 9, 6, ill. 8 kísérleti alanyon. Itt az egyetlen tényező a diéta, annak 3 szintje van. A diéta hatására a súlyvesztéseket a fenti táblázat mutatja. Vizsgáljuk meg azt a nullhipotézist, hogy a diéták által okozott súlyvesztések várható értékei egyenlők! \square

A megfigyelések: Y_{ij} , $j = 1, \dots, n_i$, $i = 1, \dots, p$. Y_{ij} az i -edik szinten végzett j -edik megfigyelés. Az egyes szinteken nem feltétlen kell azonos számú mérést végezni.

Feltesszük hogy,

$$Y_{ij} \sim \mathcal{N}(\mu + \alpha_i, \sigma^2) \text{ és } Y_{ij}\text{-k függetlenek.}$$

Elérhető, hogy $\sum_{i=1}^p n_i \alpha_i = 0$ legyen. Vezessük be az $n = n_1 + \dots + n_p$ jelölést.

Vizsgáljuk a

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p = 0$$

nullhipotézis teljesülését! H_0 azt jelenti, hogy az egyes szinteknek nincs hatása.

A Steiner-formula alapján az n db Y_{ij} négyzetösszege előáll

$$\sum_{i=1}^p \sum_{j=1}^{n_i} Y_{ij}^2 = \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 + n\bar{Y}_{..}^2$$

alakban, ahol $\bar{Y}_{..} = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^{n_i} Y_{ij}$ a **teljes átlag**. A fenti felbontásban szereplő első négyzetösszeg jelölése Q , elnevezése **teljes négyzetösszeg**. Q előáll

$$Q = Q_1 + Q_2 = \sum_{i=1}^p \sum_{j=1}^{n_i} (\bar{Y}_i - \bar{Y}_{..})^2 + \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$$

alakban, ahol $\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$ az i -edik **szint átlaga**.

Q_1 méri a **szintek közötti szóródást**, Q_2 pedig a **szinteken belüli szóródást** (azaz a véletlen hibát). H_0 -at akkor vetjük el, ha Q_1 túlságosan nagy Q_2 -höz képest.

2.2. A Fisher-Cochran-tétel alkalmazása az egyszeres osztályozásra

A próba konstruálására használjuk a Fisher-Cochran-tételt. Tekintsük a teljes felbontást:

$$\sum_{i=1}^p \sum_{j=1}^{n_i} Y_{ij}^2 = Q_1 + Q_2 + n\bar{Y}_{..}^2.$$

2.1. állítás. Q_1 szabadsági foka $p - 1$, Q_2 szabadsági foka $n - p$, $n\bar{Y}_{..}^2$ szabadsági foka 1.

BIZONYÍTÁS. A kvadratikus formák itt lineáris formák négyzetösszegeiként vannak előállítva. Rendezzük a lineáris formák együtthatóit

$$Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{p1}, \dots, Y_{pn_p}$$

szerint.

Ekkor a Q_1 -hez tartozó együtthatók mátrixa (valójában az azonos sorok közül 1-1-et megtartva):

$$\begin{pmatrix} \frac{1}{n_1} - \frac{1}{n} & \cdots & \frac{1}{n_1} - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} & \cdots & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & \cdots & -\frac{1}{n} & \frac{1}{n_2} - \frac{1}{n} & \cdots & \frac{1}{n_2} - \frac{1}{n} & \cdots & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ -\frac{1}{n} & \cdots & -\frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} & \cdots & \frac{1}{n_p} - \frac{1}{n} & \cdots & \frac{1}{n_p} - \frac{1}{n} \end{pmatrix}.$$

Ezen p sor lineárisan függő, hisz n_1, \dots, n_p együtthatókkal belőlük 0 kombinálható. De az elsőt kivonva az összes többiből, $p - 1$ lineárisan független sort kapunk. Azaz Q_1 szabadsági foka $p - 1$.

A Q_2 -ben szereplő lineáris formák együtthatói olyan

$$D = \begin{pmatrix} D_1 & & \\ & \ddots & \\ & & D_p \end{pmatrix}$$

mátrixot alkotnak, melyben a diagonálison fekvő D_1, \dots, D_p blokkokon kívüli elemek 0-k. A D_i blokk:

$$D_i = \begin{pmatrix} 1 - \frac{1}{n_i} & -\frac{1}{n_i} & \dots & -\frac{1}{n_i} \\ -\frac{1}{n_i} & 1 - \frac{1}{n_i} & \dots & -\frac{1}{n_i} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n_i} & -\frac{1}{n_i} & \dots & 1 - \frac{1}{n_i} \end{pmatrix}.$$

Ez az $n_i \times n_i$ méretű mátrix $n_i - 1$ rangú. Ugyanis sorainak összege 0, de az első sort az összes többiből kivonva, $n_i - 1$ db lineárisan független sort kapunk. Ezek alapján Q_2 szabadsági foka $(n_1 - 1) + \dots + (n_p - 1) = n - p$.

$n\bar{Y}_{..}^2$ szabadsági foka nyilván 1. □

2.2. állítás. Q_1 és Q_2 függetlenek. $\frac{1}{\sigma^2}Q_1 \sim \chi_{p-1}^2(\lambda_1)$, $\frac{1}{\sigma^2}Q_2 \sim \chi_{n-p}^2(0)$, ahol $\lambda_1 = \frac{1}{\sigma^2} \sum_{i=1}^p n_i \alpha_i^2$.

$\frac{1}{\sigma^2}Q_1$ akkor és csak akkor centrált χ^2 -eloszlású, ha a

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p = 0$$

nullhipotézis teljesül.

BIZONYÍTÁS. Az előző állítás miatt a Fisher-Cochran-tétel feltétele teljesül, hiszen a rangok összege $(p - 1) + (n - p) + 1 = n$. Így a kvadratikus formák függetlenek és χ^2 -eloszlásúak. A nem-centralitási paramétereket úgy kapjuk, hogy a várható értékeket beírjuk a kvadratikus formákba. Mivel $\mathbb{E}(Y_{ij} - \bar{Y}_{i.}) = 0$, így Q_2 nem-centralitási paramétere 0. Viszont $\mathbb{E}(\bar{Y}_{i.} - \bar{Y}_{..}) = \alpha_i$, ezért $\frac{1}{\sigma^2}Q_1$ nem-centralitási paramétere: $\lambda_1 = \frac{1}{\sigma^2} \sum_{i=1}^p n_i \alpha_i^2$. □

2.3. A szórásfelbontó táblázat

A 2.2. állítás alapján a Q_2 véletlen hiba (skalárszorosa) mindig centrált χ^2 -eloszlású, Q_1 (skalárszorosa) viszont pontosan akkor, ha H_0 teljesül.

2.3. állítás. Az

$$F = \frac{Q_1}{p - 1} \bigg/ \frac{Q_2}{n - p}$$

statisztika pontosan akkor $p - 1$ és $n - p$ szabadsági fokú centrált F -eloszlású, ha a

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p = 0$$

nullhipotézis teljesül.

BIZONYÍTÁS. Az 2.2. állítás szerint F két független, a szabadsági fokával elosztott, (H_0 esetén centrált) χ^2 -eloszlású valószínűségi változó hányadosa. Ezért az F -eloszlás definíciója alapján ennek eloszlása (H_0 esetén centrált) F -eloszlás $p - 1$ és $n - p$ szabadsági fokokkal. □

Az eddigiek alapján az alábbi szórásfelbontó táblázatot adhatjuk meg az egytényezős osztályozásra.

A szóródás forrásai	Szabadsági fokok	Négyzet-összegek	Négyzet-átlagok	F -hányados
Részsokaságok (szintek) közötti eltérések	$p - 1$	Q_1	$\frac{Q_1}{p-1}$	$F = \frac{n-p}{p-1} \cdot \frac{Q_1}{Q_2}$
Szinteken belüli eltérések (véletlen hibák)	$n - p$	Q_2	$\frac{Q_2}{n-p}$	
Teljes	$n - 1$	$Q_1 + Q_2$		

Szórásfelbontó táblázat az egytényezős osztályozásra

H_0 -at $1-\alpha$ szinten elvetjük, ha a kapott F -statisztika értéke nagyobb, mint $F_{[p-1, n-p; \alpha]}$, azaz a megfelelő szabadsági fokú centrált F -eloszlás táblázatából kikeresett (felső) kritikus érték.

2.2. példa. (A 2.1. példa folytatása.)

A (számítógépes) eredményt a szórásfelbontó tábla tartalmazza:

ANOVA Table				
Source	df	SS	MS	F
Columns	2	313.9	156.90	13.31
Error	20	235.8	11.79	
Total	22	549.7		

Az elnevezések magyarázata. Source = a szóródás forrása; Columns = oszlophatás (szintek közötti eltérések); Error = véletlen hiba; Total = teljes négyzetösszeg; df (degree of freedom) = szabadsági fok; SS (Sum of Squares) = négyzetösszeg; MS (Mean Square) = tapasztalati szórásnégyzet (négyzet átlag), F = F -statisztika.

Annak kérdéséről, hogy a diéta három szintjének van-e hatása, az F alatti mennyiség alapján döntünk. Amennyiben H_0 : a tényező szintjeinek nincs hatása nullhipotézis teljesül, az F alatti statisztika centrált F -eloszlású (jelenleg (2, 20) szabadsági fokkal). Ez alapján határozható meg a próba pontos terjedelme: p . Példánkban $p = 0.00021$ érték adódott, azaz minden használatos szinten elvetjük a diéták egyforma hatását.

A hagyományos (táblázatos) kiértékelés ugyanerre a következtetésre vezet. F értékét összehasonlítva a (2, 20) szabadsági fokú centrált F -eloszlás $F_{[2, 20; 0.05]} = 3.49$ kritikus értékével, azt kapjuk, hogy a H_0 nullhipotézist 95%-os szinten el kell vetni. Ez azt jelenti, hogy a diéta tényező különböző szintjeinek van hatásuk a súlycsökkenésre.

Megjegyezzük, hogy az eljárást formálisan végrehajtottuk, azonban az alapfeltevések nem teljesülnek. Példánkban sem a szórások nem egyenlőek, sem a normalitás nem igaz (ez utóbbi grafikus eljárások, azaz hisztogram és Gauss-papír alapján adódott). Transzformációkkal (logaritmus, illetve törtkitevős hatvány vétele) részleges javulást sikerült elérni, a transzformáció elvégzését az olvasóra bízuk. \square

2.4. Konfidencia intervallum a várható értékre

Amennyiben a

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p = 0$$

nullhipotézist elvetjük, úgy konfidencia intervallumot kell szerkesztenünk a (különböző) várható értékekre. Megmutatjuk, hogy az egyszeres osztályozás esetén hogyan lehet konfidencia intervallumot szerkeszteni a várható értékekre, ill. azok különbségeire.

A teljes négyzetösszeg felbontásában szereplő lineáris formák függetlenek. Így

$$(\bar{Y}_i - \bar{Y}_., \bar{Y}_j - \bar{Y}_.), \quad Q_2, \quad \bar{Y}_.$$

független (kettő-, illetve egy-egydimenziós) valószínűségi változók ($i \neq j$).

Emiatt \bar{Y}_i és Q_2 függetlenek. Ezek eloszlása:

$$\bar{Y}_i \sim \mathcal{N}\left(\mu + \alpha_i, \frac{\sigma^2}{n_i}\right), \quad \frac{1}{\sigma^2}Q_2 \sim \chi_{n-p}^2(0).$$

Így

$$\sqrt{n_i}(\bar{Y}_i - (\mu + \alpha_i)) / \sqrt{Q_2/(n-p)} \sim t_{n-p}.$$

Ez alapján (ha H_0 -at elvetjük) lehet konfidencia intervallumot szerkeszteni $(\mu + \alpha_i)$ -re.

Másrészt, a fentiek alapján,

$$\bar{Y}_i - \bar{Y}_. \sim \mathcal{N}\left(\alpha_i, \left(\frac{1}{n_i} - \frac{1}{n}\right)\sigma^2\right),$$

ahol $\bar{Y}_i - \bar{Y}_.$ szórása a lineáris forma együtthatói alapján adódott. Tehát ebből a lineáris formából és a Q_2 kvadratikus formából lehet α_i -re konfidencia intervallumot szerkeszteni a t_{n-p} -eloszlás alapján.

Most szerkesszünk konfidencia intervallumot a várható értékek különbségeire.

$$\bar{Y}_i - \bar{Y}_j \quad \text{és} \quad Q_2$$

függetlenek. Továbbá

$$\bar{Y}_i - \bar{Y}_j \sim \mathcal{N}\left(\alpha_i - \alpha_j, \sigma^2\left(\frac{1}{n_i} + \frac{1}{n_j}\right)\right).$$

Ebből és Q_2 -ből $(\alpha_i - \alpha_j)$ -re lehet konfidencia intervallumot szerkeszteni, hiszen

$$\sqrt{\frac{n_i n_j}{n_i + n_j}} (\bar{Y}_i - \bar{Y}_j - (\alpha_i - \alpha_j)) / \sqrt{Q_2/(n-p)} \sim t_{n-p}.$$

3. Kétszeres osztályozás interakció nélkül

3.1. A kétszeres osztályozás feladata

3.1. példa. Egy üzemben 5 gépen, és 3-féle anyagból gyártják ugyanazt a terméket. Jelölje Y a termék egy jellemzőjét. Azonosnak tekinthető-e Y a különböző gépeken és különböző anyagokból készített termékek esetén? \square

A kétszeres osztályozás azt jelenti, hogy két tényező befolyásolja a kísérlet eredményét. Ezeket jelölje A és B . Tegyük fel, hogy az A tényezőnek I , a B tényezőnek pedig J különböző szintje van. Fenti példánkban A a gép, ennek 5 szintje az 5 gép maga. B pedig az anyag, ezen tényező 3 szintje a 3-féle anyag.

A két tényezőnek minden szintkombinációjára végzünk megfigyeléseket, a mérési eredményeket mátrixba rendezzük. Az egyik tényező szintjei a mátrix sorait, a másik tényező szintjei az oszlopait jelöli. Ezért az egyik tényezőt **sorhatásnak**, a másikat **oszlophatásnak** nevezzük. (Egy szántóföldi kísérletnél a tényleges parcellák is így helyezkedhetnek el.)

A két tényező szintkombinációinak száma, azaz a cellák száma IJ . Az elsőként tárgyalandó modellünkben feltételezzük, hogy a tényezőknek együttes hatása (interakciója) nincs. Ekkor elegendő cellánként egy megfigyelést végezni. Legyen az (i, j) -edik cellában végzett kísérlet eredménye Y_{ij} , ($i = 1, 2, \dots, I$; $j = 1, 2, \dots, J$). Tegyük fel, hogy

$$(3.1) \quad Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad i = 1, 2, \dots, I, \quad j = 1, 2, \dots, J,$$

ahol α_i az i -edik sor hatása, β_j pedig a j -edik oszlop hatása; továbbá az ε_{ij} , $i = 1, 2, \dots, I$, $j = 1, 2, \dots, J$, változók függetlenek, normális eloszlásúak:

$$\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2).$$

Fontos (és a tényleges adatok esetén ellenőrizendő), hogy a megfigyelések független, azonos szórású normálisak. Az egyes tényezők szintjei csak a megfigyelések várható értékét befolyásolják.

Elérhető, hogy a paraméterek kielégítsék a következő feltételeket:

$$(3.2) \quad \sum_{i=1}^I \alpha_i = 0, \quad \sum_{j=1}^J \beta_j = 0.$$

A vizsgálandó hipotézisek:

$$\begin{aligned} H_A : \alpha_1 = \alpha_2 = \dots = \alpha_I = 0, \\ H_B : \beta_1 = \beta_2 = \dots = \beta_J = 0, \end{aligned}$$

azaz nincs A -hatás, ill. nincs B -hatás.

Vezessük be a szokásos jelöléseket:

$$\begin{aligned}\bar{Y}_i &= \frac{1}{J} \sum_{j=1}^J Y_{ij}, & i = 1, 2, \dots, I; \\ \bar{Y}_j &= \frac{1}{I} \sum_{i=1}^I Y_{ij}, & j = 1, 2, \dots, J; \\ \bar{Y}_{..} &= \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J Y_{ij}, \\ Q_1 &= \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_i - \bar{Y}_{..})^2, & Q_2 = \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_j - \bar{Y}_{..})^2, \\ Q_3 &= \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_i - \bar{Y}_j + \bar{Y}_{..})^2, & Q = \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_{..})^2.\end{aligned}$$

Itt Q a teljes négyzetösszeg, Q_1 az A tényező szintjei közötti eltérések négyzetösszege, Q_2 a B tényező szintjei közötti eltérések négyzetösszege, végül Q_3 a hibatag. Elemi számolással megmutatható, hogy Q felbontható Q_1 , Q_2 és Q_3 összegére.

3.2. A Fisher-Cochran-tétel alkalmazása a kétszeres osztályozásra

A Fisher-Cochran-tétel alkalmazása szempontjából praktikusabb az alábbi felbontásból kiindulni.

$$\sum_{i=1}^I \sum_{j=1}^J Y_{ij}^2 = Q_1 + Q_2 + Q_3 + IJ\bar{Y}_{..}^2.$$

3.1. állítás. Q_1 szabadsági foka $I - 1$, Q_2 szabadsági foka $J - 1$, Q_3 szabadsági foka $(I - 1)(J - 1)$, $IJ\bar{Y}_{..}^2$ szabadsági foka pedig 1. \square

3.2. állítás. (a) Q_1 , Q_2 és Q_3 függetlenek;
 (b) Q_3/σ^2 mindig centrált χ^2 -eloszlású $(I - 1)(J - 1)$ szabadsági fokkal;
 (c) $\frac{1}{\sigma^2}Q_1 \sim \chi_{I-1}^2(\lambda_1)$, ahol $\lambda_1 = \frac{J}{\sigma^2} \sum_{i=1}^I \alpha_i^2$;
 (d) $\frac{1}{\sigma^2}Q_2 \sim \chi_{J-1}^2(\lambda_2)$, ahol $\lambda_2 = \frac{I}{\sigma^2} \sum_{j=1}^J \beta_j^2$;
 (e) $\frac{1}{\sigma^2}Q_1$ akkor és csak akkor centrált χ^2 -eloszlású, ha H_A teljesül;
 (f) $\frac{1}{\sigma^2}Q_2$ akkor és csak akkor centrált χ^2 -eloszlású, ha H_B teljesül.

BIZONYÍTÁS. Az előző állítás miatt a Fisher-Cochran-tétel feltétele teljesül, hiszen a rangok összege $(I - 1) + (J - 1) + (I - 1)(J - 1) + 1 = IJ$. Ezért a kvadratikus formák függetlenek és χ^2 -eloszlásúak. A nem-centralitási paramétereket úgy kapjuk, hogy a várható értékeket beírjuk a kvadratikus formákba. \square

3.3. A szórásfelbontó táblázat

A 3.2. állítás alapján a Q_3 véletlen hiba (skalárszorosa) mindig centrált χ^2 -eloszlású, Q_1 és Q_2 (skalárszorosa) viszont pontosan akkor, ha H_A , ill. H_B teljesül.

3.3. állítás. Az

$$F_A = \frac{Q_1}{I-1} \bigg/ \frac{Q_3}{(I-1)(J-1)}$$

statisztika pontosan akkor $I-1$ és $(I-1)(J-1)$ szabadsági fokú centrált F -eloszlású, ha a

$$H_A : \alpha_1 = \alpha_2 = \dots = \alpha_I = 0$$

nullhipotézis teljesül.

Az

$$F_B = \frac{Q_2}{J-1} \bigg/ \frac{Q_3}{(I-1)(J-1)}$$

statisztika pontosan akkor $J-1$ és $(I-1)(J-1)$ szabadsági fokú centrált F -eloszlású, ha a

$$H_B : \beta_1 = \beta_2 = \dots = \beta_J = 0$$

nullhipotézis teljesül.

