

Ökonometria

Ferenci Tamás, tamas.ferenci@medstat.hu

2025. szeptember 4.

Tartalom

| | |
|--|-----------|
| Előszó | 5 |
| 1 Út a többváltozós regresszióhoz | 7 |
| 1.1 Történetünk első szála: néhány motiváló példa | 7 |
| 1.2 A példák tanulságai: az empirikus adatok elemzésének legnagyobb problémája | 15 |
| 1.3 A confounding megoldásai: kísérlet és megfigyelés | 17 |
| 1.4 Történetünk második szála: a regressziós modellek | 18 |
| 1.5 Regresszió a sokaságban | 20 |
| 1.6 A szálak összeérnek | 33 |
| 2 Regresszió a mintában: következtetés | 35 |
| 2.1 A hagyományos legkisebb négyzetek (OLS) elve | 35 |
| 2.2 Lineáris regresszió becslése OLS-elven | 37 |
| 3 Kategorialis magyarázó változók | 43 |
| 3.1 Regresszió csak minőségi változóval (ANOVA) | 43 |
| 3.2 Regresszió minőségi és mennyiségi magyarázó változóval (ANCOVA) | 47 |
| 4 Nemlineáris modellek | 53 |
| 4.1 Elöljáróban: a marginális hatás általánosabb értelmezése | 53 |
| 4.2 A linearitás feloldása | 54 |
| 4.3 Néhány nevezetes, paraméterében nemlineáris modell | 61 |
| 4.4 Specifikációs tesztek | 63 |

Előszó

Az ökonometria a társadalmi-gazdasági jelenségek számszerűsített, empirikus – azaz tapasztalati, tényadatokon alapuló – vizsgálatának, modellezésének a tudománya. Szemben azzal, amit esetleg a név sugallhat, közel sem csak közgazdászoknak fontos: szociológusok, társadalomkutatók számára ugyanúgy alapvető az ökonometria ismerete. Sőt, maguk a módszerei még ennél is szélesebb körben, ahol empirikus adatok kezelésére szükség van, biostatistikától a pszichometrián át az agrometriáig, felhasználhatóak. (Ahogy szokták is mondani: a “statisztika egy”.)

Ez a jegyzet ezeket a módszereket tárgyalja, az alapoktól kezdve. Nem célja mély matematikai részletek tárgyalása (noha a korszerű ökonometriára ez nagyon is jellemző), inkább a módszerek, alkalmazási területek, és eszközök sokféleségét kívánja bemutatni. Az elméleti szigorból azonban nem enged.

Manapság ökonometria elképzelhetetlen számítógépes támogatás nélkül. Bár a jegyzetnek nem kifejezett célja ennek megtanítása, de igyekszik hozzá minden segítséget megadni: valamennyi eredmény előállítását is bemutatja a manapság egyre népszerűbb R statisztikai programcsomag alatt. (Az R általános statisztikai programcsomag, és bár klasszikus orientációja nem kimondottan ökonometria, erre a célra is egyre jobban használható, kitűnő általános tulajdonságai és a megjelenő kiegészítő csomagok sokaságának köszönhetően.)

A jegyzettel kapcsolatban minden visszajelzést, véleményt, kritikát a lehető legnagyobb örömmel veszek a tamas.ferenci@medstat.hu email-címen!

A jegyzet weboldala (oktatási segédanyagokkal, technikai információkkal) a https://github.com/tamas-ferenci/FerenciTamas_Okonometria címen érhető el.

1. fejezet

Út a többváltozós regresszióhoz

Egy ilyen tudomány esetén az első feladat annak tisztázása, hogy egyáltalán mi az az ökonometria, mire szolgál, és mi szükség van rá a társadalmi-gazdasági jelenségek elemzése során. A kérdést két történetszálon fogjuk végigvezetni (persze mint minden valamirevaló kortás norvég regényben, a szálak végül össze fognak érni).

1.1 Történetünk első szála: néhány motiváló példa

Elsőként, ahelyett, hogy rögtön a definíciókra térnénk, talán érdekesebb pár példát átbeszélni, melyek már mutatni fogják az ökonometriai vizsgálatok fő problémáit.

1.1.1 Hogyan hat az osztálylétszám a tanulók teljesítményére?

A közoktatásokkal kapcsolatos vizsgálatok egyik klasszikus kérdése, hogy az osztálylétszám hogyan hat a tanulók teljesítményére. Sokan amellett érvelnek, hogy a kisebb létszámú osztályokban több tanári figyelem jut egy diákra, így a tanulók teljesítménye jobb lesz. De vajon tényleg így van?

E kérdésre számos módon válaszolhatunk: felállíthatunk elméleti modelleket, papíron és ceruzával, készíthetünk interjúkat szakértőkkel, vizsgálhatunk analóg helyzeteket más területekről stb.

