

Statisztika Jegyzet  
az üzleti informatika szakirány számára  
(kézirat gyanánt)

Telcs András

December 16, 2005



## CONTENTS

0.1	Előszó . . . . .	5
0.2	Bevezetés . . . . .	6
<b>1</b>	<b>A leíró statisztika elemei</b>	<b>9</b>
<b>2</b>	<b>Valószínűségszámítási alapfogalmak</b>	<b>13</b>
2.1	A valószínűségi mező . . . . .	13
<b>3</b>	<b>Statisztikai alapfogalmak</b>	<b>17</b>
3.1	Bevezető, sokaság, minta . . . . .	17
3.2	Alapstatisztikák . . . . .	18
3.3	Határeloszlástételek avagy a valóság megismerhetősége . . . . .	20
<b>4</b>	<b>Becslélmélet</b>	<b>27</b>
<b>5</b>	<b>A legnagyobb valószínűség elve</b>	<b>33</b>
5.1	További példák, feladatok . . . . .	35
<b>6</b>	<b>Hipotézis vizsgálat</b>	<b>39</b>
6.1	Intervallum becslés . . . . .	39
6.1.1	$t$ - eloszlásra épített konfidencia intervallum . . . . .	41
6.1.2	Konfidencia intervallum az ismeretlen szórásra . . . . .	42
6.1.3	A mintaméret megválasztása . . . . .	43
6.2	Hipotézis vizsgálat . . . . .	43
6.2.1	A hipotézis vizsgálat menete . . . . .	45
6.2.2	Paraméteres próbák . . . . .	47
6.2.3	Az egymintás próbák további esetei . . . . .	47
6.2.4	Kétmintás próbák . . . . .	48
6.3	Próbák a szórásra vonatkozóan . . . . .	51
6.3.1	Egymintás próba . . . . .	51
6.3.2	Kétmintás próba . . . . .	52
6.3.3	A másodfajú hiba . . . . .	53

<b>7 Nem paraméteres próbák</b>	<b>55</b>
7.0.4 Khinégyzet próbák . . . . .	55
7.0.5 Illeszkedés, normalitás vizsgálat . . . . .	58
7.0.6 Próbák helyzeti paraméterek vizsgálatára . . . . .	58
7.0.7 Man-Whitney próba . . . . .	60
<b>8 Szórásanalízis</b>	<b>63</b>
8.0.8 Kétrészes osztályozás . . . . .	66
<b>9 Lineáris regresszió</b>	<b>67</b>
<b>10 Főkomponens analízis</b>	<b>73</b>
10.1 A lineáris algebra néhány eleme . . . . .	73
10.2 Véletlen vektorok elforgatása . . . . .	75
10.3 A vektóváltozó elemi statisztikai viselkedése . . . . .	76
10.4 A tapasztalati főkomponens . . . . .	78
<b>11 Osztályozás, klaszterezés</b>	<b>81</b>
11.1 Osztályozás . . . . .	81
11.1.1 A legközelebbi társ módszer . . . . .	83
11.2 Klaszter analízis . . . . .	83
11.2.1 $K$ -közép módszer . . . . .	84
11.2.2 Hierarchikus eljárások . . . . .	84
<b>12 Idősorok</b>	<b>87</b>
12.1 Alapfogalmak, definíciók . . . . .	87
12.1.1 Összefüggőségi struktúrák . . . . .	87
12.1.2 Az autokovariancia függvény tulajdonságai . . . . .	88
12.2 Idősorok transzformációja . . . . .	89
12.2.1 Nincs periodikus komponens . . . . .	89
12.2.2 Trend és szezonális . . . . .	90
12.3 Tapasztalati autokovariancia és autokorreláció . . . . .	90
12.4 Parciális autokovariancia függvény . . . . .	91
12.5 Fehér zaj . . . . .	91
12.6 Mozgóátlag (MA) folyamatok . . . . .	91
12.7 Autoregresszív (AR) folyamatok . . . . .	92
12.7.1 Példa, AR(1) folyamat . . . . .	93
12.7.2 Yule-Walker egyenletek . . . . .	94
12.8 Autoregresszív - mozgóátlag (ARMA) folyamatok . . . . .	94
12.8.1 A kauzalitás szükséges és elégséges feltétele . . . . .	94
12.9 Az átlag és az autokovariancia becslései . . . . .	95
12.9.1 A spektrálfüggvény és az autokovariancia kapcsolata . . . . .	95
12.9.2 Aszimptotikus normalitás . . . . .	96
12.9.3 $\gamma(n)$ becslése . . . . .	96
12.9.4 Az autokorrelációk mikor különböznek szignifikánsan 0-tól? . . . . .	97
12.10 ARMA modellek becslései . . . . .	97

12.10.1 Ismert $p$ és $q$ . . . . .	98
12.10.2 Ismeretlen $p$ . . . . .	98
12.10.3 A Durbin-Levinson algoritmus . . . . .	98
12.10.4 Az innovációs algoritmus . . . . .	99
12.10.5 Mozgóátlag folyamatok becslései . . . . .	101
12.10.6 Aszimptotikus viselkedés ARMA folyamatok esetén . . . . .	101
12.10.7 Maximum likelihood becslések . . . . .	102

## 0.1 Előszó

E jegyzet informatikus hallgatók számára íródott. Ezen belül az Üzleti informatika szakirány keretében kerül felhasználásra. Szerkesztésekor e két szempont alapján arra törekedtünk, hogy az informatikus képzés során szerzett ismeretekre, informatikus és mérnöki szemléletre alapozzunk. Ez egyben azt is jelenti, hogy lehetőség szerint a gyakorlati igényekkel fellépő olvasó számára íródott. Ezt erősíti a szakirány szabta feladat. A keretek szabta korlátok között az üzleti életben felmerülő problémákon keresztül kerülnek elő egyes statisztikai kérdések. Bizonyos mértékig az üzleti életben szokásos zsargont is becsempesztük a szövegbe, hogy a szerzett ismeretek későbbi alkalmazását ezzel is könnyítsük.

A kurzus több szinten sajátítható el. Mindenki maga dönti el mit céloz meg.

A jegyzet anyaga lehetővé teszi, hogy minimális üzleti, statisztikai szókinccsre tegyen szert az olvasó. Képes legyen egy üzleti, gazdasági problémában felismerni a statisztikai feladatot és azonosítani a feladat megoldásához szükséges módszert. Végül megértse a mások által megoldott statisztikai elemzés főbb üzenetét, interpretációját.

Az aki ennél többre törekszik, az alaposabb elsajátítás révén képessé válhat a statisztikusokkal együttműködni a feladat korrekt specifikálásában, mediálni a vállalat és a statisztikusok között biztosítva, hogy a válasz valóban arra a kérdésre adatik amit az üzleti élet felvetett.

A kitűzhető maximális cél, hogy a tárgyhoz kapcsolódó laborra is támaszkodva képessé válik a feladat meghatározásától a teljes megvalósításig minden lépést megoldani. Azaz végighaladnia következő lépéseken:

1. probléma specifikálás
2. munkafelvételek megfogalmazása
3. módszer kiválasztás
4. adatigény felmérése
5. adatgyűjtés megtervezése, kivitelezése, adatfeldolgozás
6. adatjavítás, szűrés, tesztelés
7. előzetes statisztikai feldolgozás
8. a statisztikai módszer alkalmazása
9. egybevetés a munkafelvételekkel, egyéb szempontokkal

10. esetleg a 2-8 ciklus részleges vagy teljes ismétlése
11. következtetések levonása, interpretálás, az eredmények visszafordítása az üzleti probléma nyelvére

Aki sikerrel megbirkózott a kurzus ismereteinek ilyen magas szintű elsajátításával hasznos és erős fegyvertár birtokába jutott. Természetesen nem rendelkezik a profi statisztikus teljes arzenáljával, de egye-egy irányban már képes lehet önállóan is fejleszteni ezirányú ismereteit illetve, ha kedvet érez haladó statisztikai stúdiumokra is vállalkozhat mint például a napjainkban oly fontos nemparaméteres eljárások vagy adatbányászati módszerek.

Annak érdekében, hogy a kitűzött feladatoknak megfeleljen, jegyzet felépítése a következő. Az első fejezetben a valószínűségi számítás alapfogalmai kerülnek röviden összefoglalásra majd a másodikban statisztika alapfogalmai kerülnek ismertetésre. A következő fejezetek egyre összetettebb módszereket ismertetnek. Minden fejezet végén külön papíron illetve számítógép segítségével megoldható feladatok találhatóak. Ezek megoldása biztosítja az anyag elsajátításának második illetve harmadik szintjét.

### **Köszönetnyilvánítás**

A szerző köszönetét szeretné kifejezni Maricza Istvánnak, azért, hogy jegyzetének idősor fejezete áttemeléséhez hozzájárult. Ugyanőt illeti még a közönet a baráti és szakmai beszélgetésekért, amellyel a szerzőt messzemenően segítette.

## **0.2 Bevezetés**

A statisztika a rendszerezett számbavétel igényéből fejlődött ki az évszázadok során. Az első számbavételi feladatok az időszámításunk előtti 4000 évre nyúlnak vissza. Kínában már ekkor összeírták a lakosságot, házakat, birtokokat. Ezen adatok a hatalom két nagyon fontos célját szolgálták az adókvetést és a katonai szolgálatot. Hasonló számbavételi igénye volt az egyiptomi uralkodóknak is a termény és a munkáerő felmérése kapcsán. Mózes is számbaveszi nemzettségét (Lásd Mózes IV. könyve) és nem kevesebb mint 603.550 felnőtt férfit tud magáénak. A gyermekek a magas halandóság miatt, illetve mert még sem munkára sem hadra nem foghatóak, nem számítottak, ahogy a nők sem. Másutt több nemzettségfő részletes vagyonsorsorolását találhatjuk a Bibliában mennyi nény illetve kétlábú jószágga rendelkeztek.

Az első hivatalos összeírásról is a Biblia számol be (Lukács 2.) "Augusztus császár rendeletet adott ki, hogy az egész földkerekséget összeírják össze. Ezt az első összeírást Cirinus, Szíria helytartója, bonyoltította le. Mindenki elment a maga városába, hogy összeírják." Adat Rómáról maradt fenn, fénykorában mintegy 1 millió lakosa volt. Ezt megelőző időből származó adatok szerint Athénban és Khorinthoszban 400 illetve 460 ezer rabszolga volt, amely feltehetően nagyobb a valóságnál.

A következő nagyszabású vállalkozásra ezer évet kell várni. A XI. századba születik az un. Doomsday Book, ami az akkori anglia fölbirtok és hűbéri viszonyainka felmérését szolgáltatta, megint csak adó és katonai szolgálat céljából. E hatalmas mű egyben számos hely és kortörténeti leírással volt kiegészítve.

Az első modernnek tekinthető statisztika John Graunt és William Petty nevéhez fűződik 1650-ből. Ők készítették az első születési és halálozási statisztikát.

S leíró statisztika, ezen belül a grafikus ábrázolás úttörője volt Florence Nightingale (1820-1910) aki adatsorokra támaszkodó ábrákkal győzte meg Nagybritannia katonai vezetését, hogy a Krími háború sebesültjei közül többen pusztultak el a hadikórházak rossz rökülményeinek következtében mint a sérülésekben a csatatéren. Ő az első statisztikus nő, egyben az ápolási szakma megteremtője is.

E régi időkre visszanyúló példák is mutatják a statisztika gyakorlati jelentőségét. Napjaink üzleti döntései pedig végképp elképzelhetetlenek azok nélkül. Hogy csak egy példát említsünk, az egyes termékek szenozális fogyasztási szokásai ismeretében gyártanak a termelők, rendelnek és készleteznek a kereskedők, például sört. A későbbiekben számos ilyen jellegű példát founk még ismertetni, konkrét módszerek kapcsán.





## Chapter 1

### A LEÍRÓ STATISZTIKA ELEMELI

Ebben a fejezetben a statisztika legrégebbi ágával a leíró statisztikával ismerkedünk meg.

A leíró statisztika a vizsgálandó kérdés kapcsán szóba jöhető összes objektum megfigyelésén alapul. Vegyünk egy példát. Csoportkirándulást szervezzünk. Az étkezések megrendeléséhez néhány kérdést kell tisztázni. Van-e a résztvevők között vegetáriánus, cukorbeteg, tej- vagy lisztérzékeny. Ki milyen italt fogyaszt reggelire, kávé, teát, tejet, kakaót. Természetesen ezeket az információkat egyszerűen összegyűjthetjük a (tegyük fel) 36 résztvevőtől. Ennek eredményeképpen 36 válaszlap lesz a kezünkben. Természetesen az étkezést nem személyenként fogjuk megrendelni. A megrendelésen az kell, hogy szerepeljen, hogy a fenti kérdésekre hány igen válasz érkezett. Elkezdünk tehát strigulázni. Az eredmény már a leíró statisztika egyik alapeleme.

**Definíció 1** *Alapsokaság a vizsgált objektumok, egyedek összesége.*

Esetünkben a kiránduló csoport tagjai.

**Definíció 2** *Oszály vagy kategória a megfigyelt objektumok valamely ismerv szerinti felosztása.*

Esetünkben például a reggeli ital mint az objektumok egy attribútuma négy lehetséges dolog lehet: kávé, tea, tej, kakaó. Ennek következtében a megfigyelt egyedek halmazának egy partícióját kapjuk, aszerint, hogy ki mit választ (kizárva a többszörös választást és a nem választ).

**Definíció 3** *A partíció osztályainak elemszámát abszolút gyakoriságnak nevezzük.*

Fenti kérdéseink kissé egyoldalúak voltak, a kapott válaszok mind kategoriális ismérveket tartalmaztak. Természetesen numerikus adatokra is szükségünk lehet, mint például, egy vagy két kávé iszik reggel. Általában a megfigyelt objektumok attribútumai, hasonlóan egy adatbázi egy rekordjának mezőjéhez különböző jellegűek lehetnek. Az attribútumok két nagy osztályát különböztetjük meg.

**Definíció 4** *Beszélhetünk kvantitatív és kvalitatív ismérvekről. Ezen belül különböző skálákról szokás beszélni. Kvalitatív skála lehet nominális, ha az attribútum csak egy név, címke az objektumok valamely osztályba sorolását jelöli. Ilyen például a hajszín vagy a személyi azonosító szám. Lehet a skála ordinális, avagy rendezett amikor az osztályok között valamilyen értelem szerű sorrend van, de továbbra is csak osztálycímkek, nevek az adataink. Ilyen például a fizetési kategóriák, A7, B4, J11 vagy a labdarúgó csapatok liga besorolása.*

Vigyázat a nominális címke is lehet szám, mint a személyi azonosító szám, mégsem érdemes például átlagot számítani belőle. Ennél is óvatosabban kezelendő, hogy gyakran a számítógépes feldolgozás céljára az osztályok neve helyett számokat használnak, például a hajszín kódok fekete=1, barna=2 és így tovább, de e számok szintén csak címkék, műveleteket végezni nincs értelme velük. Ez ugyanakkor nem jelenti azt, hogy az előfordulások számával ne lehetne műveleteket végezni és a segítségükkel statisztikai következtetésekre jutni.

**Definíció 5** *A kvantitatív adatok a következő csoportokba sorolhatóak. **Intervallum skáláról** lehet beszélni, ha az ismerv számszerű és elemei egy intervallumból kerülnek ki. Ilyen például az őszi félév oktatási napjainka dátuma.*

**Hányados skálával** van dolgunk, ha a megfigyelt objektumok kérdéses numerikus ismérveinek hányadosa ismert. Például nem tudjuk, ki mennyi cukrot tesz az italába, de azt igen, hogy  $X$  másfélszer annyit mint  $Y$ .  $Y$  pedig 0.9-szer annyit mint  $Z$ .

A leíró statisztika célja, hogy a rendelkezésre álló adatok alapján általános képet kapjunk a vizsgált objektumok összeségéről és egyben az adatok minőségéről. Ezt gyakran grafikus ábrázolással lehet elérni. Kiindulásul a megfigyelt gyakoriságok szolgálnak.

Egy ideig figyelmünket olyan vizsgálatokra korlátozzuk, amelyekben az objektumok egyetlen paraméterével foglalkozunk. Természetesen később több paraméter vizsgálatára sor kerül, hiszen gyakran igazán azok az izgalmas kérdések, hogy az egyik jól megfigyelhető paraméter viselkedéséből következtessünk a másik esetleg kevésbé megfigyelhetőre.

**Definíció 6 Gyakorisági tábla.** *Legyen a megfigyelt összes objektum egy adott ismerv szerint osztályokba van sorolva. Az egyes osztályokba eső elemek száma a gyakoriság, az osztálycímkék és e gyakoriságok alkotják a gyakorisági táblát. Például egy csoportban a hajszín gyakorisági táblája*

<i>fekete</i>	4
<i>barna</i>	12
<i>vörös</i>	1
<i>szőke</i>	3
<i>összesen</i>	20

Jelölje  $f_i$  az  $i$ -edik osztály gyakoriságát, frekvenciáját.

**Definíció 7** *A relatív gyakoriság*

$$r_i = \frac{f_i}{N}$$

ahol  $N$  az összes megfigyelt elemek száma,  $i = 1, 2, \dots, K$  pedig az osztályok sorszáma,  $r_i$  az osztályok részarányát fejezi ki 0 és 1 között.

**Definíció 8** *A gyakorisági diagramm avagy hisztogramm a gyakorisági táblát ábrázolja grafikusan.*

```

1  * * * *
2  * * * * * * * * * *
3  *
4  **

```

Abban az esetben is lehet gyakoriságról, relatív gyakoriságról beszélni, illetve ábrát készíteni. Ilyenkor jól megválasztott (általában egyenlő) szélességű intervallumok alkotják az osztályokat és a gyakoriság az abba eső elemek száma. Az intervallum szélességének megválasztásakor úgy célszerű eljárni, hogy a kapott gyakoriságok jól tükrözzék a vizsgált kérdést, egy-egy intervallumba az összes elemek közül se túl sok se túl kevés ne essen.

Szokás még a relatív gyakoriságokat kördiagrammon is ábrázolni. **[ide abra]**

A gyakoriságok elemzését néhány egyszerű statisztika segíti.

**Definíció 9** *Módusz az az osztály vagy osztályok, amelyek a maximális elemszámot tartalmazzák.*

Fenti példánkban a barna a modális osztály.

A továbbiakban numerikus adatok, kvalitatív ismérvekre szorítkozunk.

**Definíció 10** *Medián a vizsgált mennyiség skáláján egy olyan érték, amely két egyenlő számú részre vágja a megfigyelt elemeket halmazát. Páratlan elemszám esetén ez a sorba rendezett értékek közül a középső, páros számú elem esetén a középső kettő átlaga.*

Vegyünk egy példát. Az alábbi adatok

7.57, 9.55, 8.82, 8.72, 6.96, 6.83, 11.42, 16.08, 6.83, 13.05, 11.36, 2.72, 8.55, 10.79, 9.97, 9.85, 7.86, 7.72, 12.90, 18.17, 7.72, 14.75, 12.84, 3.08, 9.67, 12.19, 11.26, 11.13, 8.89, 2.41

30 háztartás papírhulladék "termelésének" értékei. A sorbarendezés után a középső két elem, azaz a 15 és 16. átlaga 9.61. Azaz ez a medián.

**Definíció 11** *Hasonlóan definiáljuk a percentilist. a p-percentilis az a skálaérték, amely alatt a p százaléka található az elemeknek.*

**Definíció 12** *Nevezetes percentilis az alsó és felső kvartilis,  $Q_1$  illetve  $Q_3$ , amelyek a 25 illetve 75 százalékos percentilisen felelnek meg.*

Példánkban a 30% percentilis 8.21,  $Q_1 = 7.72$ ,  $Q_3 = 11.81$ .

**Definíció 13** *A fenti statisztikákat szokás **helyzeti paramétereknek** is nevezni, mert az adatok elhelyezkedéséről adnak felvilágosítást. Talán a legfontosabb helyzeti paraméter az **átlag**, azaz a mért  $x_i, i = 1, 2, \dots, N$  értékek számtani közepe.*

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i.$$

Az adatok szóródásának jellemzésére is több statisztikát lehet használni.

**Definíció 14** *A **terjedelem** nem más mint a maximális és minimális mért érték különbsége.*

$$range = \max_{1 \leq i \leq N} x_i - \min_{1 \leq i \leq N} x_i.$$

**Definíció 15** Az *interkvartilis terjedelem*

$$IQR = Q_3 - Q_1.$$

**Definíció 16** A leghasznosabb szóródási mérőszám a *variancia*

$$VAR = \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

illetve az eredeti skálára visszatérve a *szórás*

$$\bar{\sigma} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

**Gyakorlat 1** Számoljuk ki példánkban az átlagot és a szórást.

**Definíció 17** A *korrigált tapasztalati szórás* kis elemszámok esetén játszik szerepet, jelentőségéről még később lesz szó.

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$$

Szokás még az adatok lapultsgára illetve jobb vagy baloldalra hízására vonatkozó mérőszámokta is bevezetni. Ezekre itt nem térünk ki.

## Chapter 2

### VALÓSZÍNŰSÉGSZÁMÍTÁSI ALAPFOGALMAK

#### 2.1 A valószínűségi mező

Korábbi tanulmányaink során megismerkedtünk a Kolmogorov féle valószínűségi mező és azon értelmezett valószínűségi változó fogalmával. Az alábbiakban röviden áttekintjük az ide vonatkozó és a statisztikában nélkülözhetetlen fogalmakat és összefüggéseket. Bővebb bevezetőért, példákért, gyakorlatokért lásd [?].

**Definíció 18** Legyen  $\Omega$  tetszőleges halmaz. Ez az összes események halmaza. Szoktuk ezt **eseménytérnek** is nevezni. Abban az esetben, ha  $\Omega$  véges vagy megszámlálhatóan végtelen, akkor **atomos eseménytér**ről beszélünk. Az  $\omega \in \Omega$  eseményeket szokás **elemi eseményeknek** nevezni. Ezek már nem bonthatóak további eseményekre.

Véges eseménytér esetén könnyű elképzelni az eseményeket. Ilyen például, hogy a tanulócsoporthból véletlenül kiválasztott diák hajszíne fekete.

**Definíció 19** Összetett esemény minden nem egyelemű  $A \subset \Omega$  részhalmaza  $\Omega$ -nak.

**Definíció 20** Bevezetünk műveleteket az események illetve az azokkal azonosított halmazok között.

$$\begin{aligned} A + B &= \{\omega : \omega \in A \text{ vagy } \omega \in B\}. \\ AB &= \{\omega : \omega \in A \text{ és } \omega \in B\}. \end{aligned} \tag{2.1}$$

E definíciók az eseményalgebrában szokásos jelöléseket használják, természetesen ugyanakkor egybeesnek a halmazok közötti műveletekkel.  $A + B = A \cup B$ ,  $AB = A \cap B$ . Ha adott  $A_1, \dots, A_n, \dots$  megszámlálhatóan végtelen esemény ezek összegét jelölje

$$\sum_{i=1}^{\infty} A_i = A_1 + A_2 + \dots + A_n + \dots \tag{2.2}$$

**Definíció 21** Az  $A \subset \Omega$  esemény komplementere

$$\bar{A} = \{\omega \in \Omega : \omega \notin A\}. \tag{2.3}$$

**Definíció 22** A lehetetlen esemény az üres halmaz, jele  $\emptyset$ .

**Definíció 23**  $\mathcal{F}$  az  $\Omega$  részhalmazainak egy családja szigmaalgebra, ha nem vezet ki belőle a (2.2)összeadás, (2.1) szorzás és (2.3) komplementer képzés.

**Definíció 24** Az  $(\Omega, \mathcal{F}, \mathbb{P})$  hármast valószínűségi mező, ha  $\Omega$  valamely alaphalmaz,  $\mathcal{F}$  e fölött egy szigmaalgebra,  $\mathbb{P}$  pedig valószínűségi mérték, ami azt jelenti, hogy az alábbi axiómákat teljesíti. Legyen  $A, B \subset \Omega$ .

1.  $P(A) \geq 0$ .

2.  $P(\Omega) = 1$

3. ha  $AB = \emptyset$  akkor

$$P(A + B) = P(A) + P(B)$$

**Definíció 25** Feltételes valószínűséget definiálhatunk tetszőleges  $B \subset \Omega, P(B) \neq 0$  eseményre

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

segítségével.

**Gyakorlat 2** Lássuk be, hogy a feltételes valószínűség is kielégíti az 1.-3. axiómákat.

**Definíció 26** A  $\mathbb{B}$  valós halmazok családja Borel szigmaalgebrát alkot, ha tartalmazza az összes  $[x, y)$  intervallumot és szigmaalgebra. A  $B \in \mathbb{B}$  halmazokat Borel halmazoknak nevezzük.

**Definíció 27** Az  $(\Omega, \mathcal{F}, \mathbb{P})$ -n egy  $X$  leképezés,  $X : \Omega \rightarrow \mathbb{R}$ , valószínűségi változó, ha minden  $B$  Borel halmazra annak ősképe, azaz

$$X^{-1}B = \{\omega : X\omega \in B\} \in \mathcal{F}$$

azaz  $\mathcal{F}$  beli. Az ilyen halmazokat szoktuk még  $\mathcal{F}$ -mérhetőnek is nevezni.

Ha  $\Omega$  véges vagy megszámlálhatóan végtelen, azaz

$$\Omega = \{\omega_1, \dots, \omega_n, \dots\}$$

akkor diszkrét valószínűségi mezőről beszélünk, egyébként folytonosról. A várható érték és szórás fogalmát először diszkrét valószínűségi mezőn definiáljuk.

A diszkrét valószínűségi változó valószínűség eloszlását meghatározzák a

$$p_i = \mathbb{P}(\omega_i)$$

valószínűségek.

**Definíció 28** Legyen a diszkrét  $(\Omega, \mathcal{F}, \mathbb{P})$ -n egy  $X$  valószínűségi változó akkor ennek várható értéke, amennyiben az alábbi összeg létezik

$$E(X) = \sum_{i=1}^{\infty} X(\omega_i) \mathbb{P}(\omega_i) =: \sum_{i=1}^{\infty} x_i p_i.$$

Jelölje mint általában ezt röviden  $\mu = E(X)$ .

**Definíció 29** Legyen a diszkrét  $(\Omega, \mathcal{F}, \mathbb{P})$ -n egy  $X$  valószínűségi változó akkor ennek szórásnégyzete, avagy varianciája

$$V(X) = \sum_{i=1}^{\infty} (x_i - \mu)^2 p_i$$

amennyiben az összeg létezik. Szórása ez esetben

$$\sigma(X) = \sqrt{\sum_{i=1}^{\infty} (x_i - \mu)^2 p_i}.$$

Folytonos valószínűségi mező esetén az általánosság enyhe szűkítésével definiáljuk a várható értéket és szórást.

**Definíció 30** Legyen az  $X$  valószínűségi változó az  $(\Omega, \mathcal{F}, \mathbb{P})$ -n.  $X$  eloszlásfüggvénye  $F(x)$ , ha minden  $x \in \text{Im } X = \{y \in \mathbb{R} : y = X(\omega)\}$ -ra

$$0 \leq F(x) = \mathbb{P}(X < x) \leq 1$$

valószínűség sűrűségfüggvénye pedig  $f(x) \geq 0$ , ha  $F(x)$  abszolút folytonos és minden  $x$ -re, ahol  $F$  értelmezett.

$$F(x) = \int_{-\infty}^x f(y) dy$$

**Definíció 31** Adott  $(\Omega, \mathcal{F}, \mathbb{P})$ -n az  $A, B \subset \Omega$ , eseményekre a feltétele valószínűséget a

$$P(A|B) = \frac{P(AB)}{P(B)}$$

összefüggés definiálja, ahol feltesszük, hogy  $P(B) > 0$ .

**Definíció 32** Adott  $(\Omega, \mathcal{F}, \mathbb{P})$ -n az  $A, B \subset \Omega$ , események függetlenek, ha

$$P(A|B) = P(A).$$

Ez ekvivalens azzal, hogy

$$P(AB) = P(A)P(B).$$

**Definíció 33** Adott  $(\Omega, \mathcal{F}, \mathbb{P})$ -n, az  $X_i, i = 1, 2, \dots$  valószínűségi változók páronként függetlenek, ha

$$\mathbb{P}(\{X_i < x\} \{X_j < y\}) = \mathbb{P}(X_i < x) \mathbb{P}(X_j < y)$$

teljesen függetlenek, ha minden  $k$ -ra és minden  $i_1, \dots, i_k$  index  $k$ -asra

$$\begin{aligned} & \mathbb{P}(\{X_{i_1} < x_{i_1}\} \{X_{i_2} < x_{i_2}\} \dots \{X_{i_k} < x_{i_k}\}) \\ &= \mathbb{P}(X_{i_1} < x_{i_1}) \mathbb{P}(X_{i_2} < x_{i_2}) \dots \mathbb{P}(X_{i_k} < x_{i_k}). \end{aligned}$$

**Definíció 34** Adott  $(\Omega, \mathcal{F}, \mathbb{P})$ -n,  $A \subset \Omega, \mathbb{P}(A) > 0$  eseményre értelmeztük a feltételes valószínűség fogalmát. Ez egy újabb valószínűségi mezőt is definiál  $(\Omega, \mathcal{F}, \mathbb{P}(\cdot|A))$ -t, amin az  $A$ -ra vonatkozó feltételes várható érték valamely  $X$ -re a fenti definícióból adódik, jele  $E(X|A)$ , diszkrét esetben

$$E(X|A) = \sum_{i=0}^{\infty} x_i P(X = x_i|A) \mathbb{I}(A)$$

folytonos esetben, tegyük fel, hogy létezik a  $(\Omega, \mathcal{F}, \mathbb{P}(\cdot|A))$ -n az  $X$  feltételes sűrűség függvénye  $f(x|A)$  ekkor

$$E(X|A) = \int x f(x|A) dx.$$

A feltételes várható érték fogalmát kissé szűkített értelmezéssel definiáljuk, elkerülendő bizonyos fogalmi nehézségeket. Diszkrét esetben a definíció viszonylag egyszerű.