BIZONYÍTÁS. A 3.2. állítás szerint F_A két független, a szabadsági fokával elosztott, (H_0 esetén centrált) χ^2 -eloszlású valószínűségi változó hányadosa. Ezért az F -eloszlás definíciója alapján ennek eloszlása (H_0 esetén centrált) F -eloszlás $I-1$ és $(I-1)(J-1)$ szabadsági fokokkal. \square

A fentiek alapján a kéttényezős (interakció nélküli) modell szórásfelbontó táblázata az alábbi.

A szóródás forrásai	Szabadsági fokok	Négyzet-összegek	Négyzet-átlagok	F -hányadosok
Sorok közötti eltérések (A tényező szintjei)	$I-1$	Q_1	$S_1 = \frac{Q_1}{I-1}$	$\frac{S_1}{S_3}$
Oszlopok közötti eltérések (B tényező szintjei)	$J-1$	Q_2	$S_2 = \frac{Q_2}{J-1}$	$\frac{S_2}{S_3}$
Véletlen hiba	$(I-1)(J-1)$	Q_3	$S_3 = \frac{Q_3}{(I-1)(J-1)}$	
Teljes	$IJ-1$	Q		

A kéttényezős (interakció nélküli) modell szórásfelbontó táblázata

A H_A nullhipotézist $1-\alpha$ szinten elvetjük, ha

$$S_1/S_3 \geq F_{[I-1, (I-1)(J-1); \alpha]},$$

míg a H_B nullhipotézist $1-\alpha$ szinten elvetjük, ha

$$S_2/S_3 \geq F_{[J-1, (I-1)(J-1); \alpha]},$$

ahol az $F_{[l,m;\alpha]}$ az (l, m) szabadsági fokú centrált F -eloszlás $1 - \alpha$ szinthez tartozó (felső) kritikus értéke.

3.4. Konfidencia intervallum a várható értékre

A nullhipotézisek elutasítása esetén az egytényezős modellnél mondottakhoz hasonlóan becsülhetjük az ismeretlen α_i, β_j paramétereket, különbségeiket és a σ^2 -et.

Mivel Q_3/σ^2 mindig centrált χ^2 -eloszlású $(I-1)(J-1)$ szabadsági fokkal, ezért σ^2 torzítatlan becslése $S_3 = \frac{Q_3}{(I-1)(J-1)}$.

Belátható, hogy $\bar{Y}_i. - \bar{Y}_{i'}. \sim \mathcal{N}\left(\alpha_i - \alpha_{i'}, \frac{2\sigma^2}{J}\right)$. Ezért $1 - \alpha$ megbízhatósági szintű konfidencia intervallum az $\alpha_i - \alpha_{i'}$ különbségre (ahol $i \neq i'$)

$$(3.3) \quad (\bar{Y}_i. - \bar{Y}_{i'}.) \mp t_{[(I-1)(J-1); \alpha/2]} \sqrt{2S_3/J},$$

ahol $t_{[(I-1)(J-1); \alpha/2]}$ az $1 - \alpha$ szintű kétoldali t -próba kritikus értéke.

Továbbá, $\bar{Y}_i. - \bar{Y}..$ eloszlása

$$(3.4) \quad \mathcal{N}\left(\alpha_i, \frac{\sigma^2}{J} \left(1 - \frac{1}{I}\right)\right),$$

ahol a szórás a lineáris forma együtthatói alapján adódott. Ezért az $1 - \alpha$ szintű konfidencia intervallum α_i -re a következő:

$$(\bar{Y}_i. - \bar{Y}..) \mp t_{[(I-1)(J-1); \alpha/2]} \sqrt{\frac{S_3}{J} \left(1 - \frac{1}{I}\right)}.$$

Hasonló módon adhatunk konfidencia intervallumot a $\beta_j - \beta_{j'}$ különbségre és β_j -re.

3.5. Feladatok

1. A (3.1)-ben szereplő mennyiségekkel definiáljuk az alábbi új mennyiségeket. $\tilde{\mu} = \mu + \bar{\alpha} + \bar{\beta}$, $\tilde{\alpha}_i = \alpha_i - \bar{\alpha}$, $\tilde{\beta}_j = \beta_j - \bar{\beta}$, ahol $\bar{\alpha} = (1/I) \sum_{i=1}^I \alpha_i$, $\bar{\beta} = (1/J) \sum_{j=1}^J \beta_j$. Lássuk be, hogy minden i és j esetén $\tilde{\mu} + \tilde{\alpha}_i + \tilde{\beta}_j = \mu + \alpha_i + \beta_j$, és a $\tilde{\cdot}$ -mal jelölt új mennyiségek már teljesítik a (3.2) feltételeket!

2. Lássuk be, hogy a (3.1) modell esetén, ha $i \neq i'$,

$$\bar{Y}_i. - \bar{Y}_{i'}. \sim \mathcal{N}\left(\alpha_i - \alpha_{i'}, \frac{2\sigma^2}{J}\right),$$

$\frac{1}{\sigma^2} Q_3 \sim \chi_{(I-1)(J-1)}^2$, és ezek egymástól függetlenek. Ez alapján igazoljuk (3.3)-at!

3. Igazoljuk (3.4)-et!

4. A kétszeres osztályozás az interakció figyelembevételével

4.1. A modell

Most a kétszeres osztályozás realiztikusabb esetét tekintjük: kétszeres osztályozást interakcióval.

Legyen az A tényezőnek I szintje, a B tényezőnek J szintje. Tételezzük fel tehát azt, hogy a két tényező együttes hatása (interakciója) is befolyásolhatja a kísérlet eredményét. Ilyenkor cellánként több megfigyelést kell végeznünk.

Végezzünk a két tényező IJ számú szintkombinációjának mindegyikére $K (> 1)$ számú független megfigyelést. A cellánkénti egyenlő számú megfigyeléses kísérleti elrendezést **kiegyensúlyozott elrendezésnek** nevezzük. K -t ismétlésszámnak mondjuk.

Jelölje Y_{ijk} az (i, j) -edik cellában a k -edik megfigyelés eredményét ($k = 1, \dots, K$). Tegyük fel, hogy

$$(4.1) \quad Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk},$$

ahol az ε_{ijk} véletlen hibák független, 0 és σ^2 paraméterű normális eloszlású valószínűségi változók.

A vizsgálandó hipotézisek:

$$\begin{aligned} H_A : \alpha_1 = \alpha_2 = \dots = \alpha_I = 0, \\ H_B : \beta_1 = \beta_2 = \dots = \beta_J = 0, \\ H_{AB} : \text{minden } i = 1, 2, \dots, I, j = 1, 2, \dots, J\text{-re } \gamma_{ij} = 0. \end{aligned}$$

Azaz nincs A -hatás, nincs B -hatás, illetve nincs AB interakció.

Elérhető, hogy a paraméterek kielégítsék a következő feltételeket:

$$(4.2) \quad \begin{aligned} \sum_{i=1}^I \alpha_i = 0, \quad \sum_{j=1}^J \beta_j = 0, \\ \sum_{i=1}^I \gamma_{ij} = 0 \quad (j = 1, 2, \dots, J), \quad \sum_{j=1}^J \gamma_{ij} = 0 \quad (i = 1, 2, \dots, I). \end{aligned}$$

Ekkor az ismeretlen μ , α_i , β_j és γ_{ij} paraméterek legkisebb négyzetes becslései az alábbi alakúak. (A maximum likelihood becslések is ugyanezek.)

$$(4.3) \quad \begin{aligned} \hat{\mu} &= \bar{Y}_{\dots}, \\ \hat{\alpha}_i &= \bar{Y}_{i..} - \bar{Y}_{\dots}, \\ \hat{\beta}_j &= \bar{Y}_{.j.} - \bar{Y}_{\dots}, \\ \hat{\gamma}_{ij} &= \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{\dots}, \end{aligned}$$

ahol

$$\begin{aligned}\bar{Y}_{...} &= \frac{1}{IJK} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K Y_{ijk}, \\ \bar{Y}_{i..} &= \frac{1}{JK} \sum_{j=1}^J \sum_{k=1}^K Y_{ijk}, \\ \bar{Y}_{.j.} &= \frac{1}{IK} \sum_{i=1}^I \sum_{k=1}^K Y_{ijk}, \\ \bar{Y}_{ij.} &= \frac{1}{K} \sum_{k=1}^K Y_{ijk}\end{aligned}$$

minden $i = 1, 2, \dots, I$, $j = 1, 2, \dots, J$ esetén. Ezek segítségével adhatjuk meg teljes négyzetösszeg felbontását. Az elemzés az alábbi négyzetösszegeken fog alapulni.

$$\begin{aligned}Q_1 &= JK \sum_{i=1}^I (\bar{Y}_{i..} - \bar{Y}_{...})^2, \\ Q_2 &= IK \sum_{j=1}^J (\bar{Y}_{.j.} - \bar{Y}_{...})^2, \\ Q_3 &= K \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...})^2, \\ Q_4 &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij.})^2, \\ Q &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{...})^2.\end{aligned}$$

Q_1 az első tényező szintjei közötti eltéréseket, Q_2 a második tényező szintjei közötti eltéréseket, Q_3 pedig a két tényező együttes hatását írja le. Q_4 a véletlen hibák négyzetösszege, míg Q a teljes négyzetösszeg. A teljes négyzetösszeg felbontása:

$$Q = Q_1 + Q_2 + Q_3 + Q_4.$$

Ez a formula a kétszeres szorzatok 0 volta miatt igaz.

4.2. A Fisher-Cochran-tétel alkalmazása

A Fisher-Cochran-tétel alkalmazása szempontjából praktikusabb az alábbi felbontásból kiindulni.

$$(4.4) \quad \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K Y_{ijk}^2 = Q_1 + Q_2 + Q_3 + Q_4 + IJK \bar{Y}_{...}^2.$$

Q_1 szabadsági foka $I - 1$, Q_2 szabadsági foka $J - 1$, Q_3 szabadsági foka $(I - 1)(J - 1)$, Q_4 szabadsági foka $IJ(K - 1)$, $IJK\bar{Y}_{...}^2$ szabadsági foka pedig 1.

4.1. állítás. Az (4.1) modell esetén

- (a) Q_1, Q_2, Q_3 és Q_4 függetlenek és χ^2 -eloszlásúak, de nem feltétlenül centráltak;
- (b) $\frac{1}{\sigma^2}Q_4$ mindig centrált χ^2 -eloszlású $IJ(K - 1)$ szabadsági fokkal;
- (c) $\frac{1}{\sigma^2}Q_1$ akkor és csak akkor centrált χ_{I-1}^2 -eloszlású, ha H_A teljesül;
- (d) $\frac{1}{\sigma^2}Q_2$ akkor és csak akkor centrált χ_{J-1}^2 -eloszlású, ha H_B teljesül;
- (e) $\frac{1}{\sigma^2}Q_3$ akkor és csak akkor centrált $\chi_{(I-1)(J-1)}^2$ -eloszlású, ha H_{AB} teljesül.

BIZONYÍTÁS. A (4.4) előállításban a Fisher-Cochran-tétel feltétele teljesül, hiszen a rangok összege $(I - 1) + (J - 1) + (I - 1)(J - 1) + IJ(K - 1) + 1 = IJK$. Ezért a kvadratikus formák függetlenek és χ^2 -eloszlásúak. A nem-centralitási paramétereket úgy kapjuk, hogy a várható értékeket beírjuk a kvadratikus formákba. □

4.3. A szórásfelbontó táblázat

A fentiek alapján a kéttényezős (interakciót figyelembe vevő) modell szórásfelbontó táblázata az alábbi.

A szóródás forrásai	Szabadsági fokok	Négyzet-összegek	Négyzet-átlagok	F-hányadosok
A tényező	$I - 1$	Q_1	$S_1 = \frac{Q_1}{I-1}$	S_1/S_4
B tényező	$J - 1$	Q_2	$S_2 = \frac{Q_2}{J-1}$	S_2/S_4
AB interakció	$(I - 1)(J - 1)$	Q_3	$S_3 = \frac{Q_3}{(I-1)(J-1)}$	S_3/S_4
Véletlen hiba	$IJ(K - 1)$	Q_4	$S_4 = \frac{Q_4}{IJ(K-1)}$	
Teljes	$IJK - 1$	Q		

Szórásfelbontó táblázat kéttényezős (interakciót figyelembe vevő) modelle

A különböző nullhipotézisek F -próbáit ugyanúgy hajtjuk végre, mint az előző modelleknél. Elsőként a H_{AB} hipotézist vizsgáljuk.

4.1. példa. Kétféle vas (Fe^{2+} és Fe^{3+}) szervezetből való kiürülését vizsgáljuk. Mindkét vas fajtából 10.2, 1.2 és 0.3 millimólos koncentrációt adunk, 18-18 egyednek.

A lenti mátrix első 3 oszlopa a Fe^{2+} három különböző mennyisége esetén a szervezetben maradt százalékos arányt mutatja a 18-18 egyed esetén. Hasonló a mátrix további 3 oszlopa. A 3 különböző vas-mennyiség az oszlophatás. A vas 2 minősége a sorhatás. Végül 18 az ismétlésszám. Tehát egy kiegyensúlyozott elrendezésünk van.

Miért nem végeztük el mégsem az ANOVA eljárást ezekre az adatokra? Mert a mintának normális eloszlásúnak kell lennie azonos szórással! Viszont az empirikus szórásnégyzetek az mutatták, hogy cellánként vett 18-18 elemű minta szórásai nagyon eltérnek. Ugyanezen mintákra a hisztogram a haranggörbétől nagy eltérés mutatott, ugyanígy a Gauss-papír az egyenestől, tehát a normalitás sem áll fenn. Próbálkozzunk a minta transzformálásával! A logaritmus vétel jelentős javulást hozott mind a normalitásban, mind a szórások egyenlőségében. (Megjegyezzük, hogy normalitást 18 elemű

mintából nem igazán lehet jól tesztelni. A módszert tekintve: a normalitást Shapiro-Wilk-próbával, a szórások azonosságát egy Bartlett-próbával kellene ellenőrizni.)

Fe ²⁺			Fe ³⁺		
10.2	1.2	0.3	10.2	1.2	0.3
0.71	2.20	2.25	2.20	4.04	2.71
1.66	2.93	3.93	2.69	4.16	5.43
2.01	3.08	5.08	3.54	4.42	6.38
2.16	3.49	5.82	3.75	4.93	6.38
2.42	4.11	5.84	3.83	5.49	8.32
2.42	4.95	6.89	4.08	5.77	9.04
2.56	5.16	8.50	4.27	5.86	9.56
2.60	5.54	8.56	4.53	6.28	10.01
3.31	5.68	9.44	5.32	6.97	10.08
3.64	6.25	10.52	6.18	7.06	10.62
3.74	7.25	13.46	6.22	7.78	13.80
3.74	7.90	13.57	6.33	9.23	15.99
4.39	8.85	14.56	6.97	9.34	17.90
4.50	11.96	16.41	6.97	9.91	18.25
5.07	15.54	16.96	7.52	13.46	19.32
5.26	15.89	17.56	8.36	18.4	19.87
8.15	18.3	22.82	11.65	23.89	21.60
8.24	18.59	29.13	12.45	26.39	22.25

Az adatok logaritmusára az alábbi szórásfelbontó táblázatot kaptuk:

ANOVA Table				
Source	SS	df	MS	F
Columns	15.58	2	7.79	22.52
Rows	2.079	1	2.079	6.01
Interaction	0.8089	2	0.4045	1.62
Error	35.39	102	0.346	
Total	53.76	107		

A szokásos H_A , H_B , H_{AB} nullhipotézisek esetén a megfelelő F -próbák pontos terjedelmére a következőket kaptuk: $7.9 \cdot 10^{-9}$, 0.015, 0.314. Ezen p értékek alapján el kell vetni, hogy H_A : *nincs oszlophatás*, H_B : *nincs sorhatás*, de el kell fogadni, hogy H_{AB} : *nincs interakció*.

Magyarázat. SS oszlopában a teljes négyzetösszeg felbontása van oszlophatásra, sorhatásra, interakcióra és véletlen hibára. A df alatt ezek szabadsági foka áll. MS alatt SS és df hányadosa. F pedig a megfelelő hatás MS értékének és a véletlen hiba MS értékének hányadosa. Ha valamelyik hatás a véletlen hibához képest túl nagy, akkor elvetjük az azon hatásra vonatkozó nullhipotézist. A $(p_1, p_2, p_3) = (7.9 \cdot 10^{-9}, 0.015, 0.314)$ vektor i -edik komponense annak az értékét adja, hogy a nullhipotézis igaz volta esetén milyen valószínűséggel lehetne a centrált F -eloszlású valószínűségi

változó legalább olyan nagy, mint a kapott F érték. Így kicsi p_i esetén elvetjük a nullhipotézist. A p_i értékét F -eloszlás alapján határozzuk meg: az F oszlopban lévő statisztikák mindig F -eloszlásúak, a megfelelő nullhipotézis fennállása esetén centráltak, ha viszont a nullhipotézis nem igaz, akkor nem centráltak. (Megjegyezzük, hogy a két fenti szórásanalízisbeli példában a MATLAB és a SAS futási eredménye teljesen megegyezett.) \square

4.4. Feladatok

1. A (4.1)-ben szereplő mennyiségekkel definiáljuk az alábbi új mennyiségeket. $\tilde{\mu} = \mu + \bar{\alpha} + \bar{\beta} + \bar{\gamma}_{..}$, $\tilde{\alpha}_i = \alpha_i - \bar{\alpha} + \bar{\gamma}_{i.} - \bar{\gamma}_{..}$, $\tilde{\beta}_j = \beta_j - \bar{\beta} + \bar{\gamma}_{.j} - \bar{\gamma}_{..}$, $\tilde{\gamma}_{ij} = \gamma_{ij} - \bar{\gamma}_{i.} - \bar{\gamma}_{.j} + \bar{\gamma}_{..}$, ahol $\bar{\alpha} = (1/I) \sum_{i=1}^I \alpha_i$, $\bar{\beta} = (1/J) \sum_{j=1}^J \beta_j$, $\bar{\gamma}_{i.} = (1/J) \sum_{j=1}^J \gamma_{ij}$, $\bar{\gamma}_{.j} = (1/I) \sum_{i=1}^I \gamma_{ij}$, $\bar{\gamma}_{..} = (1/IJ) \sum_{i=1}^I \sum_{j=1}^J \gamma_{ij}$. Lássuk be, hogy minden i és j esetén $\tilde{\mu} + \tilde{\alpha}_i + \tilde{\beta}_j + \tilde{\gamma}_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$, és a $\tilde{\cdot}$ -mal jelölt új mennyiségek már teljesítik a (4.2) feltételeket!

2. A (4.1) modellben a legkisebb négyzetes becslés, azaz a hiba négyzetösszeg minimalizálása a

$$(4.5) \quad \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left(Y_{ijk} - (\mu + \alpha_i + \beta_j + \gamma_{ij}) \right)^2$$

függvény minimumának meghatározását jelenti μ , α_i , β_j és γ_{ij} szerint. Lássuk be, hogy a minimum a (4.3) képletbeni értékeknél adódik! Ehhez használjuk a minimum deriváltakal történő meghatározását. A kapott lineáris egyenletrendszer megoldásakor vegyük figyelembe (4.2)-ot.

II. fejezet

Regresszióanalízis és a lineáris modell

1. A lineáris modell

1.1. A lineáris modell definíciója

$$(1.1) \quad \boxed{\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}}$$

a lineáris modell, ahol

\mathbf{Y} n -dimenziós megfigyelés vektor,

X $n \times p$ méretű, nem véletlen, megfigyelt mátrix (a magyarázó változók mátrixa),

$\boldsymbol{\beta}$ p -dimenziós ismeretlen paraméter,

$\boldsymbol{\varepsilon}$ nem megfigyelhető n -dimenziós véletlen vektor (hiba).