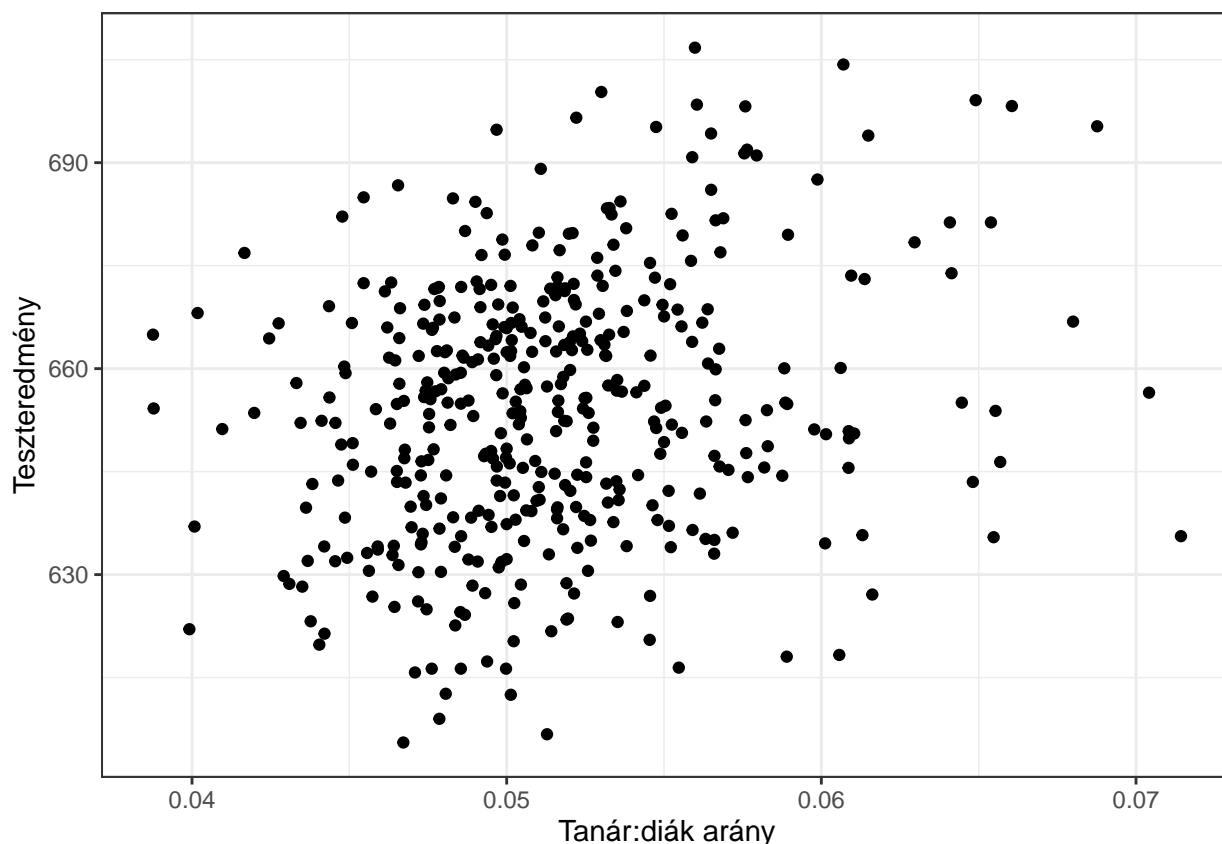
Mi azonban a továbbiakban egy módszerrel fogunk foglalkozni: ha empirikusan igyekszünk válaszolni a kérdésre. Empirikusan, annyi mint a tapasztalatok alap-

ján, tehát való életbeli tényadatok begyűjtével. Elvégre az osztálylétszámokra csak van valamilyen adatgyűjtés, ha az országban futnak standardizált képesség-mérő felmérő-programok, akkor a tanulók teljesítményére is van adatunk – mi lenne, ha

1999-ben Kalifornia állam pontosan ezzel a kérdéssel szembesült. 420 iskolai körzetből gyűjtött adatokat, melyek – számos egyéb mellett – tartalmazták a tanulók és tanárok létszámát, valamint az elért teszteredményeket¹. Az AER csomag CASchools néven tartalmazza a tényleges adatokat. Lássuk is akkor az eredményt! Íme a tanár:diák arányok és a teszteredmények szóródási diagramok szemléltetve:

```
data("CASchools", package = "AER")
CASchools$tsratio <- with(CASchools, teachers/students)
CASchools$score <- with(CASchools, (math + read)/2)
ggplot(CASchools, aes(x = tsratio, y = score)) + geom_point() +
  labs(x = "Tanár:diák arány", y = "Teszteredmény")
```

¹Az adatok igazából nem osztály-szintűek, hanem körzetenkénti átlagok, de ez minket, a mostani kérdésünk szempontjából nem érint, így a továbbiakban az egyszerűség kedvéért osztályt fogok mondani.



Az eredmények első ránézésre megerősítik a sejtésünket: ha több tanár jut egy diákra (kisebbség az osztályok), akkor az jobb teszteredménnyel jár együtt. Nem túl erős az összefüggés, de azért egyértelműen létezik (vájtfülűek kedvéért: $r = 0.23$, $p < 0.001$) és némi munkával még az is kihozható, hogy ezen eredmények szerint ha egy századdal megnöveljük a tanár:diák arányt, akkor várhatóan 8.38 ponttal fog javulni a tesztponyszám.

Ezek az eredmények húsbavágóak. Ha lecsökkentjük az osztálylétszámokat, akkor több tanárt kell alkalmazni, több osztályt kell indítani, adott esetben több osztályteremre lesz szükség, de még az sem kizárt, hogy új iskolaépületre. Pontosán tudni kell tehát, hogy tényleg elérünk-e ezzel valamit. Sőt, valójában ennél többről van szó, azt is tudni kell, hogy *mennyit* érünk el ezzel, egész egyszerűen azért, hogy költség-haszon mérleget lehessen csinálni: a tanárok bérét meg az iskolafelújítások árát megmondják a kontrollerek, de mit rakunk a mérleg másik serpenyőjébe? Ehhez kell tudni a fenti számot, hogy a kettőből együtt meg tudjuk mondani: hány millió dollárt kell költenünk 1 pont teljesítményjavításra. Hogy aztán ez megéri-e, az természetesen már nem statisztikai kérdés, függ a rendelkezésre álló büdzsétől, az egyéb feladatoktól, de még a kormány