**Definíció 35** Legyen,  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mező,  $X, Y$  két valószínűségi változó. legyenek értékészleteik rendre  $x_i, y_i$ , azt az eseményt pedig, hogy  $X = x_i, A_i$  illetve  $Y = y_j$  jelölje  $B_j$ . Ekkor

$$\begin{aligned} E(X|Y) &= \sum_{j=0}^{\infty} E(X|B_j) \mathbb{I}(B_j) \\ &= \sum_{j=0}^{\infty} \sum_{i=0}^{\infty} x_i \mathbb{P}(A_i|B_j) \mathbb{I}(B_j). \end{aligned}$$

**Definíció 36** Folytonos  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mező és  $X, Y$  változók esetén az  $f(x|Y)$  sűrűség függvényt a  $h(x, y)$  együttes sűrűségfüggvényen keresztül definiáljuk, feltéve annak létezését, valamint, hogy  $Y$  sűrűségfüggvénye  $f_Y(y) > 0$ :

$$f(x|y) = \frac{h(x, y)}{f_Y(y)}.$$

**Definíció 37** Folytonos  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mező és  $X, Y$  változók esetén tegyük fel, hogy létezik az  $f(x|Y)$  sűrűség függvény, ekkor a feltételes várható érték

$$E(X|Y) = \int x f(x|Y) dx.$$



## Chapter 3

### STATISZTIKAI ALAPFOGALMAK

#### 3.1 Bevezető, sokaság, minta

A statisztika mint azt a bevezetőben említettük két fő ágra bomlik, leíró statisztikára és matematikai statisztikára. Utóbbi módszerei matematikai alapon nyugszanak, lényegüket tekintve viszont olyan praktikus alkalmazható módszereket foglal össze, amelyek segítségével korlátozott ismeretek alapján, következtetéseket lehet levonni, ezek alapján például üzleti döntéseket lehet hozni mégpedig úgy, hogy közben a következtetés, döntés megbízhatóságára vonatkozóan is vannak ismereteink. Az élet számos területén találkozunk ilyen feladatokkal, az üzleti élet különösen sok ilyen állít elénk. Állandóan döntéseket kell hoznunk, kockázatvállalást, esélylatolgatást kell végeznünk. Álljon itt csak egyetlen példa. Új terméket szereténk piacra dobni. A termék fogadtatása nagyban függ a bevezető reklám sikerétől. A kampány kialakítására több alternatívát is kidolgozunk. Melyik lesz igazán eredményes? Melyiket válasszuk? A szubjektív vezetői döntés, a vállalati tapasztalatok felhasználása éppúgy elképzelhető, mint a megcélzott fogyasztói csoportonból kiválasztott kis csoporton végzett kísérlet, mérés. Azaz ún peer csoportot hívunk meg mintát választunk ki. Ezt több alcsoportra bontjuk és az egyes csoportokon megfelelő módszerekkel lemérjük a reklám hatékonyságát. Mit jegyeznek meg, mivel társítják s.t.b. . Az így kapott eredményeket mintegy kivetítjük a teljes fogyasztói körre, célközönségre, feltételezzük, hogy (kellően nagy és gondosan választott minta és gondosan kivitelezett mérés eredményeképpen) az alternatív kampányok közül az lesz a leghatékonyab a piacon, amelyik a kísérleti körülmények között az volt.

Általánosan tehát egy sokaság, azaz egyedek objektumok összességének bizonyos, teljes egészében nem megfigyelhető tulajdonságairól szeretnénk tudomást szerezni. Ehhez mintát veszünk (megfelelő módon) a sokaságból. A minta minden elemének kérdéses tulajdonságait megvizsgáljuk, majd a kapott eredmény segítségével visszakövetkeztetünk a teljes sokaság ugyanezen tulajdonságaira. Szokás ezt statisztikai következtetésnek nevezni. (lásd ?? ábrát).

Ismrekedésünket a statisztikai módszerekkel az egyváltozós statisztika körében kezdjük, azaz amikor a sokaság egyedeinek egyetlen tulajdonságát vizsgáljuk. Ilyen például az az egyszerű adat, hogy ki hány percet beszél telefonon, hányas cipőt visel, hány éves. Ha egy értékpapírt vizsgálunk, a sokaság lehet a papír kereskedésének napjai, a tulajdonság pedig a napi árváltozás.

Általánosságban egy  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mezőt vizsgálunk, ahol  $\mathbb{P}$  nem ismert. A  $\mathbb{P}$  valószínűségi mérték gyakran egy ismert mértékcsalád eleme, amelyet (egy vagy több) paraméter jellemez. Ilyen lehet például a normális eloszlás vagy a Poisson eloszlás. Ez esetben paraméteres problémát vizsgálunk.

Később nem paraméteres problémák vizsgálatára is szép módszerekkel ismerkedünk majd meg, de először a paraméteres problémákkal foglalkozunk, mivel történetileg is ezek alakultak ki előbb és talán az első ismerkedés is velük a könnyebb.

Fenti példáink is számszerűsíthető tulajdonságra vonatkoztak, tehát alapvetően egy  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mezőn értelmezett  $X$  valószínűségi változót vizsgálunk. Miről is van szó? Miért ne a sokaság összes eleme a vizsgálódás tárgya. Általában nincs mód megmérni a sokaság minden elemére vonatkozóan a kérdéses mennyiséget, ez eleve lehetetlen, vagy igen költséges, esetleg más okból kerülendő. Például egy új termék bevezetése előtt természetesen nincs információnk a fogadtatásról, ha meg már bevezettünk a kérdés már eldőlt. Más esetekben a teljes körű mérés bár lehetséges lenne, de a konkurensok előtt titkolni szeretnénk milyen döntésre készülünk, ezért ezt el kell kerülni.

Igy tehát modellt alkotunk. A sokaság bármely eleme lehetne a megfigyelésünk tárgya, ezért azt tesszük fel, hogy a megfigyelt objektum véletlen, annak tulajdonságát méri  $X$ , ez pedig egy valószínűségi változó. Érdeklődésünket erre az  $X$ -re összpontosítjuk. Ahhoz hogy az  $X$  eloszlásáról képet alkassunk mintát vezetünk a sokaságból, megmérjük a minta elemein a kérdéses mennyiséget, a kapott értékeket jelölje  $X_1, X_2, \dots, X_n$ .

A statisztika készítés gyakorlatának egyik sarkalatos lépése a minta kiválasztása. Később röviden majd kitérünk majd a mintavételi eljárások alapjaira, de ezek alapvetően meghaladják e jegyzet kereteit.

Mindvégig feltesszük, hogy az  $X_i$  valószínűségi változók teljesen függetlenek és azonos eloszlásúak, ha ettől a konvenciótól eltérünk, azt expliciten ott megjegyezzük.

### 3.2 Alapstatisztikák

Legyenek  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változók. Azaz az  $X$  valószínűségi változó független replikátumai. Azt a tényt, hogy az  $X, Y$  változó azonos eloszlású az  $X \sim Y$ -al fogjuk röviden jelölni. A fenti feltevés tehát azt jelenti, hogy  $X_i \sim X$  minden  $i = 1, 2, \dots, n$ -re.

**Definíció 38** *Az*

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

*értéket mintaátlagnak nevezzük. Ha konkrét realizációról beszélünk akkor ezt néha a*

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

*jelöléssel hangsúlyozzuk.*

**Állítás 1** *Ha létezik a  $\mu = E(X)$  várható érték, akkor*

$$E(\bar{X}) = \mu.$$

A bizonyítást az olvasóra bizzuk.

**Definíció 39** A minta tapasztalati szórásnégyzete, avagy varianciája ( $n > 1$  esetén)

$$V = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

és a tapasztalati szórás

$$S = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}.$$

Ha a konkrét realizációt akarjuk hangsúlyozni, akkor a

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

jelölést fogjuk használni.

**Tétel 3 (Steiner)** Tetszőleges  $x_i$  és  $c$  valósakra

$$\frac{1}{n} \sum_{i=1}^n (x_i - c)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - c)^2.$$

**Bizonyítás.**

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (x_i - c)^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - c)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{2}{n} \sum_{i=1}^n (x_i - \bar{x})(\bar{x} - c) + (\bar{x} - c)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - c)^2 \end{aligned}$$

mert a  $\sum_{i=1}^n (x_i - \bar{x}) = 0$ . ■

**Következmény 1**

$$\min_c \frac{1}{n} \sum_{i=1}^n (x_i - c)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

**Állítás 2**

$$S^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \quad (3.1)$$

**Bizonyítás.** Következik a Steiner Tételből  $c = 0$ -val. ■

A tapasztalait szórásnégyzet és szórás mellett bevezetjük a korrigált tapasztalati szórásnégyzetet szórást.

**Definíció 40** *A korrigált tapasztalati szórásnégyzet*

$$V^* = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

*és korrigált tapasztalati szórás*

$$s^* = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}.$$

### 3.3 Határeloszlástételek avagy a valóság megismerhetősége

A cím esetleg talányosnak vagy fellengzősnek tűnik, de mindjárt látni fogjuk, hogy a határeloszlás tételek valóban a valóság megismerésének egyik kulcsa. A formális állítások előtt hagy utaljunk vissza a legegyszerűbb, mondhatni ma már közhely számba menő összfüggésre. Ha egy pénzdarabbal sok dobást végzünk, a fejel részaránya egyre pontosabban közelíti az  $1/2$  értéket. Ez kicsit általánosabban is így van, ha a pénz nem szabályos, vagy más függelenül ismételhető azonos módon lezajló kísérletet ismételgetünk, amelynek több véletlen kimenetele van, akkor egy adott kimenetel relatív gyakorisága egyre pontosabban közelíti azt az értéket, amit axiomaként mint valószínűség ahhoz rendelünk. Példa lehet akár a kocka dobás esetén a 6-s kimenetele. De lehet az a megfigyelt jelenség, hogy milyen valószínűséggel választanak a vevők a jobb illetve a bal kezük ügyébe eső ugyanazon termékből.

Az előző részben láttuk, hogy a tapasztalati átlag átlaga maga a sokaság átlaga (lásd 1. Állítás), azaz a mintaátlag jól céloz. Ennél jóval több is igaz. Először is a mintaátlag szórására vonatkozó alapvető észrevétel következik.

**Állítás 3** *Ha  $X$ -nek létezik várható értéke és szórása  $\sigma$ , akkor az  $X$ -ből vett  $n$  elemű minta átlagára igaz, hogy*

$$\sigma(\bar{X}) = \frac{\sigma}{\sqrt{n}}. \quad (3.2)$$

**Bizonyítás.** A függetlenség alapján

$$\sigma^2(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2(X_i) = \frac{n\sigma}{n^2}$$

amiből az állítás már adódik. ■

**Definíció 41** *Ha egy  $A$  esemény bekövetkezésére vonatkozóan ismételt független kísérleteket végzünk, kiszámolhatjuk a tapasztalati részarányt a bekövetkezések  $k_A$  száma és a megfigyelések  $n$  számának hányadosaként.*

$$\bar{p} = \frac{k_A}{n}.$$

Természetesen azt várjuk, hogy  $\bar{p}$  jól közelíti az ismeretlen  $p = \bar{P}(A)$  értéket. Valóban, igaz a következő.

**Tétel 4** (*A nagy számok Bernoulli-féle törvénye*)

Tekintsünk egy  $A$  eseményt, amelynek valószínűsége  $p = \mathbb{P}(A)$  ( $0 < p < 1$ ). Legyen  $n$  független megfigyelésből a relatív gyakoriság  $\bar{p} = \bar{p}_n(A)$ , ekkor minden  $\varepsilon, \delta > 0$ -ra létezik  $N = N(\varepsilon, \delta)$ , hogy minden  $n > N$ -re

$$\mathbb{P}(|\bar{p}_n - p| < \varepsilon) > 1 - \delta.$$

**Definíció 42** Legyen adott egy  $F$  eloszlásfüggvényű valószínűségi változó  $X$ . Az ennől vett  $n$  elemű minta egy tapasztalati eloszlásfüggvényt határoz meg. Legyen  $k_n(x)$  az  $n$  elemű mintában az  $x$  érték alá eső  $x_i$  megfigyelések száma. Ekkor az  $F_n(x)$  tapasztalati eloszlásfüggvény a következő

$$F_n(x) = \frac{k_n(x)}{n}.$$

**Következmény 2** Ha speciálisan  $A = \{X < x\}$  valamely valószínűségi változóra, akkor

$$\mathbb{P}(|F_n(x) - F(x)| < \varepsilon) > 1 - \delta$$

ha  $n > N$ .

**Tétel 5** (*Nagy számok erős törvénye*)

Legyen egy esemény  $A, p = \mathbb{P}(A)$  ( $0 < p < 1$ ) valószínűséggel. Legyen  $n$  független megfigyelésből a relatív gyakorisága  $\bar{p} = \bar{p}_n(A)$ , ekkor minden  $\varepsilon > 0$ -ra létezik  $N = N(\varepsilon)$ , hogy minden  $n > N$ -re

$$\mathbb{P}(|F_n(x) - F(x)| < \varepsilon) = 1.$$

**Tétel 6** (*Centrális határelasztlás tétel*) Ha  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változók  $\mu$  várható értékkel és  $\sigma$  szórással, akkor

$$\mathbb{P}\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < y\right) \xrightarrow{n \rightarrow \infty} \Phi(y)$$

ahol  $\Phi$  a standard normális eloszlás eloszlásfüggvénye.

**Megjegyzés 1** Az állítás kissé pongyolán azt jelenti, hogy elég nagy  $n$  esetén ( $n > 30$  használható mint ökölszabály) az  $\bar{X}$  mint valószínűségi változó közelítőleg normális eloszlású, pontosabban  $N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$  eloszlású.

**Tétel 7** Ha  $X$  normális eloszlású valószínűségi változó, akkor

$$\mathbb{P}\left(\frac{\bar{X} - \mu}{s/\sqrt{n}} < y\right) = t^{(n-1)}(y)$$

ahol  $t^{(n-1)}(y)$  az  $n - 1$  szabadságfokú Student eloszlás eloszlásfüggvénye.

**Megjegyzés 2** Ha  $n > 30$  a Student eloszlás igen jól közelíti a standard normális eloszlást. Vegyük észre, hogy a második tétel sokkal erősebb feltevésen alapszik, miszerint az alapsokaság normális eloszlású, cserébe viszont nem csak közelítő állítás adható, hanem az eloszlás, a Student eloszlás akkor is adódik, ha a szórás értékének ismerete hiányában  $s/\sqrt{n}$ -el normáljuk a tapasztalati várható értéket.

**Megjegyzés 3** Megjegyezzük, hogy normális eloszlású valószínűségi változók összege is normális, azaz a tétel feltételei mellett  $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}$  maga standard normális eloszlású. A gyakorlatban az okoz nehézséget, hogy a sokaság szórása, sőt esetleg várható értéke sem ismert. Többek között a küzdelem ezen paraméterek meghatározásáért folyik és mindkét tétel fő mondanivalója éppen az, hogy a minta elemszámának növelésével a minta és a sokaság átlaga kis, kontrolált valószínűséggel tér csak el egymástól.

A centrális határeloszlás tétel speciális esete a bevezetőben említett feladat egzakt megoldása.

### Tétel 8

$$\mathbb{P} \left( \frac{\bar{p} - p}{\sqrt{\frac{p(1-p)}{n}}} < y \right) \xrightarrow{n \rightarrow \infty} \Phi(y).$$

**Megjegyzés 4** Azaz megint csak kissé egyszerűsítve, nagy  $n$ -re  $\bar{p}$  közelítőleg  $N \left( p, \sqrt{\frac{p(1-p)}{n}} \right)$  eloszlású. Igaz továbbá az is, hogy a nem ismert  $p$  helyett annak közelítését helyettesítve a képletbe az továbbra is fennáll, azaz:

$$\mathbb{P} \left( \frac{\bar{p} - p}{\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}} < y \right) \xrightarrow{n \rightarrow \infty} \Phi(y).$$

Az 6,7,7 Tételeket nem igazoljuk (lásd [?]).

Ezek után a megismerhetőség még teljesebb voltát igazoló tételt, a statisztika alaptételét mondjuk ki. Ez lényegében azt állítja, hogy az  $F$  eloszlásfüggvényű valószínűségi változóból vett egyre nagyobb mintából az  $F$  tetszőlegesen pontosan meghatározható.

### Tétel 9 (A statisztika alaptétele)

$$\sup_x |F_n(x) - F(x)| \xrightarrow{n \rightarrow \infty} 0$$

$I$  valószínűséggel.

A bizonyítástól eltekintünk. Az egyetlen nehézséget a szuprémum kezelése jelenti, hiszen minden fix  $x$ -re az  $A = \{X < x\}$  eseményre alkalmazható a fenti tétel, hiszen  $F(x) = \mathbb{P}(A)$ ,  $F_n(x) = \bar{p}(A)$ .

Az alábbiakban sokkal erősebb, úgynevezett próbákat is biztosító élesebb eredményeket ismertetünk.

**Tétel 10** (Szmirnov)

$$\lim_{n \rightarrow \infty} \mathbb{P} [\sqrt{n} (F_n(x) - F(x)) < y] = S(y)$$

ahol

$$S(y) = \begin{cases} 0 & \text{ha } y \leq 0 \\ 1 - e^{-2y^2} & \text{ha } y > 0 \end{cases}.$$

**Tétel 11** (Kolmogorov)

$$\lim_{n \rightarrow \infty} \mathbb{P} [\sqrt{n} |F_n(x) - F(x)| < y] = K(y)$$

ahol

$$K(y) = \begin{cases} 0 & \text{ha } y \leq 0 \\ \sum_{i=-\infty}^{\infty} (-1)^i e^{-2i^2 y^2} & \text{ha } y > 0 \end{cases}.$$

**Tétel 12** (Gnyegyenko) Legyen  $c = \lceil y\sqrt{2n} \rceil$ , ha  $F_n$  és  $G_n$  két tapasztalati eloszlásfüggvény, valamely  $F, G$  eloszlásfüggvény  $n$  elemű mintáira vonatkozóan, akkor  $F = G$ -ből következik, hogy

$$\mathbb{P} \left( \sqrt{\frac{n}{2}} \sup_x (F_n(x) - G_n(x)) < y \right) = \begin{cases} 0 & \text{ha } y \leq 0 \\ 1 - \frac{\binom{2n}{n+c}}{\binom{2n}{n}} & \text{ha } 0 < y < \sqrt{\frac{n}{2}} \\ 1 & \text{egyébként} \end{cases}$$

A tétel igazlásához bevezetõül egy klasszikus kombinatorikus gondolatmenetre van szükségünk.

**Probléma 13** Az alább ismertetendõ balott-tétel a szavazatszámolás lefolyására vonatkozó egy érdekes összefüggést mutat be. Ha egy polgármesteri posztért két jelölt  $A$  és  $B$  küzdött, a lehetséges  $n = a + b$  szavazatból  $a - t$  szerzett meg  $A$ ,  $b < a - t$  pedig  $B$ , felfet[dik, hogy mi annak a valószínűsége, hogy a szavazatok összeszámálása során mindvégig  $A$  vezet, azaz az összes részeredmény is az  $\tilde{o}$  gyõzelmét jelenti? Az összeszámolás folyamatát jól lehet ábrázolni. Tekintsük a szokásos síkbeli koordináta rendszert. Induljunk ki az origóból és rajzoljunk egy  $(1, 1)$  vektort, ha az elsõ szavazatot  $A$  kapta, ellenkezõ esetben az  $(1, -1)$  vektort rajzoljuk. Ezt az eljárást folytatjuk mindig az elõbbi két vektor egyikének lerajzolásával "megtoldva" az addig lerajzolt törött vonalat. (lásd a ?? ábrát.) Világos, hogy a törött vonal  $n$  hosszúságú és az  $y = a - b$  magasságban ér véget, azaz a  $(0, 0)$  és  $(n, a - b)$  pontokat köti össze. Hívjuk röviden az ilyen szabállyal rajzolt törött vonalakat utaknak.

**Tétel 14** (ballot-tétel) Legyen  $n \geq a > 0, b = n - a$ . Azon utak hossza amelyek azorigóból indulnak és  $(n, a - b)$ -ben úgy végzõdnek, hogy közben végig a felsõ félsíkban haladnak az  $x$  tengely érintése nélkül egyenlõ

$$\frac{a - b}{n} \binom{n}{a}.$$

A bizonyítás az igen sok helyen alkalmazható tükrözési elven alapul. Jelölje egy síkbeli  $C$  pont  $x$  tengelyre vonatkozó tükörképét  $C'$ .

**Tétel 15** (tükrözési elv) *Az  $x$  tengelyt érintő vagy metsző  $C$  pontból  $D$ -be vezető utak száma megegyezik az összes  $C'$  és  $D$  közötti utak számával.*

**Bizonyítás.** A bizonyítás a leszámolások köréből ismert módon történik, úgy, hogy a kérdéses utak között egy egy-egy értelmű megfeleltetést adunk, amiből persze következik, hogy az utak száma is egyenlő. Vegyünk szemre egy adott olyan utat ami  $C$ -ből  $D$ -be halad és érinti vagy metszi az  $x$  tengelyt. Ezen metszések közül van első, azaz balról jobbra (a lerajzolás sorrendje szerint, azaz időben) az első pont ahol az út érinti vagy metszi az  $x$  tengelyt. Legyen ez a pont  $(k, 0)$ , azaz a metszés a  $k$ -edik szakasz lerajzolásakor. Tükrözzük az út  $(0, 0)$  és  $(k, 0)$  közötti részét az  $x$  tengelyre (lásd ?? Ábra). Ezzel egy  $C'$ -ből  $D$ -be tartó úthoz jutunk. Az is világos, hogy minden ilyen  $C'$ -ből  $D$ -be tartó út valahol metszi az  $x$  tengelyt és mivel az első érintést illetve metszést választottuk a  $C, D$  úton, ezért a  $(k, 0)$  pont lesz az első metszés a  $C', D$  úton is. Ez az első (érintés illetve) metszéspont az utak mindkét szoban forgó halmazát diszjunkt részhalmazokra osztja. Az egy-egy értelmű megfeleltetést ezen részhalmazokon hozzuk létre. A tükrözés kölcsönösen egyértelmű megfeleltetés, ezért az utak között így kapott megfeleltetés is egy-egy értelmű. Ebből viszont következik, hogy a megfelelő utak száma is egyenlő. ■

**Következmény 3** *Annak a valószínűsége, hogy a fenti szavazatszámológási problémában végig  $A$  vezet  $\frac{a-b}{a+b}$ .*

**Bizonyítás.** Világos, hogyha az összes szavazatsorrendek száma  $\binom{a+b}{a}$ , ha ezek mind egyenlően valószínűek, akkor egy adott sorrend valószínűsége  $\frac{1}{\binom{a+b}{a}}$ . A 14 tételből ezért következik, hogy azon utak száma amik végig az  $x$  tengely felett haladnak az összes lehetséges utak számából kivéve azokat, amelyek érintenek vagy metszenek. Természetesen az első lépés felfelé kell, hogy történjen, ezért a vizsgált összes út  $(1, 1)$ -ből halad  $(n, a - b)$ -be, ahol  $n = a + b$ , vagyis egyel balra és letolva  $(0, 0)$ -ből halad  $(n - 1, a - b - 1)$ -be, ezek száma  $\binom{n-1}{a-1}$ , ugyanakkor a ballot tétel szerint az érintő vagy metsző utak száma egyenlő a  $(0, -1)$ ,  $(n, a)$  utak számával, azaz  $\binom{n-1}{a}$ -val, a keresett nem érintő utak száma ezért

$$\begin{aligned} \binom{n-1}{a-1} - \binom{n-1}{a} &= \frac{(n-1)!}{(a-1)!b!} - \frac{(n-1)!}{a!(b-1)!} \\ &= \frac{a}{n} \frac{n!}{a!b!} - \frac{b}{n} \frac{n!}{a!b!} = \frac{a-b}{a+b} \binom{n}{a}. \end{aligned}$$

Mivel az egyes utak (szavazócédula sorrendek) valószínűsége azonos, ezek összes száma pedig  $\binom{n}{a}$ , ezzel igazoltuk, hogy annak a valószínűsége, hogy a számlálás során végig  $A$  vezet  $\frac{a-b}{a+b}$ . ■

A ballot tétel gondolatmenete segítségével igazolhatjuk Gnyegyenko tételét.

Először egy újabb egyszerű kombinatorikus gondolatra van szükségünk. Legyen  $X_1, X_2, \dots, X_n$  az  $F$  illetve  $Y_1, Y_2, \dots, Y_n$  a  $G$ -ből választott  $n$  elemű minta. Keverjük össze és



rendezzük nagyság szerint őket. Így a  $Z_1 \dots Z_{2n}$  sorozathoz jutunk. Feltesszük, hogy ezen értékek mind különbözőek. Legyen

$$\varepsilon_i = \begin{cases} 1 & \text{ha } Z_i \text{ az } X_k \text{ sorozat eleme} \\ -1 & \text{egyébként} \end{cases}.$$

Ezzel megint egy  $2n$  hosszú töröttvonalat is kapunk, ha az  $(1, \varepsilon_i)$  vektorokat összefűzzük.

**Lemma 1** *A fenti jelölések mellett*

$$\sup_x (F_n(x) - G_n(x)) = \frac{1}{n} \max_{0 \leq i \leq 2n} S_i.$$

**Bizonyítás.** Az  $n(F_n(x) - G_n(x))$  kifejezés az  $x$ -nél kisebb  $X_j$  belüli és  $Y_k$  belüli elemek számának különbsége. Az  $x$  növekedtével ez pontosan akkor változik, mégpedig  $\varepsilon_i$ -vel, ha  $X_i$  illetve  $Y_i$  éppen  $x$ . ■

**A 12 Tétel bizonyítása.** Mivel a feltevés szerint  $F = G$  ezért az  $X_i$  és  $Y_i$  sorozatok elemei azonos eloszlásúak és teljesen függetlenek, amiből következik, hogy azok összes sorrendje azonos valószínűségű. Azaz bármelyik sorozat valószínűsége  $\frac{1}{\binom{2n}{n}}$ . Azon sorozatok száma pedig amelyekre  $\max_{0 \leq i \leq 2n} S_i < z\sqrt{2n}$  azon utak száma, amelyek a  $c = \lceil z\sqrt{2n} \rceil$  egyenes alatt maradnak. Alkalmazzuk a tükrözési elvet most az  $y = c$  egyenesre (lásd ?? Ábra). Világos, hogy  $(0, 0)$ -ból a  $(2n, 0)$ -ba haladó  $c - t$  nem érintő utak száma  $\binom{2n}{n} - \binom{2n}{n-c}$  lesz. Osztva az összes utak számával adódik az állítás. ■

**Megjegyzés 5** *Határátmenet képzésével Gnyegyenkó 12 tételéből következik Szmirnov 10 tétele. A tükrözési elv ismételt alkalmazásával és határátmenet képzésével lehet igazolni Kolmogorov 11 tételét is.*



## Chapter 4

### BECSLÉSELMÉLET

A vizsgált  $(\Omega, \mathcal{F}, \mathbb{P})$  valószínűségi mezőt gyakran jellemzi egy a  $\mathbb{P}$  valószínűségi métrékhez kapcsolódó  $\vartheta$  paraméter. Ezt néha a  $\mathbb{P}_\vartheta$  jelöléssel is hangsúlyozzák. Ilyen paraméter mondjuk az exponenciális eloszlás  $\lambda$ -ja, vagy a normális eloszlás várható értéke,  $\mu$ , vagy szórása  $\sigma$ . De ha például egy pénz feldobásánál a fej valószínűségét vizsgáljuk, lehet a  $\vartheta$  paraméter az  $1/2$ -től való eltérés is. Igen gyakori feladat, hogy egy  $X_1, X_2, \dots, X_n$  mintából megbecsüljük  $\vartheta - t$ .