Általában $n \gg p$, ezt szükség esetén fel fogjuk tenni. A gyakorlatban p a magyarázó változók száma, n pedig a megfigyelt objektumok száma, tehát $n \gg p$ ésszerű feltétel.

1.1. példa.

1.2. A legkisebb négyzetek módszere

Ha $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{0}$ és $\text{var } \boldsymbol{\varepsilon} = \sigma^2 I$ (σ^2 ismeretlen paraméter), akkor **homoszkedasztikus** esetről beszélünk. Ekkor a **legkisebb négyzetes becslést** (OLS=Ordinary Least Squares) alkalmazzuk $\boldsymbol{\beta}$ -ra: ez lesz $\hat{\boldsymbol{\beta}}$.

Legyen tehát $\hat{\boldsymbol{\beta}}$ az $\|\mathbf{Y} - X\boldsymbol{\beta}\|^2$ -et minimalizáló vektor. (Itt $\|\cdot\|$ a norma \mathbb{R}^n -ben.) Jelölje P_F az X oszlopai által generált F altérre való merőleges vetítést.

1.1. állítás. 1. $\hat{\boldsymbol{\beta}}$ legkisebb négyzetes becslés \iff

$$X\hat{\boldsymbol{\beta}} = P_F \mathbf{Y}.$$

2. $\hat{\boldsymbol{\beta}}$ legkisebb négyzetes becslés \iff $\hat{\boldsymbol{\beta}}$ az

$$\boxed{X^\top \mathbf{Y} = X^\top X \boldsymbol{\beta}}$$

normálegyenlet megoldása.

BIZONYÍTÁS. 1. \mathbf{Y} -hoz az X oszlopai által generált altér mely $X\boldsymbol{\beta}$ eleme lesz legközelebb? Éppen az \mathbf{Y} vetülete, azaz $P_F \mathbf{Y}$. Így $X\hat{\boldsymbol{\beta}} = P_F \mathbf{Y}$.

2. $\|\mathbf{Y} - X\boldsymbol{\beta}\|^2$ mikor a legkisebb? Ha $\mathbf{Y} - X\boldsymbol{\beta}$ éppen az \mathbf{Y} ortogonális komplementere az F altérre vonatkozóan. Azaz $\mathbf{Y} - X\boldsymbol{\beta}$ merőleges X minden oszlopára, tehát

$$X^\top \mathbf{Y} - X^\top X \boldsymbol{\beta} = \mathbf{0},$$

vagyis

$$X^\top X \hat{\boldsymbol{\beta}} = X^\top \mathbf{Y}. \quad \square$$

1.1. megjegyzés. $X^\top X$ invertálható \iff $\text{rang } X = p$.

BIZONYÍTÁS. 1. Ha $X^\top X$ invertálható, akkor $X^\top X \mathbf{a} = \mathbf{0}$ esetén $\mathbf{a} = \mathbf{0}$. Ezért $X \mathbf{a}$ is csak akkor lehet $\mathbf{0}$, ha $\mathbf{a} = \mathbf{0}$. Azaz X oszlopainak csak 0 együtthatókkal való lineáris kombinációi egyenlők $\mathbf{0}$ -val, tehát X rangja egyenlő az oszlopai számával, vagyis p -vel.

2. Ha $\text{rang } X = p$, akkor $X \mathbf{a} \neq \mathbf{0}$ tetszőleges $\mathbf{a} \neq \mathbf{0}$ esetén; így $\mathbf{a}^\top X^\top X \mathbf{a} = \|X \mathbf{a}\|^2 \neq 0$. Viszont ha $X^\top X$ nem lenne invertálható, akkor a vele képzett kvadratikus forma 0 lenne valamely $\mathbf{a} \neq \mathbf{0}$ vektorra. \square

1.2. megjegyzés. Ha $\text{rang } X = p$, akkor

$$\boxed{\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{Y}}.$$

Ez éppen a normálegyenlettel ekvivalens, ha $(X^\top X)$ invertálható. \square

1.2. állítás. Legyen $\mathbb{E} \boldsymbol{\varepsilon} = \mathbf{0}$, $\text{var } \boldsymbol{\varepsilon} = \sigma^2 I$ és $\text{rang } X = p$. Ekkor $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{Y}$ torzítatlan becslése $\boldsymbol{\beta}$ -nak, továbbá $\text{var } \hat{\boldsymbol{\beta}} = \sigma^2 (X^\top X)^{-1}$.

Ha $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I)$, akkor $\hat{\boldsymbol{\beta}} \sim \mathcal{N}_p(\boldsymbol{\beta}, \sigma^2 (X^\top X)^{-1})$.

BIZONYÍTÁS. Ha $\text{rang } X = p$, akkor $(X^\top X)$ invertálható. Ekkor

$$\mathbb{E} \hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbb{E} \mathbf{Y} = \boldsymbol{\beta},$$

hiszen $\mathbb{E} \mathbf{Y} = X \boldsymbol{\beta}$. Másrészt

$$\text{var}(\hat{\boldsymbol{\beta}}) = (X^\top X)^{-1} X^\top (\text{var}(\mathbf{Y})) X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1},$$

ugyanis $\text{var}(\mathbf{Y}) = \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 I$.

Ha $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I)$, akkor $\mathbf{Y} \sim \mathcal{N}_n(X \boldsymbol{\beta}, \sigma^2 I)$, így $\hat{\boldsymbol{\beta}}$ — lévén \mathbf{Y} lineáris függvénye — maga is normális eloszlású. \square

1.3. A Gauss-Markov-tétel legegyszerűbb alakja

A homoszkedasztikus esetben $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{Y}$ legkisebb négyzetes becslés a **legjobb lineáris torzítatlan becslés** (BLUE=Best Linear Unbiased Estimator). Ezt mondja ki a **Gauss-Markov-tétel** legegyszerűbb alakja.

1.1. tétel. (Gauss-Markov.) Ha $\mathbb{E} \boldsymbol{\varepsilon} = \mathbf{0}$, $\text{var } \boldsymbol{\varepsilon} = \sigma^2 I$ és $\text{rang } X = p$, akkor $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{Y}$ a $\boldsymbol{\beta}$ paraméter vektor legjobb lineáris torzítatlan becslése.

BIZONYÍTÁS. $\hat{\boldsymbol{\beta}}$ az \mathbf{Y} vektor lineáris függvénye $((X^\top X)^{-1} X^\top$ mátrixszal való szorzata), azaz $\hat{\boldsymbol{\beta}}$ lineáris becslés. Továbbá $\mathbb{E} \hat{\boldsymbol{\beta}} = \boldsymbol{\beta}$, tehát torzítatlan.

Legyen $C \mathbf{Y}$ a $\boldsymbol{\beta}$ -nak egy másik lineáris torzítatlan becslése. Be kell látnunk, hogy $C \mathbf{Y}$ szórásmatrixa „nagyobb” $\hat{\boldsymbol{\beta}}$ szórásmatrixánál, azaz $\text{var}(C \mathbf{Y}) - \text{var}(\hat{\boldsymbol{\beta}})$ pozitív szemidefinit. $C \mathbf{Y}$ torzítatlanságából

$$\boldsymbol{\beta} = \mathbb{E} C \mathbf{Y} = C \mathbb{E}(X \boldsymbol{\beta} + \boldsymbol{\varepsilon}) = C X \boldsymbol{\beta} \quad \forall \boldsymbol{\beta}.$$

Tehát $CX = I$.

$$\begin{aligned} \text{var}(CY) &= \mathbb{E}(CY - \beta)(CY - \beta)^\top = \\ &= \mathbb{E}(CY - \hat{\beta} + \hat{\beta} - \beta)(CY - \hat{\beta} + \hat{\beta} - \beta)^\top = \\ &= \mathbb{E}(CY - \hat{\beta})(CY - \hat{\beta})^\top + \mathbb{E}(CY - \hat{\beta})(\hat{\beta} - \beta)^\top + \mathbb{E}(\hat{\beta} - \beta)(CY - \hat{\beta})^\top \\ &\quad + \mathbb{E}(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top \geq \mathbb{E}(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top = \text{var}(\hat{\beta}), \end{aligned}$$

hiszen

$$\mathbb{E}(CY - \hat{\beta})(CY - \hat{\beta})^\top \quad \text{pozitív szemidefinit, és}$$

$$\mathbb{E}(CY - \hat{\beta})(\hat{\beta} - \beta)^\top = 0$$

(ezt utóbbit alább megmutatjuk). Így $\text{var}(CY) \geq \text{var}(\hat{\beta})$, azaz $\hat{\beta}$ a legjobb (a legkisebb szórású).

Figyeljük meg, hogy

$$\begin{aligned} CY - \hat{\beta} &= \underbrace{CX}_I \beta + C\varepsilon - (X^\top X)^{-1} X^\top X \beta - (X^\top X)^{-1} X^\top \varepsilon \\ &= [C - (X^\top X)^{-1} X^\top] \varepsilon \end{aligned}$$

és $\hat{\beta} - \beta = (X^\top X)^{-1} X^\top \varepsilon$, ahonnan

$$\begin{aligned} \mathbb{E}(CY - \hat{\beta})(\hat{\beta} - \beta)^\top &= [C - (X^\top X)^{-1} X^\top] \underbrace{\mathbb{E}(\varepsilon \varepsilon^\top)}_{\sigma^2 I} X (X^\top X)^{-1} = \\ &= \sigma^2 \left[\underbrace{CX}_I (X^\top X)^{-1} - \underbrace{(X^\top X)^{-1} X^\top X}_I (X^\top X)^{-1} \right] = 0. \quad \square \end{aligned}$$

1.4. Mátrixok általánosított inverze

1.1. definíció. Az A $n \times n$ -es mátrix **általánosított inverze** az az A^- $n \times n$ -es mátrix, melyre

$$\boxed{AA^-A = A}. \quad \square$$

1.3. megjegyzés. Ha A invertálható, akkor A^- egyértelműen létezik, és $A^- = A^{-1}$.

Ugyanis ekkor $AA^{-1}A = A$, azaz A^{-1} kielégíti A^- definícióját. Másrészt az $AA^-A = A$ egyenletet jobbról és balról is megszorozva A^{-1} -gyel, $A^- = A^{-1}$ -et kapjuk. \square

1.3. állítás. Ha A szimmetrikus, akkor létezik általánosított inverze.

BIZONYÍTÁS. Legyenek $\lambda_1, \dots, \lambda_k$ az A nem zérus sajátértékei, $\mathbf{u}_1, \dots, \mathbf{u}_n$ pedig az A sajátvektorainak ortonormált rendszere. Legyen U az $\mathbf{u}_1, \dots, \mathbf{u}_n$ -et tartalmazó ortogonális mátrix,

$$\Lambda^- = \begin{pmatrix} \lambda_1^{-1} & & & & \\ & \ddots & & & \\ & & \lambda_k^{-1} & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_1 & & & & \\ & \ddots & & & \\ & & \lambda_k & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{pmatrix}$$

$n \times n$ -es diagonális mátrixok. Ekkor $A^- = U\Lambda^-U^\top$ teljesíti a feltételeket:

$$AA^-A = (U\Lambda U^\top)(U\Lambda^-U^\top)(U\Lambda U^\top) = U\Lambda U^\top = A. \quad \square$$

1.4. megjegyzés. Erre a speciális A^- -ra $A^-AA^- = A^-$ és $(A^-)^\top = A^-$ is teljesül. \square

1.4. állítás. $X(X^\top X)^-X^\top$ ortogonális projekció az X oszlopvektorai által generált F altérre.

BIZONYÍTÁS. Ha $\mathbf{w} \perp F$, akkor $X^\top \mathbf{w} = \mathbf{0}$, így $X(X^\top X)^-X^\top \mathbf{w} = \mathbf{0}$.

Ahhoz, hogy $X(X^\top X)^-X^\top$ az F -beli vektorokat önmagukba viszi, elég belátni, hogy $X(X^\top X)^-X^\top X = X$. Ez viszont igaz, mert tetszőleges \mathbf{v} vektor $\mathbf{v} = \mathbf{v}_1 + X\mathbf{v}_2$ alakba írható, ahol $\mathbf{v}_1 \perp F$. Ezért

$$\begin{aligned} \mathbf{v}^\top X(X^\top X)^-X^\top X &= \mathbf{v}_1^\top X(X^\top X)^-X^\top X + \mathbf{v}_2^\top X^\top X(X^\top X)^-X^\top X = \\ &= \mathbf{0} + \mathbf{v}_2^\top (X^\top X) = \mathbf{v}^\top X. \end{aligned} \quad \square$$

1.5. A Gauss-Markov-tétel általános alakja

1.5. megjegyzés. $\hat{\boldsymbol{\beta}} = (X^\top X)^-X^\top \mathbf{Y}$ mindig. Hiszen ekkor az $X\hat{\boldsymbol{\beta}} = P_F \mathbf{Y}$ egyenlet teljesül. Ugyanis $X(X^\top X)^-X^\top$ éppen a P_F vetítéssel egyenlő. \square

1.6. megjegyzés. Ha X oszlopvektorai lineárisan függetlenek, akkor azt mondjuk, hogy a magyarázó változók között **kollinearitás** áll fenn. Ekkor — azaz ha $\text{rang } X < p$ — a $\boldsymbol{\beta}$ vektornak csupán bizonyos koordinátáit, illetve azok lineáris függvényeit tudjuk becsülni, a teljes $\boldsymbol{\beta}$ -t nem. \square

1.2. definíció. A $\boldsymbol{\beta}$ paraméter $\mathbf{c}^\top \boldsymbol{\beta}$ lineáris függvényét (lineárisan és torzítatlanul) **becsülhetőnek** nevezzük, ha van \mathbf{Y} -nak olyan $\mathbf{b}^\top \mathbf{Y}$ lineáris függvénye, melyre

$$\mathbb{E} \mathbf{b}^\top \mathbf{Y} = \mathbf{c}^\top \boldsymbol{\beta} \quad \forall \boldsymbol{\beta} \in \mathbb{R}^p.$$

(Itt $\mathbf{c} \in \mathbb{R}^p$, $\mathbf{b} \in \mathbb{R}^n$.) Ilyenkor röviden azt is mondjuk, hogy \mathbf{c}^\top becsülhető. \square

Amikor $\text{rang } X = p$, akkor a teljes $\boldsymbol{\beta}$ vektor becsülhető.

Az előzőekben az, hogy $\hat{\boldsymbol{\beta}}$ a legjobb lineáris torzítatlan becslés, azt jelentette, hogy $\hat{\boldsymbol{\beta}}$ szórás mátrixa a legkisebb a lineáris torzítatlan becslések szórás mátrixai között. (A rendezést a „ $D_1 \geq D_2 \iff D_1 - D_2$ pozitív szemidefinit” alapján értelmezzük.) Tehát $\mathbf{c}^\top \hat{\boldsymbol{\beta}}$ szórása a legkisebb a $\mathbf{c}^\top \boldsymbol{\beta}$ lineáris torzítatlan becslései szórásai között ($\mathbf{c} \in \mathbb{R}^p$). Ez a tény a Gauss-Markov-tétel általánosabb változatának is a lényege.

1.2. tétel. (Gauss-Markov.) 1. $\mathbf{c}^\top \boldsymbol{\beta}$ becsülhető $\iff \mathbf{c} = X^\top \mathbf{b}$ valamely $\mathbf{b} \in \mathbb{R}^n$ -re.

2. A $\mathbf{c}^\top \boldsymbol{\beta}$ minden lineáris torzítatlan $\mathbf{b}^\top \mathbf{Y}$ becslésére $P_F \mathbf{b} = \mathbf{b}^*$, ahol \mathbf{b}^* (a \mathbf{b} -től nem függő) állandó vektor. Itt $P_F = X(X^\top X)^-X^\top$ az X oszlopai alterére való vetítés.

3. $\mathbf{b}^{*\top} \mathbf{Y}$ adja a legkisebb szórású lineáris torzítatlan becslést.

4. $\mathbf{b}^{*\top} \mathbf{Y} = \mathbf{c}^\top \hat{\boldsymbol{\beta}}$.

BIZONYÍTÁS. **1.** Legyen $\mathbf{c}^\top \boldsymbol{\beta}$ becslhető. Ekkor $\mathbf{c}^\top \boldsymbol{\beta} = \mathbb{E}(\mathbf{b}^\top \mathbf{Y})$ valamely \mathbf{b} -re. Azaz $\mathbf{c}^\top \boldsymbol{\beta} = \mathbf{b}^\top X \boldsymbol{\beta} \quad \forall \boldsymbol{\beta}$. Tehát $\mathbf{c}^\top = \mathbf{b}^\top X$.

Megfordítva, $\mathbf{c}^\top \boldsymbol{\beta} = \mathbf{b}^\top X \boldsymbol{\beta}$ becslhető $\mathbf{b}^\top \mathbf{Y}$ -nal.

2. Legyen $\mathbf{b}^\top \mathbf{Y}$ a $\mathbf{c}^\top \boldsymbol{\beta}$ lineáris torzítatlan becslése. Ekkor tekintsük a

$$\mathbf{b} = (\mathbf{b} - P_F \mathbf{b}) + P_F \mathbf{b}$$

ortogonális felbontást. Ezt \mathbf{Y} -nal jobbról szorozva és várható értéket véve:

$$\mathbb{E} \mathbf{b}^\top \mathbf{Y} = (\mathbf{b} - P_F \mathbf{b})^\top X \boldsymbol{\beta} + \mathbb{E}(P_F \mathbf{b})^\top \mathbf{Y}.$$

Mivel $\mathbf{b} - P_F \mathbf{b}$ és $X \boldsymbol{\beta}$ egymásra merőleges és $\mathbb{E} \mathbf{b}^\top \mathbf{Y} = \mathbf{c}^\top \boldsymbol{\beta}$, így

$$\mathbf{c}^\top \boldsymbol{\beta} = \mathbb{E} \mathbf{b}^{*\top} \mathbf{Y}$$

adódik, ahol $\mathbf{b}^* = P_F \mathbf{b}$.

Legyen most $\mathbf{b}_1^\top \mathbf{Y}$ a $\mathbf{c}^\top \boldsymbol{\beta}$ egy másik lineáris torzítatlan becslése. Ekkor

$$0 = \mathbf{c}^\top \boldsymbol{\beta} - \mathbf{c}^\top \boldsymbol{\beta} = \mathbb{E} \mathbf{b}^{*\top} \mathbf{Y} - \mathbb{E} \mathbf{b}_1^\top \mathbf{Y} = (\mathbf{b}^{*\top} - \mathbf{b}_1^\top) X \boldsymbol{\beta} \quad \forall \boldsymbol{\beta}.$$

Azaz $\mathbf{b}^* - \mathbf{b}_1$ merőleges X oszlopaira. Tehát $\mathbf{b}_1 = \mathbf{b}^* + (\mathbf{b}_1 - \mathbf{b}^*)$ a merőleges felbontás, azaz \mathbf{b}^* a \mathbf{b}_1 -nek is vetülete.

3. Az előzőekben kaptuk, hogy $\mathbb{E} \mathbf{b}^{*\top} \mathbf{Y} = \mathbf{c}^\top \boldsymbol{\beta}$, azaz $\mathbf{b}^{*\top} \mathbf{Y}$ lineáris torzítatlan becslés. Az ortogonalitás alapján:

$$\mathbb{D}^2 \mathbf{b}^\top \mathbf{Y} = \sigma^2 \mathbf{b}^\top \mathbf{b} = \sigma^2 [\mathbf{b}^{*\top} \mathbf{b}^* + (\mathbf{b} - \mathbf{b}^*)^\top (\mathbf{b} - \mathbf{b}^*)] \geq \sigma^2 \mathbf{b}^{*\top} \mathbf{b}^* = \mathbb{D}^2 (\mathbf{b}^{*\top} \mathbf{Y}).$$

Tehát $\mathbf{b}^{*\top} \mathbf{Y}$ szórása a legkisebb.