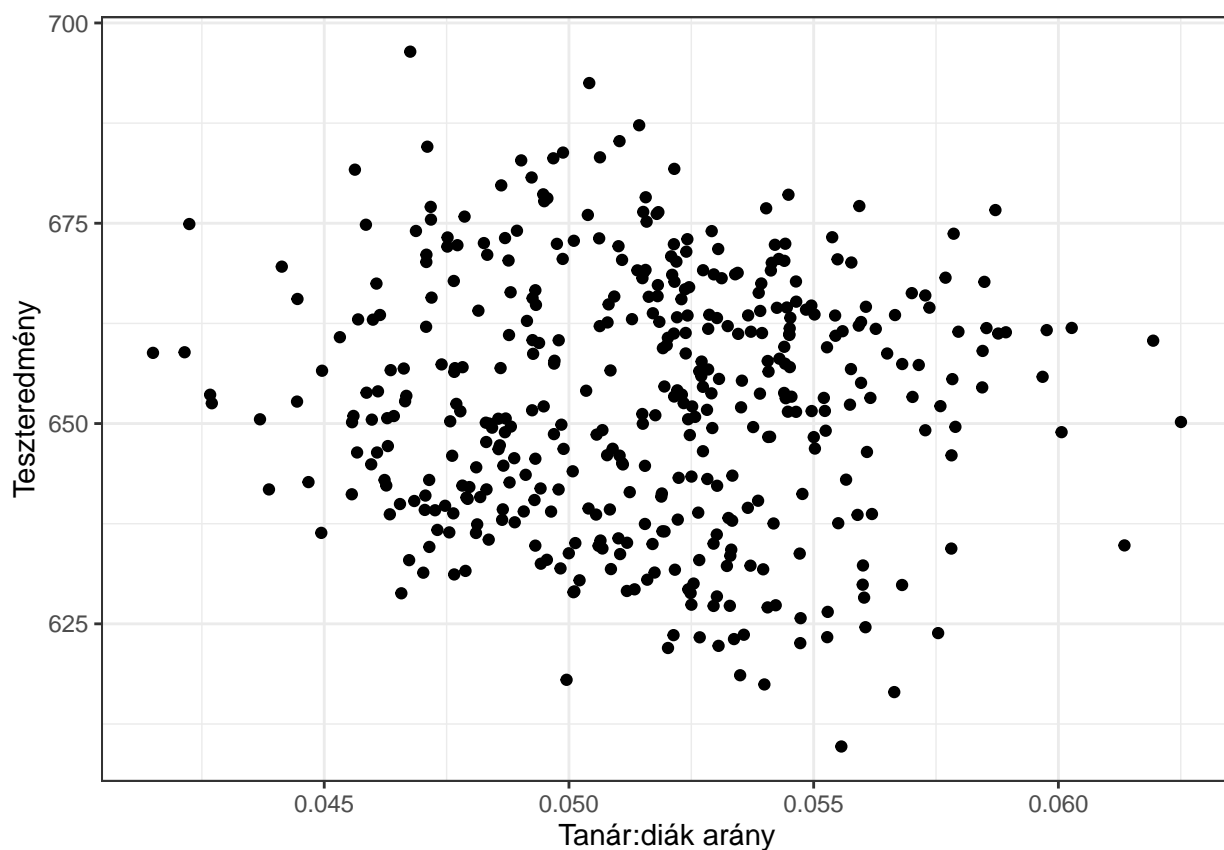
értékválasztásától is, ám a statisztikának kell ezt, mint inputadatot szolgáltatni a döntéshez.

De vajon biztos jól van így minden? Ha az ember elkezd jobban nézni a problémát, esetleg “szociológiai” szemmel is igyekszik ránézni, akkor hamar szöveget üthet a fejében valami. És nem is kell Kaliforniáig menni, magyar, vagy akár még konkrétabban budapesti viszonylatban is ugyanúgy érzékelhető a probléma: ha veszünk kis átlaglétszámú osztályokat (sztereotipikusan mondjuk 2. kerület) és nagy átlaglétszámú osztályokat (sztereotipikusan mondjuk 8. kerület), akkor hihető, hogy az előbbiek teljesítménye jobb, na de álljunk meg egy pillanatra! Csak és kizárólag az osztálylétszám nagyságában térnek el ezek az osztályok egymástól?! Dehogy! Akkor meg honnan tudjuk, hogy a tapasztalt különbség tényleg az osztálylétszámbeli eltérés miatt van...?

A rövid válasz – sajnos – az, hogy sehonnan! És itt van a bökkenő: igaz, hogy a 2. kerületi osztályok kisebbek mint a 8. kerületiek, de *együttal* másban is eltérnek, az oda járó gyerekek szocioökonómiai háttere tendenciájában jobb, tanulásra motiválódóbb otthoni környezetből érkeztek, a szülők anyagilag is megengedhetik maguknak, hogy a gyermekeiket különóra járassák stb. E ponton viszont nagy baj van: innen kezdve fogalmunk sem lehet, hogy a tapasztalt különbség *tényleg* a kisebb osztálylétszám miatt van, vagy esetleg az osztálylétszámnak a világon semmi hatása nincs, csak egyszerűen a kisebb osztályokba jobb szocioökonómiai helyzetű diákok járnak és *ez* a valódi oka az ott tapasztalt jobb tesztpontszámoknak?