**Definíció 43** Az  $X_1, X_2, \dots, X_n$  mintából készített statisztika

$$\bar{\vartheta} = f(X_1, X_2, \dots, X_n)$$

az aminek segítségével vissza kívánunk következtetni a sokaságot jellemző  $\vartheta$ -ra.

A feladat megfogalmazása már magában rejti azt, hogy valamilyen előfeltevéssel élünk (általában) a sokaságot jellemző  $\mathbb{P}$  valószínűségi mértéket illetően. Várható értékről sok esetben lehet például beszélni, de a Cauchy eloszlásnak nincs várható értéke, helyette helyzeti paramétréről szokás beszélni. Hasonlóan  $\lambda$ -ról beszélünk az exponenciális eloszlás esetében, tehát például, ha egymást követő telefonhívások között eltelt időt tekintünk, ilyenkor jó okkal feltesszük, hogy valamilyen ismeretlen  $\lambda$  paraméterű, de exponenciális eloszlással van dolgunk. Ha viszont a jelenségről tudható, hogy normális eloszlást követ, semmi értelme  $\lambda$  becsléséről beszélni.

Megközelítésünk tehát előzetes feltevésre, a jelenség valamilyen szintű ismeretére alapoz. Ez általában jellemzi a paraméteres statisztikai vizsgálatokat. A következő fejezetekben erről lesz szó.

Elkészítve egy  $\bar{\vartheta} = f(X_1, X_2, \dots, X_n)$  statisztikát, illetve becslést felmerül a kérdés, meyyire "jó" ez a becslés. Például megéri-e  $n$  elemű elvégezni a mérést. Ez a kérdés különösen akkor fontos, ha a vizsgálat tönkreteszi annak tárgyát, például a sorozatgyártásból kiemelt villanyégőket tartóssági tesztnek vetik alá. Több ezer órán át égnek, mérik a teljes élettartamukat. Ezt pedig csak akkor állapíthatják meg, ha a lámpa végül kiég. Természetesen ez az eljárás hosszadalmas és költséges, ha a teljes legyártott mennyiségen hajtánánk végre, akkor remek becslésünk lenne az átlagos élettartamra, csak nem lenne egy eladható égőnk sem. Így tisztázni kell, hogy hány elemű mintát érdemes használni. Ehhez egyrészt azt kell tisztázni, milyen pontos becslésre, milyen megbízható becslésre van szükségünk, másrészt azt kell valahogy megállapítanunk, hogy maga a mérési módszer, esetünkben a  $\bar{\vartheta} = f(X_1, X_2, \dots, X_n)$  statisztika milyen hibát "hordoz" magában. Az előbbi példára visszatérve, az égők élettartamára elég  $\pm 10$  óra pontos becslést adni. Jó

lenne ugyanakkor, ha az adott becslésünk megbízható lenne, azaz mondjuk, ha 1000 számítványra alkalmazzuk, akkor lehetőleg csak egy-két esetben forduljon elő, hogy a valódi élettartam és az általunk ígért (mérés alapján becsült) érték jobban eltér mint  $\pm 10$  óra. Szintén jogos elvárás lehet, hogy a költség, azaz a minta elemszámának növekedtével javuljon a becslés valamelyik fenti értelemben.

A becslés elmélet ezen kérdések ekzakt megközelítésére szolgál, melynek alapai kerülnek ebben a fejezetben bemutatásra.

Az alábbiakban tehát egy ismeretlen  $\vartheta$  paraméter  $\bar{\vartheta}$  becslésének jóságát leíró fogalmakat vezetünk be, majd néhány egyszerű statisztikára alkalmazzuk e fogalmakat.

**Definíció 44** Azt mondjuk, hogy  $\bar{\vartheta} = f(X_1, X_2, \dots, X_n)$  torzítatlan becslése  $\vartheta$ -nek, ha minden  $\vartheta$  esetén

$$\mathbb{E}_{\vartheta}(\bar{\vartheta}) = \vartheta.$$

Itt a várható érték a  $\mathbb{P}_{\vartheta}$  valószínűségi mértékből vett független  $X_i$  mintaelemekre vonatkozik.

**Állítás 4** Legyen

$$T(X_1, X_2, \dots, X_n) = \sum_{i=1}^n c_i X_i$$

Tegyük fel, hogy  $X$ -nek létezik  $\mu$  várható értéke. Ekkor  $T$  akkor és csak akkor torzítatlan, ha

$$\sum_{i=1}^n c_i = 1.$$

**Bizonyítás.** A várható érték linearitása miatt

$$\mathbb{E}_{\mu}(T) = \sum_{i=1}^n c_i E(X_i) = \mu \sum_{i=1}^n c_i$$

amiből következik az állítás. ■

**Következmény 4** Ebből következik  $c_i = \frac{1}{n}$ -et tekintve az is, hogy a mintaátlag a sokaság átlagának torzítatlan becslése, azaz

$$\mathbb{E}(\bar{X}) = \mu.$$

**Definíció 45** Azt mondjuk, hogy egy  $\bar{\vartheta}$  becslés hatásosabb mint egy  $\bar{\vartheta}'$ , ha torzítatlan becslések és

$$\sigma(\bar{\vartheta}) \leq \sigma(\bar{\vartheta}'),$$

feltéve persze a szórások létezését.

**Definíció 46** Azt mondjuk, hogy  $\bar{\vartheta}$  hatásos, ha minden más  $\bar{\vartheta}'$ -nél hatásosabb.

**Definíció 47** Egy  $\bar{\vartheta}'$  becslés hatásfokáról akkor beszélhetünk, ha létezik egy hatásos  $\bar{\vartheta}$  becslés, ekkor a hatásfok (efficiencia)

$$e(\bar{\vartheta}') = \frac{\sigma(\bar{\vartheta}')}{\sigma(\bar{\vartheta})}.$$

**Állítás 5** Ha a vizsgált eloszlás normális akkor  $\bar{X}$  hatásos.

Ezt az állítást nem igazoljuk, de a következőt viszont igen.

**Állítás 6** A várható érték lineáris becslései közül a mintaátlag a leghatékonyabb, ha létezik második momentum..

**Bizonyítás.** Tekintsünk egy

$$T = \sum_{i=1}^n c_i X_i$$

lineáris becslést. Ekkor a függetlenség miatt

$$\sigma^2(T) = \sigma^2\left(\sum_{i=1}^n c_i X_i\right) = \sum_{i=1}^n c_i^2 \sigma^2.$$

De a számtani és mértani közép közötti összefüggésből tudjuk hogy

$$\sum_{i=1}^n c_i^2 \geq \frac{1}{n} \left(\sum_{i=1}^n c_i\right)^2 = \frac{1}{n}$$

és egyenlőség akkor és csak akkor áll fenn ha minden  $c_i$  azonos, egyenlő  $\frac{1}{n}$ -el. Így tehát

$$\sigma^2(T) \geq \sigma^2(\bar{X}).$$

■

**Definíció 48** Egy  $\bar{\vartheta}_n = f(X_1, X_2, \dots, X_n)$  becslés asszimptotikusan torzítatlan, ha

$$\mathbb{E}(\bar{\vartheta}_n) \xrightarrow{n \rightarrow \infty} \vartheta.$$

**Definíció 49** Egy  $\bar{\vartheta}$  becslés gyengén konzisztens, ha minden  $\varepsilon > 0$ -ra létezik  $\delta > 0$  és  $N > 0$ , hogy minden  $n > N$ -re

$$\mathbb{P}_\vartheta(|\bar{\vartheta} - \vartheta| < \varepsilon) \geq 1 - \delta. \quad (4.1)$$

Erősen konzisztens, ha

$$\mathbb{P}_\vartheta\left(\lim_{n \rightarrow \infty} \bar{\vartheta} = \vartheta\right) = 1 \quad (4.2)$$

valamint  $n$ ;gyzetes  $k$ öz;  $p$ ben, ha

$$\mathbb{E}_\vartheta\left(|\bar{\vartheta} - \vartheta|^2\right) \xrightarrow{n \rightarrow \infty} 0. \quad (4.3)$$

Mint azt a valószínűségszámítás alapjainál megismertük az erős (4.2)és négyzetes középben (4.3) vett konvergencia (esetünkben konzisztencia) egyaránt maga után vonja a gyenge értelemben vett konvergenciát (4.1)(konzisztenciát).

**Tétel 16** *Ha létezik második momentum, akkor  $\bar{X}$  erősen konzisztens becslés.*

Az állítás következi a nagy számok erős törvényéből.

**Tétel 17** *A tapasztalati szórásnégyzet  $V$  és a korrigált tapasztalati szórásnégyzet  $V^*$  egyaránt erősen konzisztens.*

**Bizonyítás.**

$$V = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \rightarrow \mathbb{E}(\bar{X}^2) - \mathbb{E}(\bar{X})^2$$

igaz egy valószínűséggel a nagy számok erős törvénye szerint. Hasonlóan a korrigált tapasztalati szórásnégyzetre is  $\frac{n-1}{n} \rightarrow 1$  miatt. ■

**Tétel 18** *A  $V$  tapasztalati szóráss nem torzítatlan becslése  $\sigma^2$ -nek,  $S^*$  viszont az.*

**Bizonyítás.**

$$\begin{aligned} \mathbb{E}(V) &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \mu + \mu - \bar{x})^2 \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \{ (x_i - \mu)^2 - 2(x_i - \mu)(\bar{x} - \mu) + (\bar{x} - \mu)^2 \} \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - 2n(\bar{x} - \mu)^2 + n(\bar{x} - \mu)^2 \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - n(\bar{x} - \mu)^2 \right] = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right] - \mathbb{E}[(\bar{x} - \mu)^2] \end{aligned}$$

De mint láttuk  $\sigma^2(\bar{X}) = \frac{\sigma^2}{n}$  (lásd 3.2) ezért az utolsó kifejezés egyenlő

$$\sigma^2 - \mathbb{E}[(\bar{x} - \mu)^2] = \sigma^2 - \frac{1}{n}\sigma^2 = \frac{n-1}{n}\sigma^2 \neq \sigma^2.$$

Természetesen a korrigálás éppen kioltja az  $\frac{n-1}{n}$  tényezőt, így a korrigált tapasztalati szórásnégyzet már torzítatlan. ■

**Definíció 50** Egy  $\bar{\vartheta}$  statisztikát *elégésesnek* nevezünk, ha a becslendő  $\vartheta$  paraméterre vonatkozó minden információt tartalmaz, azaz ha  $F_{\alpha,n}$  jelöli az  $X_1, X_2, \dots, X_n$  együttes eloszlásfüggvényét, akkor

$$F_{\alpha,n}(x|\bar{\vartheta} = t) = F_n(x|\bar{\vartheta} = t),$$

azaz a  $\bar{\vartheta} = t$  feltevés mellett a jobboldal mint formula sem tartalmazza a  $\alpha$  paramétert.

E fogalom megértéséhez az alábbi példa a dhat segítséget.

**Állítás 7** Egy Poisson eloszlású sokaság paraméterének becslésére a mintaközép elégéses becslés.

**Bizonyítás.** Tegyük fel, hogy a paraméter  $\lambda$ . Mint tudjuk, Poisson eloszlású változók összege is ilyen és a paramétereik összegződnek, ezért speciálisan  $n\bar{X}$  is Poisson eloszlású  $n\lambda$  paraméterrel. Tekintsük a minta együttes eloszlásfüggvényét és használjuk ki a függetlenséget. Jelölje  $\bar{\lambda} = \frac{N}{n}$ , ahol  $N = \sum_{i=1}^n x_i$

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \bar{X} = \bar{\lambda}) \\ &= P(X_1 = x_1 | \bar{X} = \bar{\lambda}) P(X_2 = x_2 | \bar{X} = \bar{\lambda}) \dots P(X_n = x_n | \bar{X} = \bar{\lambda}) P^{-1}(n\bar{X} = N) \\ &= \frac{\lambda^{x_1}}{x_1!} e^{-x_1\lambda} \frac{\lambda^{x_2}}{x_2!} e^{-x_2\lambda} \dots \frac{\lambda^{x_n}}{x_n!} e^{-x_n\lambda} \frac{N!}{(n\lambda)^N} e^{\lambda N} = \frac{1}{n^N} \binom{N}{x_1, x_2, \dots, x_n}. \end{aligned}$$

Azaz a kapott képlet valóben nem tartalmazza a  $\lambda$  paramétert, tehát  $\bar{X}$  elégéses statisztika.

■





## Chapter 5

### A LEGNAGYOBB VALÓSZÍNŰSÉG ELVE

Ha egy doboz franciadrazs e k ozott z old, piros  s k ek sz in ek t at alhat ok, sz amuk pedig 20,50,33, akkor ha arra a k erd ésre kell v alaszt adnunk, hogy milyen sz in  lesz a v eletlen l h ih uzott drazs e, mindenki azt feleli, hogy piros, hiszen ezekb ol van a legt obb. Ugy is mondhatjuk, hogy az adott esetben ennek a val osz in s ege  $\frac{50}{103}$  a legnagyobb, mindenki azt v arja, hogy a legnagyobb val osz in s eg u esem eny k ovetkezik be. Nos ezt a természetes gondolatmenetet sok esetben alkalmazza a val osz in s eg sz am its a  s a statisztika is. Ebben a fejezetben a legnagyobb val osz in s eg elv en alapul o m odszerrel idegen sz oval a maximum likelihood m odszerrel ismerked unk meg,  ujra az egyv altoz os egyparam eteres probl em ak k or ere szor itkozva.

**P elda 19** *Legyen egy legy artott sz eri aban a hib as gy artm anyok r esz ar anya ismeretlen  $p$ . Azaz ha v eletlen szer uen egyet kih uzunk a sz eri ab ol, annak a val osz in s ege, hogy selejtet v alatszunk  $p$ . Szeretn enk mintav etel seg its eg evel meghat arozni  $p$ -t. Tegy uk fel, hogy egy  $n$  elem u mint at vett unk ki  s abb ol  $k$  bizonyult selejtesnek. Sz amoljunk form alisan, mi ennek a val osz in s ege.*

$$\mathbb{P}_p(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

*Keress uk azt a  $p$ -t amire ez az  ert ek maxim alis. Tegy uk fel, hogy  $0 < k < n$ .*

$$\frac{d}{dp} \mathbb{P}_p(X = k) = \binom{n}{k} k p^{k-1} (1 - p)^{n-k} - \binom{n}{k} p^k (n - k) (1 - p)^{n-k-1}$$

*A jobboldalt null aval egyenl ov e t eve*

$$0 = k(1 - p) - p(n - k)$$

*ad odik, amib ol  atrendez essel a nem t ul meglep o*

$$p = \frac{k}{n}$$

*ad odik. Term eszletesen meg kell gy oz odn unk arról, hogy a kapott  ert ek val oban sz els o ert ek  s maximum hely-e, de ez ebben az esetben j ol l athat o.*

A fenti p elda a modellje a maximum likelihood m odszernek.

**Definíció 51** Egy ismeretlen  $\vartheta$  paraméterű  $\mathbb{P}_\vartheta$  valószínűségi mértékből vett elemű minta maximum likelihood függvénye diszkrét valószínűségi változó esetében

$$L(x_1, x_2, \dots, x_n | \vartheta) = \mathbb{P}_\vartheta(X_1 = x_1) \mathbb{P}_\vartheta(X_2 = x_2) \dots \mathbb{P}_\vartheta(X_n = x_n)$$

illetve folytonos valószínűségi változó esetében

$$L(x_1, x_2, \dots, x_n | \vartheta) = f_\vartheta(x_1) f_\vartheta(x_2) \dots f_\vartheta(x_n)$$

ahol  $f_\vartheta(x)$  a  $\vartheta$  paraméterű valószínűségi változó sűrűségfüggvénye.

**Példa 20** Határozzuk meg a maximum likelihood módszer segítségével a normális eloszlású sokaság várható értékének és szórásának becslését. Folytonos esetben felhasználva a logaritmus függvény monotonicitását célszerű a likelihood függvény logaritmusát tekinteni. Mivel a normális eloszlás sűrűségfüggvénye

$$f_{\mu, \sigma}(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

ezért

$$\ln L(x_1, x_2, \dots, x_n | \vartheta) = \sum_{i=1}^n \left[ -\ln(\sigma) - \frac{1}{2} \ln 2\pi - \frac{(x_i - \mu)^2}{2\sigma^2} \right].$$

Véve ennek  $\mu$  illetve  $\sigma$  szerinti deriváltját

$$\frac{\partial}{\partial \mu} \ln L(x_1, x_2, \dots, x_n | \vartheta) = \sum_{i=1}^n \frac{2(x_i - \mu)}{2\sigma^2}, \quad (5.1)$$

$$\frac{\partial}{\partial \sigma} \ln L(x_1, x_2, \dots, x_n | \vartheta) = \sum_{i=1}^n \left[ -\frac{1}{\sigma} + \frac{(x_i - \mu)^2}{\sigma^3} \right]. \quad (5.2)$$

Keresve a gyököket

$$\sum_{i=1}^n \frac{2(x_i - \mu)}{2\sigma^2} = 0$$

egyenletből

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad (5.3)$$

adódik, illetve (5.2)-ből abba behelyettesítve (5.3) – t

$$\sum_{i=1}^n \left[ -\frac{1}{\sigma} + \frac{(x_i - \mu)^2}{\sigma^3} \right] = 0$$

$$\sum_{i=1}^n (x_i - \bar{\mu})^2 = n\bar{\sigma}^2$$

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mu})^2$$

adódik.

## 5.1 További példák, feladatok

**Gyakorlat 21** Legyen most a sokaság ismeretlen  $\lambda$  paraméterű Poisson eloszlású. ekkor a maximum likelihood függvény logaritmus az

$$f_{\lambda}(x) = \frac{\lambda^k}{k!} e^{-\lambda k}$$

diszkrét eloszlásfüggvény alapján

$$\ln L(x_1, x_2, \dots, x_n | \lambda) = \log \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \log(x_i!) - n\lambda$$

amiből

$$\frac{1}{\lambda} \sum_{i=1}^n x_i - n = 0$$

alján a maximum likelihood becslés  $\lambda$ -ra  $\bar{\lambda}_n = \bar{X}$ . Ehhez természetesen még szükséges annak ellenőrzése, hogy ez maximum hely, ami következik a második derivált negativitásából:

$$\frac{\partial^2}{\partial \lambda^2} \ln L(x_1, x_2, \dots, x_n | \lambda) = -\frac{1}{\lambda^2} \sum_{i=1}^n x_i < 0$$

ha  $x_i \neq 0$ .

**Gyakorlat 22** Ha a sokaság  $\lambda$  paraméterű exponenciális eloszlású, az okoskodás igen hasonló. A sűrűség függvény

$$f_{\lambda}(x) = \lambda e^{-x\lambda},$$

ezért a maximum likelihood függvény logaritmus az

$$\ln L(x_1, x_2, \dots, x_n | \lambda) = n \ln \lambda - \lambda \sum_{i=1}^n x_i.$$

Innen deriválásall a

$$\frac{n}{\lambda} - \sum_{i=1}^n x_i = 0$$

egyenlethez jutunk. Tehát  $\bar{\lambda} = \bar{X}$ , ha a maximum globális, ami megint következik a

$$\frac{\partial^2}{\partial \lambda^2} \ln L(x_1, x_2, \dots, x_n | \lambda) = -\frac{n}{\lambda^2} < 0$$

összefüggésből.

**Gyakorlat 23** Legyen most egy ismeretlen  $\lambda$  paraméterű Poisson eloszlásunk. Igazoljuk, hogy a mintaátlag elégséges statisztika  $\lambda$ -ra nézve.

A fenti példákban szinte automatikusan lehet alkalmazni a maximum likelihood módszert, az alábbiakban néhány trükkösebb példa következik.

**Gyakorlat 24** Legyen most az  $X \in [\alpha, 2\alpha]$  valószínűségi változó, amelynek sűrűségfüggvénye

$$f(x) = \frac{2x}{3\alpha^2},$$

$x \in [\alpha, 2\alpha]$ -n értelmezve. Adjunk maximum likelihood becslést  $\alpha$ -ra.

**Gyakorlat 25** Az visszatevéses mintavétel módszere. Szeretnénk megszámolni, hány kék-bálna él egy tengerezszakaszon. Legyen a keresett szám  $N$ . Most ez az ismeretlen paraméter. A következőképpen járunk el. hétig "vadászva" sárga festéklövedékekkel jelöljük meg a bálnákat. Legyen a megjelölt bálnák száma  $M$ . Ezek után a következő héten  $n$  darabot láttunk, ezek közül pedig  $s$  darab volt sárga festékekkel megjelölve. Mi az  $N$  értékének maximum likelihood becslése? Határozzuk meg először az  $L(s) = L(s|N)$  maximum likelihood függvényt.

$$L(s|N) = \frac{\binom{M}{s} \binom{N-M}{n-s}}{\binom{N}{n}}$$

$n$ .

Vizsgáljuk a

$$\frac{L(s+1)}{L(s)}$$

hányadost, mely értékekre  $>$  illetve  $<$  mint 1.

$$\begin{aligned} \frac{L(s+1)}{L(s)} &= \frac{\frac{M!}{s!(M-s)!} \frac{(N-M)!}{(n-s)!(N-M-n+s)!} \frac{n!(N-n)!}{N!}}{\frac{M!}{(s-1)!(M-s+1)!} \frac{(N-M)!}{(n-s+1)!(N-M-n+s-1)!} \frac{n!(N-n)!}{N!}} \\ &= \frac{\frac{1}{s!(M-s)!} \frac{1}{(n-s)!(N-M-n+s)!}}{\frac{1}{(s-1)!(M-s+1)!} \frac{1}{(n-s+1)!(N-M-n+s-1)!}} \\ &= \frac{(M-s+1)(n-s+1)}{s(N-M-n+s)} \geq 1 \end{aligned}$$

ha

$$N \leq n \frac{M}{s} - 1$$

azaz  $L$  ilyen értékekre nő, nagyobbakra csökken. ellenőrizni!

**Gyakorlat 26** Legyen most az  $X \in [\alpha, 2\alpha]$  egyenletes eloszlású valószínűségi változó. Adjunk maximum likelihood becslést  $\alpha$ -ra.

A problémához egy paradoxon is tartozik (lásd még [?]). Ennek ismertetése előtt oldjuk meg az eredeti feladatot. Nyilván a sűrűség függvény

$$f_\alpha(x) = \frac{x}{\alpha} : x \in [\alpha, 2\alpha].$$

*A maximum likelihood függvény logaritmus*

$$\ln L(x_1, x_2, \dots, x_n | \alpha) = \sum_{i=1}^n \ln \frac{x_i}{\alpha}$$

és ennek

$$\frac{\partial}{\partial \alpha} \ln L(x_1, x_2, \dots, x_n | \alpha) = -\frac{n}{\alpha},$$

ahol  $X_i \in [\alpha, 2\alpha]$  miatt  $-\frac{n}{\alpha}$  a maximumát

$$\bar{\alpha} = X_n^* = \max_{1 \leq i \leq n} X_i$$

helyen veszi fel. Ez tehát a maximum likelihood becslés. Tekintsük ezek után a

$$\beta = \frac{1}{2} \frac{n+1}{n+2} \bar{\alpha}$$

statisztikát és igazoljuk, hogy ez torzítatlan becslése  $\alpha$ -nak. Igazoljuk továbbá, hogy

$$\sigma^2(\beta) = \frac{1}{4n^2}.$$

Ezután készítsük el a

$$\gamma = \frac{n+1}{5n+4} \left( \min_{1 \leq i \leq n} X_i + 2 \max_{1 \leq i \leq n} X_i \right)$$

statisztikát és lássuk be, hogy ez hatásosabb mint  $\beta$ . Igaz ugyanis, hogy

$$\sigma^2(\gamma) \simeq \frac{1}{5n^2}.$$

Hogyan lehetséges ez? A válasz összefügg az elégségesség fogalmával. Szemléletesen is érthető, hogy  $\gamma$  több információt hordoz  $\alpha$ -ról mint  $\beta$ . Sőt az is igaz, hogy  $\gamma$  elégséges statisztika  $\alpha$ -ra nézve, ezért aztán nem meglepő, hogy az kisebb hibával közelíti. Ezzel elárultuk a választ, de természetesen az olvasónak még maradt feladata, igazolni kell, hogy a szóban forgó becslések torzítatlanok, fennállnak a szórásokra vonatkozó állítások, és érdemes megpróbálkozni a  $\gamma$  elégségességének igazolásával is.

**Gyakorlat 27** A következő szintén sokaságméretre vonatkozó példa igazán történelmi. A II. Világháború alatt az angol katonai hírszerzés meg kívánta becsülni a német ipar tankgyártási kapacitását. A hagyományos hírszerzési módszerek alapján a havi termelésre vonatkozóan 1550 darabos becslést adtak. Statisztikusok a következő módszert javasolták. Írják össze a kilőtt tankok gyártási sorozatszámait. Ennek alapján adnak majd becslést. Mint kiderült a precíz németek, minden hónapban más betűjellel kezdődő sorozatszámmal látták el a legyártott tankokat. Az 1941 júniusában gyártott tankok közül a kilőtt tankok közül legnagyobb sorozatszám a 244 volt. Adjunk maximum likelihood becslést ha feltesszük, hogy kilőtt tankok sorszáma egyenletesen oszlik el az összes sorszám között. Használjuk az előző gyakorlat becsléseit. (A legkisebb leolvasott sorozatszám 31 volt).

**Megjegyzés 6** Itt nem részletezendő szintén a populáció méretére vonatkozó becsléssel állapították meg angol statisztikusok (Bradley és, Efron [?]), hogy Shakespeare aktív szókincse 31534 szóból állt, passzívan további 35.000-t ismert. Érdekes ezt összevetni azzal, hogy az átlagember 2000 szót használ aktívan, 3-5000 szó a passzív szókincse, míg az igen választékos irodalmi közlés hozzávetőleg 8000-10.000 szó aktív ismeretét tételezi fel.

**Megjegyzés 7** Az átlag és a módusz (illetve a legvalószínűbb osztály) használata némi óvatosságot igényel. Nem igaz például, hogy az átlagos a leggyakoribb. Ez egyszerűen azért igaz, mert ferde eloszlások esetén az átlag és a módusz nem esik egybe. **[ábra!]** Például egy társadalomban sok szegény ebre él és igen kevés gazdag. Átlagos vagyoni helyzetű igen kevés van.

J. Reynolds vélte úgy, hogy az átlagot látjuk „szépnek” Nem világos, hogy e kijelentést hogyan kell értelmezni. Legyen ugyanis  $h$  az átlagmagasság,  $w$  az átlagsúly. Ha valaki arányos testfelépítésű akkor  $X$  magassághoz

$$Y = c_w X^3 \quad (5.4)$$

súly tartozik valamilyen  $c_w$  arányossági tényezővel. Mi következik akkor az átlagokra?

$$w = \mathbb{E}(Y)$$

ugyanakkor

$$h = \mathbb{E}(X).$$

Véve (5.4) mindkét oldalán a várható értéket

$$w = c_w \mathbb{E}(X^3).$$

Viszont nyilvánvaló, hogy általában  $\mathbb{E}(X^3) \neq \mathbb{E}(X)^3 = h^3$ , azaz az átlag magassághoz nem átlagos súly tartozik, vagyis, Átlag Polgár, akinek  $h$  magasságot tulajdonítunk és  $w$  súlyt nem lesz arányos, nem lesz "szép".

## Chapter 6

### HIPOTÉZIS VIZSGÁLAT

Ebben a fejezetben igen hasznos, jól alkalmazható módszereket ismertetünk amelyek egy-egy üzleti döntéshez adhatnak megbízható támpontot. Mintapéldánk lehet egy új termék, szolgáltatás bevezetése. Egy új szappanopera sikeréhez mondjuk az a minimum feltétel, hogy az egyidoben sugárzott műsorokkal szemben szerezzon 25% nézettséget, azaz egy-két hét után a piaci részesedése legyen 25%. A sorozat sorsa felől megint minta alapján döntünk. Zártkörű vetítést tartunk nézők egy csoportjának, akik mint otthon szabadon választhatnak a csatornák között. (Esetleg nem is tudják a vizsgálat valódi tárgyát, mondjuk chips-et és üdítőt kapnak, választhatnak a TV nézéshez mit fogyasztanak és ezekről kérjük a véleményüket.)

#### 6.1 Intervallum becslés

A  $p$  piaci részarány megállapítása nem jelent mást mint annak a  $p = \mathbb{P}(A)$  valószínűségnek a becslése, hogy egy véletlenül kiválasztott néző sorozatot választja-e. Mint má láttuk, a  $\bar{p}$  relatív gyakoriság remek pontbecslése  $p$ -nek. Ugyanakkor mint az az átlagmagasság kapcsán megjegyeztük, sem egyetlen ember esetében, sem egy  $n > 1$  elemű minta esetében sem fog pontosan fennállni, hogy  $p = \bar{p}_n$ . Mit mondhatunk akkor e helyett?