4. Mivel $X \widehat{\boldsymbol{\beta}} = P_F \mathbf{Y}$, így $\mathbf{Y} - X \widehat{\boldsymbol{\beta}}$ merőleges X oszlopai alterére. Viszont \mathbf{b}^* benne van ebben az alterben. Tehát $\mathbf{b}^{*\top} (\mathbf{Y} - X \widehat{\boldsymbol{\beta}}) = 0$. Azaz

$$\mathbf{b}^{*\top} \mathbf{Y} = \mathbf{b}^{*\top} X \widehat{\boldsymbol{\beta}} = \mathbf{c}^\top \widehat{\boldsymbol{\beta}},$$

mert $\mathbf{c} = X^\top \mathbf{b}$ és $X^\top \mathbf{b} = X^\top \mathbf{b}^*$, lévén \mathbf{b}^* a \mathbf{b} vetülete X oszlopai alterére. \square

2. A lineáris modell normális eloszlás esetén

2.1. Maximum-likelihood becslés

2.1. tétel. Ha $\varepsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I)$, akkor β maximum-likelihood becslése $\hat{\beta}$, σ^2 maximum-likelihood becslése pedig $\frac{1}{n} \|\mathbf{Y} - X\hat{\beta}\|^2$.

BIZONYÍTÁS. Mivel $\mathbf{Y} \sim \mathcal{N}_n(X\beta, \sigma^2 I)$, így \mathbf{Y} sűrűségfüggvénye:

$$f(\mathbf{y}) = (2\pi)^{-n/2} (\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{y} - X\beta\|^2 \right\}.$$

Ennek β szerinti maximuma az exponensben lévő kifejezés, nevezetesen $\|\mathbf{Y} - X\beta\|^2$ minimuma helyén van, azaz $\hat{\beta}$ -ban.

A σ^2 szerinti minimumához

$$0 = \frac{\partial \log f}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2} \frac{1}{(\sigma^2)^2} \|\mathbf{Y} - X\hat{\beta}\|^2.$$

Innen $\hat{\sigma}^2 = \frac{1}{n} \|\mathbf{Y} - X\hat{\beta}\|^2$. □

2.1. megjegyzés. $\hat{\sigma}^2 = \frac{1}{n} \|\mathbf{Y} - X\hat{\beta}\|^2$ a becsült β esetén a **maradék négyzetösszeg**, (RSS=Residual Sum of Squares) azaz $\|\mathbf{Y} - X\hat{\beta}\|^2$ n -ed része. □

2.1. lemma. Legyen V egy s -dimenziós altérre való merőleges vetítés mátrixa, $\eta \sim \mathcal{N}_n(\mathbf{m}, I)$. Ekkor $V\eta$ és $(I - V)\eta$ függetlenek és normális eloszlásúak. Továbbá

$$\|V(\eta - \mathbf{m})\|^2 \sim \chi_s^2.$$

BIZONYÍTÁS. Nyilván $\begin{pmatrix} V\eta \\ (I - V)\eta \end{pmatrix} = \begin{pmatrix} V \\ I - V \end{pmatrix} \eta$, azaz normális eloszlású vektor lineáris transzformáltja, tehát normális eloszlású. Másrészt

$$\text{cov}(V\eta, (I - V)\eta) = VI(I - V) = 0,$$

azaz korrelálatlanok, tehát függetlenek.

Legyen $\mathbf{e}_1, \dots, \mathbf{e}_s$ az s -dimenziós alterünk ortonormált bázisa. Ezt egészítsük ki ortonormált bázissá: $\mathbf{e}_1, \dots, \mathbf{e}_n$. Ezen vektorokat foglaljuk az U ortogonális mátrixba. Ekkor $V = U\Lambda U^\top$, ahol Λ olyan diagonális mátrix, melynek első s főátlós eleme 1, a többi 0.

$$\|V(\eta - \mathbf{m})\|^2 = \|U\Lambda U^\top(\eta - \mathbf{m})\|^2 = \|U\Lambda\xi\|^2 = \xi^\top \Lambda U^\top U \Lambda \xi,$$

ahol $\xi \sim \mathcal{N}_n(\mathbf{0}, I)$. De $U^\top U = I$ és $\Lambda\Lambda = \Lambda$, tehát

$$\|V(\eta - \mathbf{m})\|^2 = \xi_1^2 + \dots + \xi_s^2 \sim \chi_s^2. \quad \square$$

2.2. tétel. Legyen $\mathbf{Y} \sim \mathcal{N}_n(X\beta, \sigma^2 I)$, $r = \text{rang } X$. Ekkor

$$\frac{1}{\sigma^2} \|\mathbf{Y} - X\hat{\beta}\|^2 \sim \chi_{n-r}^2,$$

$\frac{\|\mathbf{Y} - X\hat{\beta}\|^2}{n-r}$ a σ^2 paraméter torzítatlan becslése.

BIZONYÍTÁS. Emlékeztetünk, hogy F az X oszlopvektorai által kifeszített altér, P_F pedig az F -re merőleges vetítés.

$\mathbf{Y} - X\hat{\boldsymbol{\beta}} = \mathbf{Y} - P_F\mathbf{Y} = (I - P_F)\mathbf{Y}$. Mivel $I - P_F$ az F ortogonális komplementerére való merőleges vetítés, így $\mathbb{E}(\mathbf{Y} - X\hat{\boldsymbol{\beta}}) = (I - P_F)\mathbb{E}\mathbf{Y} = (I - P_F)X\boldsymbol{\beta} = \mathbf{0}$. Tehát $\mathbf{Y} - X\hat{\boldsymbol{\beta}} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2(I - P_F))$.

Tudjuk, hogy $(I - P_F)(\mathbf{Y} - \mathbb{E}\mathbf{Y}) = (I - P_F)\mathbf{Y} = \mathbf{Y} - X\hat{\boldsymbol{\beta}}$. Az 2.1. lemma miatt ezek közül a legelső hossz négyzetének $1/\sigma^2$ -szerese χ_{n-r}^2 eloszlású, így $\frac{1}{\sigma^2}\|\mathbf{Y} - X\hat{\boldsymbol{\beta}}\|^2 \sim \chi_{n-r}^2$.

Mivel $\mathbb{E}\chi_{n-r}^2 = n - r$, ezért $\mathbb{E}\frac{\|\mathbf{Y} - X\hat{\boldsymbol{\beta}}\|^2}{n-r} = \sigma^2$. \square

2.2. megjegyzés. $\mathbf{Y} - X\hat{\boldsymbol{\beta}}$ és $X\hat{\boldsymbol{\beta}}$ függetlenek.

Ugyanis $\mathbf{Y} - X\hat{\boldsymbol{\beta}} = (I - P_F)\mathbf{Y}$ és $X\hat{\boldsymbol{\beta}} = P_F\mathbf{Y}$. \square

2.2. Hipotézisvizsgálat a lineáris modellben

Legyen $Y = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, ahol $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 I)$. A nullhipotézis:

$$H_0 : B\boldsymbol{\beta} = \mathbf{0},$$

ahol a B mátrix ismert és sorai becsülhetőek.

2.3. megjegyzés. Legyen az X oszlopai által generált altér F . Legyenek B sorai becsülhetőek. Ekkor létezik olyan F_0 altére F -nek, hogy $B\boldsymbol{\beta} = \mathbf{0} \iff X\boldsymbol{\beta} \in F_0$.

Ennek bizonyítására megjegyezzük, hogy B sorai becsülhetőek $\iff B = B^*X$ valamely B^* mátrixra. Ezért $B\boldsymbol{\beta} = \mathbf{0} \iff B^*X\boldsymbol{\beta} = \mathbf{0} \iff X\boldsymbol{\beta}$ a B^* mátrix nullterében van, azaz $X\boldsymbol{\beta} \in F_0$, ahol F_0 a $B^*|_F$ nulltere. \square

2.3. tétel. Legyen

$$S_B^2 = \min_{\boldsymbol{\beta}: B\boldsymbol{\beta}=\mathbf{0}} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2, \quad S_e^2 = \min_{\boldsymbol{\beta}} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2$$

a hiba négyzetösszeg feltételes, illetve feltétel nélküli minimuma. Legyen $\text{rang } X = r$, $\text{rang}(F - F_0) = q$. (Itt $F - F_0$ az F_0 altérnek az F altérre vett ortogonális komplementerét jelöli.)

Ekkor

$$\frac{S_B^2 - S_e^2}{q} : \frac{S_e}{n-r} \sim F_{q, n-r}$$

akkor és csak akkor, ha H_0 igaz.

BIZONYÍTÁS. Tekintsük az

$$\mathbf{Y} - P_{F_0}\mathbf{Y} = (\mathbf{Y} - P_F\mathbf{Y}) + (P_F\mathbf{Y} - P_{F_0}\mathbf{Y})$$

felbontást. Mivel $F_0 \subseteq F$, így $P_F\mathbf{Y} - P_{F_0}\mathbf{Y}$ benne van F -ben. $\mathbf{Y} - P_F\mathbf{Y}$ viszont merőleges F -re. Tehát a fenti felbontás ortogonális felbontás. Ezért

$$\|\mathbf{Y} - P_{F_0}\mathbf{Y}\|^2 = \|\mathbf{Y} - P_F\mathbf{Y}\|^2 + \|P_F\mathbf{Y} - P_{F_0}\mathbf{Y}\|^2.$$

Másrészt

$$S_B^2 = \min_{\boldsymbol{\beta}: B\boldsymbol{\beta}=\mathbf{0}} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2 = \min_{\boldsymbol{\beta}: X\boldsymbol{\beta} \in F_0} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2 = \|\mathbf{Y} - P_{F_0}\mathbf{Y}\|^2$$

és

$$S_e^2 = \min_{\boldsymbol{\beta}} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2 = \|\mathbf{Y} - P_F \mathbf{Y}\|^2 = \|(I - P_F)\mathbf{Y}\|^2.$$

Így

$$S_B^2 - S_e^2 = \|P_F \mathbf{Y} - P_{F_0} \mathbf{Y}\|^2 = \|(P_F - P_{F_0})\mathbf{Y}\|^2.$$

Tehát az $S_B^2 - S_e^2$ -ben szereplő lineáris kifejezés $(P_F - P_{F_0})\mathbf{Y}$, az S_e^2 -ben szereplő pedig $(I - P_F)\mathbf{Y}$. Ezek együttesen normális eloszlásúak. Korrelálatlanok, mert $(I - P_F)(P_F - P_{F_0}) = 0$. Ezért függetlenek.

$$\mathbb{E}(I - P_F)\mathbf{Y} = (I - P_F)X\boldsymbol{\beta} = \mathbf{0}$$

mindig teljesül. Viszont

$$\mathbb{E}(P_F - P_{F_0})\mathbf{Y} = P_F X\boldsymbol{\beta} - P_{F_0} X\boldsymbol{\beta} = X\boldsymbol{\beta} - P_{F_0} X\boldsymbol{\beta} = \mathbf{0} \iff B\boldsymbol{\beta} = \mathbf{0},$$

mivel ez utóbbi $X\boldsymbol{\beta} \in F_0$ -lal ekvivalens.

Tehát $\frac{1}{\sigma^2} S_e^2 \sim \chi_{n-r}^2$ (ahol r az X rangja), hiszen $I - P_F$ $(n - r)$ -dimenziós altérbe vetít.

Ha $H_0 : B\boldsymbol{\beta} = \mathbf{0}$ teljesül, akkor $\frac{1}{\sigma^2} (S_B^2 - S_e^2) \sim \chi_q^2$, hiszen $P_F - P_{F_0}$ q -dimenziós altérbe vetít. \square

2.3. Az egyszeres osztályozás mint speciális lineáris modell

A lineáris modell alapján elvégezzük az egyszeres osztályozás elemzését.

2.1. példa. Legyen $Y_{ij} = a_i + \varepsilon_{ij}$, $j = 1, \dots, n_i$, $i = 1, \dots, p$, ahol $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$ és a ε_{ij} változók függetlenek. Ekkor

$$\mathbf{Y} = X\mathbf{a} + \boldsymbol{\varepsilon}$$

a modell, ahol

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ Y_{pn_p} \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & & & 1 \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \\ \vdots \\ \varepsilon_{pn_p} \end{pmatrix}, \quad \mathbf{a} = \begin{pmatrix} a_1 \\ \vdots \\ a_p \end{pmatrix}.$$

Itt \mathbf{Y} és $\boldsymbol{\varepsilon}$ n -dimenziós vektorok, X pedig $n \times p$ méretű mátrix, ahol $n = n_1 + \dots + n_p$.

A

$$H : a_1 = a_2 = \dots = a_p$$

nullhipotézis akkor és csak akkor teljesül, ha $B\mathbf{a} = \mathbf{0}$, ahol

$$B = \begin{pmatrix} 1 & & 0 & -1 \\ & \ddots & & \vdots \\ 0 & & 1 & -1 \end{pmatrix}$$

$(p-1) \times p$ típusú mátrix. B sorai becsülhetőek, hiszen $B = B^*X$, ahol

$$B^* = \begin{pmatrix} 1 & 0 & \dots & \dots & 0 & -1 \\ 0 & \dots & 0 & 1 & 0 & \dots & 0 & -1 \\ \vdots & & & & & & & \end{pmatrix}.$$

$$S_e^2 = \min_{\mathbf{a}} \|Y - X\mathbf{a}\|^2 = \min_{\mathbf{a}} \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - a_i)^2.$$

Ennek minimumhelye (deriválással meghatározva) $\hat{a}_i = \bar{Y}_{i.}$, azaz az i -edik szinten vett átlag. Így $S_e^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{i.})^2$. Ez éppen a Q_2 hiba négyzetösszeg. Ennek szabadsági foka a 2.2. tétel alapján $n - p$, mivel X rangja p .

$$S_B^2 = \min_{\mathbf{a}: B\mathbf{a}=\mathbf{0}} \sum_i \sum_j (Y_{ij} - a_i)^2.$$

Multiplikátorokkal dolgozunk. $\lambda_1, \dots, \lambda_p$ a multiplikátorok. A Lagrange-függvény:

$$\sum_i \sum_j (Y_{ij} - a_i)^2 + \sum_{i=1}^{p-1} \lambda_i (a_i - a_p) \quad \text{és} \quad a_i = a_p, \quad i = 1, \dots, p-1.$$

a_i szerint deriválva a fenti függvényt:

$$-2 \sum_j (Y_{ij} - a_i) + \lambda_i = 0, \quad i = 1, \dots, p-1,$$

és

$$-2 \sum_j (Y_{pj} - a_p) - \sum_{i=1}^{p-1} \lambda_i = 0.$$

Ezeket összeadva:

$$\sum_{i=1}^{p-1} \sum_{j=1}^{n_i} (Y_{ij} - a_i) + \sum_{j=1}^{n_p} (Y_{pj} - a_p) = 0.$$

Felhasználva az $a_i = a_p$ ($i = 1, \dots, p-1$) feltételt: $a_p = \frac{1}{n} \sum_i \sum_j Y_{ij} = \bar{Y}_{..}$, azaz minden a_i becslése most az $\bar{Y}_{..}$ teljes átlag. Tehát

$$S_B^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2 = Q$$

teljes négyzetösszeg.

$$\begin{aligned} S_B^2 - S_e^2 &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2 - \sum_i \sum_j (Y_{ij} - \bar{Y}_{i.})^2 = \\ &= Q - Q_2 = Q_1 = \sum_i \sum_j (\bar{Y}_{i.} - \bar{Y}_{..})^2. \end{aligned}$$

Azaz $S_B^2 - S_e^2 = Q_1$, a szintek közötti eltérést mérő négyzetösszeg.

$B^*|_F$ nulltere az azonos koordinátájú vektorokból áll, ami egydimenziós. F az X oszlopai által generált altér. Ezért F p -dimenziós. Így $\text{rang}(F - F_0) = p - 1$.

Tehát az egyszeres osztályozás elemezhető a lineáris modell segítségével. A 2.3. tétel alapján adódó F -próba ugyanaz lesz, mint a korábban a Fischer-Cochran-tétel alapján kapott F -próba. \square

III. fejezet

A maximum-likelihood módszer

1. A Rao-Cramér-féle egyenlőtlenség

1.1. A likelihood- és a loglikelihood-függvény. A regularitási feltételek

Felsoroljuk azokat a regularitási feltételeket, amelyek esetén a Rao-Cramér-egyenlőtlenséget bizonyítani fogjuk.

Legyen a Θ paramétertér nyílt halmaz \mathbb{R}^k -ban. Legyen $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_k)^\top \in \Theta$.

Legyen X minta (általános értelemben).

Legyen $f(x, \boldsymbol{\vartheta})$ a minta sűrűségfüggvénye (általános értelemben, azaz egy ν domináló mértékre nézve).

Tegyük fel, hogy az $f(x, \boldsymbol{\vartheta})$ függvény $\boldsymbol{\vartheta}$ szerint kétszer differenciálható, és az

$$\int f(x, \boldsymbol{\vartheta}) d\nu(x) = 1$$

kifejezésben a $\boldsymbol{\vartheta}$ szerinti első és második deriválás és az integrálás sorrendje fölcserélhető.

Ekkor

$$(1.1) \quad \int \frac{\partial f(x, \boldsymbol{\vartheta})}{\partial \vartheta_i} d\nu(x) = 0, \quad i = 1, \dots, k,$$

és

$$(1.2) \quad \int \frac{\partial^2 f(x, \boldsymbol{\vartheta})}{\partial \vartheta_i \partial \vartheta_j} d\nu(x) = 0, \quad i, j = 1, \dots, k.$$

Tegyük fel, hogy

$$(1.3) \quad \mathbb{E} \left(\frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_i} \right)^2 < \infty, \quad i = 1, \dots, k.$$

Ekkor

$$(1.4) \quad I_{ij}(\boldsymbol{\vartheta}) = \mathbb{E} \left[\frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_i}, \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_j} \right]$$

véges minden $i, j = 1, \dots, k$ esetén.

1.1. definíció. Az

$$I(\boldsymbol{\vartheta}) = \left(I_{ij}(\boldsymbol{\vartheta}) \right)_{i,j=1}^k$$

mátrixot **Fisher-féle információs mátrixnak** nevezzük.

Az $f(x, \boldsymbol{\vartheta})$ függvényt **likelihood-függvénynek** hívjuk, az $\ln f(x, \boldsymbol{\vartheta})$ függvényt pedig **loglikelihood-függvénynek**. A

$$(1.5) \quad \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} = \left(\frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_1}, \dots, \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_k} \right)^\top$$

valószínűségi vektorváltozó az ún. **score vektor**. □

A fenti feltételek esetén a score vektor várható értéke $\mathbf{0}$. Ugyanis

$$(1.6) \quad \mathbb{E} \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} = \int \frac{1}{f(x, \boldsymbol{\vartheta})} \frac{\partial f(x, \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} f(x, \boldsymbol{\vartheta}) d\nu(x) = \mathbf{0}$$

az (1.1) feltétel miatt. (Ne tévesszük szem elől, hogy a várható érték \mathbb{E} jele után az X véletlen mennyiség szerepel az aktuális függvényben, amikor viszont az \int integrálra térünk, akkor pedig az x változó szerepel ugyanabban a függvényben. Természetesen az utóbbi esetben az f sűrűségfüggvénnyel szoroznunk kell, mielőtt a ν mérték szerint integrálnánk.)