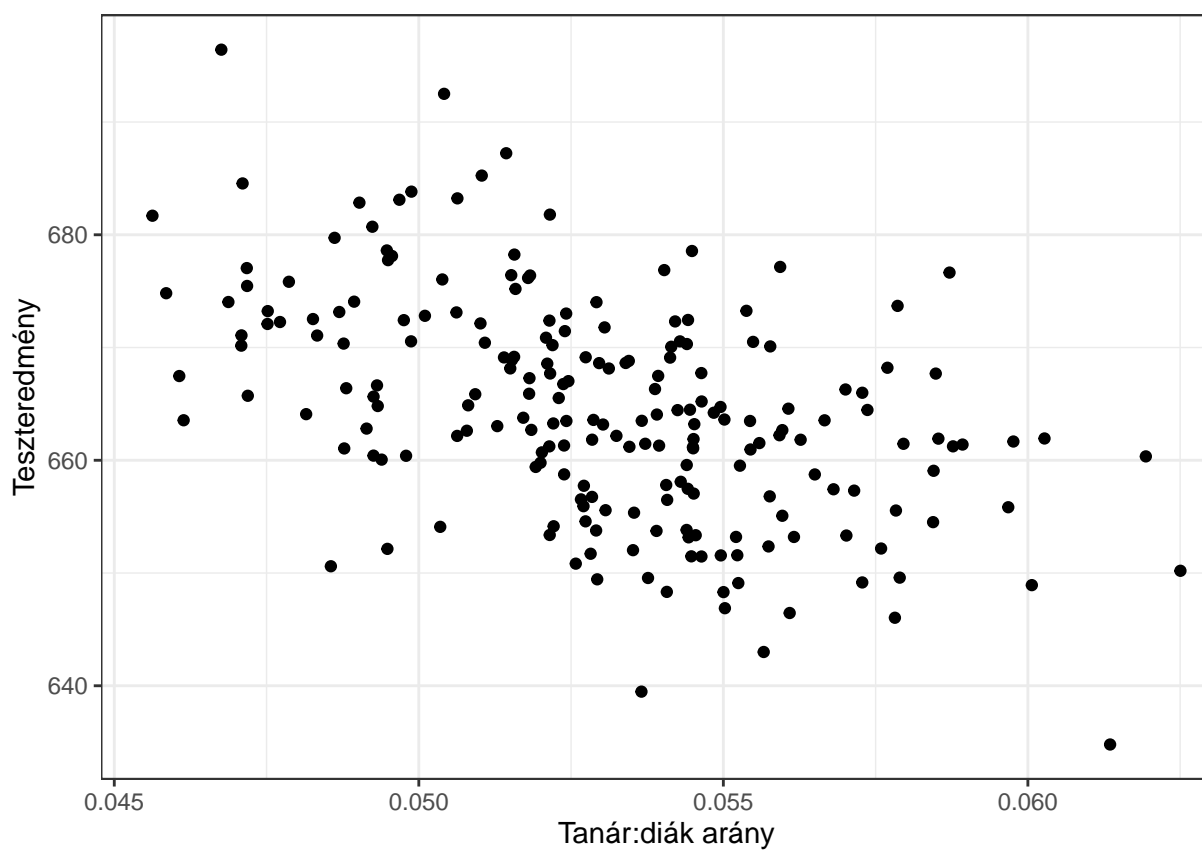
Sőt! Innentől kezdve még akár az is elképzelhető, hogy a kisebb osztálylétszám igazából kifejezetten *ront* a teljesítményen *önmagában*, csak épp a kisebb osztályokba annyival jobb szociális helyzetű diákok járnak, hogy az átfordítja a helyzetet.

Valaki nem hiszi el, hogy ez még is lehetséges? Nos, gyártsunk egy egyszerű szimulációt! Egyelőre nem árulom el, hogy hogyan készítettem, de íme a végeredménye:

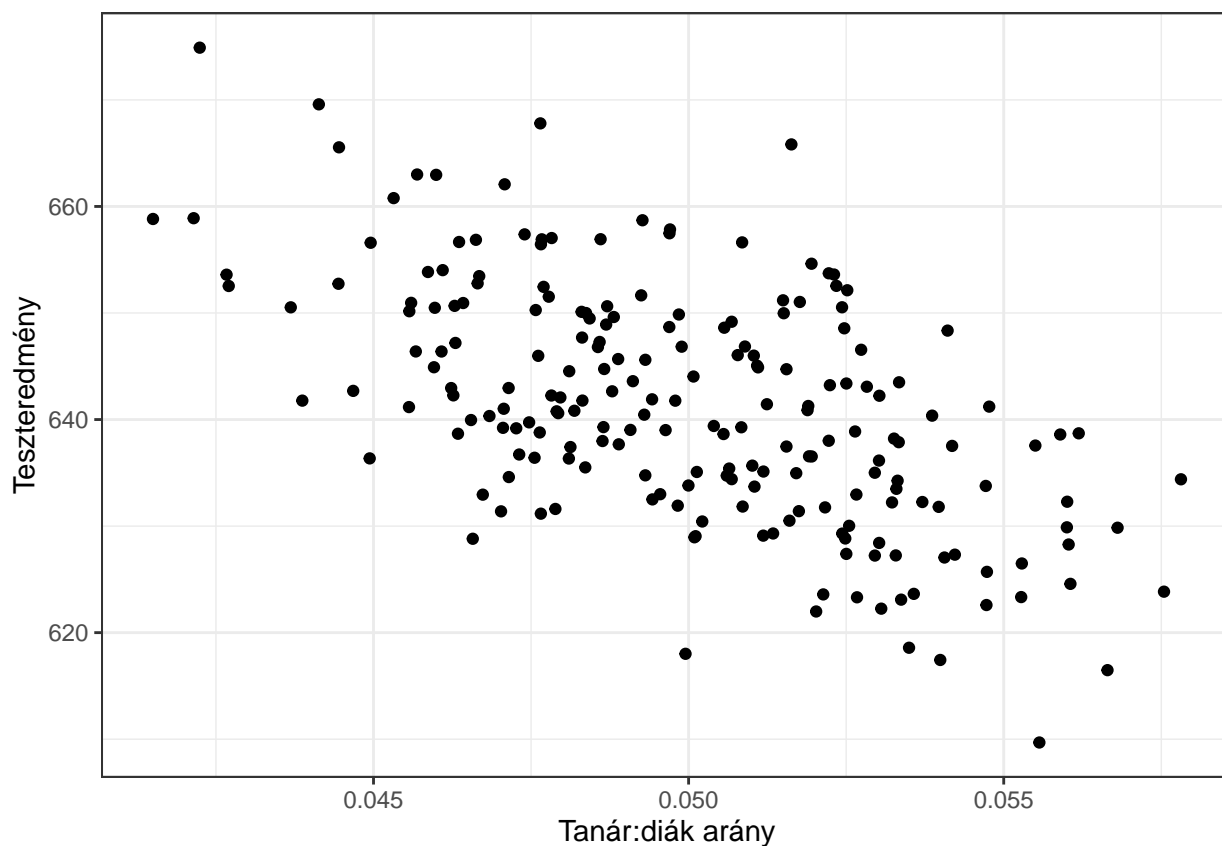


Első ránézésre nagyjából megfelel a korábbi képnek. De nézzük csak meg jobban mi történik itt! (Mivel ezt saját kezűleg generáltuk, így megtehetjük, hiszen tudjuk mi van a valóságban, az adatok háttérében). Az egyszerűség kedvéért mondjuk, hogy a szocioökonómiai státusz egy bináris változó, “jó” és “rossz” a két lehetséges értéke.

Nézzük a tanár:diák arány és a pontszám összefüggését a jó szocioökonómiai státuszú osztályok körében:

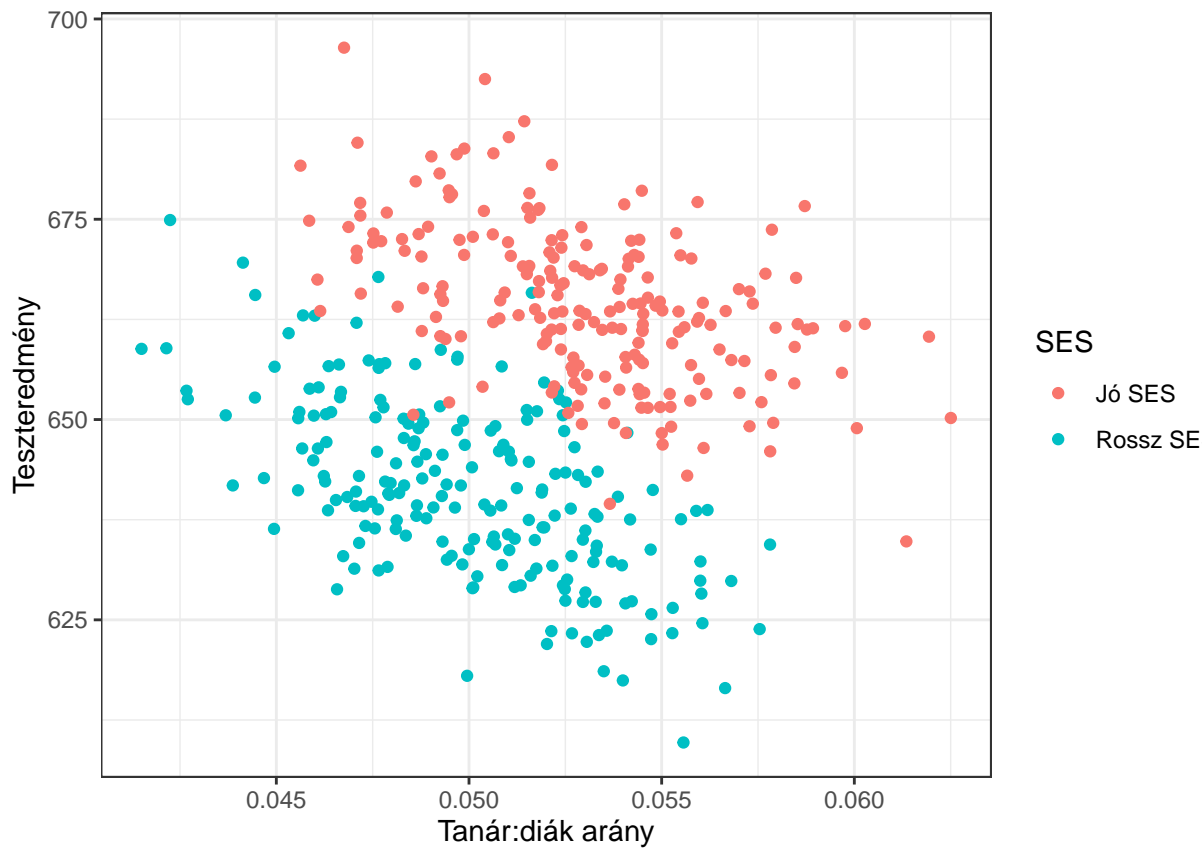


Érdekes! Itt fordított a kapcsolat. Na de mi a helyzet a rossz szocioökonómiai státuszú osztályokban? Íme:



Itt is negatív a kapcsolat!