1.  $p$  és  $\bar{p}_n$  közel van egymáshoz. Mennyire?
2.  $p$  és  $\bar{p}_n$  távolsága kisebb mint  $\delta$ . Ez biztos, vagy csak valamilyen valószínűséggel igaz;

A centrális határeloszlás tétel segítségével pontos válasz adható ezekre a kérdésekre.

Vegyünk egy tipikus példát. Egy gyártmány számos paraméterrel rendelkezik, hogy csak a legegyszerűbb esetet vegyük egy üdítő palackra az van rá írva, hogy 1 litert tartalmaz. A vevő ezt el is várja. De valóban annyi van az üvegekben? Vegyünk mintát az üvegek közül és mérjük meg a tartalmukat. Legyen az  $n$  elemű minta alapján számolt tapasztalati átlag

$$\bar{X} = \bar{X}_n.$$

A centrális határeloszlás tétel alapján  $\bar{X}$  igen jó közelítéssel normális eloszlású, ha a vizsgált  $X$  valószínűségi változónak létezik várható értéke és szórása, a mintaelemek függetlenek és  $n > 30$ . Még pontosabban várható értéke  $\mu$  lesz és szórása  $\frac{\sigma}{\sqrt{n}}$ , azaz

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right),$$

**Megjegyzés 8** Vegyük észre, hogy a szórás  $\sigma$ -ről  $\frac{\sigma}{\sqrt{n}}$ -re változott. Mint mindjárt látni fogjuk, ez biztosítja azt, hogy a becslésünk egyre jobb lesz a mintaelemszám növekedtével.

Ebből következik, hogy

$$\mathbb{P}(\bar{X} < x) = \Phi\left(\frac{x - \mu}{\frac{\sigma}{\sqrt{n}}}\right)$$

vagy másképpen

$$\mathbb{P}\left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < x\right) = \Phi(x).$$

A normális eloszlás szimmetriája miatt, akkor az is igaz, hogy

$$\mathbb{P}\left(\left|\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}\right| < x\right) = 2\Phi(x) - 1. \quad (6.1)$$

A kényelem kedvéért szokás a jobb oldalon található valószínűséget valamilyen „nevezetes” közmegegyezéssel elfogadott értéknek választani, pl. 0.95, 0.975, 0.99. Azaz a standard normális eloszlás inverzét használjuk. Tekintsük először a

$$\mathbb{P}(X < x) = \Phi(x) = 1 - \alpha$$

adott  $\alpha$ -hoz tartozó  $x = z_\alpha$  értéket. (lásd ábra) Ez a  $z_\alpha$  érték vág ki a standard normális eloszlás „farkából”  $\alpha$  valószínűséget. Természetesen akkor

$$\mathbb{P}(-z_\alpha < X < z_\alpha) = 1 - 2\alpha.$$

Ha tehát a kívánalom mondjuk, hogy (6.1) jobb oldalán .95 valószínűség álljon, akkor a kivágott össz valószínűség legyen  $1 - \alpha = .95$ , akkor a  $z$  kritikus értékek az  $\alpha/2$  helyekhez tartozóak. Tehát a standard normális eloszlás esetén

$$\mathbb{P}(-z_{\alpha/2} < X < z_{\alpha/2}) = 1 - \alpha.$$

amiből feladatunkban

$$\mathbb{P}\left(-z_{\alpha/2} < \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < z_{\alpha/2}\right) = 1 - \alpha,$$

ezt átrendezve

$$\mathbb{P}\left(-z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \bar{X} - \mu < z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha,$$

illetve

$$\mathbb{P}\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

A kapott összefüggés azt jelenti, hogy az ismeretlen  $\mu$  sokaság átlag, az

$$\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$

intervallumba esik  $1 - \alpha$  valószínűséggel. Ezzel megszerkesztettük a  $\mu$ -re vonatkozó konfidencia intervallumot.



**Definíció 52** Másképp azt szokták mondani, hogy az  $\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$  intervallum  $1 - \alpha$  megbízhatósággal tartalmazza  $\mu - t$ , vagyis ez a  $\mu$ -re vonatkozó  $1 - \alpha$  szintű **megbízhatósági (vagy konfidencia) intervallum**. Ezzel úgynevezett **intervallum becslést** adunk  $\mu$ -re. Szokás  $\alpha - t$  a szignifikancia szintjének nevezni.

Általában kétoldali konfidencia intervallumokat szokás használni, de speciális esetekben (pl. az üvegbe eleve nem fér 1 liternél több) lehet egyoldalú konfidencia intervallumot is szerkeszteni, ilyenkor persze  $\alpha$  nem feleződik meg:

$$\mathbb{P}\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu\right) = 1 - \alpha$$

vagy

$$\mathbb{P}\left(\mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

### 6.1.1 $t$ -eloszlásra épített konfidencia intervallum

Az előző szakaszban láttuk, hogy amennyiben a tapasztalati átlag normális illetve közel normális eloszlású, akkor ennek a tudásunknak a birtokában konfidencia intervallumot lehet szerkeszteni az ismeretlen sokaság, azaz populáció átlagra. Abban az esetben, amikor a vizsgály valószínűségi változó  $X$  nem normális eloszlású, továbbá a populáció  $\sigma$  szórása ismeretlen, akkor viszonylag kis minták esetén a tapasztalati szórás segítségével normált statisztika

$$\frac{\bar{X} - \mu}{s}$$

nem standard normális eloszlású.

**Állítás 8** Ha  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változók, akkor a mintaátlag alábbi standardizáltja

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim t^{(n-1)}$$

$n - 1$  szabadságfokú Student, avagy  $t$ -eloszlás követ. A jelölés némi keverésével jelölje  $t^{(n-1)}(x)$  ennek az eloszlásnak a sűrűségfüggvényét.

A  $t$ -eloszlás értékeit célszerű táblázatból kikeresni vagy számítógép segítségével meghatározni. A 8 Állítás segítségével a megint lehetőség van konfidencia intervallumot szerkeszteni.

**Állítás 9** Legyen  $t_{\alpha/2}$  az a kritikus érték, amely a  $t$ -eloszlás jobboldali „farkából”  $\alpha/2$  területet vág ki, azaz, ha  $Y$   $t$ -eloszlású, akkor

$$\mathbb{P}(Y > t_{\alpha/2}) = \frac{\alpha}{2}$$

ekkor az eloszlás szimmetriája miatt

$$\mathbb{P}(t_{\alpha/2} < Y < t_{\alpha/2}) = 1 - \alpha. \quad (6.2)$$

**Következmény 5**

$$\mathbb{P} \left( t_{\alpha/2} < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < t_{\alpha/2} \right) = 1 - \alpha$$

és

$$\mathbb{P} \left( \bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right) = 1 - \alpha.$$

Az állítás nyilván következik a 8 Állításból és (6.2)-ből.

**Gyakorlat 28** Készítsünk „döntési diagrammot mikor kell normális és mikor kell  $t$  eloszláson alapuló intervallumot készíteni.

**Állítás 10** Részarány konfidencia intervalluma mindig normális eloszlásra épül, azaz

$$\mathbb{P} \left( \bar{p} - z_{\alpha/2} \sqrt{\frac{\bar{p}(1-\bar{p})}{n}} < p < \bar{p} + z_{\alpha/2} \sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \right) = 1 - \alpha. \quad (6.3)$$

**Bizonyítás.** A centrális határeloszlás tételből tudjuk, hogy

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

ezért a részarányra vonatkozóan is igaz, hogy

$$\frac{\bar{p} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1). \quad (6.4)$$

Az (6.3) összefüggés ezután abból következik, hogy a binomiális eloszlás szórása  $\sigma = \sqrt{p(1-p)}$ , ezért  $\bar{p}$  szórása  $\sqrt{\frac{p(1-p)}{n}}$ , amibe ha  $\bar{p}$  helyettesítünk az elkövetett hiba „másodrendben” kicsis(6. továbbra is igaz marad. ■

**Gyakorlat 29** Készítsünk „döntési diagrammot” mikor kell normális és mikor kell  $t$  eloszláson alapuló intervallumot készíteni.

### 6.1.2 Konfidencia intervallum az ismeretlen szórásra

Az előző rész sémája a következő módon foglalható össze. Tegyük fel, hogy  $Y$  valószínűségi változó, statisztika szolgál egy  $\eta$  paraméter becslésére. Ha tudjuk, hogy  $Y - \eta$  vagy  $\frac{Y}{\eta}$  milyen eloszlást követ, akkor ebből konfidencia intervallum szerkeszthető  $\eta$ -ra.

**Definíció 53** A  $\chi^2$  eloszlást implicit módon definiáljuk. Legyen

$$Y = \sum_{i=1}^n X_i^2$$

ahol  $X_1, X_2, \dots, X_n$  független standard normális eloszlás valószínűségi változók. Ekkor azt mondjuk, hogy  $Y$  eloszlása  $(n-1)$  szabadságfokú  $\chi^2$  (olvassd khi-négyszet) eloszlás.

**Állítás 11** Ha egy normális populációból vettünk  $X_1, X_2, \dots, X_n$  független mintát, akkor

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2.$$

**Bizonyítás.** Az állítás közvetlenkövetkezménye a (53) Definíciónak. ■

**Állítás 12** Az ismeretlen  $\sigma^2$  szórásnégyzet becslésére normális eloszlású populáció esetén a

$$\mathbb{P} \left( \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{\alpha/2}^2} \right) = 1 - \alpha$$

intervallum szerkeszthető. Azaz  $\sigma^2$ -t  $1-\alpha$  valószínűséggel tartalmazza a  $\left( \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}, \frac{(n-1)s^2}{\chi_{\alpha/2}^2} \right)$  intervallum.

### 6.1.3 A mintaméret megválasztása

A gyakorlati életben gyakran az a feladat, hogy a populáció átlagot bizonyos hibahatáron belül tartsuk, Mint az üvegtöltés példájában a térfogat  $1 \pm 0,01$  liter kell, hogy legyen. Ehhez először a populációátlagot kell a fenti módon jól becsülnünk, ezután lehet az egyes üvegek töltésére következtetni az átlag becse mellett a szórás becslését is felhasználva. Foglalkozunk a legegyszerűbb konfidencia intervallummal. A

$$\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

egyenőtlenség  $1 - \alpha$  valószínűséggel fenáll. Szeretnénk elérni, hogy  $\mu$ -re  $\bar{X} \pm B$  intervallumot kapjunk ugyanezen  $1 - \alpha$  megbízhatósággal. Adott, rögzített  $B$  mellett. Világos, hogy ehhez a

$$z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = B$$

összefüggésnek kell teljesülnie. Ha  $\alpha$  rögzített, ez meghatározza  $z_{\alpha/2}$ -t, ezért az  $n$  mintaméret az egyetlen választható paraméter. A legkisebb megfelelő  $n$

$$n = \left\lceil \frac{z_{\alpha/2}^2 \sigma^2}{B^2} \right\rceil.$$

**Gyakorlat 30** Adjuk meg az analóg mintaméretre vonatkozó összefüggéseket  $t$ -eloszlás esetére és részaránybecslés esetére is.

## 6.2 Hipotézis vizsgálat

A mindennapi életben, döntési szituációban gyakran nem az érdekel minket, hogy konfidencia intervallumot szerkesszünk, inkább az a kérdés bent van-e az ismeretlen átlag egy adott intervallumban vagy sem. Például a palackok töltési átlaga 1 liter. Ha igen, minden rendben, de ha nem akkor állítani kell a töltőberendezésen, le kell állítani a gépsort. Ez persze termelés, profit kiesést jelent, ezért ha lehet, csak akkor döntünk így, ha eléggé biztosak vagyunk abban, hogy a töltési átlag lényegesen eltér az előírttól. A mintavétel,

a mérés és statisztikai vizsgálat alapján egy dichotóm döntést hozunk, egy liter az töltési átlag, vagy nem az. Ennek a példának a végigvitelével mutatjuk be a hipotézis vizsgálat tipikus menetét.

A döntési folyamat első lépése a munkahipotézis felállítása. Ez az egyetlen „kényes” lépés, a többi rutin feladat.

Legyen  $\mu_0 = 1$ . Elvárásunk, hogy a populáció várható értéke  $\mu$  ezzel egybe essék. Ezt fogalmazzuk meg mint **null hipotézist**.

$$H_0 : \mu = \mu_0$$

Ennek logikai ellenéte az **alternatív hipotézis**

$$H_1 : \mu \neq \mu_0.$$

Mielőtt mintát veszünk, számításokat végzünk el kell, döntenünk, hogy mennyire „fontos”, hogy jól döntsünk. A döntést minta alapján fogjuk meghozni. A minta véletlentől függő, ezért több módon is felléphet hiba. Ezeket mutatja az alábbi táblázat.

		a valóság	
		$H_0$ igaz	$H_0$ hamis
hogyan döntünk	elfogadjuk $H_0$ -t	jó a döntés	másodfajú hibát vétünk
	elvetjük $H_0$ -t	elsőfajú hibát vétünk	jó a döntés

Mielőtt a hiba valószínűségének kiszámításához látunk, két jelölést vezetünk be.

$$\alpha = \mathbb{P}_{H_0}(\text{elvetjük } H_0\text{-t})$$

ahol a  $\mathbb{P}_{H_0}$  azon valószínűségi mérték amely  $H_0$  fennállása esetén jellemzi a sokaságot.

$$\beta = \mathbb{P}_{H_1}(\text{elfogadjuk } H_0\text{-t})$$

itt pedig  $\mathbb{P}_{H_1}$  azon valószínűségi mérték amely  $H_1$  fennállása esetén jellemzi a sokaságot. Ez utóbbi nem feltétlenül jól definiált, erre később visszatérünk.

Ha a feltételek biztosítják, hogy a mintaátlag jó közelítéssel normális akkor ez igaz a

$$z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

úgynevezett **próbafüggvényre** is, pontosabban  $z$  közel standard normális eloszlású lesz. Ebből következik, hogy

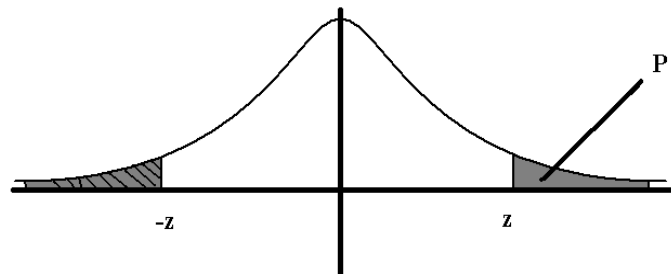
$$\mathbb{P}(|z| < z_{\alpha/2}) = 1 - \alpha$$

azaz

$$\mathbb{P}\left(\left|\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}\right| < z_{\alpha/2}\right) = 1 - \alpha$$

illetve annak a valószínűsége, hogy a próbafüggvény abszolút értéke nagyobb legyen mint  $z_{\alpha/2}$  egyenlő  $\alpha$ . Megint azt az okoskodást, használjuk, hogy egy kísérletben a nagy valószínűségű esemény bekövetkeztét tételezzük fel. Ha tehát feltesszük, hogy

$$\mu = \mu_0$$

Figure 1  $p$  az eloszlás farkában

akkor

$$\mathbb{P} \left( \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} < -z_{\alpha/2} \text{ vagy } \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} > z_{\alpha/2} \right) = \alpha.$$

Az (??) ábrára tekintve ez azt jelenti, hogy a  $z$  próbafüggvény a kis,  $\alpha/2 - \alpha/2$  valószínűségű alsó vagy felső farkba mutat. Ez nincs összhangban azzal a feltevésünkkel, hogy az  $1 - \alpha \gg \alpha$  valószínűségű esemény következik be, ezért a  $\mu = \mu_0$  feltevést vetjük el. Természetesen valamilyen más  $\mu$  lehet igaz, aminek megfelelően netrálva  $\bar{X} - t$  a próbafüggvény már a  $(-z_{\alpha/2}, z_{\alpha/2})$  intervallumba mutat, (lásd (1)Ábrát).

**Definíció 54** A  $(-z_{\alpha/2}, z_{\alpha/2})$  intervallumot szokás **elfogadási tartomány**nak nevezni, a  $(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$  pedig **elutasítási tartomány**nak. A **próba szintje**, vagy **szignifikancia szintje**  $\alpha$ , a **kritikus érték(ek)**  $\alpha$  illetve  $\alpha/2$ .

Azt mondjuk, hogy  $\alpha$  szinten elutasítjuk  $H_0$ -t ha  $z$  az elutasítási tartományba esik, illetve, hogy nem utasítjuk el, ha az ellenkezője igaz. Kicsit pongyolán lehet azt is mondani, hogy elfogadjuk  $H_0$ -t, de mint látni fogjuk ez nem igazán szerencsés, kissé félrevezető.

**Megjegyzés 9** Vegyük észre, hogy  $\alpha$ -t tetszés szerint megválaszthatjuk a döntési eljárás elején. Ebből pedig az is következik, hogy az elsőfajú hiba valószínűségét meg tudjuk szabni, azaz annak a valószínűségét, hogy elutasítjuk  $H_0$ -t pedig igaz. Hiszen ez a situáció akkor áll elő, amikor  $\mu = \mu_0$  és a próbafüggvény mégis a  $(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$  elutasítási tartományba esik, aminek a valószínűsége éppen  $\alpha$ .

### 6.2.1 A hipotézis vizsgálat menete

A hipotézis vizsgálat lépései a következőkben foglalhatóak össze.

1. A null hipotézis  $H_0$  megfogalmazása.
2. A próbafüggvény kiválasztása a feltételek alapján.
3. A  $\alpha$  szignifikancia szint meghatározása a probléma természetétől függően.
4. A próbafüggvény eloszlásának alapján az elutasítási tartomány meghatározása.
5. Mítavételezés, a próbafüggvény kiszámítása.

## 6. Döntés attól függően, hogy a próbafüggvény hova esik.

Mint említettük az egyetlen kényes feladat a nullhipotézis felállítása (és talán mint minden statisztikai feladatnál a jó mintavételezés). A gyakorlati életben a  $\mu = \mu_0$  null hipotézis mellett igen gyakori a  $\mu \geq \mu_0$  illetve a  $\mu \leq \mu_0$  alakú feltevés. Mint a palacktöltő sor működtetője a  $\mu = \mu_0$  feltevés volt fontos számunkra. A vevő szempontja más. Őt az érdekli, hogy elég sokat vagy elég jót kapjon a pénzéért. A hirdetések is igen gyakran azzal kínálnak egy portékát, hogy az azt jellemző valamilyen mennyiség kisebb vagy nagyobb mint egy adott érték. Erre lehet példa a gyümölcslé, üdítő. Nagy betűkkel a csomagolás oldalán hirdeti, hogy több mint 50% gyümölcstartalommal, vagy „kevesebb mint 1% cukor tartalom”. A fogyasztó vagy a konkurencia esetleg gyanút fog, hogy a gyártó nem teljesíti azt, amit a hirdetésben állít. kritikus esetben ez pert vonhat maga után, a bíróságnak kell döntenie a felek között. Vegyük ebből a nézőponból szemügyre a problémát.

A gyártó azt állítja, hogy

$$\mu \geq \mu_0$$

ahol mondjuk  $\mu_0 = 50\%$ . A konkurencia szerint ez nem igaz. Az ártatlanság vélelmének jogelvét követve azt mondjuk, hogy *a gyártó ártatlan mindaddig amíg bűnössége minden kétséget kizáróan bebizonyosodik*. Abban az esetben amikor nem lehet minden egyes terméket, vitás objektumot megvizsgálni, természetesen a bíróság is csak statisztikai módszerek alapján dönthet. A minta véletlen természetéből fakadóan a döntés sem lehet 100% bizonyosságú, ezért az ártatlanság vélelme azt kívánja, hogy jól kontrolálni tudjuk annak a valószínűségét, hogy bár a gyártó ártatlan, a szerencsétlen véletlen mintaválasztás azt eredményezi, hogy statisztikai vizsgálat bűnösnek találja a gyártót. Az előbb vázolt hipotézis vizsgálat keretei között ez azt jelenti, hogy az  $\alpha$  elsőfajú hiba valószínűségét tudjuk kontrolálni, a vitázó feleknek meg kell állapodnia ebben az  $\alpha$ -ban, ezt közlik a bírósággal (természetesen ez a minta nagyságát az eljárás költségét is befolyásolja). A null hipotézist ennek a logikának az értelmében úgy kell fölállítani, hogy az első fajú hiba a leírt legyen, azaz a null hipotézis a gyártó ártatlanságát, a hirdetésben közölt állítást kell, hogy tartalmazza, azaz

$$H_0 : \mu \geq \mu_0$$

a helyesen választott null hipotézis. Természetesen nem minden vizsgálandó probléma fordítható le a jog nyelvére, általában célszerű a nullhipotézis irányát úgy megválasztani, hogy az elkövethető hibák közül az legyen az elsőfajú aminek bekövetkezése a nagyobb kárt okozza, mivel ennek a valószínűségét tudjuk  $\alpha$ -val szabályozni, míg a másodfajúét nem.

Vegyük erre még egy példát. Fémlemezket vonunk be rozsdamentesítő festékkel. Ha túl sok festéket viszünk fel, akkor fölösleges költséget hozunk létre, ha viszont nem elég vastag a védőréteg, a lemezt hamar kikezdi a korrózió, tönkremegy a lemez, esetleg ennek következtében a lemez értékénél is nagyobb kár keletkezik. Ezért tehát a minőségellenőrzés során a  $\mu$  átlagos festékvastagságra vonatkozóan célszerű a

$$H_0 : \mu \leq \mu_0$$

nullhipotézist alkalmazni, hiszen akkor lesz igaz az hogy a két lehetséges hiba közül a nagyobb kárt okozó lesz az elsőfajú hiba. Konkrétan azt a valószínűséget kontrolálja  $\alpha$ , hogy a

festékréteg vékonyabb mint szükséges, de ezt a próba nem mutatta ki és ezért elvetjük a null hipotézist.

### 6.2.2 Paraméteres próbák

A hipotézis vizsgálat 5.6. lépése az úgynevezett próba. A körülményektől függően különböző próbák ismeretesek, azoknak pedig további variációi. Az alábbiakban a legfontosabakat ismertetjük. Először egy sokaság valamely paraméterére vonatkozó próbákat mutatjuk be, ezek az egymintás próbák, majd két sokaság valamilyen paraméterének az összehasonlítására vonatkozó kétmintás próbák következnek.

A kétoldali  $z$  próba (néha nevezik  $u$  próbának is) került a fejezet elején ismertetésre.

1. Ekkor a null hipotézis

$$H_0 : \mu = \mu_0$$

az alternatív hipotézis

$$H_1 : \mu \neq \mu_0.$$

2. Ha ismert a sokaság szórása  $\sigma$  vagy a minta elemszáma elég nagy ( $n > 30$ ), akkor próbafüggvény

$$z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

ami közel standard normális eloszlású.

3. A hétköznapi, üzleti életben általában a  $\alpha = 0.05$  megfelelő választás.
4. Az elutasítási tartomány  $(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$ , ahol  $z_{\alpha/2} = 1.96$
5. Legyen a példa kedvéért  $\bar{X} = 1.002\%$  liter a mintában a palackokba töltött folyadék átlagos mennyisége. Legyen  $\sigma = 0.01\%$  és  $n = 100$ . Ekkor

$$z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{0.002}{0.01} 10 = 2.$$

6. Viszont  $z_{\alpha/2} = 1.96 < z = 2$ , ezért a null hipotézist el kell vetnünk. Le kell állítani a gépsort és újraszabályozni a töltést.

### 6.2.3 Az egymintás próbák további esetei

Ebben a szakaszban a hipotézis vizsgálat variánsait mutatjuk be amelyet a feltételek változása hoz létre.

Abban az esetben, ha a populáció szórása ismeretlen azt a tapasztalati szórással helyettesítjük amikor a próbafüggvényt kiszámoljuk, azaz a mintaátlagot standardizáljuk:

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}.$$

Ilyen esetben a standardizált kifejezés normális eloszlást követ, ha a populáció maga is normális eloszlású volt vagy ha  $n > 30$ .

**Állítás 13** *Ellenkező esetben azaz ha a populáció nem normális, illetve ismeretlen szórást követ, létezik a várható érték és szórás ( de ismeretlenek) akkor a*

$$t = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

*próbafüggvény Student eloszlású és ennek megfelelően alakul a 3. lépés, illetve az 5. lépésben a kritikus értékeket  $t_{\alpha/2}$  illetve egyoldali hipotézis esetében  $t_{\alpha}$  adja.*

**Állítás 14** *Ha egy  $p$  valószínűség, illetve részarány becslése a feladat akkor mindig a*

$$z = \frac{\bar{p} - \mu}{\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}}$$

*próbafüggvény alkalmazható, ami standard normális eloszlású.*

#### 6.2.4 Kétmintás próbák

Gyakori és természetes kérdések merülnek fel, két populáció összehasonlítása során. Igaz-e, hogy a villányi szőlő cukorfoka magasabb mint a spononié? Az ilyen típusú kérdések vizsgálatát is két példán mutatjuk először be, majd tömören ismertetjük a technikai részleteket.

Legyen  $\mu_1$  a Villányi szőlő átlagos cukorfoka,  $\mu_2$  a Spononié. Vegyünk  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változókat, mintát a Villányon termelt mustokból és  $Y_1, Y_2, \dots, Y_m$  a Spononi.

*Az  $X_i$  és  $Y_j$  mintaelemekről feltesszük, hogy teljesen függetlenek.*

Tegyük fel, hogy  $n, m > 30$ , azaz a minták „nagyok”. Mi legyen a nullhipotézis? Min korábban hangsúlyoztuk a null hipotézis megválasztása függ attól, hogy hibás döntés esetén milyen kár keletkezik. Most egy olyan szituációt vázolunk, amiben a korábbival éppen ellentétes módon célszerű a null hipotézist választani. Ha egy hirdetés állítaná, hogy a villányi szőlő cukorfoka magasabb mint a spononié akkor a független döntést a  $H_0 : \mu_1 \geq \mu_2$  hipotézisből kiindulva kell hozni. Ha viszont az a helyzet, hogy a pincészetünknek régóta spononi mustot veszünk, akkor természetesen mindaddig kitartunk e mellett míg ennek az ellenkezője alaposan be nem igazolódik, azaz a statisztika nyelvén, ki nem derül, hogy a villányi cukorfoka **szignifikánsan** magasabb. Ekkor tehát a

$$H_0 : \mu_1 \leq \mu_2$$

null hipotézisből indulunk ki. Érdemes a null hipotézist kissé átalakítani,

$$H_0 : \mu_1 - \mu_2 \leq D_0 = 0.$$

Esetünkben az állítás annyiról szól, hogy az egyik átlag nagyobb mint a másik, ilyenkor  $(\mu_1 - \mu_2)_0 = D_0 = 0$ , de lehetne a kezdeti állítás pl az, hogy a villányi szőlő átlagos cukorfoka 2 fokkal magasabb mint a spononi, ekkor  $D_0 = 2$  lenne.

Világos, hogy ezt akkor fogjuk elvetni, ha

$$z = \frac{(\bar{X} - \bar{Y}) - D_0}{s_{\mu_1 - \mu_2}} > z_{\alpha} \quad (6.5)$$



ahol

$$s_{\mu_1 - \mu_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

A (6.5) egyenőtlenséget átírva

$$\mu_1 > \mu_2 + z_\alpha \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

összefüggést kapjuk, amiből azonnal világos, hogy mit értünk szignifikánsan nagyobb alatt.  $\mu_1$  nem egyszerűen nagyobb kell, hogy legyen mint  $\mu_2$  hanem  $\mu_2 + z_\alpha \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ -nél is azaz a véletlen szóródás  $z_\alpha$ -szorosánál is nagyobb az eltérés. A próba kialakításakor fontos, hogy  $\alpha$  a szignifikancia szintje,  $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$  pedig a mintaátlagok különbségének szórása. A próba további menete már megegyezik az egymintás próbáéval.

Abban az esetben, ha a hipotézis kétoldali, azaz a null hipotézis

$$H_0 : \mu_1 = \mu_2$$

illetve

$$H_0 : \mu_1 - \mu_2 = D_0 = 0$$

akkor is (6.5) a próbafüggvény, de a kritikus tartomány kétoldali, azaz  $-z_{\alpha/2}, z_{\alpha/2}$  a két kritikus érték.

### Kis mintás próba két átlag különbségére vonatkozóan

A fenti példa kis minták esetén a

$$t = \frac{(\bar{X} - \bar{Y}) - D_0}{s_{\mu_1 - \mu_2}} \quad (6.6)$$

Student eloszlású próbafüggvény segítségével oldható meg. A **szabadságfokot** a

$$df = \frac{\frac{s_1}{n_1} + \frac{s_2}{n_2}}{\frac{\left(\frac{s_1}{n_1}\right)^2}{n_1 - 1} + \frac{\left(\frac{s_2}{n_2}\right)^2}{n_2 - 1}}$$

képlet határozza meg.