1.1. megjegyzés. Az előzőből azonnal adódik, hogy az $I(\boldsymbol{\vartheta})$ Fisher-féle információs mátrix éppen a score vektor szórásátrixa:

$$(1.7) \quad I(\boldsymbol{\vartheta}) = \text{var} \left(\frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} \right).$$

Másrészt a Fisher-féle információs mátrix a score vektor segítségével előáll az

$$(1.8) \quad I_{ij}(\boldsymbol{\vartheta}) = -\mathbb{E} \left[\frac{\partial^2 \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_i \partial \vartheta_j} \right]$$

szerint is. Valóban, (1.2) felhasználásával

$$\begin{aligned} \mathbb{E} \frac{\partial^2 \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_i \partial \vartheta_j} &= \mathbb{E} \frac{\partial}{\partial \vartheta_j} \left(\frac{1}{f}, \frac{\partial f}{\partial \vartheta_i} \right) = \mathbb{E} \left[\frac{1}{f} \frac{\partial^2 f}{\partial \vartheta_i \partial \vartheta_j} - \frac{1}{f^2} \frac{\partial f}{\partial \vartheta_i} \frac{\partial f}{\partial \vartheta_j} \right] = \\ &= \int \frac{\partial^2 f(x, \boldsymbol{\vartheta})}{\partial \vartheta_i \partial \vartheta_j} d\nu(x) - \mathbb{E} \frac{1}{f} \frac{\partial f}{\partial \vartheta_i} \cdot \frac{1}{f} \frac{\partial f}{\partial \vartheta_j} = 0 - I_{ij}(\boldsymbol{\vartheta}). \quad \square \end{aligned}$$

Folytatjuk a Rao-Cramér-egyenlőtlenség feltételeinek felsorolását.

Feltesszük, hogy $I(\boldsymbol{\vartheta})$ invertálható.

Legyen $\mathbf{g} : \Theta \rightarrow \mathbb{R}^r$ differenciálható függvény.

Legyen a \mathbf{T} r -dimenziós statisztika a $\mathbf{g}(\boldsymbol{\vartheta})$ torzítatlan becslése:

$$(1.9) \quad \int \mathbf{T}(x) f(x, \boldsymbol{\vartheta}) d\nu(x) = \mathbf{g}(\boldsymbol{\vartheta}).$$

Tegyük fel, hogy ebben az egyenletben a $\boldsymbol{\vartheta}$ szerinti differenciálás és az integrálás sorrendje felcserélhető:

$$(1.10) \quad \int \mathbf{T}(x) \frac{\partial f(x, \boldsymbol{\vartheta})}{\partial \vartheta_i} d\nu(x) = \frac{\partial \mathbf{g}(\boldsymbol{\vartheta})}{\partial \vartheta_i}, \quad i = 1, \dots, k.$$

1.2. A Rao-Cramér-féle egyenlőtlenség

A Rao-Cramér-féle egyenlőtlenség azt állítja, hogy – a regularitási feltételek esetén – a torzítatlan becslés szórása nem lehet akármilyen kicsi. Egy (abszolút) alsó határ adható a Fisher-féle információ segítségével.

1.1. tétel (Rao-Cramér-egyenlőtlenség). *A fenti regularitási feltételek esetén*

$$(1.11) \quad \boxed{\text{var}(\mathbf{T}) \geq \mathbf{G} I^{-1}(\boldsymbol{\vartheta}) \mathbf{G}^\top},$$

ahol $G = \left(\frac{\partial g(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} \right)$.

BIZONYÍTÁS. Tekintsük a

$$(1.12) \quad \left(T_1(X), \dots, T_r(X), \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_1}, \dots, \frac{\partial \log f(X, \boldsymbol{\vartheta})}{\partial \vartheta_k} \right)$$

$(r + k)$ dimenziós statisztika transzponáltját. Ennek szórásmatrixa az előzőek szerint

$$(1.13) \quad \begin{pmatrix} V & G \\ G^\top & I(\boldsymbol{\vartheta}) \end{pmatrix},$$

ahol \mathbf{T} szórásmatrixa V , $G = \left(\frac{\partial g(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} \right)$ pedig $r \times k$ típusú mátrix. Mivel az (1.13) szórásmatrix, így determinánsa nem negatív.

Ezért

$$\begin{aligned} 0 &\leq \det \begin{pmatrix} E & -GI^{-1}(\boldsymbol{\vartheta}) \\ 0 & I^{-1}(\boldsymbol{\vartheta}) \end{pmatrix} \det \begin{pmatrix} V & G \\ G^\top & I(\boldsymbol{\vartheta}) \end{pmatrix} = \\ &= \det \begin{pmatrix} V - GI^{-1}(\boldsymbol{\vartheta})G^\top & 0 \\ I^{-1}(\boldsymbol{\vartheta})G^\top & E \end{pmatrix} = \det(V - GI^{-1}(\boldsymbol{\vartheta})G^\top), \end{aligned}$$

ahol E és 0 alkalmas méretű egység-, illetve nullmatrixot jelöl.

Ez a gondolatmenet igaz akkor is, ha csupán a \mathbf{T} vektor néhány koordinátáját tekintjük. Ezért a $(V - GI^{-1}(\boldsymbol{\vartheta})G^\top)$ mátrix minden, a főátlóra szimmetrikusan elhelyezkedő részmatrixának nem negatív a determinánsa. Tehát $(V - GI^{-1}(\boldsymbol{\vartheta})G^\top)$ pozitív szemidefinit. Ezzel beláttuk a tételt. \square

1.1. feladat. Legyen ξ_1, \dots, ξ_n minta $\mathcal{N}(\mu, \sigma^2)$ -ből. A paraméter kétdimenziós: $\boldsymbol{\vartheta} = (\mu, \sigma^2)^\top$. Lássuk be az alábbiakat.

(a) A Fisher-féle információs mátrix:

$$(1.14) \quad I(\boldsymbol{\vartheta}) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}, \quad I^{-1}(\boldsymbol{\vartheta}) = \begin{pmatrix} \frac{\sigma^2}{n} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix}.$$

(b) $\bar{\xi}$ és s_n^{*2} torzítatlan becslés (μ, σ^2) -re. A becslés szórásmatrixa:

$$(1.15) \quad \mathbb{D}^2 \begin{pmatrix} \bar{\xi} \\ s_n^{*2} \end{pmatrix} = \begin{pmatrix} \frac{\sigma^2}{n} & 0 \\ 0 & \frac{2\sigma^4}{n-1} \end{pmatrix}.$$

(c) Bár a Rao-Cramér-féle alsó határ nem éretik el, s_n^{*2} mégis a legkisebb szórású torzítatlan becslés σ^2 -re.

MEGOLDÁS. (c) Legyen t torzítatlan becslés σ^2 -re: $\mathbb{E}t = \sigma^2$. Legyen $u = t - s_n^{*2}$, ahonnan $t = u + s_n^{*2}$.

Viszont t és s_n^{*2} várható értéke σ^2 , így u várható értéke 0 . De a minta együttes sűrűségfüggvénye alapján:

$$(1.16) \quad 0 = \mathbb{E}u = \int u(\mathbf{x}) \cdot f(\mathbf{x}) d\mathbf{x} = \int u(\mathbf{x}) \frac{1}{(\sigma^2 2\pi)^{\frac{n}{2}}} e^{-\frac{\sum (x_i - \mu)^2}{2\sigma^2}} d\mathbf{x},$$

ahol $\mathbf{x} = (x_1, \dots, x_n)^\top$. μ szerint deriválva kétszer (közben (1.16)-ot még egyszer használva):

$$(1.17) \quad \int \underbrace{\left(\sum (x_i - \mu)\right)^2}_{(n(\bar{x} - \mu))^2} u(\mathbf{x}) \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n e^{-\frac{\sum(x_i - \mu)^2}{2\sigma^2}} d\mathbf{x} = 0.$$

σ^2 szerint deriválva (1.16)-ot (közben (1.16)-ot még egyszer használva):

$$(1.18) \quad \int \sum (x_i - \mu)^2 u(\mathbf{x}) \frac{1}{(2\pi)^{n/2}} e^{-\frac{\sum(x_i - \mu)^2}{2\sigma^2}} d\mathbf{x} = 0.$$

Viszont a Steiner-formula szerint $(n-1)s_n^{*2} = \sum (x_i - \mu)^2 - n(\bar{x} - \mu)^2$, így az előző két egyenlet miatt u és s_n^{*2} korrelálatlan. Ezért

$$\mathbb{D}^2(t) = \mathbb{D}^2(u) + \mathbb{D}^2(s_n^{*2}) \geq \mathbb{D}^2(s_n^{*2}). \quad \square$$

2. A maximum-likelihood becslés

2.1. A maximum-likelihood becslés definíciója

2.1. definíció. Legyen X minta (általános értelemben), x a minta realizáció, Θ a paraméterter. Legyen ν mérték a mintatéren, tegyük fel, hogy a minta eloszlása abszolút folytonos ν -re nézve, jelölje $\ell(x, \boldsymbol{\vartheta})$ a sűrűségfüggvényt.

$\ell(x, \boldsymbol{\vartheta})$ -t – mint $\boldsymbol{\vartheta} \in \Theta$ függvényét tekintve – **likelihood-függvénynek** nevezzük. \square

2.2. definíció. A $\boldsymbol{\vartheta}$ paraméter **maximum-likelihood becslésének** azt a $\hat{\boldsymbol{\vartheta}} = \hat{\boldsymbol{\vartheta}}(x)$ statisztikát nevezzük, melyre

$$(2.1) \quad \ell(x, \hat{\boldsymbol{\vartheta}}) = \sup_{\boldsymbol{\vartheta} \in \Theta} \ell(x, \boldsymbol{\vartheta}). \quad \square$$

2.1. megjegyzés. Ha létezik $T(x)$ elégséges statisztika, akkor a maximum-likelihood becslés ennek függvénye. Valóban, a Neyman-Fisher-féle faktorizációs tétel szerint

$$\ell(x, \boldsymbol{\vartheta}) = h(x)g(T(x), \boldsymbol{\vartheta}),$$

így a maximum hely csak $T(x)$ -től függ.

A fenti állítás megfordítása nem igaz, a maximum-likelihood becslés nem mindig elégséges statisztika. \square

2.3. definíció. A likelihood-függvény logaritmusát **loglikelihood-függvénynek** nevezzük:

$$(2.2) \quad L(x, \boldsymbol{\vartheta}) = \log \ell(x, \boldsymbol{\vartheta}). \quad \square$$

Ha $\Theta \subseteq \mathbb{R}^k$ nyílt halmaz és $\ell(x, \boldsymbol{\vartheta})$ $\boldsymbol{\vartheta}$ szerint differenciálható, akkor $L(x, \boldsymbol{\vartheta})$ maximum helyén a $\boldsymbol{\vartheta}$ szerinti parciális deriváltak eltűnnek.

2.4. definíció. A

$$(2.3) \quad \frac{\partial L(x, \boldsymbol{\vartheta})}{\partial \vartheta_i} = 0, \quad i = 1, \dots, k,$$

egyenleteket **likelihood-egyenleteknek** nevezzük. \square

2.2. megjegyzés. Ha \boldsymbol{x} független valószínűségi változókból álló vektor, akkor a loglikelihood-függvény:

$$(2.4) \quad L_n = L(\boldsymbol{x}, \boldsymbol{\vartheta}) = \sum_{i=1}^n \log f_i(x_i, \boldsymbol{\vartheta}),$$

ahol $f_i(x_i, \boldsymbol{\vartheta})$ az \boldsymbol{x} vektor i -edik koordinátájának a sűrűségfüggvénye. \square

2.2. A maximum-likelihood módszer korlátai

2.3. megjegyzés. Lehetséges, hogy a likelihood-függvény nem korlátos felülről, azaz nem létezik maximum-likelihood becslés. \square

2.4. megjegyzés. Lehetséges, hogy a likelihood-egyenleteknek van olyan gyöke, mely nem konzisztens becslése a paraméternek. Ezt az alábbi példa mutatja. \square

2.1. példa. Legyen X_1, \dots, X_n minta $\mathcal{N}(\vartheta, c^2\vartheta^2)$ eloszlásból, ahol $c \neq 0$ ismert, ϑ ismeretlen paraméter ($\vartheta \neq 0$). Ekkor a loglikelihood függvény:

$$L_n(\vartheta) = n \log \frac{1}{\sqrt{2\pi}(c\vartheta)} - \sum_{i=1}^n \frac{(x_i - \vartheta)^2}{2c^2\vartheta^2}.$$

$$\frac{\partial L_n(\vartheta)}{\partial \vartheta} = -n\vartheta^{-1} + \frac{1}{2c^2} \sum_{i=1}^n [2(x_i - \vartheta)\vartheta^{-2} + (x_i - \vartheta)^2 2\vartheta^{-3}].$$

$$\frac{\partial L_n(\vartheta)}{\partial \vartheta} = 0$$

alapján

$$nc^2\vartheta^2 + \vartheta \sum_{i=1}^n x_i - \sum_{i=1}^n x_i^2 = 0.$$

Innen

$$\hat{\vartheta}_{12} = \frac{-\sum_{i=1}^n X_i \pm \sqrt{(\sum_{i=1}^n X_i)^2 + 4nc^2 \sum_{i=1}^n X_i^2}}{2nc^2}$$

a megoldás.

Ezen statisztikák határértékét $n \rightarrow \infty$ estén a nagy számok törvénye alapján határozzuk meg:

$$\frac{\sum_{i=1}^n X_i}{n} \rightarrow \vartheta \quad \text{és} \quad \frac{\sum_{i=1}^n X_i^2}{n} \rightarrow \vartheta^2(c^2 + 1).$$

Így $\hat{\vartheta}_1 \rightarrow \vartheta$ és $\hat{\vartheta}_2 \rightarrow -\vartheta(1 + \frac{1}{c^2})$.

Tehát ϑ_1 a likelihood-egyenletnek konzisztens gyöke, ϑ_2 pedig nem konzisztens gyöke. \square

2.5. megjegyzés. Lehetséges, hogy a maximum-likelihood becslés nem konzisztens. Erre utal az alábbi példa. \square

2.2. példa. Legyenek $\begin{pmatrix} Y_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ Y_n \end{pmatrix}$ független, kétdimenziós véletlen vektorok,

$$\begin{pmatrix} X_i \\ Y_i \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_i \\ \mu_i \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix} \right), \quad i = 1, \dots, n.$$

Ekkor a likelihood-függvény:

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu_i)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}}.$$

A loglikelihood-függvény:

$$L = -2n \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n [(x_i - \mu_i)^2 + (y_i - \mu_i)^2].$$

$$\frac{\partial L}{\partial \mu_i} = \frac{1}{\sigma^2} [(x_i - \mu_i) + (y_i - \mu_i)],$$

ahonnan $\hat{\mu}_i = \frac{X_i + Y_i}{2}$.

$$\frac{\partial L}{\partial \sigma^2} = -n(\sigma^2)^{-1} + (\sigma^2)^{-2} \frac{1}{2} \sum_{i=1}^n [(x_i - \mu_i)^2 + (y_i - \mu_i)^2].$$

Innen

$$\hat{\sigma}^2 = \frac{1}{2n} \sum_{i=1}^n [(X_i - \hat{\mu}_i)^2 + (Y_i - \hat{\mu}_i)^2] = \frac{1}{n} \sum_{i=1}^n \left(\frac{X_i - Y_i}{2} \right)^2.$$

Itt $\mathbb{E} \left(\frac{X_i - Y_i}{2} \right)^2 = \frac{\sigma^2}{2}$, ahonnan a nagy számok törvénye értelmében $\hat{\sigma} \rightarrow \frac{\sigma}{\sqrt{2}}$ ($n \rightarrow \infty$), azaz a becslés nem konzisztens.

Megjegyezzük, hogy a példában a paraméterek száma végtelenhez tart, ha $n \rightarrow \infty$. □

2.6. megjegyzés. Lehetséges, hogy a likelihood-egyenlet gyökei nem a likelihood-függvény maximumhelyét szolgáltatják. □

2.3. Egy konzisztencia tétel

Az alábbiakban feltesszük, hogy a mintaelemek függetlenek és azonos eloszlásúak. Tegyük fel, hogy különböző paraméter értékekhez különböző eloszlások tartoznak. Jelölje P_ϑ , illetve \mathbb{E}_ϑ a ϑ paraméternek megfelelő eloszlást, illetve várható értéket. Jelölje $f(x, \vartheta)$ egy mintaelem sűrűségfüggvényét. Jelölje ϑ^* a paraméter valódi értékét.

2.1. tétel. Legyen $\Theta \subseteq \mathbb{R}$ és legyen $f(x, \vartheta)$ deriválható ϑ szerint egy olyan intervallumon, mely ϑ^* -ot tartalmazza. Ekkor a likelihood-egyenletnek létezik olyan gyöke, mely $n \rightarrow \infty$ esetén ϑ^* -hoz tart (P_{ϑ^*} szerint) 1 valószínűséggel.

BIZONYÍTÁS. Tegyük fel, hogy $\log f(x, \vartheta)$ deriválható a $(\vartheta^* - \delta, \vartheta^* + \delta)$ intervallumon. Ismeretes, hogy $\log t \leq t - 1$. Itt t helyére $\frac{f(x, \vartheta^* \mp \delta)}{f(x, \vartheta^*)}$ -ot helyettesítve:

$$\int f(x, \vartheta^*) \log \frac{f(x, \vartheta^* - \delta)}{f(x, \vartheta^*)} d\nu(x) < 0,$$

$$\int f(x, \vartheta^*) \log \frac{f(x, \vartheta^* + \delta)}{f(x, \vartheta^*)} d\nu(x) < 0.$$

Azaz

$$\mathbb{E}_{\vartheta^*} \log f(X, \vartheta^* - \delta) < \mathbb{E}_{\vartheta^*} \log f(X, \vartheta^*),$$

$$\mathbb{E}_{\vartheta^*} \log f(X, \vartheta^* + \delta) < \mathbb{E}_{\vartheta^*} \log f(X, \vartheta^*).$$

Mivel L_n független, azonos eloszlású valószínűségi változók összege, azaz

$$L_n(X_1, \dots, X_n, \vartheta) = \sum_{i=1}^n \log f(X_i, \vartheta),$$

így a nagy számok törvénye alkalmazható. Tehát

$$\lim_{n \rightarrow \infty} \frac{1}{n} L_n(X_1, \dots, X_n, \vartheta^* + \delta) < \lim_{n \rightarrow \infty} \frac{1}{n} L_n(X_1, \dots, X_n, \vartheta^*),$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} L_n(X_1, \dots, X_n, \vartheta^* - \delta) < \lim_{n \rightarrow \infty} \frac{1}{n} L_n(X_1, \dots, X_n, \vartheta^*)$$

P_{ϑ^*} szerint 1 valószínűséggel.

Ezért elég nagy n -re L_n értéke a $(\vartheta^* - \delta, \vartheta^* + \delta)$ intervallum végpontjaiban kisebb, mint ϑ^* -ban. Így L_n -nek van a $(\vartheta^* - \delta, \vartheta^* + \delta)$ -ban lokális maximuma. Mivel $\delta > 0$ tetszőleges, így van 1 valószínűséggel konvergencia gyök. \square

2.4. A maximum-likelihood becslés aszimptotikája heurisztikus magyarázattal

Legyen $\vartheta \in \Theta \subseteq \mathbb{R}^k$ a paraméter, ϑ^* pedig ennek a valódi értéke. Az X_1, \dots, X_n független, azonos eloszlású mintaelemek esetén a loglikelihood függvény is független, azonos eloszlású tagok összege:

$$(2.5) \quad L_n(\vartheta) = \sum_{i=1}^n \log f(X_i, \vartheta),$$

ahol f a közös sűrűségfüggvény. Így a score vektor:

$$(2.6) \quad \frac{\partial L_n(\vartheta)}{\partial \vartheta} = \sum_{i=1}^n \frac{\partial \log f(X_i, \vartheta)}{\partial \vartheta}.$$

Itt az összeadandók független, azonos eloszlású, $\mathbf{0}$ várható értékű és $I(\vartheta)$ szórásmatrixú változók. (Itt $I(\vartheta)$ egyetlen megfigyelés esetén a Fisher-féle információs mátrix.) A központi határeloszlás tétel miatt $n \rightarrow \infty$ esetén

$$(2.7) \quad \frac{1}{\sqrt{n}} \frac{\partial L_n(\vartheta)}{\partial \vartheta} \implies b, \quad \text{ahol } b \sim \mathcal{N}_k(\mathbf{0}, I(\vartheta)).$$

Másrészt a loglikelihood-függvény második deriváltjából álló mátrix:

$$(2.8) \quad \frac{\partial^2 L_n(\vartheta)}{\partial \vartheta \partial \vartheta} = \sum_{i=1}^n \frac{\partial^2 \log f(X_i, \vartheta)}{\partial \vartheta \partial \vartheta}.$$

Ez független, $-I(\vartheta)$ várható értékű tagokból áll, így a nagy számok törvénye miatt:

$$(2.9) \quad \frac{1}{n} \frac{\partial^2 L_n(\vartheta)}{\partial \vartheta \partial \vartheta} \rightarrow -I(\vartheta)$$

majdnem biztosan.