Ez meg hogy a viharban lehet? – kérdezhetné valaki. A rossz szociális helyzetű csoporton belül is ront a kisebb osztálylétszám, a jó helyzetűeken belül is ront, de összességében meg javít?! Egyből világosabb a helyzet, ha egy ábrán ábrázoljuk a kettőt, csak eltérő színekkel:



Azonnal érthető, hogy mi történik, ha hozzávesszük a korábban mondottakat: a jobb tanár:diák arány igazából *ront* a helyzeten, de a jobb szociális helyzetű diákok által alkotott osztályok egyszerre kisebbek és – szociális helyzetük, nem az osztálylétszám miatt! – jobb teljesítményűek, és ez olyan erős effektus, hogy ha egyben vizsgáljuk az osztályokat, akkor a kisebb létszám rontó hatását átbillenti, hogy a kisebb osztályok a jobb szociális helyzetük miatt jobb teljesítményűek. Így összességében azt fogjuk látni, ahonnan indultunk: hogy a kisebb létszámú osztályok jobb teljesítményűek.

A végeredmény tehát: a kisebb osztálylétszám rontja a teljesítményt és a kisebb létszámú osztályok jobb teljesítményűek. És ha valaki érti a fentieket, akkor azt is érti, hogy ebben a mondatban miért nincs semmi ellentmondás!

1.1.2 Csökkenti-e a korrupció mértékét a nők részvétele a politikában?

TODO

1.1.3 Csalnak-e az orosz választásokon?

TODO

1.1.4 További példák

Hosszasan sorolhatóak a további, hasonló példák a társadalmi-gazdasági elemzések világából. Kommentár nélkül még néhány kérdés, érdemes mindegyiket a fenti példákból leszűrni, kezdődő tanulságok szemüvegén keresztül végiggondolni:

- Hogyan hat a munkanélküliség a GDP-re?
- Hogyan hat az államadósság a növekedésre?
- Return on education: mekkora az oktatás haszna, tehát, ha egy évvel többet tölt valaki az iskolapadban, az mennyivel növeli a fizetését?
- Létezik-e “cigánybűnözés”?
- Az ökonometria-előadás haszna: ha többet tölt a hallgató az öko előadáson, jobb jegyet kap-e emiatt, és ha igen, mennyivel?
- Milyen tényezők hatnak arra, hogy egy országban hány terrortámadás történik?
- Hogyan hat a rendőri erők létszáma egy adott városban az ottani bűnözési rátákra?
- Cégeknek adott továbbképzési támogatás hogyan hat a termelékenységre?

(Igen, ezekre mind válaszolhatunk ökonometriai módszerekkel!)

1.2 A példák tanulságai: az empirikus adatok elemzésének legnagyobb problémája

A mintázat most már látszik. Valamilyen tényező hatására vagyunk kíváncsiak, az okozati hatására, szép szóval: a **kauzalításra**, ez mindegyik példa lelke.

Úgy döntünk, hogy a kérdést **empirikus** adatok elemzésével igyekszünk megoldani, azaz való életbeli tényadatokat gyűjtünk be. (A vizsgált tényezőről, a hatásról, esetleg egyéb fontos változókról.)

Azért, hogy eldöntsük, hogy van-e okozati hatás (illetve, hogy lemérjük mekkora), csoportokat hasonlítunk össze, melyek eltérnek a vizsgált tényezőben. Csak épp közben a vizsgált tényezőbeli eltéréssel *automatikusán együtt járnak* egyéb tényezőbeli eltérések, és innentől kezdve bajban vagyunk, mert ha találunk is különbséget a csoportok között, nem fogjuk tudni, hogy az mi miatt van: a vizsgált tényezőbeli eltérés miatt, a vele együtt járó egyéb eltérések miatt, vagy esetleg ezek valamilyen keveréke miatt. Ezt a jelenséget, mely az empirikus adatok elemzésének legnagyobb problémája, hívják **confoundingnak**. (Angolul nagyon jó szó, amire nem sikerült hasonlóan találó magyar fordítást bevezetni. A confounding azt jelenti, hogy összemosódás, és csakugyan, az a probléma, hogy a vizsgált tényezőbeli eltérés összemosódik egyéb tényezőbeli eltérésekkel.)

Vegyük észre, hogy a confounding fellépéséhez az kellett, hogy létezzen olyan tényező, amire két dolog *egyszerre* igaz: együttmozog a vizsgált tényezővel (összefügg vele) és önmagában – azaz a vizsgált tényező minden értéke mellett – hat az eredményváltozóra. Akkor van confounding, ha ez a kettő egyidejűleg fennáll. Ha bármelyik nincs jelen, akkor nincs probléma. Ha a szociális helyzet összefügg ugyan az osztálylétszámmal, de nem hat a teljesítményre, akkor nincs baj: igaz, hogy a kisebb osztályok jobb szociális helyzetűek, de ez nem befolyásolja a teljesítményt. Hasonlóképp, ha a szociális helyzet befolyásolja ugyan a teljesítményt, de nem függ össze az osztálymérettel, akkor sincs gond: a kisebb osztályoknak nem tér el a szociális helyzete a nagobbaktól. (Gondoljuk végig az összes többi példára is!) Azokat a változókat, amelyek ezt a két dolgot egyidejűleg tudják, tehát a vizsgált tényezővel együttmozognak és az eredményváltozóra is hatnak, ilyen módon okozzák a confoundingot, szokás zavaró változónak, vagy **confoundernek** nevezni.