### Kombinált szórásbecslés

Kis minták esetében javítható a szórásra vonatkozó becslés, ha egyéb körülmények alapján tudható, hogy a két populáció szórása ugyanaz, azaz  $\sigma_1 = \sigma_2$ . Ekkor  $s_1, s_2$  ugyanazt a közös  $\sigma$  értéket becsli, ezért célszerű a kevert mintaszórással számolni, ennek négyzete:

$$s_p^2 = \frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2}{n_1 + n_2 - 2}.$$

Ekkor próbafüggvényünk a

$$t = \frac{(\bar{X} - \bar{Y}) - D_0}{\sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

lesz. E két képletet, azaz a kevert szórás becslést érdemes nagy minta esetén is használni, ha  $\sigma_1 = \sigma_2$  feltehető.

### Részarányok összehasonlítása

Igaz-e hogy a diplomások körében az  $A$  párt népszerűbb mint az alacsonyabb végzettségűek körében? Erre a kétmintás részarányokra vonatkozó próba adhat választ. Legye a null hipotézis

$$H_0 : p_1 - p_2 \geq D_0 = 0$$

ahol  $p_1$  a diplomások között az  $A$ -t preferálók részaránya,  $p_2$  a másik csoportban. Nyilván egyoldali hipotézisről van szó, a próbafüggvény pedig

$$z = \frac{(p_1 - p_2) - D_0}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$$

Kétoldali

$$H_0 : p_1 - p_2 = 0$$

hipotézis esetén lehet a közös  $p$  becslését,  $\hat{p}$ -t használni, ezért a próbafüggvény

$$z = \frac{p_1 - p_2 - 0}{\sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

ahol

$$\hat{p} = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$$

### Próba párosított mintával

Az összes eddigi esetben feltettük, hogy a két minta elemei teljesen függetlenek. Bizonyos esetekben ez bár nem áll fenn, mégis igen jó teszt készíthető. Tegyük fel, hogy egy készülő felhasználói szoftverhez két kezelőfelület terv készült. El kívánjuk dönteni, hogy melyik tevének kedvezőbbek az ergonómiai tulajdonságai, melyiken hatékonyabb a munkavégzés. Ezért a következő „kísérletet” végezzük. 50 rutinos operátort teljesen azonos feltételek között megtanítunk minkét kezelőfelület használatára. Ezek után ugyanazt a feladatot elvégzik az egyik illetve a másik felületen, a szükséges időt mérjük. Legyen  $X_i$  az 1.  $Y_i$  a 2. kezelőfelületen mért idő  $i = 1, \dots, 50$ . Az esetleges fáradási tényezőt is kiküszöböljük, 25-en először az 1. kezelőfelületet használják utána a 2-at, a másik huszonöt operátor fordítva.

Számítsuk ki a mért idők  $d_i$  különbségeit

$$d_i = X_i - Y_i.$$

Ha nincs egyéb előfeltevésünk, akkor a

$$H_0 : \mu_1 = \mu_2$$

null hipotézist használjuk. Ez természetesen ekvivalens a

$$H_0 : \mu_1 - \mu_2 = 0$$

illetve ha  $\mu_d = E(X - Y) = \mu_1 - \mu_2$  jelöli a várható értékek különbségét, akkor ez szintén ekvivalens a

$$H_0 : \mu_d = 0$$

hipotézissel. Ezután a  $d_i$  mintadifferenciákra mint egymintás próba járhatunk el. A módszer kis és nagymintás variánsai, részarányra vonatkozó variánsa egyaránt egyszerűen értelmezhető.

### 6.3 Próbák a szórásra vonatkozóan

#### 6.3.1 Egymintás próba

Eddig kizárólag az átlag illetve a részarányra vonatkozó próbákról volt szó. De ahogy a szórásnégyzetre lehet konfidencia intervallumot szerkeszteni, természetesen hipotézis vizsgálat is végezhető. Legyen egy normális eloszlású sokaságunk, aminek szórása ismeretlen. Korábbi példánkban, a palacktöltősor esetén vizsgálni érdemes a szóródást. Nem kedvező, ha túl nagy (ha kicsi általában nem baj.) Ha egy kiháló félben lévő faj genetikus változatosságáról van szó, akkor a nagy szórás éppen kívánatos is lehet. Az előbbi esetnél maradva legyen a kívánt szórás  $\sigma_0$ , ekkor a null hipotézis, a gyártó nézőpontjából

$$H_0 : \sigma^2 \leq \sigma_0^2$$

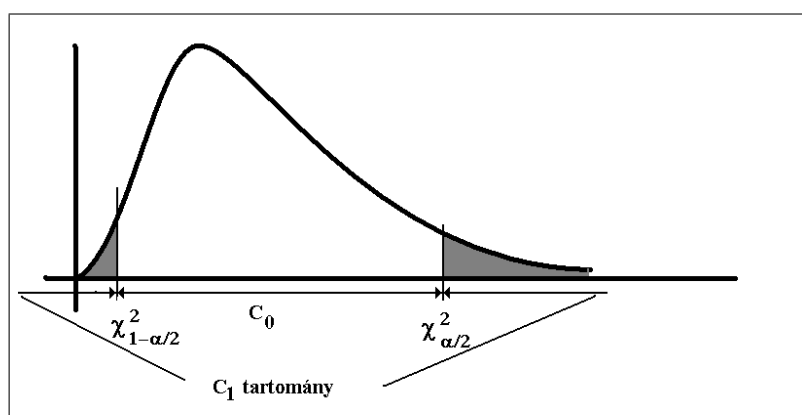
módon választandó. Az (11) Állításból tudjuk, hogy

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2$$

ahol  $\chi^2$  szabadságfokú  $n-1$ . Ezért  $\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$  próbafüggvényt alkalmazva elutasítjuk a null hipotézist, ha

$$\chi^2 > \chi_\alpha^2$$

ahol a  $\chi_\alpha^2$  az  $n-1$  szabadságfokú  $\chi^2$  eloszlás  $\alpha$ -hoz tartozó kritikus értéke.



### 6.3.2 Kétmintás próba

Az átlagokra vonatkozó kétmintás próbák között már felmerült a kérdés, egyenlő-e a két populáció szórása. Ha igen akkor a szórásra vonatkozóan igen jó becslés adható a kombinált szórással (Lásd 6.2.4 alfejezetet). Ehhez először ellenőriznünk kell, hogy a két szórás valóban egyenlő-e, azaz igaz-e a

$$H_0 : \sigma_1^2 = \sigma_2^2$$

A próba kivitelezése a következő állításon alapszik.

**Állítás 15** Ha két normális populációból származó  $n$  illetve  $m$  elemű teljesen független mintát veszünk, akkor

$$\frac{s_1^2}{s_2^2} \sim F^{(n-1, m-1)}$$

hányados  $(n-1, m-1)$  szabadságfokú  $F$  eloszlást követ.

Az  $F$  eloszlást nem definiáljuk egzakt módon, lényegében két khinégyzet eloszlású független valószínűségi változó hányadosának eloszlását írja le.

**Állítás 16**

$$F^{(n-1, m-1)}(1-x) = \frac{1}{F^{(m-1, n-1)}(x)}.$$

Az állítást nem bizonyítjuk, de intuitíve látható a definícióból

Ennek alapján a próbafüggvény

$$F = \frac{s_1^2}{s_2^2}$$

$F$  eloszlású ha igaz a null hipotézis. A kritikus érték  $F_{\alpha/2}$  illetve  $F_{1-\alpha/2}$ , de ha az általánosság megszorítása nélkül, úgy számoztuk meg a mintákat, hogy  $s_1 > s_2$ , akkor elég a felső kritikus tartománnyal foglalkozni. Abban az esetben, ha

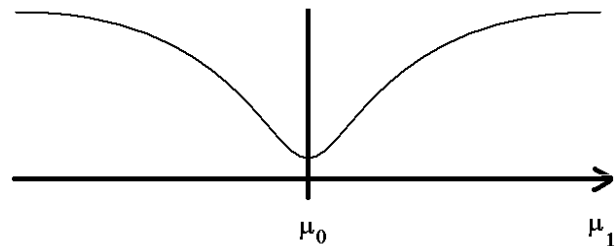
$$F > F_{\alpha/2}$$

elvetjük a null hipotézist.

### 6.3.3 A másodfajú hiba

Emlékezzünk vissza arra, hogy másodfajú hibát akkor vétünk a hipotézis vizsgálat során, ha a null hipotézis nem áll fenn, mégis a mellett döntünk. Ezzel az eseménnyel a valószínűsége (a véletlen mintavétel következtében)  $\beta$ . Az elsőfajú hibával ellentétben  $\beta$  általában nem számítható ki a valódi paraméter érték ismerete nélkül.

**Definíció 55** A próba erejének szokás nevezni  $1 - \beta$ , vagyis azt a valószínűséget, hogy helyesen ismeri fel, hogy  $H_0$  nem áll fenn. Szokás próba erőfüggvényét definiálni a



A próba erőfüggvénye a  $\mu$  változással

paraméter függvényében  $f(\vartheta) = 1 - \beta$ .

Például egy jobboldali, várható értékre vonatkozó próba erőfüggvényét mutatja a (??)Ábra. [Ide abrat]

### A másodfajú hiba kiszámítása

Képzeld el a következő szituációt. A hő ellenállók selejtes hadianyagot gyártanak a munkaszolgálat során. Természetesen ügyelni kell, hogy a szabotázs ne derüljön ki. A gyártott fegyver például űrmérete  $\mu_0 = 9.00$  kell, hogy legyen. A gyártás során ennél lényegesen kisebb nem is lehet, viszont a furat készítésekor kis trükkkel ennél nagyobbra lehet készíteni, ami selejtté teszi a fegyvert. A minőségellenőrzés a

$$H_0 : \mu \leq \mu_0$$

feltevést ellenőrizni. Mi a valószínűsége, hogy a szabotázs során gyártott  $\mu = 9.07$  átlaggal nem bukna le az ellenállók, azaz  $H_0$  nem igaz, de ezt a próbával nem mutatja ki. Tegyük fel, hogy a szórás  $\sigma = 0.02$ , a szignifikancia szintje  $\alpha = 0.05$  és a szokásos mintaméret  $n = 30$ ? Készítsük el a próbafüggvényt.

$$z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

Ekkor a próba alapján nem vetik el a null hipotézist, ha  $z < z_\alpha$ , azaz

$$\bar{x} < \mu_0 + z_\alpha \frac{\sigma}{\sqrt{n}} =: C. \quad (6.7)$$

Ha tudjuk, hogy  $\mu = 9.02$ , akkor

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

alapján adódik, hogy a másodfajú hiba  $\beta$  valószínűsége (azaz, hogy elkerülik a lebukást)

$$\mathbb{P}(\bar{x} < C) = \mathbb{P}\left(\frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < \frac{C - \mu}{\frac{\sigma}{\sqrt{n}}}\right) = \Phi\left(\frac{C - \mu}{\frac{\sigma}{\sqrt{n}}}\right).$$

Érdemes  $C$ -be behelyettesíteni annak értékét (6.7)-ből. Az átlag eredeti skáláján

$$\frac{C - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\mu_0 - \mu}{\frac{\sigma}{\sqrt{n}}} + z_\alpha$$

lesz a kritikus érték. Jól látható, hogy a kritikus tartomány nő, ha  $\mu$  nő. Másképpen a  $\frac{\sigma}{\sqrt{n}}$  szórású mintaátlag alapján jól el lehet különíteni  $\mu > \mu_0$  átlagot, ha az a szórás többszörösével haladja meg  $\mu_0$ -t, egyébként nem.

## Chapter 7

### NEM PARAMÉTERES PRÓBÁK

#### 7.0.4 Khinégyzet próbák

Ebben a fejezetben a **nem paraméteres statisztika** egy kis részét a **nem paraméteres próbákat** érintjük, ezen belül is először az úgynevezett  $\chi^2$  **próbákról** lesz szó. Közös vonásuk, mint nevük mutatja, hogy  $\chi^2$  eloszlásra vezető szellemes konstrukciókkal ragadnak meg igen összetett problémákat. A módszer mintapéldája a **multinomiális eloszlás** tesztelése.

**Példa 31** *Képzeljük el, hogy a kóla termékek piaci részesedését követjük figyelemmel. Korábban azt tapasztaltuk, hogy a  $p_1, p_2, \dots, p_k$  volt az  $1, 2, \dots, k$  termék piaci részesedése az előző időszakban. Minket első nekifutásra az érdekel, történt-e elmozdulás, vagy sem. Ezért null hipotézisként a*

$$H_0 : p'_i = p_i, i = 1, \dots, k$$

*feltevással élünk, ahol  $p'_i$  az  $i$ -edik termék ismeretlen jelenlegi piaci részaránya. Legyen egy független  $n$  elmeû mintánk, amiben azt találtuk, hogy  $x_i$  fogyasztó választotta az  $i$ -edik terméket. Világos, hogy ekkor  $x_i$  binomiális eloszlású valószínűségi változó ( $n, p_i$ ) paraméterekkel) feltéve, hogy a null hipotézis fennáll. De akkor az is igaz, hogy igen jó közelítéssel ( $x_i > 5$  feltevése már elegendő) igaz, hogy*

$$\frac{x_i - np_i}{\sqrt{np_i(1 - p_i)}} \sim N(0, 1).$$

*Tekintettel arra, hogy az osztályok száma általában elég nagy ( $1 - p_i \approx 1$  ezt el szokták hagyni. Jelölje  $e_i$  a hipotézis alapján a várt értéket,*

$$e_i = np_i$$

*ezzel a jelöléssel*

$$\chi^2 = \sum_{i=1}^k \frac{(x_i - e_i)^2}{e_i}$$

*próbafüggvény  $\chi^2$  eloszlást követ, ha  $H_0$  igaz. Ezért a hipotézis ellenőrzése innen a szokásos módon történhet.*

#### Homogenitás vizsgálat

Igen hasonló az úgynevezett **homogenitás vizsgálat**, amely két populáció rétegezethegének azonos voltát hivatott ellenőrizni. Adott  $A, B$  populációk és ezek azonos ismérv szerinti

felosztása. Az egyes osztályok részaránya legyen  $p_i, q_i$ . Példaként állhat itt mondjuk a következő  $A$  a városi lakosság,  $B$  a kistelepülési lakosság. az egyes osztályokat a szokásos legmagasabb iskolai végzettség szerint alakítjuk ki. A null hipotézis, hogy a legmagasabb végzettség szerint a városi és kistelepülési megoszlás azonos, azaz

$$H_0 : p_i = q_i, i = 1, \dots, k$$

Megint feltesszük, hogy  $1 - p_i \approx 1, 1 - q_i \approx 1$ . Legyen a két vizsgált minta elemszáma  $n, m$  a talált osztályok mérete pedig  $x_i, y_i$ . Tulajdon képpen az

$$\frac{x_i}{n} - \frac{y_i}{m} \approx 0$$

összefüggést ellenőrizzük. Ha igaz  $H_0$ , akkor

$$\frac{x_i}{n} \approx N\left(p_i, \sqrt{\frac{p_i}{n}}\right)$$

ezért  $\frac{x_i}{n} - \frac{y_i}{m}$  szórása jó közelítéssel

$$\sqrt{p_i \frac{m+n}{mn}},$$

a  $p_i$  közösbecslése pedig

$$\frac{x_i + y_i}{n + m}$$

ezért

$$z = \frac{\frac{x_i}{n} - \frac{y_i}{m}}{\sqrt{\frac{x_i + y_i}{n + m}}} \sim N(0, 1)$$

amiből

$$\chi^2 = \sum_{i=1}^k \frac{\left(\frac{x_i}{n} - \frac{y_i}{m}\right)^2}{\frac{x_i + y_i}{nm}}$$

$\chi^2$  eloszlású próbafüggvény kapható.

### Függetlenség vizsgálat

Igen gyakori kérdés, hogy két tulajdonság között van-e kapcsolat. Például a végzettség és az egyes napilapok kedveltsége között van-e összefüggés. Azaz egy sokaságot két módon is osztályba sorolunk, példánk szerint végzettség illetve a vásárolt napilap szerint. Tegyük fel, hogy  $m$  különböző választ adhatnak a megkérdezettek az újág kérdésre (nem olvas is ide tartozhat) és  $n$  féle végzettséget különböztetünk meg. Ekkor ha  $N$  volt a megkérdezettek száma készítsünk el egy úgynevezett kontingencia táblát, amelynek  $i$ -edik sorának  $j$ -edik eleme (azaz oszlopa) azon válaszadók  $k_{i,j}$  számát tartalmazza, akik az



$i$ -edik újságot olvassák és végzettségük alapján a  $j$ -edik osztályba tartoznak.

	1	2		$i$		$m$
1						
2						
$j$				$k_{i,j}$		
$n$						

A függetlenség azt jelentené, hogy ha  $A_i$  annak a valószínűsége, hogy egy véletlenül választott ember az  $i$ -edik újságot választja, illetve  $B_j$ , hogy a  $j$ -edik végzettségi kategóriába esik, akkor

$$\mathbb{P}(A_i B_j) = \mathbb{P}(A_i) \mathbb{P}(B_j).$$

Ugyanakkor a bal és jobb oldalon található valószínűségekre a kontingencia táblázat becseket tartalmaz. Világos, hogy ha

$$r_i = \sum_{j=1}^m k_{i,j}$$

akkor

$$\frac{r_i}{N}$$

becsli  $\mathbb{P}(A_i)$ -t, illetve ha

$$c_j = \sum_{i=1}^n k_{i,j}$$

akkor

$$\frac{c_j}{N}$$

becsli  $\mathbb{P}(B_j)$ , végül pedig

$$\frac{k_{i,j}}{N}$$

becsli  $\mathbb{P}(A_i B_j)$ -t. Ennek alapján a következő próba függvény konstruálható.

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^m \frac{(k_{i,j} - e_{i,j})^2}{e_{i,j}} \quad (7.1)$$

ahol

$$e_{i,j} = \frac{c_i o_j}{N}.$$

**Állítás 17** A (7.1)beli próbafüggvény jó közelítéssel  $\chi^2$  eloszlású  $(n-1)(m-1)$  szabadságfokkal, ha az egyes cellákban kapott elvárt  $e_{i,j}$  értékek mindegyike nagyobb mint 5.

Az állítást nem bizonyítjuk, helyessége a korábbi gondolatmenetből intuitíven látható.

**Megjegyzés 10** *A leírt függetlenség vizsgálat valószínűségi változók függetlenségének eldöntésére is alkalmazható oly módon, hogy elkészítjük a kétdimenziós tapasztalati együttes eloszlás alapján megfelelő intervallumok választásával a kontingencia táblázatot.*

**Gyakorlat 32** *Készítsük el a megjegyzés szerinti függetlenségvizsgálat teljes menetét.*

### 7.0.5 Illeszkedés, normalitás vizsgálat

A multinomiális eloszlásnál látott módon tetszőleges eloszláshoz való illeszkedést is lehet tesztelni.

Legyen  $\Psi$  egy tetszőleges eloszlás függvény. Tegyük fel, hogy adott  $x_i, x_{i+1}$   $i = 1..k$  értékekhez tartoznak a

$$p_i = \Psi(x_{i+1}) - \Psi(x_i)$$

értékek. Ekkor az  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változók akkor származnak a  $\Psi$  eloszlásból ( $\alpha$  szinten), ha a

$$\chi^2 = \sum_{i=1}^k \frac{(k_i - Np_i)^2}{Np_i}$$

$\chi^2$  eloszlású próbafüggvény az  $1 - \alpha/2, \alpha/2$  kritikus értékek köz; esik.

Az eljárással ellenőrizhető az a feltevés, hogy  $X$  normális eloszlású-e. Pontosabban célszerű a  $z = \frac{X - \bar{X}}{s}$  standardizáltat vizsgálni, hogy az illeszkedik-e a standard normális eloszláshoz.

**Gyakorlat 33** *Készítsük el a standard normális eloszlás esetén azt az  $x_i$  sorozatot, amelyre  $p_i = 0.15$ .*

**Megjegyzés 11** *Az eljárás nem jól alkalmazható olyan paraméteres eloszlások esetén, ahol nem áll rendelkezésre a standardizáláshoz hasonló az ismeretlen paramétert egyszerű transzformációval kiküszöbölő módszer. Ilyenkor ugyanis bár lehet, hogy az eloszlás  $\Psi_a - tköveti$ , de ha azt egy  $b \neq a$  paraméterű  $\Psi_b$ -hez illesztjük, akkor a próbafüggvény esetleg igen nagy lehet. Erre a legegyszerűbb példa, ha a  $\mu$  várható értékű normális eloszlást egy  $\mu'$ -vel centráljuk természetesen nagyon rossz illeszkedést kapunk (lásd (??)Ábra).*

### 7.0.6 Próbák helyzeti paraméterek vizsgálatára

Az előző (11) Megjegyzés is mutatja, hogy a helyzeti paraméterek próbái milyen fontosak. Az alábbiakban olyan tesztek ismertetünk amelyek többek között a helyzeti paraméterekről szolgálnak információval.

## Előjel próba

Ha egy eloszlásnak nem ismert az  $m$  mediánja, azaz az a  $m$  érték, melyre  $\mathbb{P}(X < m) = 1/2$ , az alábbi állítás segítségével lehet a

$$H_0 : m = m_0 \tag{7.2}$$

hipotézist ellenőrizni.

**Állítás 18** Legyen  $X_1, X_2, \dots, X_n$  független azonos eloszlású valószínűségi változók

$$Y_i = \begin{cases} 1 & \text{ha } X_i > m \\ 0 & \text{egyébként} \end{cases} .$$

$$Y = \sum_{i=1}^n Y_i$$

Ekkor, ha  $n \geq 20$

$$Y \sim N\left(\frac{1}{2}n, \sqrt{\frac{1}{4}n}\right)$$

**Következmény 6**

$$z = \frac{Y - \frac{1}{2}n}{\sqrt{\frac{1}{4}n}}$$

próbafüggvény standard normális eloszlású, alkalmas (7.2) ellenőrzésére.

Minkét állítás evidens a binomiális eloszlásra vonatkozó centrális határeloszlás tételből.

**Wilcoxon féle előjeles rang test**

A most ismertetésre kerülő módszer párosított (azaz nem független minta esetén) az  $X$  ill.  $Y$  változók  $F, G$  eloszlásának azonosságát ellenőrzni. Legyen  $X_1, X_2, \dots, X_n$  illetve  $Y_1, Y_2, \dots, Y_n$  független azonos eloszlású valószínűségi változók ( $X$  illetve  $Y$  példányai).

$X$	$Y$	$d$	$ d $	rang $ d $	előjeles rang
$x_1$	$y_1$	$d_1$	$ d_1 $	$r_1$	$\text{sign}(d_1) r_1$
$x_2$	$y_2$	$d_2$	$ d_2 $	$r_2$	$\text{sign}(d_2) r_2$
$x_n$	$y_n$	$d_n$	$ d_n $	$r_n$	$\text{sign}(d_n) r_n$
					$T = \sum_{i=1}^n \text{sign}(d_i) r_i$

**Állítás 19** Ha  $F = G$  akkor a fenti  $T$  próbafüggvény jó közelítéssel normális eloszlású, és

$$\mathbb{E}(T) = 0$$

$$\sigma_T = \sqrt{\frac{n(n+1)(2n+1)}{6}}$$

$$z = \frac{T}{\sigma_T}$$

alkalmas próbafüggvény.

Az állítás a centrális határeloszlás tétel egy élesítéséből következik, nem bizonyítjuk.

### 7.0.7 Man-Whitney próba

Az alábbi próba újra a két eloszlás azonossága, az

$$H_0 : F = G$$

feltevés ellenőrzésére szolgál.

Legyen mint előbb  $X_1, X_2, \dots, X_n$  független  $F$  eloszlású valószínűségi változók illetve  $Y_1, Y_2, \dots, Y_m$  független  $G$  eloszlású valószínűségi változók. Keverjük össze a két mintát és rendezzük nagyság szerint. Legyenek az  $X_i$  elemek rangjai  $r_{\pi(i)}$  ezek összege pedig  $T = sr_X$ .

**Állítás 20** *Ha  $F = G$ , akkor a fenti  $T$  statisztikák igazak az alábbiak.*

$$E(T) = \frac{1}{2}n(n+m+1)$$

$$\sigma_T = \sqrt{\frac{1}{2}nm(n+m+1)}$$

és

$$z = \frac{T - \frac{1}{2}n(n+m+1)}{\sigma_T}$$

standard normális eloszlású valószínűségi változó.

### A Kruskal-Wallis teszt

A következő teszt az összetett

$$H_0 : F_1 = F_2 = \dots = F_k$$

hipotézis ellenőrzésére szolgál, azaz  $k$  minta aláapán  $k$  eloszlás azonosságát teszteli. Legyen  $X_{1,j}, X_{2,j}, \dots, X_{n,j}$  független azonos eloszlású valószínűségi változók  $j = 1, 2, \dots, k$ -ra. Hasonlóan mint a Man-Whitney próba esetén keverjük össze a mintákat minden elem kapja meg a neki megfelelő rangot. Jelölje  $R_j$  a  $j$ -edik minta elemeinek rangösszegét,  $n_T$  pedig az összes elemek számát.

**Állítás 21** *Ha igaz a*

$$H_0 : F_1 = F_2 = \dots = F_k$$

*feltevés és minden  $j$ -re  $n_j \geq 5$ , akkor*

$$W = \frac{12}{n_T(n_T+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n_T+1)$$

*próbafüggvény  $k-1$  szabadságfokú  $\chi^2$  eloszlást követ.*

**A Spreman féle rangkorreláció**

Párosított minták esetén gyakori kérdés, hogy mi a kapcsolat, mi a korreláció a két változó között. Időnként célszerű ezt a kérdést is nem paraméteres módon megközelíteni. Erre szolgál a következő fogalom.

**Definíció 56** *A Spreman féle rangkorreláció. Legyenek a két párosított minta rangjai  $r_i, s_i$ . Legye  $d_i = r_i - s_i$ . Ekkor*

$$r_S = 1 - \frac{\sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

*a tapasztalati rangkorreláció,*

$$\rho_S = \mathbb{E}(r_S)$$

*pedig a Spreman féle rangkorreláció.*

**Állítás 22** *Ha*

$$H_0 : \rho_S = 0$$

*akkor persze*

$$\begin{aligned} \mathbb{E}(r_S) &= 0 \\ \sigma(r_S) &= \sqrt{\frac{1}{n-1}} \end{aligned}$$

*és*

$$z = \frac{r_S - 0}{\sigma(r_S)}$$

*közel standard normális eloszlás eloszlású, ha  $n \geq 30$ .*



## Chapter 8

### SZÓRÁSANALÍZIS

A szórás analízis (analysis of variance ANOVA) a paraméteres próbák egy érdekes családja, amely egy közös modellre épít. A legegyszerűbb az egyváltozós, egy faktoros eset. A módszer megismerését egy ilyen példával kezdjük.

**Példa 34** Egy étteremlánc ugyanazt a hamburgerét sok étteremben kínálja. A hamburger népszerűsítésére akciót kíván indítani. Először azt vizsgálja országonként 30-30 étteremben, hogy mennyit nőtt a nyeresége ezen a terméken. Azaz rendelkezésre állnak az  $X_{i,j}$  eredmények, ahol  $i$  az ország felsorolás,  $j = 1..30$  pedig az éttermek sorszáma. Ez azt jelenti, hogy  $X_{i,j}$   $j = 1..30$  egy független 30 elemű minta egy  $Y_i$  változóból. Első feltevésük, hogy az akció eredménye független az országtól, azaz

$$\mu_i = E(Y_i) = \mu$$

minden  $i = 1..k$ -ra, azaz

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

összetett hipotézist kell ellenőrizni.

A statisztika hagyományos kifejezésével, szokás az  $i$ -edik kezeléstről beszélni, mert a módszert a növénytermesztésben, az egyes földterületek eltérő kezelésének összevetésére alkalmazták először. Az azonos kezelésnek alávetett egyedek alkotnak egy csoportot. Szokás szerint olyan statisztikát keresünk, ami a  $H_0$  feltevés mellett jól viselkedik. Az egye mintaelemekre a modell szerinti feltevés a következő:

$$X_i(j) = \mu_i + \varepsilon_i(j)$$

ahol

$$\varepsilon_i(j) \sim N(0, \sigma)$$

független valószínűségi változók ismeretlen közös  $\sigma$ -el.

Szokás a modellt a

$$X_i(j) = \mu + \alpha_i + \varepsilon_i(j)$$

alakban megfogalmazni, ahol  $\alpha_i$  az  $i$ -edik kezelés egyedi hatása. Ha  $H_0$  igaz, akkor a sokaságban tapasztalt szórásnégyzetre a teljes minta alapján becslés is adható. Ugyanakkor az egyes csoportokon belül is adható variancia-becslés, ami a feltevés szerint független becsléseket ad, ezek átlaga szintén becsli a teljes sokaság szórásnégyzetét. Így kétféle becslést lehet készíteni a szórásnégyzetre. Ehhez először a tapasztalati átlagokra vezessünk be jelölést.