Fejtsük Taylor-sorba a loglikelihood-függvényt a $\boldsymbol{\vartheta}^*$ igazi paraméter körül:

$$(2.10) \quad \begin{aligned} L_n(\boldsymbol{\vartheta}) \approx & L_n(\boldsymbol{\vartheta}^*) + \frac{1}{\sqrt{n}} \frac{\partial L_n(\boldsymbol{\vartheta}^*)^\top}{\partial \boldsymbol{\vartheta}} \sqrt{n}(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*) + \\ & + \frac{1}{2} \sqrt{n}(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*)^\top \left[\frac{1}{n} \frac{\partial^2 L_n(\boldsymbol{\vartheta}^*)}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} \right] \sqrt{n}(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*). \end{aligned}$$

A maradék tagot elhagytuk, az n és a \sqrt{n} szorzókat utólag írtuk be. Az előzőek alapján a központi határeloszlás tétel és a nagy számok törvénye alkalmazható. Ezért (2.10) aszimptotikusan

$$(2.11) \quad L_n(\boldsymbol{\vartheta}) \approx c + \mathbf{b}^\top \mathbf{y} - \frac{1}{2} \mathbf{y}^\top I(\boldsymbol{\vartheta}^*) \mathbf{y},$$

ahol $\mathbf{y} = \sqrt{n}(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*)$, $\mathbf{b} \sim \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*))$, $c = L_n(\boldsymbol{\vartheta}^*)$.

A (2.11) kifejezés más alakban:

$$L_n(\boldsymbol{\vartheta}) \approx c + \frac{1}{2} \mathbf{b}^\top I^{-1}(\boldsymbol{\vartheta}^*) \mathbf{b} - \frac{1}{2} \left(I^{\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{y} - I^{-\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{b} \right)^\top \left(I^{\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{y} - I^{-\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{b} \right).$$

Ezen kifejezés maximuma az $\mathbf{y} = I^{-1}(\boldsymbol{\vartheta}^*) \mathbf{b}$ pontban van, a maximum értéke pedig $c + \frac{1}{2} \mathbf{b}^\top I^{-1}(\boldsymbol{\vartheta}^*) \mathbf{b}$.

Tehát a $\hat{\boldsymbol{\vartheta}}_n$ maximum-likelihood becslésre

$$(2.12) \quad \sqrt{n}(\hat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \approx I^{-1}(\boldsymbol{\vartheta}^*) \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*)) = \mathcal{N}_k(\mathbf{0}, I^{-1}(\boldsymbol{\vartheta}^*)),$$

azaz $\hat{\boldsymbol{\vartheta}}_n$ **aszimptotikusan normális eloszlású**. Innen

$$(2.13) \quad \hat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^* \approx \mathcal{N}_k(\mathbf{0}, (nI(\boldsymbol{\vartheta}^*))^{-1}).$$

Ebből látszik, hogy $\hat{\boldsymbol{\vartheta}}_n$ szórása 0-hoz tart, várható értéke viszont $\boldsymbol{\vartheta}^*$ -hoz, tehát **konzisztens becslés**. Másrészt $\hat{\boldsymbol{\vartheta}}_n$ **aszimptotikusan legkisebb szórású**, hisz a Rao-Cramér-egyenlőtlenség miatt a (torzítatlan) becslések szórásának alsó határa $(nI(\boldsymbol{\vartheta}^*))^{-1}$. (Itt kihasználtuk, hogy n db független, azonos eloszlású megfigyeléshez tartozó Fisher-információ az 1 db megfigyeléshez tartozónak n -szerese.)

A loglikelihood-függvény aszimptotikája:

$$(2.14) \quad \max_{\boldsymbol{\vartheta}} 2 [L_n(\boldsymbol{\vartheta}) - L_n(\boldsymbol{\vartheta}^*)] \approx \mathbf{b}^\top I^{-1}(\boldsymbol{\vartheta}^*) \mathbf{b} = \left[I^{-\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{b} \right]^\top \left[I^{-\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{b} \right].$$

Mivel $I^{-\frac{1}{2}}(\boldsymbol{\vartheta}^*) \mathbf{b} \sim \mathcal{N}_k(\mathbf{0}, E)$, tehát a fenti kifejezés határeloszlása χ_k^2 . Ezen fog alapulni a **likelihood-hányados próba** határeloszlása.

2.7. megjegyzés. A likelihood egyenletből is adódik Taylor-sorfejtéssel az aszimptotika:

$$(2.15) \quad \mathbf{0} = \frac{1}{\sqrt{n}} \frac{\partial L_n(\hat{\boldsymbol{\vartheta}}_n)}{\partial \boldsymbol{\vartheta}} \approx \underbrace{\frac{1}{\sqrt{n}} \frac{\partial L_n(\boldsymbol{\vartheta}^*)}{\partial \boldsymbol{\vartheta}}}_{\downarrow \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*))} + \underbrace{\frac{1}{n} \frac{\partial^2 L_n(\boldsymbol{\vartheta}^*)}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}}}_{\downarrow -I(\boldsymbol{\vartheta}^*)} \sqrt{n}(\hat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*).$$

Tehát $n \rightarrow \infty$ esetén $\sqrt{n}(\hat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \approx \mathcal{N}_k(\mathbf{0}, I^{-1}(\boldsymbol{\vartheta}^*))$. □

2.5. A maximum-likelihood becslés aszimptotikája

Bebizonyítjuk, hogy független, X -szel azonos eloszlású, X_1, \dots, X_n mintaelemek esetén a likelihood-egyenletnek van olyan $\widehat{\boldsymbol{\vartheta}}_n$ gyöke, mely erősen konzisztens, aszimptotikusan normális eloszlású és aszimptotikusan legkisebb szórású. Legyen $f(x, \boldsymbol{\vartheta})$ az X sűrűségfüggvénye (általános értelemben, tehát pl. f jelölhet diszkrét eloszlást is).

2.2. tétel. *Tegyük fel, hogy*

I. $\Theta \subseteq \mathbb{R}^k$;

II. A $\boldsymbol{\vartheta}^*$ valódi paraméter egy ω környezetében léteznek az $f(x, \boldsymbol{\vartheta})$ sűrűségfüggvénynek $\boldsymbol{\vartheta}$ szerinti harmadik parciális deriváltjai;

III. $\mathbb{E}_{\boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \vartheta_j} \log f(X, \boldsymbol{\vartheta}) \right)^2 < \infty, \quad j = 1, \dots, k$;

IV. $\mathbb{E}_{\boldsymbol{\vartheta}} \frac{\partial}{\partial \vartheta_j} \log f(X, \boldsymbol{\vartheta}) = 0, \quad j = 1, \dots, k$;

V. az $I(\boldsymbol{\vartheta}) = (I_{ij}(\boldsymbol{\vartheta}))_{i,j=1}^k$ Fisher-féle információs mátrix pozitív definit $\boldsymbol{\vartheta} \in \omega$ esetén, ahol

$$I_{ij}(\boldsymbol{\vartheta}) = \mathbb{E}_{\boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \vartheta_i} \log f(X, \boldsymbol{\vartheta}) \frac{\partial}{\partial \vartheta_j} \log f(X, \boldsymbol{\vartheta}) \right);$$

VI. léteznek olyan $M_{ij\ell}$ függvények, hogy

$$\left| \frac{\partial^3}{\partial \vartheta_i \partial \vartheta_j \partial \vartheta_\ell} \log f(x, \boldsymbol{\vartheta}) \right| \leq M_{ij\ell}(x),$$

és $m_{ij\ell} = \mathbb{E}_{\boldsymbol{\vartheta}^*} M_{ij\ell}(X) < \infty, \quad i, j, \ell = 1, \dots, k$.

Ekkor a $\frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}) = 0$ likelihood-egyenletnek létezik olyan $\widehat{\boldsymbol{\vartheta}}_n$ gyöke, melyre $n \rightarrow \infty$ esetén:

(a) $\boxed{\widehat{\boldsymbol{\vartheta}}_n \rightarrow \boldsymbol{\vartheta}^* \quad 1 \text{ valószínűséggel } P_{\boldsymbol{\vartheta}^*} \text{ szerint}};$

(b) $\boxed{\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \text{ aszimptotikusan } \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*)^{-1}) \text{ eloszlású}}.$

BIZONYÍTÁS. (a) Fejtsük $L_n(\boldsymbol{\vartheta})$ -t Taylor-sorba $\boldsymbol{\vartheta}^*$ körül.

$$\begin{aligned} \frac{1}{n} L_n(\boldsymbol{\vartheta}) - \frac{1}{n} L_n(\boldsymbol{\vartheta}^*) &= \frac{1}{n} \mathbf{A}^\top (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*) + \frac{1}{2n} (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*)^\top B (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*) + \\ &+ \frac{1}{6n} \sum_i \sum_j \sum_\ell (\vartheta_i - \vartheta_i^*) (\vartheta_j - \vartheta_j^*) (\vartheta_\ell - \vartheta_\ell^*) \sum_{s=1}^n \gamma_{ij\ell}(X_s) M_{ij\ell}(X_s) = \\ &= S_1 + S_2 + S_3, \end{aligned}$$

ahol

$$\mathbf{A} = \frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta})|_{\boldsymbol{\vartheta}=\boldsymbol{\vartheta}^*} = \left(\frac{\partial L_n(\boldsymbol{\vartheta}^*)}{\partial \vartheta_1}, \dots, \frac{\partial L_n(\boldsymbol{\vartheta}^*)}{\partial \vartheta_k} \right),$$

$$B = \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta})|_{\boldsymbol{\vartheta}=\boldsymbol{\vartheta}^*} = \left(\frac{\partial^2 L_n(\boldsymbol{\vartheta}^*)}{\partial \vartheta_i \partial \vartheta_j} \right)_{i,j=1}^k,$$

$$|\gamma_{ij\ell}(X)| \leq 1.$$

Mivel

$$\mathbf{A} = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\vartheta}} \log f(X_i, \boldsymbol{\vartheta})|_{\boldsymbol{\vartheta}=\boldsymbol{\vartheta}^*}$$

és

$$B = \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} \log f(X_i, \boldsymbol{\vartheta})|_{\boldsymbol{\vartheta}=\boldsymbol{\vartheta}^*},$$

így a nagy számok törvénye miatt $\frac{1}{n}\mathbf{A} \rightarrow \mathbf{0}$ és $\frac{1}{n}B \rightarrow -I(\boldsymbol{\vartheta}^*)$ 1 valószínűséggel $P_{\boldsymbol{\vartheta}^*}$ szerint. Hasonlóan, $\frac{1}{n} \sum_{s=1}^n M_{ij\ell}(X_s) \rightarrow m_{ij\ell}$ $P_{\boldsymbol{\vartheta}^*}$ szerint 1 valószínűséggel. Azon esemény, melyre mindhárom konvergencia teljesül, 1 valószínűségű. Az alábbiak ezen esemény egy elemi eseményére vonatkoznak. Jelölje c az $I(\boldsymbol{\vartheta}^*)$ legkisebb sajátértékét. Legyen $\varepsilon > 0$ és $b = \frac{k^3}{3} \sum_i \sum_j \sum_\ell m_{ij\ell}$. Elég nagy n -re $\|\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^*\| = \varepsilon$ esetén igaz, hogy $|S_1| < \varepsilon^3$, $S_2 < -\frac{c}{4}\varepsilon^2$, $|S_3| < b\varepsilon^3$.

Így a $\boldsymbol{\vartheta}^*$ középpontú, ε sugarú gömb felületén

$$\frac{1}{n}L_n(\boldsymbol{\vartheta}) - \frac{1}{n}L_n(\boldsymbol{\vartheta}^*) = S_1 + S_2 + S_3 < \varepsilon^3 - \frac{c}{4}\varepsilon^2 + b\varepsilon^3 < 0,$$

ha $\varepsilon < \frac{c}{4(1+b)}$.

Tehát ezen a gömbön belül van lokális maximuma $L_n(\boldsymbol{\vartheta})$ -nak. Mivel ε tetszőleges pozitív szám, ez a lokális maximum a likelihood-egyenlet $\boldsymbol{\vartheta}^*$ -hoz konvergáló gyöke. Jelölje ezt a gyököt $\widehat{\boldsymbol{\vartheta}}_n$.

(b) Fejtsük Taylor-sorba $\frac{\partial}{\partial \vartheta_j} L_n(\widehat{\boldsymbol{\vartheta}}_n)$ -et $\boldsymbol{\vartheta}^*$ körül:

$$(2.16) \quad \begin{aligned} \frac{\partial}{\partial \vartheta_j} L_n(\widehat{\boldsymbol{\vartheta}}_n) - \frac{\partial}{\partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*) &= \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*)^\top (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \\ &+ \frac{1}{2} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\widetilde{\boldsymbol{\vartheta}}) (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*), \end{aligned}$$

ahol $\widetilde{\boldsymbol{\vartheta}}$ a $\boldsymbol{\vartheta}^*$ és a $\widehat{\boldsymbol{\vartheta}}_n$ pontokat összekötő szakasz eleme. Mivel $\widehat{\boldsymbol{\vartheta}}_n$ maximum hely, így $\frac{\partial}{\partial \vartheta_j} L_n(\widehat{\boldsymbol{\vartheta}}_n) = 0$.

Minden j -re:

$$\left[\underbrace{\frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*)^\top}_{-I(\boldsymbol{\vartheta}^*)} + \frac{1}{2n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \underbrace{\frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\tilde{\boldsymbol{\vartheta}})}_0 \right] \sqrt{n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) =$$

1 valószínűséggel 1 valószínűséggel

$$= \underbrace{-\frac{1}{\sqrt{n}} \frac{\partial}{\partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*)}_{\downarrow}.$$

$\mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*))$
eloszlásban

A fenti konvergenciák úgy értendők, hogy az egyenletet minden j -re kiírva, a kapott vektorok, ill. mátrixok konvergálnak $n \rightarrow \infty$ esetén. Az első konvergencia a nagy számok törvénye miatt áll fenn. A második a $\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^* \rightarrow \mathbf{0}$ miatt, és mert $\frac{1}{n} \frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\tilde{\boldsymbol{\vartheta}})$ majorálható egy – a nagy számok törvénye miatt – konvergens szorzattal. A harmadik a centrális határeloszlás tétel következménye.

Tehát a fenti egyenletrendszer bal oldalán a $[\]$ -ben $-I(\boldsymbol{\vartheta}^*) + o(1)$ áll, ahol $o(1)$ egy 1 valószínűséggel 0-hoz tartó sorozatot jelöl. Ennek inverze $-I^{-1}(\boldsymbol{\vartheta}^*) + o(1)$. Ezzel szorozva mindkét oldalt, Szluckij lemmája alapján adódik:

$$\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \Rightarrow \mathcal{N}_k(\mathbf{0}, I^{-1}(\boldsymbol{\vartheta}^*)). \quad \square$$

2.8. megjegyzés. Az előző tétel feltételei mellett $\widehat{\boldsymbol{\vartheta}}_n$ lokális maximum (1-hez tartó valószínűséggel). □

2.9. megjegyzés.

$$(2.17) \quad (a) \quad \frac{1}{\sqrt{n}} \left[\frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*) \right] + \frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*) \sqrt{n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)$$

és

$$(2.18) \quad (b) \quad \frac{1}{\sqrt{n}} \left[\frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*) \right] - I(\boldsymbol{\vartheta}^*) \sqrt{n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)$$

sztochasztikusan $\mathbf{0}$ -hoz tart.

BIZONYÍTÁS. (a) (2.16) alapján

$$\begin{aligned} \frac{1}{\sqrt{n}} \frac{\partial}{\partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*) + \frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*)^\top \sqrt{n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) &= \\ &= -\frac{1}{2n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\tilde{\boldsymbol{\vartheta}}) \sqrt{n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*). \end{aligned}$$

Mivel $\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^* \rightarrow \mathbf{0}$ sztochasztikusan, $\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)$ eloszlásban konvergens, tehát sztochasztikusan korlátos, továbbá $\frac{1}{n} \frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\tilde{\boldsymbol{\vartheta}})$ sztochasztikusan korlátos, így a jobb oldal 0-hoz tart sztochasztikusan.

(b) (2.16) alapján

$$\begin{aligned} \frac{1}{\sqrt{n}} \frac{\partial}{\partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*) - I_j(\boldsymbol{\vartheta}^*)^\top \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) &= \\ &= - \left[\frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*)^\top + I_j(\boldsymbol{\vartheta}^*)^\top \right] \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) - \\ &\quad - \frac{1}{2n} (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \frac{\partial^3}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\tilde{\boldsymbol{\vartheta}}) \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*). \end{aligned}$$

Itt $I_j(\boldsymbol{\vartheta}^*)^\top$ az $I(\boldsymbol{\vartheta}^*)$ mátrix j -edik sorvektora. Az előzőekhez még figyelembe kell venni, hogy a nagy számok törvénye miatt $\frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \vartheta_j} L_n(\boldsymbol{\vartheta}^*) \rightarrow I_j(\boldsymbol{\vartheta}^*)$. \square

3. A likelihood-hányados próba

3.1. A loglikelihood-függvény aszimptotikája

Ahhoz, hogy a likelihood-függvényen alapuló próbák statisztikáinak határeloszlását megkapjuk, szükségünk van az L_n függvény aszimptotikájának ismeretére. Legyen X_1, \dots, X_n (független, azonos eloszlású elemekből álló) minta, $L_n(\boldsymbol{\vartheta})$ a megfelelő loglikelihood-függvény, $\widehat{\boldsymbol{\vartheta}}_n$ pedig a likelihood-egyenlet konzisztens gyöke.

3.1. tétel. A 2.2. tétel feltételeinek teljesülése esetén, ha $n \rightarrow \infty$,

$$(3.1) \quad (A) \quad \left[L_n(\widehat{\boldsymbol{\vartheta}}_n) - L_n(\boldsymbol{\vartheta}^*) \right] - \frac{1}{2}n(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top I(\boldsymbol{\vartheta}^*)(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \rightarrow 0$$

sztochasztikusan $P_{\boldsymbol{\vartheta}^*}$ szerint;

$$(3.2) \quad (B) \quad 2 \left[L_n(\widehat{\boldsymbol{\vartheta}}_n) - L_n(\boldsymbol{\vartheta}^*) \right] \implies \chi_k^2$$

eloszlásban $P_{\boldsymbol{\vartheta}^*}$ szerint.