A későbbiek szempontjából hasznos lesz ezt még másképp átfogalmazni. A probléma, hogy nem az érdekel minket, hogy egy osztály abban tér el, hogy kisebb a létszám, akkor ott jobb-e a teljesítmény, hanem az, hogy ha *csak* abban tér el, hogy kisebb a létszám, akkor ott jobb-e a teljesítmény. Ezt szokás **ceteris paribus elvnek** (“minden mást változatlanul tartva”) nevezni, ez a kulcs a kauzalitáshoz: az érdekel minket, hogy ha minden mást változatlanul tartva *csak* az osztálylétszám változik, akkor mi történik. A naiv elemzésben az osztálylétszám változásával együtt egyéb tényezők is változhatnak, így ebből nem tudunk a kauzalitásra következtetni. Figyeljünk a szóhasználatra: azt mondhatjuk, hogy a kisebb osztálylétszám jobb teljesítménnyel jár *együtt* (korreláció), de azt nem, hogy a kisebb osztálylétszám jobb teljesítményt *okoz* (kauzalitás).

Valaki esetleg azt mondhatja, hogy rendben, itt tényleg van valami módszertani gubanc, meg szép latin szavak² de igazából ez csak az ilyen módszertani kérdéseken szöszmötölő kutatóknak érdekes, a lényeg, hogy ha kisebb az osztálylétszám, akkor ott jobb a teljesítmény, ennyi a fontos, és pont. Nem! Ez az érvelés teljesen fals, az hogy mi hat mire, nem tudományos szórszálhasogatás, hanem elsőrendű gyakorlati kérdés. Miért? A beavatkozás miatt! A *valódi* okozatiság felismerése ott válik kritikussá, ha beavatkozunk a rendszerbe, ha ugyanis rosszul állapítjuk meg az okozati kapcsolatok irányát, akkor ez teljesen félremehet. Például lecsökkentjük az osztálylétszámokat, adott esetben rengeteg pénzt elkölthetve, de ha a *valódi* oka a jobb teljesítménynek nem az osztálylétszám, hanem a jobb szociális helyzet, akkor ezzel semmit nem érünk el! Sőt, mint a későbbi példa mutatja, adott esetben még kimondottan árthatunk is!

Végezetül még egy megjegyzés. A confounding felismerése nem azt jelenti, hogy akkor igazából nincs hatás, végképp nem azt, hogy bizonyítottuk, hogy ellentétes irányú hatás van. Pusztán annyit jelent, hogy a confounding-gal terhelt adatok *nagyon gyenge* bizonyítékot jelentenek a hatás léte mellett. De ettől még lehet éppenséggel hatás! – csak az ilyen adatok nagyon kevésbé támasztják ezt alá.

²Pedig még össze sem foglaltam a fentieket úgy, hogy a korreláció nem implikál kauzalitást!

1.3 A confounding megoldásai: kísérlet és megfigyelés

Most, hogy alaposan kiveséztük a confounding problémáját, természetesen adódik a kérdés: na de mit tehetünk ez ellen? Azt könnyű lenne biztosítani, hogy a csoportok egy-két általunk megadott szempont szerint ne legyenek eltérőek, de azt, hogy *egyáltalán semmilyen* szempont szerint ne térjenek el (kivéve persze az vizsgált tényezőt), olyanok szerint sem, amikről eszünkbe sem jut, hogy lehet bennük eltérés, csak egy módon lehet: ez a **randomizálás**. Ahogy a szó is sugallja, a randomizálás lényege, hogy a megfigyelési egységeket *véletlenszerűen* sorsoljuk különböző csoportokba, majd ezeket a csoportokat tesszük ki a vizsgált tényezőnek. Például pénzfeldobással döntjük el az óvoda végén minden egyes gyermekről, hogy kis vagy nagy létszámú osztályba kerüljön. Ez azért jó, mert ilyen módon a két csoport között nem lesz szisztematikus különbség szocioökonomiai státuszban, de ami még fontosabb: *semmilyen* tényezőben nem lesz szisztematikus különbség, a kék szeműek vagy a balkezesek számában sem, hiszen a pénzfeldobás nyilván ezekre is érzéketlen. Ilyen módon a csoportok összehasonlíthatóak: ha találunk köztük különbséget a tanulmányi eredményben, az *tényleg* az osztálylétszámnak lesz betudható... hiszen másban nincs szisztematikus különbség.

A randomizálásnak egy baja van: akkor alkalmazható, ha a vizsgált tényezőt tudjuk irányítani. (Hiszen nekünk kell az egyik csoportba sorsolt gyerekeket kis, a másikat nagy létszámú osztályba helyezni.) Azokat a kutatásokat, ahol a kutatást végzők tudják irányítani a vizsgált tényezőt, **kísérletes (experimentális) kutatásnak** nevezzük. És itt értünk el a bökkenőhöz: a társadalmi-gazdasági jelenségek vizsgálata az a terület, ahol tipikusan *nem* lehet kísérletet végezni. Aligha lehet gyerekeket pénzfeldobással sorsolni osztályokba, vagy országokban pénzfeldobással meghatározni, hogy mennyi nő üljön a kormányban...