Legyen az  $n_i$  elemű  $i$ -edik csoport tapasztalati átlaga

$$\bar{X}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_i(j),$$

a teljes átlag

$$\bar{X} = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^{n_i} X_i(j)$$

ahol  $n = \sum_{i=1}^p n_i$  az összes mintaelemek száma. Bevezetjük a

$$SST = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^{n_i} (X_i(j) - \bar{X})^2$$

teljes négyzetösszeget valamint a

$$SSE = \sum_{i=1}^p \sum_{j=1}^{n_i} (X_i(j) - \bar{X}_i)^2$$

csoportokon belüli négyzetösszegek összegét, végül pedig a

$$SSTR = \sum_{i=1}^p n_i (\bar{X}_i - \bar{X})^2$$

a csoportosításból fakadó négyzetösszeget.

## Lemma 2

$$SST = SSE + SSTR$$

**Bizonyítás.** Az állítás egyszerű aritmetikai átalakítással igazolható. A kidolgozást az olvasóra bizzuk. ■

A próba megalapozását a következő már nehezebb állítás biztosítja.

## Állítás 23 Jelölje

$$MSR = \frac{1}{p-1} SSTR$$

a csoportosításból fakadó átlagos tapasztalati négyzetösszeget. Ekkor

$$E(MSR) = \sigma^2 + \frac{1}{p-1} \sum_{i=1}^p n_i (\mu_i - \mu)^2.$$

**Bizonyítás.** Kezdjük a definíció szerinti zárójel bővítésével és felbontásával.

$$\begin{aligned} (\bar{x}_i - \bar{\bar{x}})^2 &= (\bar{x}_i - \mu_i + \mu_i - \mu + \mu - \bar{\bar{x}})^2 \\ &= (\bar{x}_i - \mu_i)^2 + (\mu_i - \mu)^2 + (\mu - \bar{\bar{x}})^2 \\ &\quad + 2(\bar{x}_i - \mu_i)(\mu_i - \mu) + 2(\bar{x}_i - \mu_i)(\mu - \bar{\bar{x}}) + 2(\mu_i - \mu)(\mu - \bar{\bar{x}}) \end{aligned}$$



Vegyük észre, hogy a cetrálás miatt, illetve mert a szorzat másik tagja konstans

$$E [2 (\bar{x}_i - \mu_i) (\mu_i - \mu)] = E [2 (\mu_i - \mu) (\mu - \bar{x})] = 0.$$

Az alábbi két tagban mintaátlagok szórásai szerepelnek, ezért

$$\sum_{i=1}^r n_i E [(\bar{x}_i - \mu_i)^2 + (\mu - \bar{x})^2] = \sum_{i=1}^r n_i \left( \frac{\sigma^2}{n_i} + \frac{\sigma^2}{n} \right) = r + 1. \quad (8.1)$$

A harmadik négyzetes tag a jobboldalon kívánt összege a  $\mu_i - k$  és a  $\mu$  négyzetes eltérésösszegének:

$$\frac{1}{r-1} \sum_{i=1}^r n_i (\mu_i - \mu)^2 \quad (8.2)$$

Foglalkozzunk a megmaradt egyetlen keresztszorzattal.

$$\begin{aligned} & E \left[ \sum_{i=1}^r 2n_i (\bar{x}_i - \mu_i) (\mu - \bar{x}) \right] \\ &= E \left[ \sum_{i=1}^r \frac{2n_i}{n} (\bar{x}_i - \mu_i) (n\mu - n\bar{x}) \right] \\ &= E \left[ \sum_{i=1}^r \frac{2n_i}{n} (\bar{x}_i - \mu_i) \left( n\mu - n_i\mu_i + n_i\mu_i - \sum_{j=1}^{n_i} x_{i,j} + \sum_{j=1}^{n_i} x_{i,j} - n\bar{x} \right) \right] \\ &= E \left[ \sum_{i=1}^r \frac{2n_i}{n} (\bar{x}_i - \mu_i) \left( \left( n_i\mu_i - \sum_{j=1}^{n_i} x_{i,j} \right) + \left( n\mu - n_i\mu_i - \sum_{k \neq i} \sum_{j=1}^{n_i} x_{k,j} \right) \right) \right] \\ &= E \left[ \sum_{i=1}^r \frac{2n_i^2}{n} (\bar{x}_i - \mu_i) (\mu_i - \bar{x}_i) + \sum_{i=1}^r \frac{n_i}{n} (\bar{x}_i - \mu_i) \left( n\mu - n_i\mu_i - \sum_{k \neq i} \sum_{j=1}^{n_i} x_{k,j} \right) \right] \\ &= E \left[ -2 \sum_{i=1}^r \frac{n_i^2}{n} (\bar{x}_i - \mu_i)^2 + 2 \sum_{i=1}^r \frac{n_i}{n} (\bar{x}_i - \mu_i) \left( n\mu - n_i\mu_i - \sum_{k \neq i} \sum_{j=1}^{n_i} x_{k,j} \right) \right] \end{aligned}$$

Vegyük észre, hogy a második tagban a szorzat két tényezője független, hiszen az utóbbiban pont az  $i$ -edik csoport elemei nem kerülnek összegzésre. A várható érték így a várható értékek szorzata és mindkettő centrált, mindkettő nulla. Az első tag pedig megint a mintaátlagok szórása, így

$$E \left[ -2 \sum_{i=1}^r \frac{n_i^2}{n} (\bar{x}_i - \mu_i)^2 \right] = -2 \sum_{i=1}^r \frac{n_i^2 \sigma^2}{n n_i} = -2\sigma^2. \quad (8.3)$$

Összevetve (8.1), (8.2) és (8.3)-t kapjuk az állítást. ■

Ez tehát azt jelent, hogy  $MSR$  pontosan akkor torzítatlan becslése  $\sigma^2$ -nek ha az összes várható érték egyenlő, azaz a null hipotézis teljesül. Ellenkező esetben  $\sigma^2$ -től felfelé tér el.

**Állítás 24** *A model feltevései mellett igaz, hogy*

$$E\left(\frac{1}{n-p}SSE\right) = \sigma^2$$

*továbbá*

$$E\left(\frac{1}{n-1}SST\right) = \sigma^2$$

*és  $MSE = \frac{1}{n-p}SSE$  illetve  $MSRT = \frac{1}{p-1}SSTR$  függetlenek. Igaz továbbá, hogy*

$$\frac{SSE}{\sigma^2}$$

*$n - p$  szabadságfokú  $\chi^2$  eloszlású valószínűségi változó továbbá, ha igaz a null hipotézis akkor*

$$\frac{SSTR}{\sigma^2}$$

*$p - 1$  szabadságfokú  $\chi^2$  valószínűségi változó.*

*Az állítások evidensek, kivéve a függetlenséget, ezt nem bizonyítjuk. Mindezek alapján a következő próba végezhető.*

**Tétel 35** *Ha a model feltevései fennállnak és igaz a null hipotézis, akkor a*

$$F = \frac{MSTR}{MSE}$$

*próbafüggvény  $(p - 1, n - p)$  szabadságfokú  $F$  eloszlást követ. Azaz  $F$  értéke közel van 1 - hez. Ellenkező esetben  $F$  értéke nagyobb.*

### 8.0.8 Kétrészes osztályozás

## Chapter 9

### LINEÁRIS REGRESSZIÓ

A modell:

$$Y = \alpha X + \beta + \varepsilon$$

ahol  $\varepsilon \sim N(0, \sigma)$  független az  $X$ -től. Keressük az

$$F(a, b) = \sum_{i=1}^n (y_i - ax_i - b)^2$$

minimumát  $a, b$ -ben. Keressük a  $\frac{\partial}{\partial b} F(a, b) = 0, \frac{\partial}{\partial a} F(a, b) = 0$  megoldásokat.

$$\frac{\partial}{\partial b} F(a, b) = -2 \sum_{i=1}^n (y_i - ax_i) - nb$$

amiből

$$\begin{aligned} nb &= \sum_{i=1}^n (y_i - ax_i) \\ b &= \bar{y} - a\bar{x}. \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial a} F(a, b) &= 2 \sum_{i=1}^n (y_i - ax_i - b) x_i \\ &= 2 \sum_{i=1}^n (y_i - ax_i - \bar{y} - a\bar{x}) x_i \\ &= 2 \sum_{i=1}^n [y_i x_i - ax_i^2] - n\bar{x}\bar{y} - an(\bar{x})^2. \end{aligned}$$

Ebből az  $s_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$  jelöléssel

$$\begin{aligned} ans_x^2 &= a \left( \sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right) \\ &= \sum_{i=1}^n y_i x_i - n\bar{x}\bar{y} \end{aligned}$$

azaz

$$a = \frac{\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}}{n s_x^2},$$

vagy másképpen bevezetve az  $s_{x,y} = \frac{1}{n} \sum_{i=1}^n y_i x_i$  jelölést

$$\begin{aligned} a &= \frac{\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}}{n s_x^2} = \frac{n \left( \frac{1}{n} \sum_{i=1}^n y_i x_i - \bar{x}\bar{y} \right)}{n s_x^2} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n y_i x_i - \bar{x}\bar{y}}{s_x^2}, \end{aligned}$$

azt kapjuk, hogy

$$a = \frac{s_{x,y} - \bar{x}\bar{y}}{s_x^2}.$$

Összefoglalva:

$$\begin{aligned} a &= \frac{\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}}{n s_x^2}, \\ b &= \bar{y} - a\bar{x}. \end{aligned}$$

Ebből, ha a tapasztalati covarianciát illetve korrelációs együtthatót  $Cov_{x,y}$  illetve  $r_{x,y}$  jelöli, azt nyerjük, hogy

$$a = \frac{Cov_{x,y}}{s_x^2} = \frac{r_{x,y} s_x s_y}{s_x^2} = r_{x,y} \frac{s_y}{s_x}.$$

**Definíció 57** Jelölje  $\hat{y}_i = ax_i + b$

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

*a lineáris illeszkedés négyzetes hibájának átlagát, vagy másképpen az egyenes körüli )korrigált tapasztalati szórásnégyzetet.*

**Definíció 58** Vezessük be a következő rövidítéseket

$$\begin{aligned} SX &= \sum_{i=1}^n x_i & SX^2 &= \sum_{i=1}^n x_i^2 \\ SY &= \sum_{i=1}^n y_i & SY^2 &= \sum_{i=1}^n y_i^2 \\ SXY &= \sum_{i=1}^n x_i y_i \end{aligned}$$

**Megjegyzés 12** Természetesen

$$\begin{aligned} \bar{x} &= \frac{1}{n} SX, \\ s_x^2 &= \frac{1}{n} SX^2 - (\bar{x})^2 \end{aligned}$$

**Tétel 36** Ha  $\varepsilon_i \sim N(0, \sigma)$  és korrelálatlanok, akkor

$$\begin{aligned} E(a) &= \alpha, \\ E(b) &= \beta, \\ \sigma^2(a) &= \frac{s^2}{s_x}, \\ \sigma^2(b) &= s^2 \frac{SX^2}{ns_x^2}. \end{aligned}$$

Az állítást nem bizonyítjuk.

**Definíció 59**

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2, \\ SSR &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \\ SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2. \end{aligned}$$

$SST$  az  $y$  ingadozását méri,  $SSE$  az egyenes és a méréspontok közötti hiba négyzetes összegét.

**Megjegyzés 13**

$$\begin{aligned} SSR &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (ax_i + b - a\bar{x} - b)^2 \\ &= \sum_{i=1}^n (ax_i - a\bar{x})^2 \\ &= a^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= a^2 ns_x^2 = \\ &= a^2 \left[ SX^2 - \frac{1}{n} (SX)^2 \right]. \end{aligned}$$

Ha figyelembe vesszük, hogy

$$a = \frac{\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}}{ns_x^2} = \frac{SXY - \frac{1}{n} SXSY}{ns_x^2}$$

akkor egyrészt azt kapjuk, hogy

$$SSR = a^2 n s_x^2 = \frac{[SXY - \frac{1}{n} SXSY]^2}{n s_x^2} \quad (9.1)$$

azaz

$$SSR = \frac{[SXY - \frac{1}{n} SXSY]^2}{SX^2 - \frac{1}{n} (SX)^2}.$$

Ugyanakkor a korrelációs együtthatóval kifejezve:

$$\begin{aligned} SSR &= a^2 n s_x^2 \\ &= \left[ r(x, y) \frac{s_y}{s_x} \right]^2 n s_x^2 \\ &= r(x, y)^2 n s_y s_x = n Cov(x, y). \end{aligned}$$

Tehát az SSR a lineáris kapcsolat erősségét fejezi ki.

### Lemma 3

$$\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) = \sum_{i=1}^n y_i x_i - n \bar{x} \bar{y}$$

### Bizonyítás.

$$\begin{aligned} &\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) \\ &= \sum_{i=1}^n [y_i x_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}] \\ &= \sum_{i=1}^n y_i x_i - \sum_{i=1}^n x_i \bar{y} - \bar{x} \sum_{i=1}^n y_i + n \bar{x} \bar{y} \\ &= \sum_{i=1}^n y_i x_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y} \\ &= \sum_{i=1}^n y_i x_i - n \bar{x} \bar{y}. \end{aligned}$$

■

### Tétel 37

$$SST = SSE + SSR.$$

**Bizonyítás.**

$$\begin{aligned}
SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - ax_i - (\bar{y} - a\bar{x}))^2 \\
&= \sum_{i=1}^n [(y_i - \bar{y}) - a(x_i - \bar{x})]^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 - 2a \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + a^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\
&= SST - 2a \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + a^2 ns_x^2.
\end{aligned}$$

Megint helyettesítsünk  $a$ -ba.

$$\begin{aligned}
a^2 ns_x^2 &= \left[ \frac{\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}}{ns_x^2} \right]^2 ns_x^2 \\
&= \frac{[\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}]^2}{ns_x^2}
\end{aligned}$$

amiből  $\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) = \sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}$  miatt

$$\begin{aligned}
SSE &= SST - 2 \frac{[\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}]^2}{ns_x^2} \\
&\quad + \frac{[\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}]^2}{ns_x^2} \\
&= SST - \frac{[\sum_{i=1}^n y_i x_i - n\bar{x}\bar{y}]^2}{ns_x^2} \\
&= SST - a^2 ns_x^2.
\end{aligned}$$

Ugyanakkor (9.1)-ben láttuk, hogy

$$SSR = a^2 ns_x^2,$$

azaz

$$SSE = SST - SSR.$$

■

**Definíció 60**

$$r^2 = \frac{SSR}{SST} = \frac{SXY - \frac{1}{n}SXSY}{SY^2 - \frac{1}{n}(SY)^2}.$$

**Megjegyzés 14** Tudjuk, hogy

$$\begin{aligned}
SSR &= ans_x^2, \\
SST &= ns_y^2
\end{aligned}$$

és

$$a = r_{x,y} \frac{s_y}{s_x}$$

$$a^2 = r_{x,y}^2 \left( \frac{s_y}{s_x} \right)^2,$$

ezért

$$r_{x,y}^2 = \left( \frac{s_x}{s_y} \right)^2 a^2 = \frac{ns_x^2}{ns_y^2} \frac{SSR}{ns_x^2}$$

$$= \frac{SSR}{ns_y^2} = \frac{SSR}{SST} = r^2.$$

Az  $r_{x,y}$  előjelét pedig az egyenes  $a$  meredekségének előjele határozza meg, ezért

$$r_{x,y} = \text{sign}(a) \sqrt{r^2}.$$

**Megjegyzés 15** Vegyük észre, hogy

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST},$$

$$|r| = \sqrt{1 - \frac{SSE}{SST}}.$$



## Chapter 10 FŐKOMPONENS ANALÍZIS

### 10.1 A lineáris algebra néhány eleme

**Definíció 61** Ha  $a, b \in \mathbb{R}^p$  akkor skaláris szorzatuk

$$a^T b = (a, b) = \sum_{i=1}^p a_i b_i.$$

**Definíció 62** Ha  $a \in \mathbb{R}^n, b \in \mathbb{R}^p$  akkor diadikus szorzatuk

$$ab^T = (a_i b_j)_{n \times p}$$

$n \times p$ -s mátrix.

**Definíció 63** Ha  $a \in \mathbb{R}^p$  akkor az  $a$  vektor  $l_2$  normája, avagy hossza:

$$\|a\|_2 = \|a\| = a^T a = (a, a) = \sum_{i=1}^p a_i^2.$$

**Definíció 64** Egy  $u$  vektort normálnak nevezünk, ha  $\|u\| = 1$ , azaz hossza 1.

**Definíció 65** Ha  $A$   $p \times p$ -s mátrix, akkor azt mondjuk, hogy az  $u \in \mathbb{R}^p$  vektor az  $A$  mátrix sajátvektora a  $\lambda$  sajátértékkel, ha

$$Au = \lambda u.$$

**Állítás 25** Ha  $A$  szimmetrikus mátrix akkor minden sajátértéke valós, ezek száma  $p$ . Legyenek  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  a sajátértékek  $\phi_i$  a hozzájuk tartozó (normált) sajátvektorok.

$$\lambda_1 = \max_{\|u\|=1} (u, Au),$$

$$\lambda_i = \max_{\|u\|=1, (u, \phi_j)=0: j=1 \dots i-1} (u, Au). \quad (10.1)$$

Azaz az  $i$ -edik sajátvektor az első  $i - 1$  feszítette altérre merőleges.

**Állítás 26** Legyen  $U = (\phi_1, \dots, \phi_p)$  a sajátvektorokból alkotott mátrix. Ekkor

$$A = U^T \Lambda U$$

ahol

$$\Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ & & \ddots \\ 0 & & & 0 & \lambda_p \end{pmatrix}$$

diagonális mátrix. Ez az  $A$  mátrix spektrál felbontása.

**Definíció 66** Egy  $U$  mátrix ortonormált, ha

$$U^T U = I,$$

ami egyben azt is jelenti, hogy

$$U^{-1} = U^T.$$

**Definíció 67** Egy  $A$  mátrix nyoma (trace-e)

$$\text{tr}(A) = \sum_{i=1}^p a_{i,i}$$

a diagonális elemeinek összege.

**Lemma 4** Tetszőleges  $A, B$   $k \times l$ -s mátrixokra (ahol  $k, l \geq 1$ )

$$\text{tr}(AB) = \text{tr}(BA)$$

**Bizonyítás.**

$$\begin{aligned} \text{tr}(AB) &= \sum_{i=1}^p (AB)_{i,i} = \sum_{i=1}^p \left( \sum_{k=1}^p a_{i,k} b_{k,i} \right) \\ &= \sum_{i=1}^p \left( \sum_{k=1}^p a_{i,k} b_{k,i} \right) = \sum_{i=1}^p \left( \sum_{k=1}^p b_{k,i} a_{i,k} \right) \\ &= \sum_{k=1}^p \left( \sum_{i=1}^p b_{k,i} a_{i,k} \right) = \sum_{k=1}^p (BA)_{k,k} = \text{tr}(BA). \end{aligned}$$

■

**Lemma 5**

$$u^T A v = \text{tr}(A v u^T).$$

**Bizonyítás.** Mivel  $u^T A v \in \mathbb{R}$  skalár, ezért  $u^T A v = \text{tr}(u^T A v)$ . Viszont a 4 Lemma miatt  $\text{tr}(u^T A v) = \text{tr}(A v u^T)$ . ■

**Lemma 6** Ha  $U$  ortonormált akkor

$$\text{tr}(U^T A U) = \text{tr}(A).$$

**Bizonyítás.** Alkalmazzuk a 4 Lemmát  $A = U^T, B = AU$  szereposztással.

$$\operatorname{tr}(U^T AU) = \operatorname{tr}(AUU^T).$$

Viszont  $U^T U = I$ , amiből következik az állítás. ■

**Következmény 7** Ha  $A$  szimmetrikus mátrix, akkor

$$\operatorname{tr}(A) = \sum_{i=1}^p \lambda_i.$$

**Bizonyítás.** A feltétel miatt alkalmazható a 26 Állítás, azaz

$$A = U^T \Lambda U$$

és alkalmazva a 6 Lemmát  $V = U^T$ -re

$$\operatorname{tr}(V^T AV) = \operatorname{tr}(UU^T \Lambda UU^T) = \operatorname{tr}(\Lambda)$$

amiből következik az állítás. ■

## 10.2 Véletlen vektorok elforgatása

**Definíció 68** Legyen  $X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} \in \mathbb{R}^p$   $p$ -dimenziós valószínűségi vektorváltozó. Ennek várható értéke  $\mu \in \mathbb{R}^p$ :

$$\mu = E(X) = \begin{pmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{pmatrix}.$$

Az  $X$  kovariancia mátrixa:

$$\Sigma = \operatorname{Cov}(X) = (\operatorname{Cov}(X_i, X_j))_{p \times p}.$$

Vegyük észre, hogy  $\Sigma$  főátlójában  $\operatorname{Cov}(X_i, X_j) = \sigma^2(X_i)$ -k állnak.

**Lemma 7**  $\Sigma$  szimmetrikus, pozitív szemidefinit mátrix.

**Bizonyítás.** A szimmetria nyilvánvaló. Ugyanakkor  $\Sigma = E((X - \mu)(X - \mu)^T)$  ezért minden  $a \in \mathbb{R}^p$ -re

$$a^T \Sigma a = E(a^T (X - \mu)(X - \mu)^T a) = E(Y^T Y),$$

ahol  $Y = a^T (X - \mu) \in \mathbb{R}$ . Azaz

$$a^T \Sigma a = E(Y^T Y) = E(Y^2) \geq 0.$$

■

**Következmény 8** A pozitív szemidefinités egyben azt is jelenti, hogy a  $\Sigma$  sajátértékei mind nem negatívak, hiszen

$$0 \leq \phi_i^T \Sigma \phi_i = \lambda_i \phi_i^T \phi_i = \lambda_i.$$

**Következmény 9**

$$\text{tr}(\Sigma) = \sum_{i=1}^p \sigma^2(X_i)$$

ugyanakkor mivel  $\Sigma$  szimmetrikus a 7 Következmény miatt

$$\text{tr}(\Sigma) = \sum_{i=1}^p \sigma^2(X_i) = \sum_{i=1}^p \lambda_i,$$

azaz a sajátértékek összege azonos az  $X$  komponenseinek teljes szórásnégyzetével. Ezért a sajátértékek (amelyek mint tudjuk nemnegatívak) a teljes szórás egy másik felbontását adják.

Legyen  $e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$  egységvektor. Világos, hogy

$$e_i^T \Sigma e_i = \sigma^2(X_i)$$

azaz az  $X$   $e_i$  irányú szórásnégyzete  $\sigma^2(X_i)$ . Hasonlóan bármely  $u$  normált vektorra  $u^T \Sigma u$ -t tekinthetjük az  $X$   $u$  irányú szórásának. Ez egyben a sajátértékek (10.1) előállítására alapján azt is jelenti, hogy  $\phi_1$  az az irány amelyre az  $X$  legnagyobb szóródása "esik" ez pedig  $\lambda_1$ . A következő sajátvektor  $\phi_2$  a  $\phi_1$ -re merőleges altérben az az irány amely a maximális szórást adja, ez  $\lambda_2$ , sorra így tovább. Ezzel az  $X$  teljes szórásnégyzetét nemcsak felbontottuk a  $\lambda_i$ -k összegére, de megtaláltuk azokat az egymásra merőleges irányokat, amelyekre eső szórásnégyzetek összege kiadja a teljes szórásnégyzet összeget is.

### 10.3 A vektrováltozó elemi statisztikai viselkedése

Legyen  $X_1, \dots, X_n$  az  $X$  eloszlásából származó  $n$  elemű minta. Jelölje

$$M = (X_1, \dots, X_n)$$

a  $p \times n$  dimenziós úgynevezett minta mátrixot. Vigyázat, itt most mindegyik  $X_i$  maga egy  $p$ -dimenziós vektor, nem pedig az  $X$  vektor  $i$ -edik komponense. Ha valahol ez félreértésre ad okot, ott  $(X_i)_m$  fogja jelölni az  $X_i$  vektor  $m$ -edik komponensét. Jelölje

$$\bar{X} = \frac{1}{n} \sum X_i,$$

tapasztalati vagy mintaátlag,

$$S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}) (X_i - \bar{X})^T$$

a tapasztalati kovariancia mátrix,

$$S_n^* = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}) (X_i - \bar{X})^T$$

pedig a korrigált tapasztalati kovariancia mátrix.

Az egydimenziós Steiner tétel alábbi általánosítása is igaz.

**Tétel 38** (Steiner tétele) Minden  $a \in \mathbb{R}^p$ -re

$$nS_n = \sum_{i=1}^n (X_i - a) (X_i - a)^T - n (a - \bar{X}) (a - \bar{X})^T.$$

**Bizonyítás.** A bizonyítás lényegében azonos az egydimenziós eset igazolásával. Kivitelezését az olvasóra bízunk mint gyakorlat. ■

**Tétel 39** 1.  $\bar{X}$  torzítatlan, konzisztens becslése  $\mu$ -nek,

2.  $S_{n-1}$  torzítatlan, konzisztens becslése  $\Sigma$ -nak,

Ha igaz továbbá, hogy  $X_i \in N_p(\mu, \Sigma)$   $i = 1, 2, \dots, n$ -re, akkor

3.  $\bar{X} \in N_p(\mu, \frac{1}{n}\Sigma)$ ,

4.  $\bar{X}$  és  $S_n$  független.

5. Ha  $\Sigma$  pozitív definit, azaz nem elfajuló és  $n > p$ , akkor  $S_n$  is pozitív definit 1 valószínűséggel, azaz a tapasztalati kovariancia mátrix sem elfajuló.

**Tétel 40** Ha a kovariancia mátrix pozitív definit, és

$$\Sigma = U^T \Lambda U$$

a spektrális felbontása, akkor

$$Y = U (X - \mu)$$

valószínűségi változóra

$$\begin{aligned} E(Y) &= 0 \\ \text{Cov}(Y) &= \Lambda, \end{aligned}$$

ha továbbá  $X \in N_p(\mu, \Sigma)$ , akkor

$$Y \in N(0, \Lambda).$$

**Bizonyítás.**

$$\begin{aligned} E(Y Y^T) &= E\left(U(X - \mu)(X - \mu)^T U^T\right) \\ &= U U^T \Lambda U U^T = \Lambda. \end{aligned}$$

■

**Tétel 41** Az  $X$   $L_2$  normában (négyzetes távolságban) mért legjobb  $k$  dimenziós becslését az  $X$ -nek az  $X$  első  $k$  főkomponense által kifeszített altérre vonatkozó vetülete biztosítja.

Nem bizonyítjuk.

**Tétel 42** Legyen  $W$  ortonormált mátrix. Ekkor az  $Z$  és  $WZ$  főkomponensei azonosak.

**Bizonyítás.** Legyen most  $X = WZ$ . A  $Z$  kovariancia mátrixának  $U^T \Lambda U$  felbontásával a  $Z$  főkomponense  $UZ$ -ne adódik. Ugyanakkor  $X$ -re

$$\begin{aligned} E(X X^T) &= E(W Z Z^T W^T) = W E(Z Z^T) W^T \\ &= W \Sigma W^T = W U^T \Lambda U W^T, \end{aligned}$$

ezért az  $X$  főkomponense  $Y = (U W^T) X = (U W^T) W Z$ . Ebből viszont az következik, hogy  $Y = U W^T W Z = U Z$  amit igazolni akartunk. ■

**Megjegyzés 16** Ezzel beláttuk, hogy az ortonormált transzformáció, azaz az elforgatás nem változtatja a főkomponenst. Az egyes koordináták átskálázása viszont igen, azaz nyújtás, affinitásra a főkomponens nem invariáns.

**10.4 A tapasztalati főkomponens**

A fentiek alapján mostmár azt vizsgáljuk, hogy, ha az  $X$  véletlen vektor eloszlásából egy  $n$  elemű mintával rendelkezünk, akkor ebből, hogyan lehet jól közelíteni az  $X$  főkomponenseit tapasztalati főkomponensekkel.

Keressük azt az  $a$  irányt amire az  $X$  vetületének szórása, azaz  $a^T X$  szórása maximális. Ezt a maximumot úgy keressük, hogy az  $a$  irányában a minta szórását maximalizáljuk. Ezzel a feladatot visszavezettük a fenti gondolatmenetre, úgy hogy a tapasztalati kovariancia mátrix spektrális felbontását használjuk. Legyenek a tapasztalati kovariancia mátrix sajátértékei  $\hat{\lambda}_1 \geq \hat{\lambda}_2, \dots, \hat{\lambda}_p$ , sajátvektoraik pedig  $\hat{a}_1, \hat{a}_2, \dots, \hat{a}_p$ . Világos, hogy tapasztalati főkomponensek

$$\hat{Y}_1 = \hat{a}_1^T X,$$

$\hat{Y}_2 = \hat{a}_2^T X, \dots, \hat{Y}_p = \hat{a}_p^T X$  korrelálatlanok és rendre maximalizálják az a maradék szórásból egy-egy vektorra vetíthető szórást.