BIZONYÍTÁS. (A) Fejtsük Taylor-sorba $L_n(\boldsymbol{\vartheta})$ -t a $\boldsymbol{\vartheta}^*$ körül:

$$\begin{aligned} L_n(\widehat{\boldsymbol{\vartheta}}_n) - L_n(\boldsymbol{\vartheta}^*) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*)^\top (\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \\ &+ \frac{1}{2}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*)(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \gamma \|\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*\|^3 \sum_{s=1}^n G(X_s), \end{aligned}$$

ahol $G(x) = \sum_{i,j,\ell} M_{ij\ell}(x)$ és $|\gamma| \leq \frac{1}{6}$. Ebből

$$\begin{aligned} L_n(\widehat{\boldsymbol{\vartheta}}_n) - L_n(\boldsymbol{\vartheta}^*) &= \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top I(\boldsymbol{\vartheta}^*)\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \\ &+ \left[\frac{1}{\sqrt{n}} \frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*)^\top - \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top I(\boldsymbol{\vartheta}^*) \right] \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \\ &+ \frac{1}{2}\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top \left[\frac{1}{n} \frac{\partial^2}{\partial \boldsymbol{\vartheta} \partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}^*) + I(\boldsymbol{\vartheta}^*) \right] \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) - \\ &- \frac{1}{2}\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)^\top I(\boldsymbol{\vartheta}^*)\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) + \\ &+ \left[\gamma \|\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*\| \right] \left[\sqrt{n} \|\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*\| \right] \left[\sqrt{n} \|\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*\| \right] \left[\frac{1}{n} \sum_{s=1}^n G(X_s) \right]. \end{aligned}$$

A 2.9. megjegyzés (b) része miatt a második tag sztochasztikusan 0-hoz tart. A harmadik tag ugyanezen megjegyzés (a)–(b) része, valamint $\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*)$ sztochasztikus korlátossága miatt 0-hoz tart sztochasztikusan. Az ötödik tag is 0-hoz tart, mivel ennek az első tényezője 0-hoz tart 1 valószínűséggel, míg a következő három tényezője sztochasztikusan korlátos.

Az első és a negyedik tag összege pedig éppen az (A)-ban szereplő mennyiség.

(B) Mivel

$$\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \implies \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*)^{-1}),$$

így (A)-ból következik (B). \square

3.2. Egyszerű nullhipotézis vizsgálata

A

$$H_0 : \boldsymbol{\vartheta} = \boldsymbol{\vartheta}_0$$

egyszerű nullhipotézis

$$H_1 : \boldsymbol{\vartheta} \neq \boldsymbol{\vartheta}_0$$

alternatív hipotézissel szembeni tesztelésére a Neyman-Pearson-féle, az A. Waldtól és a C. R. Raotól származó statisztikákat ismertetjük. Ezek mindegyike a maximum-likelihood módszerrel kapcsolatos. Ezen általános statisztikák alapján lehet speciális modellek esetén konkrét próbákat konstruálni.

3.1. definíció. Legyen X_1, \dots, X_n (független, azonos eloszlású elemekből álló) minta, $\boldsymbol{\vartheta} \in \Theta \subseteq \mathbb{R}^k$ a paraméter, $\ell(X_1, \dots, X_n, \boldsymbol{\vartheta})$ a megfelelő likelihood-függvény, $L_n(\boldsymbol{\vartheta}) = L_n(X_1, \dots, X_n, \boldsymbol{\vartheta})$ a loglikelihood-függvény, $I(\boldsymbol{\vartheta})$ a Fisher-féle információs mátrix, $\widehat{\boldsymbol{\vartheta}}_n$ pedig a likelihood-egyenlet konzisztens gyöke. Vezessük be az alábbi statisztikákat.

$$\text{Neyman-Pearson: } \Lambda_n = \frac{\ell(X_1, \dots, X_n, \boldsymbol{\vartheta}_0)}{\ell(X_1, \dots, X_n, \widehat{\boldsymbol{\vartheta}}_n)}.$$

$$\text{Wald: } W_n = \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}_0)^\top I(\widehat{\boldsymbol{\vartheta}}_n) \sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}_0).$$

$$\text{Rao: } V_n = \frac{1}{n} \left[\frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}_0) \right] [I(\boldsymbol{\vartheta}_0)]^{-1} \left[\frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}_0) \right]^\top. \quad \square$$

A Neyman-Pearson-féle statisztikát szokták a

$$\Lambda_n = \frac{\ell(X_1, \dots, X_n, \boldsymbol{\vartheta}_0)}{\sup_{\boldsymbol{\vartheta} \in \Theta} \ell(X_1, \dots, X_n, \boldsymbol{\vartheta})}$$

képlettel is definiálni. Az aszimptotikus eloszlás megfogalmazásához azonban az előző definíció alkalmasabb.

A Rao-féle tesztet szokták score tesztnek is nevezni.

3.2. tétel. A 2.2. tételbeni feltételek esetén $-2 \log \Lambda_n$, W_n , valamint V_n is aszimptotikusan χ_k^2 eloszlású, ha H_0 igaz.

BIZONYÍTÁS. Az, hogy $-2 \log \Lambda_n$ aszimptotikusan χ_k^2 eloszlású, azonnal adódik a 3.1. tétel (B) részéből.

Tudjuk, hogy $\sqrt{n}(\widehat{\boldsymbol{\vartheta}}_n - \boldsymbol{\vartheta}^*) \implies \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}^*)^{-1})$. Mivel $\widehat{\boldsymbol{\vartheta}}_n \rightarrow \boldsymbol{\vartheta}_0$ a H_0 teljesülése esetén 1 valószínűséggel, és $I(\boldsymbol{\vartheta})$ $\boldsymbol{\vartheta}$ -ben folytonos, így $I(\widehat{\boldsymbol{\vartheta}}_n) \rightarrow I(\boldsymbol{\vartheta}_0)$ a H_0 teljesülése esetén 1 valószínűséggel. Tehát W_n is aszimptotikusan χ_k^2 eloszlású.

Mivel $\frac{1}{\sqrt{n}} \frac{\partial}{\partial \boldsymbol{\vartheta}} L_n(\boldsymbol{\vartheta}_0) \implies \mathcal{N}_k(\mathbf{0}, I(\boldsymbol{\vartheta}_0))$, így $V_n \implies \chi_k^2$. \square

IV. fejezet

Nem-paraméteres módszerek

1. Hoeffding-féle U -statisztikák

Ebben a részben a Hoeffding-féle U -statisztikákra bizonyítunk határeloszlás-tételt, amiből azután levezetjük a Wilcoxon-féle rangösszeg statisztika aszimptotikus normalitását.

1.1. Független változókkal való közelítés

Legyenek Z_1, Z_2, \dots, Z_N független valószínűségi változók, T pedig legyen olyan valószínűségi változó, melyre $\mathbb{E}(T^2) < \infty$ és $\mathbb{E}(T) = 0$.

Keresendők olyan k_1, k_2, \dots, k_N (Borel-)mérhető függvények, melyekre

$$\mathbb{E}(T - S)^2 \quad \text{minimális,}$$

ahol

$$(1.1) \quad S = \sum_{i=1}^N k_i(Z_i).$$

Figyeljük meg, hogy a feladat és annak alábbi megoldása analógiát mutat egy vektornak egy ortogonális vektorrendszer által generált altérre való vetítésével.

1.1. tétel. *A fenti feladat megoldása*

$$k_i(Z_i) = \mathbb{E}(T \mid Z_i), \quad i = 1, \dots, N.$$

BIZONYÍTÁS. Legyen

$$(1.2) \quad T^* = \sum_{i=1}^N \mathbb{E}(T \mid Z_i).$$

Elegendő belátni, hogy tetszőleges, (1.1)-ben megadott S -re

$$(1.3) \quad \mathbb{E}(T - S)^2 = \mathbb{E}(T - T^*)^2 + \mathbb{E}(T^* - S)^2.$$

Viszont az alábbi alapján a „kétszeres szorzat” = 0, amiből (1.3) már következik.

$$\begin{aligned} \mathbb{E}(T - T^*)(T^* - S) &= \sum_{i=1}^N \mathbb{E} \left[(T - T^*)(\mathbb{E}(T \mid Z_i) - k_i(Z_i)) \right] = \\ &= \sum_{i=1}^N \mathbb{E} \left[\mathbb{E} \left\{ (T - T^*)(\mathbb{E}(T \mid Z_i) - k_i(Z_i)) \mid Z_i \right\} \right] = \\ &= \sum_{i=1}^N \mathbb{E} \left[(\mathbb{E}(T \mid Z_i) - k_i(Z_i)) \mathbb{E} \{ T - T^* \mid Z_i \} \right] = 0. \end{aligned}$$

Itt kihasználtuk, hogy a Z_i és a Z_j ($i \neq j$) függetlensége alapján

$$\mathbb{E} \{ T - T^* \mid Z_i \} = \mathbb{E}(T \mid Z_i) - \mathbb{E}(T \mid Z_i) - \sum_{i \neq j} \mathbb{E} \{ \mathbb{E}(T \mid Z_j) \mid Z_i \} = 0. \quad \square$$

1.2. A Hoeffding-féle U -statisztika

Legyen

az X valószínűségi változó eloszlásfüggvénye F ,

az Y valószínűségi változó eloszlásfüggvénye G .

Legyen

$$\begin{aligned} X_1, \dots, X_m & \quad \text{minta } X\text{-re,} \\ Y_1, \dots, Y_n & \quad \text{minta } Y\text{-ra,} \end{aligned}$$

ahol feltesszük a két minta egymástól való függetlenségét is.

1.1. definíció. A Hoeffding-féle U -statisztika

$$(1.4) \quad U = \frac{1}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n \varphi(X_i, Y_j),$$

ahol φ adott kétváltozós mérhető függvény. \square

Az a célunk, hogy a Hoeffding-féle U -statisztika (alkalmas normalizáltjának) aszimptotikus normalitását belássuk. Ennek érdekében — a fenti 1.1. tétel segítségével — független valószínűségi változók összegével közelítjük meg U -t. A független összegre már alkalmazhatjuk a központi határeloszlás-tételt. A pontos eljárás az alábbi.

Legyen $\mathbb{E}\varphi^2(X, Y) < \infty$ és legyen

$$\theta = \mathbb{E}U = \mathbb{E}\varphi(X, Y),$$

$$\psi(x, y) = \varphi(x, y) - \theta.$$

Keresendő

$$mn(U - \mathbb{E}U) = \sum_{i=1}^m \sum_{j=1}^n \psi(X_i, Y_j)$$

legjobb közelítése

$$S = \sum_{i=1}^m a_i(X_i) + \sum_{j=1}^n b_j(Y_j)$$

alakú függvényekkel, ahol az a_i és b_j mérhető függvényeket kell meghatározni.

A feltételes várható érték ismert tulajdonsága, hogy ha X és Y függetlenek, akkor $\mathbb{E}\{f(X, Y) \mid X = x\} = \mathbb{E}f(x, Y)$. Ezért

$$\mathbb{E}\{\psi(X_i, Y_j) \mid X_k = x\} = \begin{cases} \mathbb{E}\psi(x, Y_j), & i = k; \\ \mathbb{E}\psi(X_i, Y_j) = 0, & i \neq k. \end{cases}$$

Vezessük be az alábbi jelöléseket:

$$\psi_{10}(x) = \mathbb{E}\psi(x, Y),$$

$$\psi_{01}(y) = \mathbb{E}\psi(X, y).$$

A fenti megjegyzés és az 1.1. tétel alapján $mn(U - \mathbb{E}U)$ keresett közelítése:

$$(1.5) \quad \begin{aligned} & \sum_{k=1}^m \mathbb{E} \left\{ \sum_{i=1}^m \sum_{j=1}^n \psi(X_i, Y_j) \mid X_k \right\} + \sum_{\ell=1}^n \mathbb{E} \left\{ \sum_{i=1}^m \sum_{j=1}^n \psi(X_i, Y_j) \mid Y_\ell \right\} = \\ & = \sum_{k=1}^m n\psi_{10}(X_k) + \sum_{\ell=1}^n m\psi_{01}(Y_\ell). \end{aligned}$$

A következő tételben használni fogjuk az alábbi jelöléseket:

$$\sigma_{10}^2 = \mathbb{D}^2\psi_{10}(X) = \mathbb{E}\psi_{10}^2(X), \quad \sigma_{01}^2 = \mathbb{D}^2\psi_{01}(Y) = \mathbb{E}\psi_{01}^2(Y).$$

1.2. tétel. (A Hoeffding-féle U -statisztika aszimptotikus normalitása.) Tegyük fel, hogy $m \leq n$ és

$$\lim_{\substack{n \rightarrow \infty \\ m \rightarrow \infty}} m/n = \lambda,$$

ahol λ lehet 0 is. Ekkor $m, n \rightarrow \infty$ esetén

$$T = \sqrt{m}(U - \theta) \implies \mathcal{N}(0, \sigma_{10}^2 + \lambda\sigma_{01}^2).$$

(Amennyiben a határeloszlás szórása 0, akkor elfajult, pontosabban a 0-ba koncentrált a határeloszlás.)

BIZONYÍTÁS. Belátjuk, hogy $T - T^* \rightarrow 0$ L^2 -ben, és $T^* \implies \mathcal{N}(\cdot, \cdot)$, ahol T^* a tételbeni T -hez (1.2) alapján tartozik. Ezekből — Szluckij ismert lemmája alapján — már következni fog az állítás.

Emlékeztetünk Szluckij lemmájára.

Ha $X_n - Y_n \rightarrow 0$ sztochasztikusan, és $Y_n \implies Y$ eloszlásban, akkor $X_n \implies Y$ eloszlásban.

Kezdjük $T - T^* \rightarrow 0$ L^2 -ben igazolásával. (1.5)-ből:

$$(1.6) \quad T^* = \sqrt{m} \left[\frac{1}{m} \sum_{i=1}^m \psi_{10}(X_i) + \frac{1}{n} \sum_{j=1}^n \psi_{01}(Y_j) \right].$$

(1.3)-ból $S = 0$ vételével:

$$(1.7) \quad \mathbb{E}(T - T^*)^2 = \mathbb{E}T^2 - \mathbb{E}T^{*2} = \mathbb{D}^2T - \mathbb{D}^2T^*,$$

mivel $\mathbb{E}T = \mathbb{E}T^* = 0$. Kiszámítjuk $\mathbb{D}^2(T)$ -t és $\mathbb{D}^2(T^*)$ -ot.

$$(1.8) \quad \begin{aligned} \mathbb{D}^2T &= m \frac{1}{m^2 n^2} \sum_{k,\ell} \sum_{i,j} \text{cov}(\varphi(X_i, Y_j), \varphi(X_k, Y_\ell)) = \\ &= \frac{1}{mn^2} \left[\sum_{i,j} \text{cov}(\varphi(X_i, Y_j), \varphi(X_i, Y_j)) + \right. \\ &\quad \left. + \sum_{\substack{i,j,\ell \\ \ell \neq j}} \text{cov}(\varphi(X_i, Y_j), \varphi(X_i, Y_\ell)) + \sum_{\substack{i,j,k \\ i \neq k}} \text{cov}(\varphi(X_i, Y_j), \varphi(X_k, Y_j)) \right] = \\ &= \frac{1}{mn^2} [mn(\dots) + mn(n-1)(\dots) + nm(m-1)(\dots)] \rightarrow \\ &\quad \xrightarrow{m,n \rightarrow \infty} 0 + \text{cov}(\varphi(X_i, Y_j), \varphi(X_i, Y_\ell)) + \lambda \text{cov}(\varphi(X_i, Y_j), \varphi(X_k, Y_j)) = \\ &= \sigma_{10}^2 + \lambda\sigma_{01}^2. \end{aligned}$$

Ugyanis egyrészt a függetlenség miatt $\text{cov}(\varphi(X_i, Y_j), \varphi(X_k, Y_\ell)) = 0$, ha $i \neq k$ és $j \neq \ell$. Másrészt $j \neq \ell$ esetén

$$\begin{aligned} \text{cov}(\varphi(X_i, Y_j), \varphi(X_i, Y_\ell)) &= \mathbb{E}(\psi(X_i, Y_j) \cdot \psi(X_i, Y_\ell)) = \\ &= \mathbb{E}\left(\underbrace{\mathbb{E}\{\psi(X_i, Y_j) \cdot \psi(X_i, Y_\ell) \mid X_i = x\}}_{|_{x=x_i}}\right), \end{aligned}$$

valamint

$$\mathbb{E}\psi(x, Y_j) \cdot \psi(x, Y_\ell) = \mathbb{E}\psi(x, Y_j) \cdot \mathbb{E}\psi(x, Y_\ell) = \psi_{10}^2(x).$$

Innen $\sigma_{10}^2 = \mathbb{E}\psi_{10}^2(X_i)$ már adódik. Hasonlóan σ_{01}^2 is.

Rátérünk $\mathbb{D}^2 T^*$ -ra. (1.6) alapján, a függetlenséget kihasználva:

$$\mathbb{D}^2 T^* = m \left[\frac{1}{m^2} \sum_{i=1}^m \sigma_{10}^2 + \frac{1}{n^2} \sum_{j=1}^n \sigma_{01}^2 \right] = \sigma_{10}^2 + \frac{m}{n} \sigma_{01}^2 \rightarrow \sigma_{10}^2 + \lambda \sigma_{01}^2.$$

Így, (1.7) alapján, $T - T^* \rightarrow 0$ L^2 -ben.

Térjünk rá $T^* \implies \mathcal{N}(\cdot, \cdot)$ igazolására. $n, m \rightarrow \infty$ esetén

$$\begin{array}{ccc} T^* = \frac{1}{\sqrt{m}} \sum_{i=1}^m \psi_{10}(X_i) + \sqrt{\frac{m}{n}} \frac{1}{\sqrt{n}} \sum_{j=1}^n \psi_{01}(Y_j) & & \\ \underbrace{\hspace{10em}}_{\downarrow} & \downarrow & \underbrace{\hspace{10em}}_{\downarrow} \\ \mathcal{N}(0, \sigma_{10}^2) & \lambda & \mathcal{N}(0, \sigma_{01}^2) \\ \underbrace{\hspace{10em}}_{\downarrow} & & \\ \mathcal{N}(0, \sigma_{10}^2 + \lambda \sigma_{01}^2), & & \end{array}$$

ha $\lambda \neq 0$ és legalább $\sigma_{10}^2, \sigma_{01}^2$ egyike pozitív; illetve ha $\lambda = 0$ és σ_{10}^2 pozitív (ellenkező esetben T^* (határ)eloszlása a 0-ba koncentrálódik). A fenti határeloszlás meghatározásához a központi határeloszlás tételt és a két összeg egymástól való függetlenségét használtuk ki (lásd a 3. és a 4. Feladatot). \square

1.3. Feladatok

1. Az 1.1. tétel bizonyításának mely lépéseinél használtuk ki a feltételes várható érték alábbi tulajdonságait?

- $\mathbb{E}\{\xi \mid \eta\}$ az η -nak (Borel-)mérhető függvénye.
- $\mathbb{E}(\mathbb{E}\{\xi \mid \eta\}) = \mathbb{E}\xi$.
- $\mathbb{E}\{\xi\zeta \mid \eta\} = \zeta \mathbb{E}\{\xi \mid \eta\}$, ha ζ η -mérhető.
- $\mathbb{E}\{a_1\xi_1 + a_2\xi_2 \mid \eta\} = a_1\mathbb{E}\{\xi_1 \mid \eta\} + a_2\mathbb{E}\{\xi_2 \mid \eta\}$, ahol a_1 és a_2 konstansok.
- $\mathbb{E}\{\xi \mid \eta\} = \mathbb{E}\xi$, ha ξ és η függetlenek.

2. Ismeretes, hogy $g(y) = \mathbb{E}\{\xi \mid \eta = y\}$ jelöléssel: $\mathbb{E}\{\xi \mid \eta\} = g(\eta)$; vagy tömörebben: $\mathbb{E}\{\xi \mid \eta\} = (\mathbb{E}\{\xi \mid \eta = y\})_{|_{y=\eta}}$. A fenti szakaszban mely helyeken használtuk ezt az összefüggést?

3. Legyen ξ_n és η_n független minden n esetén. Tegyük fel, hogy $\xi_n \Rightarrow \xi$ és $\eta_n \Rightarrow \eta$ eloszlásban. Igazoljuk (a karakterisztikus függvények módszerével), hogy ekkor $\xi_n + \eta_n$

eloszlásban konvergál a ξ és η eloszlásának a konvolúciójához. Hogyan alkalmaztuk ezt az 1.2. tétel bizonyításának a végén?

4. Tegyük fel, hogy $\xi_n \Rightarrow \xi$ eloszlásban, $a_n \rightarrow a$. Lássuk be, hogy $a_n \xi_n \Rightarrow a\xi$ eloszlásban.