**Definíció 69** Tekintsük a sajátértékek következő hányadosát.

$$\Psi = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_p}$$

és ennek tapasztalati megfelelőjét.

$$\hat{\Psi} = \frac{\hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_k}{\hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_k + \dots + \hat{\lambda}_p}$$

**Tétel 43** Ha  $X_i \in N_p(\mu, \Sigma)$  akkor

$$\hat{\Psi} \in N_p(\Psi, \cdot).$$

*A szórásnégyzetet szándékosan nem adtuk meg. Az állítást nem bizonyítjuk.*

**Következmény 10** A fenti tétel alapján a

$$H_0 : \Psi = \Psi_0$$

hipotézist a  $\hat{\Psi}$  próbafüggvény segítségével lehet ellenőrizni.





## Chapter 11

### OSZTÁLYOZÁS, KLASZTEREZÉS

#### 11.1 Osztályozás

Az osztályozás feladata igen sok gyakorlati probléma során felmerül. Az üzleti élet mindennapi feladata a vevők, ügyfelek osztályozása. Jó példa erre a hitelbírálat. A banktisztviselő az ügyfél adatai alapján kell, hogy döntsön arról, hogy kaphat-e hitelt, vagy differenciáltabb esetben, milyen hitelkonstrukciókat érdemes ajánlani a számára aszerint, hogy melyik ügyfélosztályba esik. Ezen ügyfélosztályok előzetesen kerülnek kialakításra, ennek módjáról később esik szó. Az informatika is számos osztályozási feladattal foglalkozik. Klasszikus példa a karakterolvasó program. Ez egy pixelhalmazról kell, hogy eldöntse, melyik betűhöz hasonlít leginkább, melyik gépi jellel azonosítsa.

Feladatunk a következő képpen fogalmazható meg. Adott  $M$  osztály,  $D = \{1, 2, \dots, M\}$ . Az objektumok  $N$  attribútummal rendelkeznek, így az objektumok azonosíthatóak az  $X \subset \mathbb{R}^N$  halmaza elemeivel, ezek az objektumokat leíró vektorok. Feltesszük, hogy adott egy  $\{X, \mathcal{F}, P\}$  Kolmogorov féle valószínűségi mező továbbá egy veszteségfüggvény  $w_{i,j}$ , ami az objektumok hibás besorolásákor felmerülő veszteséget jelenti.

$$w_{i,j} = \begin{cases} \geq 0 & \text{ha } i \neq j \\ = 0 & \text{ha } i = j \end{cases}.$$

Jelöljön  $\xi \in X$  egy véletlen objektumot, amelynek eloszlása  $P$ . Legyen  $d$  diszkriminancia függvény, azaz osztályba sorolás és

$$\hat{\eta} = d(\xi)$$

a  $d$  által javasolt osztály, legyen továbbá  $\eta$  a valódi osztálya  $\xi$ -nek. A rizikófüggvényt az

$$R(w, d) = E(w(\hat{\eta}, \eta)) = E(w(\hat{\eta}(\xi), \eta(\xi))).$$

Azt az osztálybasorolást amely  $R$ -t minimalizálja Bayes féle döntésnek, osztálybasorolásnak nevezzük. Jelölje ezt  $d^*, \eta^* = d^*(\xi)$ . Ha  $d$  ismert, akkor  $X = X_1 \cup X_2 \dots X_M$  osztályozás is ismert. Jelölje  $C_i^*$  a Bayes féle döntés osztályait és  $R^*$  a Bayes féle döntés rizikófüggvényét.

**Állítás 27** Legyen  $R$  egy tetszőleges

$$R \geq R^*$$

**Bizonyítás.**

$$R = E(w(\hat{\eta}, \eta)).$$

Alkalmazzuk a feltételes várható érték alaptulajdonságát, miszerint átlaga az eredeti átlag. Tekintsük először a diszkrét eloszlású  $\xi$  *ük* esetét.

$$\begin{aligned} E(w(\hat{\eta}, \eta)) &= E[E(w(\hat{\eta}, \eta) | \xi)] \\ &= E\left[\sum w_{i,j} P(\eta = i, \hat{\eta} = j | \xi)\right] \\ &= E\left[\sum w_{i,j} I(\hat{\eta}(\xi) = j) P(\eta = i | \xi)\right] \\ &= E\left[\sum w_{i,j} I(\hat{\eta}(\xi) = j) q_i\right], \end{aligned}$$

ahol

$$q_i = P(\eta(\xi) = i) = P(\xi \in C_i),$$

annak a valószínűsége, hogy a véletlen  $\xi$  objektum az  $i$ -edik osztályba tartozik. Így

$$\begin{aligned} R(w, d) &= E\left[\sum_{i,j} w_{i,j} I(\hat{\eta}(\xi) = j) q_i\right] \\ &= E\left[\sum_j I(\hat{\eta}(\xi) = j) \sum_i w_{i,j} q_i\right] \\ &\geq E\left[\min_j \sum_i w_{i,j} q_i\right]. \end{aligned}$$

Hasonlóan, abszolút folytonos eloszlású  $\xi$  esetén, ha  $f_i(x)$  az  $i$ -edik osztályban a  $\xi$  sűrűségfüggvénye, akkor  $p_i = P(\eta = 1)$  mellett

$$f(x) = \sum_{i=1}^M p_i f_i(x)$$

és Bayes tétele alapján

$$q_i(x) = P(\eta = i | \xi = x) = p_i \frac{f_i(x)}{f(x)}$$

ezért, ha  $d^*$  a Bayes döntés, akkor

$$\sum_{i,j} w_{i,j} f_i(x) \geq \sum_i w_{i,d^*(x)} p_i f_i(x)$$

és

$$R^* = \int \min_j \sum_{i=1}^M w_{i,j} p_i(x) f_i(x) dx.$$

■

### 11.1.1 A legközelebbi társ módszer

A legközelebbi társ módszer (Nearest Neighbor) az  $X$  téren definiált valamilyen metrikára épít. Az objektumok természetéről megszerzett ismereteink ebben a metrikában öltenek testet. Annál jobb a módszer, minnél pontosabb a metrika szétválasztó képessége.

Legyen  $T = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots, (\xi_n, \eta_n)\}$   $n$  elemű "tananyag". Elemei jó besorolási döntéseket tartalmaznak, azaz  $\xi_i$  valóban a  $C_{\eta_i}$  osztályhoz tartozik. A legközelebbi társ módszer az új  $\xi$  elemet abba a  $j_0$  osztályba sorolja, amelyre igaz, hogy

$$d(\xi, \xi_{j_0}) \leq d(\xi, \xi_i) : \forall i = 1 \dots n.$$

A legközelebbi társ  $\xi_{j_0}$ .

Legyen most két osztályunk,  $M = 2$  és tegyük fel, hogy  $W = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ . Legyen

$$R^* = \int \min \{p_1 f_1(x), p_2 f_2(x)\} dx,$$

legyen továbbá

$$R = \lim R_n$$

ahol  $R_n$  az  $n$  elemű tananyag alapján adódó rizikófüggvény. Világos, hogyha  $n \rightarrow \infty$  és az elemek valamilyen véletlen eloszlás szerint kerülnek a tananyagba akkor a végtelen sok mintaelem tökéletesen letapogatja az objektumok terét és ezért az erre épített döntés optimális lesz.

**Tétel 44** (Cover&Hart)

$$R^* \leq R \leq 2R^* (1 - R^*)$$

továbbá

$$R_n \rightarrow R^* = R,$$

azaz a Bayes féle besorolás asszimptotikusan optimális.

## 11.2 Klaszter analízis

A klaszteranalízis feladata hasonló az osztályozás feladatához. A klaszterezési feladat során ismertek az osztályok, amikbe az objektumokat be kell sorolni, hanem éppen az a kérdés, hogy hogyan alakítsunk ki úgy klasztereket (csoportokat) adott objektumokból, hogy azok

1. viszonylag homogén klasztereket alkossanak
2. a klaszterek jól elkülönüljenek
3. az egyes klasztereket a feladat szellemében jól tudjuk jellemezni.

Ilyen feladattal találkozhatunk éppen akkor, amikor egy gyártó vagy forgalmazó a fogyasztókat csoportokba kívánja sorolni, piacszegmentációt kíván végezni, annak érdekében, hogy az egyes szegmenseket egyedi módon célozza meg termékkel, marketinggel. Gyakori

feladat a biológia vagy az orvostudomány területén is, hogy nagyszámú megfigyelt objektumainkat csoportokba soroljuk, aminek alapján a lényegi eltérésekre lehet koncentrálni.

Legyenek  $X(1), X(2) \dots X(n)$  az adott objektumok, jelölje halmazukat  $X$ . Mindegyik azonosítható a saját leíró vektorával:  $(X_1(i), \dots, X_N(i))^* \in \mathbb{R}^*$ . Keressük az

$$X = \cup_{i=1}^M C_i$$

diszjunk osztályozást, ahol  $M$  sem ismert.

### 11.2.1 $K$ -közép módszer

Ez a módszer előre választott  $M$  számú klaszter kialakítására szolgál. A módszer iteratív. Tegyük fel, hogy már adóttak a  $C_1^k, \dots, C_M^k$  klaszterek. A  $k+1$  generáció a következő képpen kapható. Képezzük a klaszterek  $Z_i^k$  középpontjait.

$$Z_i^k = \frac{1}{|C_i^k|} \sum_{X_j \in C_i^k} X_j.$$

Definiáljuk egy objektum és egy klaszter távolságát mint az objektum és a klaszter középpontjának távolságát:

$$d(X_j, C_i) = d(X_j, Z_i^k).$$

Sorra vesszük az  $X_j$  objektumokat. Tegyük fel, hogy

$$d(X_j, C_l) \leq d(X_j, C_i) : \forall i = 1 \dots M.$$

Ekkor az  $X_j$  elemet a  $C_l$  osztályba helyezzük és ennek megfelelően módosítjuk az osztályközepeket.

$$Z_i^{k+1} = \frac{1}{|Z_i^k| + 1} (|C_i^k| Z_i^k + X_j).$$

Megállunk az eljárás iterálásával, ha elfogynak a besorolandó elemek vagy, ha a klaszterek nem nagyon változnak már.

Legyen a veszteségfüggvény

$$w = w(C_1^k, \dots, C_M^k) = \sum_i \sum_j d(X_j, C_i).$$

**Állítás 28** Igen általános feltételek mellett  $w$ -t minimalizálja a módszer, a konvergencia kielégítő és független a kezdeti  $C^0$  klaszterek megválasztásától.

### 11.2.2 Hierarchikus eljárások

Ez a módszer az  $N$  objektum minden  $M = 1, 2, \dots, N$  klaszterbe sorolását elvégzi és az eredmény alapján választhatjuk meg az osztályok kívánt számát. Megint feltesszük, hogy az objektumok leíróvektorai között egy  $d$  távolság definiált, (nem feltétlenül Euklideszi ez a távolság). A klaszterek távolságát most a középpontjaik távolságával definiáljuk:

$$d(C_i, C_j) = d(Z_i, Z_j).$$

Az eljárás  $M^0 = N$  klaszterrel indul, minden objektum külön klasztert alkot. Ezután összevonjuk azt a két osztályt amelyek távolsága minimális. Meghatározzuk az új klaszter

középpontját. Innen az előző lépés folytatható. Alternatív módszer, amely gyorsabb algoritmust eredményez, ha minden osztályt összevonunk, amelyek távolsága egy adott küszöb alatt van. Az eredményt szokás a klaszterek egymást követő generációjának fájával, úgynevezett dendogrammal ábrázolni. Világos, hogy az első módszer esetén a fa  $N$  levéllel rendelkezik, az alattuk lévő szinteken sorra  $N - 1$ ,  $N - 2..$  csúcspont helyezkedik el, a fa magassága pedig  $N$ . A második eljárásnál a szintek elemszáma gyorsabban is csökkenhet, így a fa magassága kisebb egyenlő mint  $N$ .



## Chapter 12 IDŐSOROK

\*

### 12.1 Alapfogalmak, definíciók

A továbbiakban olyan folyamatokat vizsgálunk, amelyeknél  $X_1, \dots, X_n$  nem független, azonos eloszlású változók.

**Véges dimenziós eloszlások**  $X_t : t \in T$

$t_1, \dots, t_n \in T : (X_{t_1}, \dots, X_{t_n})$

$P(X_{t_1} < x_1, \dots, X_{t_n} < x_n)$

$P((X_{t_1}, \dots, X_{t_n}) \in B), B \in \mathfrak{R}^n$

[Kolmogorov] Ha adott véges dimenziós eloszlásoknak egy kompatibilis rendszere, akkor létezik egy valószínűségi mező és azon egy sztochasztikus folyamat, amelynek pont ezek a véges dimenziós eloszlásai (vagyis véges dimenziós eloszlásai meghatározzák a folyamatot).

**Gauss folyamat** Akkor nevezünk egy folyamatot *Gauss folyamatnak*, ha minden véges dimenziós eloszlása normális.

#### 12.1.1 Összefüggőségi struktúrák

1. véges rendű Markov-i összefüggés

2. *martingál tulajdonság*:  $E(X_{n+1} | X_1, \dots, X_n) = X_n$

3. *stacionaritás*

**Erős stacionaritás** Egy sztochasztikus folyamatot akkor nevezünk *erősen stacionáriusnak* ha

$$\forall t_1, \dots, t_n, s : (X_{t_1}, \dots, X_{t_n}) \stackrel{d}{=} (X_{t_1+s}, \dots, X_{t_n+s})$$

**Gyenge stacionaritás** Egy sztochasztikus folyamatot akkor nevezünk *gyengén stacionáriusnak* ha

1.  $\forall t : E(X_t) = E(X_1)$

2.  $\text{cov}(X_t, X_s) = \gamma(t - s)$

A  $\gamma$  függvényt a folyamat *autokovariancia* függvényének nevezzük.

---

\*Marizca István jegyzete alapján. A szerző engedélyével.

### 12.1.2 Az autokovariancia függvény tulajdonságai

1.  $\gamma(0) = D^2[X_t] \geq 0$  (ha létezik)

2.  $|\gamma(h)| \leq \gamma(0)$

Bizonyítás: Cauchy-Schwartz:  $\langle a, b \rangle^2 \leq \|a\|^2 \|b\|^2$ .

$$\gamma(h) = \text{cov}(X_t, X_{t+h})$$

$$|\text{cov}(X_t, X_{t+h})|^2 \leq D^2[X_t] D^2[X_{t+h}]$$

3.  $\gamma(h) = \gamma(-h), \forall h \in \mathbb{Z}$

4. pozitív szemidefinit

[Pozitív Szemidefinit] Az  $M \in \mathfrak{R}^{n \times n}$  mátrix pozitív szemidefinit ha  $\forall a \in \mathfrak{R}^n$  esetén  $a^T M a \geq 0$ .

Ha egy egész számokon értelmezett  $u$  függvényre igaz, az, hogy  $\forall t_1, \dots, t_n \in \mathbb{Z}$  esetén a  $[u(t_i - t_j)]_{1 \leq i \leq n, 1 \leq j \leq n}$  mátrix pozitív szemidefinit akkor a függvényt pozitív szemidefinitnek nevezzük.

Gyengén stacionárius folyamat autokovariancia függvénye pozitív szemidefinit.

**Bizonyítás.** Legyen  $Z := (X_{t_1} - E(X), X_{t_2} - E(X), \dots, X_{t_n} - E(X))^T$ . Tudjuk, hogy  $E(Z) = 0$ , ezért  $\forall a \in \mathfrak{R}^n: E(a^T Z) = 0$ . Ezért:

$$D^2[a^T Z] = E(a^T Z a^T Z).$$

Mivel  $a^T Z$  egy skalár, ezért egyenlő önmaga transzponáltjával, így:

$$D^2[a^T Z] = E(a^T Z Z^T a) = a^T E(Z Z^T) a.$$

A  $Z Z^T$  mátrix egy olyan (diadikus) mátrix melynek  $i, j$ -edik eleme  $(X_{t_i} - E(X))(X_{t_j} - E(X))$ , így az  $E(Z Z^T)$  mátrix  $i, j$ -edik eleme  $E(X_{t_i} - E(X))(X_{t_j} - E(X)) = \gamma(t_i - t_j)$ . Tudjuk azonban, hogy mivel  $D^2[a^T Z]$  egy valószínűségi változó szórását jelöli ezért nem lehet negatív, így  $\forall a \in \mathbb{Z}^n$ :

$$0 \leq D^2[a^T Z] = a^T E(Z Z^T) a.$$

Vagyis az  $E(Z Z^T)$  mátrix pozitív szemidefinit és így a  $\gamma$  függvény is. ■

[Herglotz] Legyen  $X_1, \dots, X_n, \dots$  komplex értékű stacioner folyamat  $\gamma(h) = \text{cov}(X_t, X_{t+h}^-)$  autokovariancia-függvénnyel. Ekkor

$$\gamma(h) = \int_{-\pi}^{\pi} e^{ih\nu} dF(\nu),$$

ahol  $F(\nu)$  a spektrálfüggvény, amelyre igaz, hogy  $F(-\pi) = 0$ , jobbról folytonos, korlátos.



## 12.2 Idősorok transzformációja

Klasszikus dekompozíció:

$$X_t = m_t + \Delta_t + Y_t,$$

ahol

$m_t$ : trend, lassan változó determinisztikus,

$\Delta_t$ : szezonális, periodikus függvény,

$Y_t$ : stacionárius folyamat.

### 12.2.1 Nincs periodikus komponens

Kiindulunk egy ismert trendfüggvényből:  $m_t = a + bt$ . A minták alapján  $a$ -t, és  $b$ -t úgy határozzuk meg, hogy a  $\sum_{t=1}^n [X_t - (a + bt)]^2$  négyzetes eltérés minimális legyen.

### Mozgó átlagos simítás (moving average smoothing)

$W_t = \frac{1}{2q+1} \sum_{j=-q}^q X_{t+j}$  a simított folyamat.

$$X_t \rightarrow \frac{1}{2q+1} \sum_{j=-q}^q m_{t+j} + \frac{1}{2q+1} \sum_{j=-q}^q Y_{t+j}$$

$$\frac{1}{2q+1} \sum_{j=-q}^q Y_{t+j} \approx 0$$

$W_t = \hat{m}_t$  a trend becslése, feltéve, ha az *lineáris*!  $\hat{Y}_t = X_t - \hat{m}_t$ .

### Exponenciális simítás (exponential smoothing)

Legyen  $a \in (0, 1)$ .

$$\hat{m}_1 := X_1$$

$$\hat{m}_t := aX_t + (1-a)\hat{m}_{t-1}$$

$$\hat{m}_t = \sum_{j=0}^{t-2} a(1-a)^j X_{t-j} + (1-a)^t X_1$$

### Különbségképzés (differencing)

Definiáljuk a különbségképző (differencing,  $\nabla$ ), illetve backward shift ( $B$ ) operátorokat a következőképpen:

$$B : \mathfrak{R}^{\mathbb{Z}} \mapsto \mathfrak{R}^{\mathbb{Z}}, BX_t = X_{t-1},$$

illetve

$$\nabla : \mathfrak{R}^{\mathbb{Z}} \mapsto \mathfrak{R}^{\mathbb{Z}}, \nabla = 1 - B$$

$$\nabla X_t = (1 - B)X_t = X_t - X_{t-1}$$

Ezek felhasználásával:  $\nabla(at + b) = a$ , illetve lineáris  $m_t (= at + b)$  trend esetén  $\nabla(m_t + Y_t) = a + \nabla Y_t$ .

Hasonlóképp definiálhatjuk a  $\nabla^k$  operátort is:

$$\nabla^k = (1 - B)^k = 1 - \binom{k}{1}B + \binom{k}{2}B^2 + \cdots + \binom{k}{k}(-B)^k.$$

Például:  $\nabla^2 X_h = X_h - 2X_{h-1} + X_{h-2}$ , illetve:  $\nabla^k(a_k t^k + \cdots + a_0) = k!a_k$ .

### 12.2.2 Trend és szezonális

#### Lassú trend

Legyen

$$X_{j,k} = Y_{k+12(j-1)},$$

ahol  $j$  jelentheti pl. az évet és  $k$  a hónapot. Feltesszük, hogy egy éven belül a trend konstans:  $m_j$ .

$$\hat{m}_j := \frac{1}{12} \sum_{k=1}^{12} X_{j,k}.$$

Szezonális becslése:

$$\hat{s}_k := \frac{1}{20} \sum_{j=1}^{20} (X_{j,k} - \hat{m}_j).$$

$$\sum_{k=1}^{12} \hat{s}_k = 0$$

A szezonális periódusát ismernünk kell! Ennek meghatározásához használhatjuk például a periodogram módszert.

#### Mozgó átlagos simítás

Trend ( $\hat{m}_t$ ) becslése:

$$\hat{m}_t = \begin{cases} \frac{1}{2q+1} \sum_{i=-q}^q X_{t+i} & : d = 2q + 1 \\ \frac{1}{2q} (0.5(X_{t-q} + X_{t+q}) + \sum_{i=-q+1}^{q-1} X_{t+i}) & : d = 2q \end{cases}$$

$$\nabla_d X_t := X_t - X_{t-d} = X_t - B^d X_t = (1 - B^d) X_t$$

### 12.3 Tapasztalati autokovariancia és autokorreláció

Adott egy  $n$  elemű minta. Ekkor a *tapasztalati autokovariancia függvényt* a következőképpen definiáljuk:

$$\hat{\gamma}(h) := \frac{1}{n} \sum_{j=1}^{n-h} (X_{j+h} - \bar{X})(X_j - \bar{X}), \quad (12.1)$$

ahol  $\bar{X}$  jelöli a mintaátlagot.

Bizonyítható, hogy az így definiált empirikus autokovariancia tagokból képzett  $\Sigma = [\hat{\gamma}(i - j)]_{1 \leq i, j \leq n}$  mátrix pozitív szemidefinit. Amennyiben a 12.1 egyenletben  $n$  helyett  $(n - r)$ -rel normálnánk, úgy a kapott empirikus autokovariancia mátrixra ez nem teljesülne.

Hasonlóképp értelmezhetjük az *empirikus autokorrelációs függvényt* is:

$$\hat{\rho}(h) := \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)}.$$

Box és Jenkins ökölszabálya szerint  $n \geq 50$  és  $h \leq n/4$  esetén van értelme az idősor analízisével foglalkozni.

Fontos megjegyezni, hogy az empirikus autokovariancia ill. autokorreláció függvényeket nem-stacionárius esetben is ki tudjuk számítani, ezért a kapott eredmények értelmezésénél ezt figyelembe kell venni.

Ha például  $\hat{\rho}(h)$  függvény lecsengése lassú, hatványfüggvény jellegű, az az erős függés helyett jelentheti lassú determinisztikus trend jelenlétét is.

Amennyiben a folyamatban szezonáltság van jelen ez a  $\hat{\rho}(h)$  periodicitását eredményezi.

## 12.4 Parciális autokovariancia függvény

$$\alpha(1) = \text{cov}(X_1, X_2)$$

$$\alpha(h) = \text{cov}(X_1 - \text{E}(X_1|X_2, \dots, X_h), X_{h+1} - \text{E}(X_{h+1}|X_2, \dots, X_h))$$

$$\alpha(k) = \frac{\det R_k^*}{\det R_k},$$

ahol  $R_k$  az autokorreláció mátrix,  $R_k = [\rho(i - j)]_{1 \leq i, j \leq k}$ ,  $R_k^*$ -t pedig úgy kapjuk  $R_k$ -ből, hogy annak utolsó sorát a  $[\rho_1, \dots, \rho_k]$  vektorra cseréljük.

## 12.5 Fehér zaj

**Fehér zaj** *Fehér zajnak* hívjuk, és  $WN(0, \sigma^2)$ -tel jelöljük azokat a folyamatokat, melyre

$$\gamma(k) = \begin{cases} \sigma^2, & k = 0 \\ 0, & k \neq 0 \end{cases}$$

A fehér zaj spektrálfüggvénye konstans.

$$f(\lambda) = \frac{\sigma^2}{2\pi}, \lambda \in [-\pi, \pi]$$

## 12.6 Mozgóátlag (MA) folyamatok

**Mozgó átlag folyamat** Akkor nevezünk egy  $X_t$  folyamatot *mozgó átlag folyamatnak* ha felírható

$$X_t = \Theta_0 e_t + \Theta_1 e_{t-1} + \dots + \Theta_q e_{t-q}$$

alakban, ahol  $e_t \sim WN(0, 1)$  fehér zaj,  $\Theta_i, 0 \leq i \leq q$  pedig konstansok.

$X_t$  és  $e_t$  közötti összefüggést felírhatjuk a  $B$  backshift operátor segítségével:

$$\begin{aligned} X_t &= \Theta_0 e_t + \Theta_1 B e_t + \cdots + \Theta_q B^q e_t \\ &= (\Theta_0 + \Theta_1 B + \cdots + \Theta_q B^q) e_t \end{aligned}$$

Formálisan definiálhatjuk a  $\Theta(z)$  polinomot a következőképp:  $\Theta(z) := \Theta_0 + \Theta_1 z + \cdots + \Theta_q z^q$ . Ezzel a jelöléssel:

$$X_t = \Theta(B) e_t$$

Az algebra alaptétele szerint egy  $n$  változós polinomnak pontosan  $n$  darab nem feltétlenül különböző gyöke van. Ennek alapján felírhatjuk  $\Theta(z)$  gyöktényezős alakját:

$$\Theta(z) = \Theta_q \prod_{i=1}^q (z - z_i).$$

$$E(X_r e_s) = \begin{cases} \Theta_{r-s} & : r - q \leq s \leq r \\ 0 & : \text{különben} \end{cases}$$

$$E(X_{t+k} X_t) = E\left(X_{t+k} \sum_{i=0}^q \Theta_i e_{t-i}\right) = \sum_{i=0}^q \Theta_i E[X_{t+k} e_{t-i}] = \sum_{i=0}^{q-k} \Theta_i \Theta_{k+i}$$

Mivel  $E(X_{t+k} X_t)$   $t$ -től nem, csak  $k$ -tól függ, ezért  $X_t$  gyengén stacionárius, és az autokovariancia függvénye:

$$\gamma(k) = \begin{cases} \sum_{i=0}^{q-k} \Theta_i \Theta_{k+i} & : k \leq q \\ 0 & : k > q \end{cases} \quad (12.2)$$

Az előző egyenletben kifejeztük  $\gamma(k)$  értékét a  $\Theta_i$  értékek segítségével. Felvetődik a kérdés, hogy lehet-e ugyanezt visszefelé, illetve mi annak a feltétele, hogy adott  $\gamma(k)$ ,  $k = 0, 1, 2, \dots, q$  autokovariancia függvényhez létezzenek a 12.2 egyenletet kielégítő  $\Theta_i$  értékek.

Vagyis a kérdés: adott  $\gamma(k)$ ,  $k = 0, 1, 2, \dots, q$  esetén megoldható-e  $\gamma(k) = b_0 b_k + b_1 b_{k+1} + \dots + b_{q-k} b_q$ ,  $k = 0, 1, 2, \dots, q$  egyenletrendszer.

A válasz pedig, hogy a megoldhatóság feltétele, hogy a

$$\Gamma(s) := \gamma(0) + \sum_{k=1}^q \gamma(k)(s^k + s^{-k})$$

függvénynek az az  $|s| = 1$  egységkörön csak páros multiplicitású gyökei legyenek.

## 12.7 Autoregresszív (AR) folyamatok

**Autoregresszív folyamatok** Akkor nevezünk egy  $(X_t)$  folyamatot *autoregresszív*nek, ha felírható

$$X_t = \Phi_1 X_{t-1} + \cdots + \Phi_p X_{t-p} + e_t$$

alakban, ahol  $e_t$  fehér zaj.

A mozgó átlag folyamatoknál definiált  $\Theta$  függvényhez hasonlóan definiálhatjuk a  $\Phi(z) := \Phi_0 + \Phi_1 z + \cdots + \Phi_p z^p$  polinomot. Így

$$\Phi(B) X_t = e_t.$$

## 12.7.1 Példa, AR(1) folyamat

$$X_t := \Lambda X_{t-1} + e_t = e_t + \Lambda(\Lambda X_{t-2} + e_{t-1}) = \sum_{j=0}^k \Lambda^j e_{t-j} + \Lambda^{k+1} X_{t-k-1}$$

$$\left\| X_t - \sum_{j=0}^k \Lambda^j e_{t-j} \right\|^2 = \Lambda^{2(k+1)} \|X_{t-k-1}\|^2 \rightarrow 0$$

ha  $|\Lambda| < 1$ . Ebben az esetben azt mondjuk, hogy az AR folyamat *kauzális*. Ebben az esetben felírhatjuk az  $X_t$ -t

$$X_t = \sum_{j=0}^{\infty} \Lambda^j e_{t-j}$$

alakban is, ami egy MA( $\infty$ ) folyamatnak felel meg. Ezt nevezzük az AR(1) folyamat *kauzális MA( $\infty$ ) előállításának*.

Észrevehetjük, hogy amennyiben az AR(1) folyamat *kauzális*, azaz létezik MA( $\infty$ ) előállítása, akkor a

$$\Phi(z) = 1 - \Lambda z$$

polinomnak az egyedüli gyöke az egységkörön kívül helyezkedik el. Általánosságban is igaz a következő

Egy AR(p) folyamatnak akkor és csak akkor létezik *kauzális MA( $\infty$ ) előállítása* ha a  $\Phi(z) = 0$  egyenletnek nincsen a  $|z| \leq 1$  egységkörön belül gyöke.

Tekintsük most azt az esetet amikor  $|\Lambda| > 1$ !

$$\begin{aligned} X_{t+1} &= \Lambda X_t + e_{t+1} \\ X_t &= \frac{1}{\Lambda} X_{t+1} - \frac{1}{\Lambda} e_{t+1} = \\ &= - \sum_{j=1}^k \frac{1}{\Lambda^j} e_{t+j} + \frac{1}{\Lambda} X_{t+k} \end{aligned}$$

$|\Lambda| > 1$  esetén

$$X_t = - \sum_{j=1}^{\infty} \frac{1}{\Lambda^j} e_{t+j}$$

Vagyis  $|\Lambda| > 1$  esetén is létezik MA( $\infty$ ) előállítás, ez azonban nem *kauzális*.

AR(1) folyamatok autokovariancia - függvényét könnyen kifejezhetjük az MA( $\infty$ ) előállításuk segítségével.

$$\text{cov}(X_{t+h}, X_t) = \lim_{n \rightarrow \infty} E \left( \sum_{j=0}^n \Lambda^j e_{t+k-j} \sum_{k=0}^n \Lambda^k e_{t-k} \right) = \Lambda^h \sum_{j=0}^{\infty} \Lambda^{2j} = \frac{\Lambda^h}{1 - \Lambda^2}$$

### 12.7.2 Yule-Walker egyenletek

Ezeket az egyenleteket az  $AR(p)$  folyamatokra tekintjük.

$$\begin{aligned} X_t &= \Phi_1 X_{t-1} + \dots + \Phi_p X_{t-p} + e_t \\ X_{t-k} X_t &= \Phi_1 X_{t-k} X_{t-1} + \dots + \Phi_p X_{t-k} X_{t-p} + X_{t-k} e_t \\ \gamma(k) = E(X_{t-k} X_t) &= \Phi_1 \gamma(k-1) + \dots + \Phi_p \gamma(k-p) \\ \rho(k) &= \Phi_1 \rho(k-1) + \dots + \Phi_p \rho(k-p) \end{aligned}$$

Ezen utóbbi egyenletet  $k = 1, 2, \dots, p$  értékekre felírva és mátrix alakba rendezve kapjuk a Yule-Walker egyenletrendszerét.

$$\begin{pmatrix} \rho_0 & \rho_1 & \dots & \rho_{p-1} \\ \rho_1 & \rho_0 & \dots & \rho_{p-2} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p-1} & \dots & \rho_1 & \rho_0 \end{pmatrix} \begin{pmatrix} \Phi_1 \\ \Phi_2 \\ \vdots \\ \Phi_p \end{pmatrix} = \begin{pmatrix} \rho_1 \\ \rho_2 \\ \vdots \\ \rho_p \end{pmatrix}$$

Ennek segítségével találhatunk adott  $\rho(1), \dots, \rho(p)$  értékekhez olyan  $\Phi_1, \dots, \Phi_p$  együtthatókat, hogy a kapott  $AR(p)$  folyamat autokorreláció - függvényének első  $p$  eleme az adott  $\rho(1), \dots, \rho(p)$ -val egyezzen. Ugyanez az egyenlet használható a  $\hat{\Phi}_1, \dots, \hat{\Phi}_p$  értékek becslésére is a  $\hat{\rho}_1, \dots, \hat{\rho}_p$  segítségével.

## 12.8 Autoregresszív - mozgóátlag (ARMA) folyamatok

Az  $X_t$  folyamatot  $ARMA(p, q)$  folyamatnak nevezzük, ha

$$\Phi(B)X_t = \Theta(B)e_t, \quad (12.3)$$

ahol  $\Phi(B)$   $p$ -ed és  $\Theta(B)$   $q$ -ad fokú polinomok.

Akkor nevezzük a folyamatot kauzálisnak, ha létezik  $MA(\infty)$  előállítás. Egy kauzális ARMA folyamat autokovariancia - függvényét az  $MA(\infty)$  előállításában szereplő együtthatókkal a következő tétel segítségével tudjuk kifejezni:

Amennyiben  $X_t$  stacionárius, és  $\sum_{j=-\infty}^{\infty} |\Psi_j| < \infty$  akkor az  $\eta_t := \sum_{j=-\infty}^{\infty} \Psi_j X_{t-j}$  függvény is stacionárius, és autokovariancia - függvénye:

$$\gamma_\eta(h) = \sum_{j=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} \Psi_j \Psi_k \gamma_\xi(h - j + k)$$

### 12.8.1 A kauzalitás szükséges és elégséges feltétele

Tegyük fel, hogy a  $\Phi(z)$  és  $\Theta(z)$  polinomoknak nincs közös gyöke. Ezt nyugodtan feltehetjük, mivel ellenkező esetben a (12.3) egyenletben ezzel a közös gyökkel egyszerűsíthetünk.

Az  $X_t$  ARMA(p, q) folyamat kauzalitásának szükséges és elégséges feltétele, hogy a  $\Phi(z) = 0$  egyenletnek ne legyen a  $|z| \leq 1$  egységkörön belül gyöke.

$$\sum_{j=0}^{\infty} \Psi_j z^j = \frac{\Theta(z)}{\Phi(z)}, |z| \leq 1$$

**Invertálhatóság** Az  $X_t$  folyamatot *invertálhatónak* nevezzük, ha  $\exists \Pi_j : \sum_{j=0}^{\infty} |\Pi_j| < \infty$  és  $e_t = \sum_{j=0}^{\infty} \Pi_j X_{t-j}$ .

Az  $X_t$  ARMA(p,q) folyamat invertálhatóságának szükséges és elégséges feltétele, hogy a  $\Theta(z) = 0$  egyenletnek ne legyen a  $|z| \leq 1$  egységkörön belül gyöke.

$$\sum_{j=0}^{\infty} \Pi_j z^j = \frac{\Phi(z)}{\Theta(z)}, |z| \leq 1$$

## 12.9 Az átlag és az autokovariancia becslései

$$E(X_t) = \mu, \text{Cov}(X_t, X_{t+n}) = \gamma(n)$$

Az átlag természetes becslése

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n} \text{ és erre } E(\bar{X}_n) = \mu$$

Az átlag szórásnégyzete

$$\begin{aligned} nD^2(\bar{X}_n) &= \frac{1}{n} \sum_{i,j=1}^n \text{Cov}(X_i, X_j) = \sum_{i,j=1}^n \frac{1}{n} \gamma(i-j) = \\ &= \sum_{h=-n+1}^{n-1} \sum_{j=1}^{n-|h|} \frac{1}{n} \gamma(h) = \sum_{|h|<n} \frac{n-|h|}{n} \gamma(h) \leq \sum_{|h|<n} |\gamma(h)| \end{aligned}$$

Ha  $\gamma(n) \rightarrow 0$ , akkor

$$D^2(\bar{X}_n) \leq \sum_{|h|<n} |\gamma(h)| \rightarrow 0.$$

Ha pedig még az is igaz, hogy a  $\gamma(n)$  sorozat abszolút konvergens, akkor az átlag szórásnégyzetére az alábbi aszimptotikát adhatjuk

$$\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty \Rightarrow nD^2(\bar{X}_n) \rightarrow \sum_{h=-\infty}^{\infty} \gamma(h)$$

### 12.9.1 A spektrálfüggvény és az autokovariancia kapcsolata

[Inverz Fourier transzformált]

Igaz, hogy

$$\sum_{n=-\infty}^{\infty} |K(n)| < \infty \Rightarrow K(h) = \int_{-\pi}^{\pi} e^{ih\nu} f(\nu) d\nu,$$

ahol

$$f(\lambda) = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} e^{in\lambda} K(n)$$

ugyanis

$$\int_{-\pi}^{\pi} e^{ih\nu} f(\nu) d\nu = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} K(n) \int_{-\pi}^{\pi} e^{i(h-n)\nu} = K(h).$$

Ennek egyszerű következménye az alábbi állítás. Egy  $\gamma(h)$  sorozatra

$$\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty \quad \text{és} \quad \Leftrightarrow \quad f(\lambda) = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} e^{-in\lambda} \gamma(n) \geq 0.$$

$\gamma$  autokovariancia függvény

Ennek alkalmazásaként számítsuk ki, mikor lesz az alábbi alakú  $K$  függvény autokovariancia függvény!

$$K(h) = \begin{cases} 1 & \text{ha } h = 0 \\ \rho & \text{ha } h = \pm 1 \\ 0 & \text{egyébként} \end{cases}$$

$$f(\lambda) = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} e^{-in\lambda} K(n) = \frac{1}{2\pi} (1 + 2\rho \cos \lambda) \geq 0 \Rightarrow |\rho| \leq \frac{1}{2}$$

### 12.9.2 Aszimptotikus normalitás

Legyen például

$$X_t = \mu + \sum_{j=-\infty}^{\infty} \Psi_j Z_{t-j},$$

ahol  $Z_t$  független, azonos eloszlású 0 várható értékkel és  $\sigma^2$  szórásnégyzettel, továbbá

$$\sum_{j=-\infty}^{\infty} |\Psi_j| < \infty \quad \text{de} \quad \sum_{j=-\infty}^{\infty} \Psi_j \neq 0.$$

Ekkor

$$\bar{X}_n \xrightarrow{d} N \left( \mu, \frac{1}{n} \sum_{j=-\infty}^{\infty} \Psi(n) \right)$$

### 12.9.3 $\gamma(n)$ becslése

$$\widehat{\gamma}(n) = \frac{1}{n} \sum_{t=1}^{n-h} (X_t - \bar{X}_n) (X_{t+h} - \bar{X}_n),$$

ahol  $0 \leq h \leq n-1$ .

Ez általánosságban torzított becslés (bár bizonyos további feltételek mellett aszimptotikusan torzítatlan), de viszont a  $\widehat{\Gamma}_n = \left[ \widehat{\gamma}(i-j) \right]_{1 \leq i, j \leq n}$  mátrixa pozitív szemidefinit.

Ennek bizonyításához elég, hogy  $\widehat{\Gamma}_n = \frac{1}{n} M M^T$  a következő az  $Y_i = X_i - \bar{X}_n$  jelöléssel kifejezett  $M$  mátrixszal



$$M = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 & Y_1 & Y_2 & \dots & Y_{n-1} & Y_n \\ 0 & 0 & \dots & 0 & Y_1 & Y_2 & Y_3 & \dots & Y_n & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & Y_1 & \dots & Y_{n-2} & Y_{n-1} & Y_n & 0 & \dots & 0 & 0 \end{bmatrix}.$$

Általános ökölszabályként elmondhatjuk, hogy  $\rho(h)$  becslése  $(\widehat{\rho(h)} = \frac{\widehat{\rho(h)}}{\widehat{\rho(0)}})$  akkor jó, ha  $n \geq 50$  és  $h \leq \frac{n}{4}$ .

#### 12.9.4 Az autokorrelációk mikor különböznek szignifikánsan 0-tól?

Ha a következő alakú szűrt független zajra

$$X_t - \mu = \sum_{j=-\infty}^{\infty} \Psi_j Z_{t-j}$$

és  $Z_t$  független, azonos eloszlású 0 várható értékkel és  $\sigma^2$  szórásnégyzettel, továbbá

$$\sum_{j=-\infty}^{\infty} |\Psi_j| < \infty \text{ és } E(Z_i^4) < \infty$$

akkor

$$[\widehat{\rho(1)}, \dots, \widehat{\rho(h)}] \xrightarrow{d} N\left([\rho(1), \dots, \rho(h)], \frac{1}{n}W\right),$$

ahol  $W$  az ún. Bartlett mátrix.

$$W_{i,j} = \sum_{k=1}^{\infty} [\rho(k+i) + \rho(k-i) - 2\rho(i)\rho(k)] [\rho(k+j) + \rho(k-j) - 2\rho(j)\rho(k)].$$

Például a független fehér zaj folyamatra  $\rho(l) \neq 0$ , ha  $l \neq 0$ , azaz

$$W_{i,j} = \begin{cases} 1 & \text{ha } i = j \\ 0 & \text{ha } i \neq j \end{cases}$$

azaz

$$\widehat{\rho(1)}, \dots, \widehat{\rho(h)} \approx \text{iid } N\left(0, \frac{1}{n}\right),$$

ennek konfidenciaintervalluma  $\pm 1.96 \frac{1}{\sqrt{n}}$ , amely értéket a normális eloszlás táblázatából olvashatunk ki.

### 12.10 ARMA modellek becslései

Amikor egy folyamatot ARMA modellel közelítünk, a következő lépések szerint járunk el:

1. megbecsüljük  $p$ -t és  $q$ -t, az ARMA folyamathoz tartozó két polinom fokszámát

2. megbecsüljük a polinomok együtthatóit

3. megbecsüljük a szórásnégyzetet

### 12.10.1 Ismert $p$ és $q$

Tiszta autoregresszív esetben felírhatjuk a Yule-Walker egyenleteket:

$$X_t - \Phi_1 X_{t-1} - \dots - \Phi_p X_{t-p} = e_t \text{ ahol } e_t \sim WN(0, \sigma^2)$$

$$\Gamma_p \Phi = \gamma_p,$$

ahol

$$\Gamma_p = [\gamma(i-j)]_{i,j=1}^p$$

$$\gamma_p^T = [\gamma(1), \dots, \gamma(p)]$$

$$\Phi^T = [\Phi(1), \dots, \Phi(p)].$$

Továbbá

$$\sigma^2 = D^2 e_t = D^2 (X_t - \Phi_1 X_{t-1} - \dots - \Phi_p X_{t-p}) = \gamma(0) - \Phi^T \gamma_p.$$

Így a Yule-Walker becslések a következő alakúak lesznek:

$$\widehat{\Gamma}_p \widehat{\Phi} = \widehat{\gamma}_p$$

és

$$\widehat{\sigma}^2 = \widehat{\gamma}(0) - \widehat{\Phi}^T \widehat{\gamma}_p.$$

### 12.10.2 Ismeretlen $p$

Ha  $p$  nem ismert és  $AR(m)$ -et próbálunk illeszteni, akkor azt várjuk, hogy  $\widehat{\Phi}_{m,m}$  kicsi lesz.

$$\sqrt{n} (\widehat{\Phi}_m - \Phi_m) \xrightarrow{d} N_m(0, \sigma^2 \Gamma_m^{-1}),$$

ahol  $\Phi_m$  az  $X_{m+1}$  legjobb lineáris közelítésének együtthatóvektora:

$$\|X_{m+1} - \Phi_m^T (X_1, \dots, X_m)\| \rightarrow \min.$$

### 12.10.3 A Durbin-Levinson algoritmus

A mátrixinvertálás kikerülésére a Durbin-Levinson algoritmust használjuk. Legyenek az  $m$ -ed rendű illesztés együtthatói

$$\widehat{\Phi}_m = (\widehat{\Phi}_{m,1}, \widehat{\Phi}_{m,2}, \dots, \widehat{\Phi}_{m,m}) = \widehat{R}_m^{-1} \widehat{\rho}_m$$

és

$$\widehat{v}_m = \widehat{\gamma}(0) \left[ 1 - \widehat{\rho}_m^T \widehat{R}_m^{-1} \widehat{\rho}_m \right].$$

Ekkor  $\widehat{\Phi}_{1,1} = \widehat{\rho}(1)$  és  $\widehat{v}_1 = \widehat{\gamma}(0) [1 - \widehat{\rho}^2(1)]$ . Továbbá a becsült parciális autokovariancia függvény

$$\widehat{\Phi}_{m,m} = \left[ \widehat{\gamma}(m) - \sum_{j=1}^{m-1} \widehat{\Phi}_{m-1,j} \widehat{\gamma}(m-j) \right] / \widehat{v}_{m-1}$$

$$\begin{bmatrix} \widehat{\Phi}_{m,1} \\ \vdots \\ \widehat{\Phi}_{m,m-1} \end{bmatrix} = \widehat{\Phi}_{m-1} - \widehat{\Phi}_{m,m} \begin{bmatrix} \widehat{\Phi}_{m-1,m-1} \\ \vdots \\ \widehat{\Phi}_{m-1,1} \end{bmatrix}$$

és

$$\widehat{v}_m = \widehat{v}_{m-1} (1 - \widehat{\Phi}_{m,m}^2).$$

Elméletileg  $\alpha(m) = \Phi_{m,m} = 0$ , ha  $m > p$ , gyakorlatilag  $\sqrt{n}\widehat{\Phi}_{m,m} \rightarrow N(0,1)$ , azaz  $P\left(-1.96\frac{1}{\sqrt{n}} < \widehat{\Phi}_{m,m} < 1.96\frac{1}{\sqrt{n}}\right) = 0.95$ . A rendre ezzel előzetes becslést adhatunk:  $\widehat{p} = \min \left\{ r : \forall m > r \ |\widehat{\Phi}_{m,m}| < 1.96\frac{1}{\sqrt{n}} \right\}$ .

#### 12.10.4 Az innovációs algoritmus

A Gram-Schmidt ortogonalizációs eljárással független vektorokból ortogonális rendszert készíthetünk vetítésekkel. Az eljárást idősorokra is alkalmazhatjuk a következő módon:

$$E(X)_t = 0 \text{ és } \kappa(i, j) := E(X_i X_j)$$

$$H_n := \langle X_1, \dots, X_n \rangle = \langle X_1 - \widehat{X}_1, \dots, X_n - \widehat{X}_n \rangle \text{ ahol } \widehat{X}_{n+1} := pr_{H_n} X_{n+1}.$$

$$\widehat{X}_{n+1} = \begin{cases} 0, & n = 0 \\ \sum_{j=0}^n \Theta_{n,j} (X_{n-j+1} - \widehat{X}_{n-j+1}), & n \neq 0 \end{cases}.$$

A  $\Theta$  együtthatók rekurzív kiszámítását a  $v$  segédváltozóval (szórásnégyzet) a következő rendszer adja meg:

$$v_n := \left\| X_{n+1} - \widehat{X}_{n+1} \right\|^2$$

így

$$v_0 = \kappa(1, 1)$$

$$\Theta_{n,n-k} = \frac{1}{v_k} \left[ \kappa(n+1, k+1) - \sum_{j=0}^{k-1} \Theta_{k,k-j} \Theta_{n,n-j} v_j \right] \text{ ahol } k = 0..n-1$$

$$v_n = \kappa(n+1, n+1) - \sum_{j=0}^{n-1} \Theta_{n,n-j}^2 v_j$$

Ennek bizonyítását úgy kezdjük, hogy  $0 \leq k \leq n$  esetén  $X_1 - \widehat{X}_1, \dots, X_n - \widehat{X}_n$  ortogonális, hiszen  $X_i - \widehat{X}_i \in H_{j-1}$ , ha  $i < j$ , és  $X_j - \widehat{X}_j \perp H_{j-1}$ . Tehát  $\widehat{X}_{n+1}$  definícióját használva

$$\langle \hat{X}_{n+1}, X_{k+1} - \hat{X}_{k+1} \rangle = \Theta_{n,n-k} v_k,$$

amely egyenlethez a  $X_{n+1} - \hat{X}_{n+1} \perp X_{k+1} - \hat{X}_{k+1}$  azonosságot adva kapjuk, hogy

$$\langle X_{n+1}, X_{k+1} - \hat{X}_{k+1} \rangle = \Theta_{n,n-k} v_k.$$

Ebbe  $\hat{X}_{k+1}$  definícióját írva

$$\begin{aligned} \Theta_{n,n-k} &= \frac{1}{v_k} \langle X_{n+1}, X_{k+1} - \sum_{j=0}^{k-1} \Theta_{k,k-j} (X_{j+1} - \hat{X}_j + 1) \rangle = \\ &= \frac{1}{v_k} \left[ \kappa(n+1, k+1) - \sum_{j=0}^{k-1} \Theta_{k,k-j} \Theta_{n,n-j} v_j \right] \end{aligned}$$

továbbá

$$\begin{aligned} v_n &= |X_{n+1} - \hat{X}_{n+1}|^2 = |X_{n+1}|^2 + |\hat{X}_{n+1}|^2 = \\ &= \kappa(n+1, k+1) - \sum_{j=0}^{n-1} \Theta_{n,n-j}^2 v_j \end{aligned}$$

Például a MA(1) folyamat predikcióját így adhatjuk meg:

$$X_t = Z_t + \Theta Z_{t-1} \text{ ahol } Z_t \sim WN(0, \sigma^2)$$

$$\kappa(i, j) = \begin{cases} \sigma^2(1 + \Theta^2) & i = j \\ \Theta\sigma^2 & i = j + 1 \\ 0 & \text{egyébként} \end{cases}$$

$$v_0 = \kappa(1, 1) = \sigma^2(1 + \Theta^2),$$

$$\Theta_{n,j} = \begin{cases} \frac{1}{v_{n-1}} \Theta \sigma^2 & j = 1 \\ 0 & 2 \leq j \leq n \end{cases}$$

$$v_n = (1 + \Theta^2) \sigma^2 - \frac{1}{v_{n-1}} \Theta^2 \sigma^4,$$

$$r_n := \frac{v_n}{\sigma^2} = (1 + \Theta^2) - \frac{1}{v_{n-1}} \Theta^2 \sigma^2.$$

Tehát a predikció:

$$\hat{X}_{n+1} = \frac{\Theta}{r_{n-1}} (X_n - \hat{X}_n)$$

## 12.10.5 Mozgóátlag folyamatok becslései

Az  $X_1, \dots, X_n$  adatokra a következő előfeltevést tesszük annak analógiájára, hogy az  $X\hat{X}$  mennyiségek voltak a hibák:

$$X_t = Z_t + \hat{\Theta}_{m,1}Z_{t-1} + \dots + \hat{\Theta}_{m,m}Z_{t-m} \text{ ahol } Z_t \sim WN(0, \hat{v}_m^2).$$

Ha  $\hat{\gamma}(0) > 0$ , akkor vezessük be a becsült együtthatók vektorára a  $\hat{\Theta}_m = (\hat{\Theta}_{m,1}, \dots, \hat{\Theta}_{m,m})$  jelölést! Ezekre a következő rekurzív becslés érvényes:

$$\hat{v}_0 = \hat{\gamma}(0)$$

és

$$\hat{\Theta}_{m,m-k} = \hat{v}_k^{-1} \left[ \hat{\gamma}(m-k) - \sum_{j=0}^{k-1} \hat{\Theta}_{m,m-j} \hat{\Theta}_{k,k-j} \hat{v}_j \right]$$

$$\hat{v}_m = \hat{\gamma}(0) - \sum_{j=0}^{k-1} \hat{\Theta}_{m,m-j}^2 \hat{v}_j \text{ ahol } k = 0, \dots, m-1$$

## 12.10.6 Aszimptotikus viselkedés ARMA folyamatok esetén

A jelölések rövid leírása a következő:

$$\Phi(B)X_t = \Theta(B)Z_t \text{ ahol } Z_t \sim IID(0, \sigma^2) \text{ és } E(Z_t^4) < \infty$$

$$\Psi(z) = \sum_{j=0}^{\infty} \frac{\Theta(z)}{\Phi(z)}, |z| \leq 1, \Psi_0 = 1$$

Ekkor minden  $k$ -ra

$$\sqrt{n} \left[ \hat{\Theta}_{m,1} - \Psi_1, \dots, \hat{\Theta}_{m,k} - \Psi_k \right] \xrightarrow{d} N(0, A)$$

ahol

$$A_{i,j} = \sum_{r=1}^{\min(i,j)} \Psi_{i-r} \Psi_{j-r}$$

továbbá

$$m(n) \rightarrow \infty$$

úgy hogy

$$m(n) = o(\sqrt[3]{n}) \text{ és } \hat{v}_m \xrightarrow{p} \sigma^2.$$

Itt jegyezzük meg, hogy  $AR(p)$  esetben a Durbin-Levinson algoritmus által a  $\Phi_p$ -re adott  $\hat{\Phi}_p = (\hat{\Phi}_{p,1}, \dots, \hat{\Phi}_{p,p})$  becslés konzisztens, ha  $n \rightarrow \infty$ . Viszont  $MA(k)$  esetben az innovációs algoritmus által adott  $\hat{\Theta}_q = (\hat{\Theta}_{q,1}, \dots, \hat{\Theta}_{q,q})$  becslés nem konzisztens, viszont a  $(\hat{\Theta}_{m,1}, \dots, \hat{\Theta}_{m,q})$  már az.

A gyakorlatban  $MA(q)$  esetben tudjuk, hogy  $\rho(m) = 0$ , ha  $m > q$ , és Bartlett tétele miatt

$$\hat{\rho}(m) \xrightarrow{d} N \left( 0, \frac{1}{n} \sum_{i=-q}^{n-q} \rho(i) \right).$$

### 12.10.7 Maximum likelihood becslések

$E(X_t) = 0$  tulajdonságú Gauss folyamat esetén a  $\Gamma_n = E(\underline{X}_n \underline{X}_n^T)$  jelöléssel a likelihood függvény a következő:

$$L(\Gamma_n) = \frac{1}{(2\pi)^{n/2}} \frac{1}{(\det \Gamma_n)^{1/2}} \exp \left( -\frac{1}{2} \underline{X}_n^T \Gamma_n^{-1} \underline{X}_n \right).$$

A mátrixinvertálás és a determinánsszámítás kikerülésére a következő algoritmus javasolt:  $k = 0, \dots, n-1$  esetén

$$\hat{X}_{k+1} = \sum_{j=0}^{k-1} \Theta_{k,k-j} (X_{j+1} - \hat{X}_{j+1})$$

azaz

$$\begin{pmatrix} \hat{X}_1 \\ \vdots \\ \hat{X}_n \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ \Theta_{1,1} & 0 & 0 & \dots & 0 \\ \Theta_{2,2} & \Theta_{2,1} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \Theta_{n-1,n-1} & \Theta_{n-1,n-2} & \Theta_{n-1,n-3} & \dots & 0 \end{pmatrix} \begin{pmatrix} X_1 - \hat{X}_1 \\ \vdots \\ X_n - \hat{X}_n \end{pmatrix}.$$

A képletben szereplő mátrixot jelöljük a következőképpen

$$\bar{C} = [\Theta_{i,i-j}]_{i,j=0}^{n-1} \text{ ahol } j \leq 0 \text{ esetén } \Theta_{i,j} = 0.$$

Ezt a mátrixot módosítsuk úgy, hogy a főátlóba 1-eket írunk:  $C := \bar{C} + Id$ . Ekkor

$$C (\underline{X}_n - \hat{\underline{X}}_n) = (\bar{C} + Id) (\underline{X}_n - \hat{\underline{X}}_n) = \hat{\underline{X}}_n + \underline{X}_n - \hat{\underline{X}}_n = \underline{X}_n$$

azaz

$$\Gamma_n = E(\underline{X}_n \underline{X}_n^T) = CE \left( (\underline{X}_n - \hat{\underline{X}}_n) (\underline{X}_n - \hat{\underline{X}}_n)^T \right) C^T = CDC^T$$

ahol

$$D = \text{Diag}(v_0, v_1, \dots, v_{n-1})$$

ezért a determináns egyszerűen számítható.

$$\det \Gamma_n = (\det C)^2 \det D = v_0 v_1 \dots v_{n-1}.$$

A kitevő is egyszerűbb alakra hozható:

$$\begin{aligned} \underline{X}_n^T \Gamma_n^{-1} \underline{X}_n &= \left( \underline{X}_n - \hat{\underline{X}}_n \right)^T C^T \Gamma^{-1} C \left( \underline{X}_n - \hat{\underline{X}}_n \right) = \\ & \left( \underline{X}_n - \hat{\underline{X}}_n \right)^T C^T C^{T-1} D^{-1} C^{-1} C \left( \underline{X}_n - \hat{\underline{X}}_n \right) = \\ & \left( \underline{X}_n - \hat{\underline{X}}_n \right)^T D^{-1} \left( \underline{X}_n - \hat{\underline{X}}_n \right) = \sum_{j=1}^n \frac{\left( X_j - \hat{X}_j \right)^2}{v_{j-1}} \end{aligned}$$

Tehát végeredményképpen a likelihood függvény legegyszerűbb alakja:

$$L(\Gamma_n) = (2\pi)^{-n/2} (v_0 v_1 \dots v_{n-1})^{-1/2} \exp \left[ -\frac{1}{2} \sum_{j=1}^n \frac{\left( X_j - \hat{X}_j \right)^2}{v_{j-1}} \right].$$