4. ÚTMUTATÁS. Első módszer. Jelölje φ a karakterisztikus függvényt. Egyrészt $\varphi_{a_n \xi_n}(t) = \mathbb{E}e^{it a_n \xi_n} = \varphi_{\xi_n}(a_n t)$.

Másrészt $|\varphi_{\xi_n}(a_n t) - \varphi_{\xi}(a_n t)| \leq \varepsilon$, ha $n > n_\varepsilon$, mert $\varphi_{\xi_n}(t) \rightarrow \varphi_{\xi}(t)$ minden korlátos intervallumban egyenletesen (l. Rényi (1973)).

Harmadrészt $\varphi_{\xi}(a_n t) \rightarrow \varphi_{\xi}(at) = \varphi_{a\xi}(t)$, mert φ folytonos.

Második módszer. Alkalmazzuk a Szluckij-lemmát! Ekkor kihasználjuk, hogy $a\xi_n \Rightarrow a\xi$ eloszlásban és $(a_n \xi_n - a\xi_n) \rightarrow 0$ sztochasztikusan. Ez utóbbi belátásához szükséges: a ξ_n sorozat sztochasztikusan korlátos (mert eloszlásban konvergens).

5. Érvényben marad-e a Hoeffding-féle U -statisztika aszimptotikus normalitása $m > n$ esetén?

5. MEGOLDÁS. Az 1.2. tétel alábbi változata áll fenn. Tegyük fel, hogy $m > n$ és

$$\lim_{\substack{n \rightarrow \infty \\ m \rightarrow \infty}} n/m = \tilde{\lambda},$$

ahol $\tilde{\lambda}$ lehet 0 is. Ekkor $m, n \rightarrow \infty$ esetén

$$\tilde{T} = \sqrt{n}(U - \theta) \implies \mathcal{N}(0, \tilde{\lambda}\sigma_{10}^2 + \sigma_{01}^2).$$

(Amennyiben a határeloszlás szórása 0, akkor elfajult, pontosabban a 0-ba koncentrált, a határeloszlás.)

A bizonyítás ugyanaz, mint az 1.2. tételé, csak a \sqrt{m} helyett \sqrt{n} -nel kell normálni.

6. Érvényben marad-e az 1.2. tétel, ha csupán $n \rightarrow \infty$, de m korlátos?

6. MEGOLDÁS. Igen. Ekkor $\lambda = 0$ és $\mathcal{N}(0, \sigma_{10}^2)$ a határeloszlás. Továbbá, ha $m \rightarrow \infty$, de n korlátos, akkor (lásd az 5. Feladatot) $\tilde{T} = \sqrt{n}(U - \theta) \implies \mathcal{N}(0, \sigma_{01}^2)$, hiszen $\tilde{\lambda} = 0$.

2. A Mann-Whitney-féle U -próba

2.1. A Mann-Whitney-féle U -próba bevezetése

A **Mann-Whitney-féle U -próbát** (más néven a **Wilcoxon-féle rangösszeg próbát**) mint speciális Hoeffding-féle statisztikát tekintjük. Legyen X_1, \dots, X_m minta F eloszlásfüggvényű, míg Y_1, \dots, Y_n minta G eloszlásfüggvényű sokaságból. Felteszük a két minta egymástól való függetlenségét is. Legyen az F és a G eloszlásfüggvény folytonos. F és G számunkra ismertelen, és éppen az F és G azonosságát kívánjuk vizsgálni.

Legyen

$$D_{ij} = \begin{cases} 1, & \text{ha } X_i > Y_j, \\ 0, & \text{ha } X_i \leq Y_j, \end{cases}$$

$i = 1, \dots, m, j = 1, \dots, n$. Ekkor

$$U'_{mn} = \sum_{i=1}^m \sum_{j=1}^n D_{ij} = \sum_{i=1}^m (\text{rang}(X_i) - i)$$

a Wilcoxon-féle rangösszeg statisztika. Itt $\text{rang}(X_i)$ az X_i sorszáma az X_1, \dots, X_m és az Y_1, \dots, Y_n minták együttes rendezése esetén.

2.1. tétel. (A Wilcoxon-féle rangösszeg statisztika várható értéke és szórása.)

$$\mathbb{E}U'_{mn} = mn\pi,$$

$$\mathbb{D}^2U'_{mn} = nm \left[\pi - \pi^2(n + m - 1) + (n - 1)\pi_1 + (m - 1)\pi_2 \right],$$

ahol

$$\pi = \mathbb{P}(Y < X) = \int_{-\infty}^{+\infty} G(x)dF(x) = \mathbb{E}G(X) = \mathbb{E}(1 - F(Y)),$$

$$\pi_1 = \mathbb{P}(Y_j < X_i, Y_k < X_i) = \int_{-\infty}^{+\infty} G^2(x)dF(x) = \mathbb{E}G^2(X),$$

$$\pi_2 = \mathbb{P}(Y_i < X_j, Y_i < X_k) = \int_{-\infty}^{+\infty} (1 - F(y))^2 dG(y) = \mathbb{E}(1 - F(Y))^2,$$

$j \neq k$. Speciálisan,

$$H_0 : F(x) \equiv G(x)$$

esetén

$$\mathbb{E}U'_{mn} = \frac{mn}{2},$$

$$\mathbb{D}^2U'_{mn} = \frac{nm(n + m + 1)}{12}.$$

BIZONYÍTÁS. Határozzuk meg először U'_{mn} várható értékét! Mivel X és Y független,

$$\pi = \mathbb{P}(Y < X) = \int_{-\infty}^{+\infty} \int_{-\infty}^x dG(y)dF(x) = \int_{-\infty}^{+\infty} G(x)dF(x) = \mathbb{E}G(X).$$

Ha $H_0 : F(x) \equiv G(x)$ teljesül, akkor

$$\pi = \int_{-\infty}^{+\infty} F(x)dF(x) = \mathbb{E}F(X) = \frac{1}{2},$$

lévén $F(X)$ a $[0, 1]$ -en egyenletes eloszlású (lásd az 1. Feladatot).

Másrészt nyilván $\mathbb{E}D_{ij} = \pi$, így $\mathbb{E}U'_{mn} = mn\pi$. Tehát ha H_0 teljesül, akkor $\mathbb{E}U'_{mn} = \frac{mn}{2}$.

Most $\mathbb{D}^2U'_{mn}$ összetevőit határozzuk meg. $\mathbb{D}^2D_{ij} = \pi(1 - \pi)$ adódik D_{ij} definíciójából (hiszen Bernoulli-eloszlásról van szó). A függetlenség miatt

$$\text{cov}(D_{ij}, D_{k\ell}) = 0, \quad \text{ha } i \neq k \text{ és } j \neq \ell.$$

A $\text{cov}(X, Y) = \mathbb{E}(X \cdot Y) - \mathbb{E}X \cdot \mathbb{E}Y$ képlet alapján

$$\text{cov}(D_{ij}, D_{ik}) = \pi_1 - \pi^2, \quad \text{ha } j \neq k,$$

ahol

$$\begin{aligned} \pi_1 &= \mathbb{P}(Y_j < X_i, Y_k < X_i) = \int_{-\infty}^{+\infty} \mathbb{P}(Y_j < x, Y_k < x)dF(x) = \\ &= \int_{-\infty}^{+\infty} G^2(x)dF(x) = \mathbb{E}G^2(X). \end{aligned}$$

Hasonlóan

$$\text{cov}(D_{ij}, D_{kj}) = \pi_2 - \pi^2, \quad \text{ha } i \neq k,$$

ahol

$$\pi_2 = \mathbb{P}(X_i > Y_j, X_k > Y_j) = \int_{-\infty}^{+\infty} (1 - F(x))^2 dG(x).$$

Így

$$\begin{aligned} \mathbb{D}^2U'_{mn} &= \sum_{i=1}^m \sum_{j=1}^n \mathbb{D}^2(D_{ij}) + \sum_{j \neq k} \sum \text{cov}(D_{ij}, D_{ik}) + \\ &+ \sum_{i \neq k} \sum \text{cov}(D_{ij}, D_{kj}) + \sum_{\substack{i \neq k \\ j \neq \ell}} \sum \text{cov}(D_{ij}, D_{k\ell}) = \\ &= nm\pi(1 - \pi) + mn(n - 1)(\pi_1 - \pi^2) + nm(m - 1)(\pi_2 - \pi^2) = \\ &= nm[\pi - \pi^2(n + m - 1) + (n - 1)\pi_1 + (m - 1)\pi_2]. \end{aligned}$$

Speciálisan, $H_0 : F(x) \equiv G(x)$ esetén

$$\pi_1 = \int_{-\infty}^{+\infty} F^2(x)dF(x) = \mathbb{E}F(X)^2 = \frac{1}{3}, \quad \pi_2 = \frac{1}{3},$$

ahol megint az $F(X)$ valószínűségi változó $[0, 1]$ -en való egyenletes eloszlását használtuk ki. Tehát H_0 esetén

$$\mathbb{D}^2U'_{mn} = \frac{nm(n + m + 1)}{12}. \quad \square$$

2.2. A Mann-Whitney-statisztika határeloszlása a Hoeffding-statisztika alapján

Láttuk, hogy

$$\varphi(x, y) = \begin{cases} 1, & \text{ha } x > y, \\ 0, & \text{ha } x \leq y \end{cases}$$

esetén $U'_{mn} = mnU$ éppen a Mann-Whitney-statisztika, ahol U a Hoeffding-statisztika. Ekkor (lásd a 3. Feladatot)

$$(2.1) \quad \sigma_{10}^2 = \pi_1 - \pi^2 \quad \text{és} \quad \sigma_{01}^2 = \pi_2 - \pi^2.$$

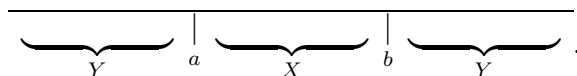
Tehát az 2.1 tétel alapján $\sigma_{10}^2 = \mathbb{D}^2 G(X)$ és $\sigma_{01}^2 = \mathbb{D}^2(1 - F(Y)) = \mathbb{D}^2 F(Y)$.

Most már az 1.2 tétel segítségével meg tudjuk vizsgálni, hogy U'_{mn} mikor lesz aszimptotikusan normális, illetve mikor elfajult. Ha $\lambda > 0$, úgy U'_{mn} akkor lesz elfajult (aszimptotikusan és véges m, n -re is), ha $\sigma_{10}^2 = \sigma_{01}^2 = 0$.

$$0 = \sigma_{10}^2 = \mathbb{D}^2 G(X) \iff \mathbb{P}(G(X) = \text{const.}) = 1 \iff$$

$$G(x) = c, \quad \text{ha } x \in (a, b] \quad \text{és} \quad \mathbb{P}(X \in (a, b]) = 1,$$

azaz X és Y eloszlása a számegyenesen az alábbi séma szerint alakul (ahol $a = \infty$, ill. $b = \infty$ is lehetséges)



$\sigma_{01}^2 = 0$ jelentése hasonló. Tehát

$$\sigma_{10}^2 = 0 \quad \text{és} \quad \sigma_{01}^2 = 0 \quad \iff \quad \mathbb{P}(X \leq Y) = 1 \quad \text{vagy} \quad \mathbb{P}(X > Y) = 1.$$

Az első esetben $U'_{mn} = 0$, a másodikban $U'_{mn} = mn$, azaz értékük determinisztikus. Tehát, ha $\lambda > 0$, ezen elfajult esetektől eltekintve U'_{mn} aszimptotikusan normális.

Ha H_0 teljesül, akkor ezek az elfajult esetek nem fordulnak elő.

2.2. tétel. (A Wilcoxon-féle rangösszeg statisztika aszimptotikus normalitása.) *Tegyük fel, hogy F és G folytonos, valamint $H_0 : F(x) \equiv G(x)$ teljesül. Ha $m, n \rightarrow \infty$ úgy, hogy $m/n \rightarrow \lambda$, akkor*

$$\frac{U'_{mn} - \frac{mn}{2}}{\sqrt{\frac{nm(n+m+1)}{12}}}$$

standardizált Wilcoxon-féle rangösszeg statisztika eloszlása a standard normális eloszláshoz konvergál.

BIZONYÍTÁS. Ekkor sem $\sigma_{10}^2 = 0$, sem $\sigma_{01}^2 = 0$ nem fordulhat elő. Tehát ekkor a 1.2 tételben a határeloszlás nem lehet elfajult. (Lásd az 5. Feladatot is!) \square

A Wilcoxon-féle rangösszeg próba leírása megtalálható az alábbi helyeken: Gibbons (1971), 142-146. o., Lehmann-D'Abrea (1975), 20. o. Az aszimptotikus normalitás

bizonyítása Lehmann–D’Abrea (1975), 362. o., míg az alkalmazások Hollander–Wolfe (1973), 68. o. és Vincze–Varbanova (1993), 101. o. lehetők.

2.3. Feladatok

1. Legyen a ξ eloszlásfüggvénye az $F(x)$ folytonos függvény. Lássuk be, hogy $F(\xi)$ egyenletes eloszlású a $[0, 1]$ -en!

1. MEGOLDÁS. Belátjuk, hogy $F(\xi)$ eloszlásfüggvénye éppen a $[0, 1]$ -en egyenletes eloszlásfüggvény. $\mathbb{P}(F(\xi) < y) = 0$, ha $y \leq 0$, és $\mathbb{P}(F(\xi) < y) = 1$, ha $y > 1$, nyilvánvaló. Legyen $0 < y \leq 1$ rögzített. Belátjuk, hogy

$$F(a) < y \iff a < x_y,$$

ahol $x_y = \sup\{x : F(x) < y\}$. Ugyanis, ha $a < x_y$, akkor $a < x_0$ legalább egy x_0 -ra az $\{x : F(x) < y\}$ halmazból. De akkor $F(a) \leq F(x_0) < y$. A másik irányhoz belátjuk, hogy $F(x_y) = y$. F monoton növekvő és folytonos voltából: $F(x_y) = \sup\{F(x) : F(x) < y\}$. De, szintén F folytonos voltából, ez utóbbi éppen y -nal egyenlő. Visszatérve a másik irányhoz: ha $F(a) < y$, akkor — az F monoton növekvő voltát is használva — $a < x_y$. Végül $F(\xi)$ eloszlásfüggvényére:

$$\mathbb{P}(\omega : F(\xi(\omega)) < y) = \mathbb{P}(\omega : \xi(\omega) < x_y) = F(x_y) = y,$$

ha $0 < y \leq 1$.

2. Legyen ξ egyenletes eloszlású a $[0, 1]$ -en. Lássuk be, hogy $1 - \xi$ is egyenletes eloszlású a $[0, 1]$ -en! Hol volt erre szükség?

2. ÚTMUTATÁS. Alkalmazzuk az eloszlásfüggvényt!

3. Lássuk be az (2.1) összefüggést!

3. MEGOLDÁS.

$$\begin{aligned} \sigma_{10}^2 &= \mathbb{D}^2 \{\psi_{10}(X)\} = \mathbb{D}^2 \left\{ (\mathbb{E}\psi(x, Y))_{|_{x=X}} \right\} = \\ &= \mathbb{D}^2 \left\{ (\mathbb{E}\varphi(x, Y))_{|_{x=X}} \right\} = \mathbb{D}^2 \left\{ (\mathbb{P}(Y < x))_{|_{x=X}} \right\} = \\ &= \mathbb{D}^2 \left\{ (G(x))_{|_{x=X}} \right\} = \mathbb{D}^2 G(X) = \mathbb{E}G^2(X) - (\mathbb{E}G(X))^2 = \\ &= \pi_1 - \pi^2. \end{aligned}$$

4. A $H_0 : F(x) \equiv G(x)$ teljesülése esetén határozzuk meg U'_{mn} várható értékét és szórását elemi eszközökkel!

4. MEGOLDÁS. Ekkor X és Y független, azonos eloszlású, folytonos eloszlásfüggvényű. Tehát $\mathbb{P}(X = Y) = 0$, $\mathbb{P}(X < Y) = \mathbb{P}(X > Y)$. Azaz $\pi = \mathbb{P}(Y < X) = 1/2$. X_i, Y_j, Y_k is független, azonos eloszlású. Bármely sorrendjük egyformán valószínű. Tehát $\pi_1 = \mathbb{P}(Y_j < X_i, Y_k < X_i) = 1/3$.

5. Érvényben marad-e az 2.2. tétel, ha csupán n és m egyike tart végtelenhez, a másik pedig korlátos?

5. MEGOLDÁS. Igen. H_0 esetén sem $\sigma_{10}^2 = 0$, sem $\sigma_{01}^2 = 0$ nem fordul elő, így a határeloszlások nem elfajultak (lásd az előző szakasz 6. Feladatát).

Tárgymutató

- F -eloszlás
 - nem-centrál, 15
- χ^2
 - eloszlás, 14, 16
 - addíciós tétel, 15
- $\chi_n^2(\lambda)$, 14
- általánosított inverz, 36
- ANOVA, 18
- becsülhető, 37
- best linear unbiased estimator, 35
- BLUE, 35
- cella, 23
- egyszerű nullhipotézis, 60
- egyszeres osztályozás, 18, 41
- elégéses statisztika, 50
- eloszlás
 - nem-centrál χ^2 -, 14
 - nem-centrál F -, 15
- feltételes várható érték, 63
- Fisher–Cochran-tétel, 10, 16, 18, 19, 24, 28
- Fisher-féle információs mátrix, 46
- Gauss–Markov-tétel, 35, 37
- hipotézisvizsgálat
 - lineáris modellben, 40
- Hoeffding-féle U -statisztika, 62, 63
 - aszimptotikus normalitása, 64
- homoszkedasztikus, 34
- interakció, 23, 27
- központi határeloszlás tétel, 65
- kétszeres osztályozás, 23, 27
 - interakció figyelembevételével, 27
 - interakció nélkül, 23
- kollinearitás, 37
- konfidencia intervallum
 - várható értékre, 22, 26
- kvadratikus forma
 - négyzetösszeg alakja, 11
 - rangja, 10
 - szabadsági foka, 10
- legjobb lineáris torzítatlan becslés, 35
- legkisebb négyzetek módszere, 34
- legkisebb négyzetes becslés, 34
- likelihood-egyenlet, 50
- likelihood-függvény, 46, 50
- likelihood-hányados próba, 54, 59
- lineáris modell, 34, 39
- loglikelihood-függvény, 46, 50
 - aszimptotikája, 59
- magyarázó változók, 34
- Mann–Whitney-féle
 - U -próba, 67
 - U -statisztika határeloszlása, 69
- maradék négyzetösszeg, 39
- maximum-likelihood becslés, 39, 50
 - aszimptotikája, 53, 55
 - aszimptotikusan legkisebb szórású, 54
 - aszimptotikusan normális, 54
 - konzisztens, 52, 54
- nem-centralitási paraméter
 - χ^2 -eloszlásé, 14, 16
 - F -eloszlásé, 15
- normálegyenlet, 34
- normális eloszlás
 - paramétereinek becslése, 48
- OLS, 34
- ordinary least squares, 34
- ortogonális projekció, 37
- oszlophatás, 23
- próba

- likelihood-hányados, 59
- Mann-Whitney-féle U -, 67
- Neyman-Pearson-féle, 60
- Rao-féle, 60
- Wald-féle, 60
- Wilcoxon-féle
 - rangösszeg, 62, 67
- Rao-Cramér-egyenlőtlenség, 47
- regularitási feltételek, 46
- residual sum of squares, 39
- RSS, 39
- score vektor, 46
- sorhatás, 23
- statisztika
 - Hoeffding-féle U -, 62
- szórásfelbontó táblázat, 20, 25, 29
- szabadsági fok, 19
 - kvadratikus formáé, 10
- Szluckij lemmája, 64
- teljes átlag, 19
- teljes négyzetösszeg, 19
- Wilcoxon-féle
 - rangösszeg próba, 62, 67
 - rangösszeg statisztika
 - aszimptotikus normalitása, 69
 - szórása, 67
 - várható értéke, 67