(Persze ez sincs kőbe vésve. Néha lehet kísérletet csinálni, ahogy a választási megfigyelők példája is mutatja. Másik oldalról, például az orvostudományban sokszor lehet kísérletet csinálni, ez a jellemző új gyógyszerek bevezetésénél, ahol a vizsgálat során véletlenszerűen kiválasztott alanyok kapnak gyógyszert, míg a többiek placebo, de ott is van olyan kérdés, ahol nem lehet kísérletet csinálni! Tipikusan ilyenek fordulnak elő az epidemiológiában: a vörös hús rákkeltő? Aligha lehet emberekkel pénzfeldobás alapján évtizedekig több vagy kevesebb vörös húst etetni... Innentől a probléma ott is ugyanaz: ha a vörös húst evők között több a rákos, az nagyon gyenge bizonyíték, mert a több vörös húst fogyasztó emberek milliónyi *egyéb* dologban is eltérnek a kevesebb vörös húst fogyasztó emberektől a vörös hús fogyasztás mértékén túl – és mi van, ha ezek közül valami növeli a rákkockázatot...?)

Azokat a vizsgálatokat, ahol a kutatást végzők nem tudják befolyásolni a vizsgált tényezőt, az alakul a maga rendje szerint, és a kutatók csak passzíve feljegyzik a történéseket külső szemlélőként, **megfigyeléses (obszervációs) vizsgálatnak** nevezzük. A társadalmi-gazdasági elemzések során tehát szinte mindig ilyenekkel

lesz dolgunk. Márpedig ezeknél mindig fejünk felett fog lebegni a confounding problémája.

1.4 Történetünk második szála: a regressziós modellek

Folytassuk most valami – látszólag – teljesen más témával.

Minden fenti példában volt egy változó, mely az eredménye volt a vizsgálatunknak, a kimenet szerepét játszotta, tehát aminek az alakulását le kívántuk írni (tesztpontszám, korrupció mértéke, szavazati arány stb.). A továbbiakban ezt **eredményváltozónak** (vagy függő változónak, angolul response) fogjuk hívni, jele Y . Az első példánkban $Y = \text{Teszteredmény}$. Másrészt voltak változók, adott esetben nem is egy, amikkel le akarjuk írni az eredményváltozó alakulását, amelyekről azt mondjuk, hogy hatnak, vagy hathatnak az eredményváltozóra; ezek neve **magyarázó változó** (vagy független változó, angolul predictor). Ezekből több is lehet (az első példában ilyen az osztálylétszám és a szocioökonómiai helyzet), jelöljük számukat k -val, és az egyes változókat X_i -vel ($i = 1, 2, \dots, k$). Az első példában $k = 2$ és $X_1 = \text{Tanár:diák arány}$, $X_2 = \text{Szocioökonómiai státusz}$. Összefoglalva, az eredményváltozó a vizsgált kimenet, a magyarázó változók az azt – potenciálisan – befolyásoló tényezők (tehát a fontos, vizsgált változók és a – potenciális – confounderek egyaránt).

Az X -ek hatnak az Y -ra, vagy fordítva megfogalmazva, az Y függ az X -ektől – ragadjuk meg most ezt matematikailag. Szerencsére arra, hogy egy változó függ más változóktól, ismerünk egy jó matematikai objektumot, ez a függvény fogalma:

$$Y = f(X_1, X_2, \dots, X_k).$$

A későbbiekben erre azt fogjuk mondani, hogy ez egy statisztikai modell. Ennek az általánosságával nehéz lenne vitatkozni, de egy baja mégis csak van.

A fő probléma, hogy a modell azt feltételezi, hogy az Y és az X -ek kapcsolata **determinisztikus**. Szinte teljesen mindegy is, hogy mi az Y és mik az X -ek, hogy mi a vizsgált probléma, a társadalmi-gazdasági jelenségek vizsgálata kapcsán lényegében általánosan kijelenthető, hogy ez irreális: bármilyen ügyesek vagyunk, soha az életben nem fogunk tudni determinisztikus modelleket alkotni társadalmi-gazdasági jelenségekre. (Aligha lehet olyan modellt alkotni, ami *pontosan*, hiba nélkül megmondja előre, hogy egy osztály milyen pontszámot fog elérni, vagy, hogy egy választáson pontosan hány szavazat érkezik egy pártra.) Ez legfeljebb középiskolás fizikában működik, a társadalmi-gazdasági jelenségekben szinte kizárt, hogy *függvényszerű* módon meghatározzák a magyarázó változók az eredményváltozót. Hiszen lesznek változók amiket nem ismerünk, rosszul mérünk, rosszul veszünk figyelembe, az, hogy egy gyerek hány pontot ír egy

1.4. TÖRTÉNETÜNK MÁSODIK SZÁLA: A REGRESSZIÓS MODELLEK¹⁹

teszten, mindig függ a mi közelítési szintünk ténylegesen véletlen dolgoktól stb. A valódi modell tehát **sztochasztikus** kell legyen:

$$Y = f(X_1, X_2, \dots, X_k) + \varepsilon$$

Itt ε jelzi a fentiekből fakadó bizonytalanságot, a neve: **hibatag**.

Rövid jelölésként az X -eket gyakran egy vektorba vonjuk össze: $Y = f(\underline{X}) + \varepsilon$.

Az így kapott modellünk már egy teljes értékű statisztikai (ökonometriai) modell! Az ilyen f -et hívjuk (sokasági) **regressziófüggvénynek**.

Már most is fontos, hogy lássuk, hogy az f -nek van egy nagyon is földhözragadt értelmezése: ezt kell használni, ha szeretném “megtippelni” Y értékét X -ek ismeretében. Hogy ez miért lesz fontos, azt majd később fogjuk látni, de a feladat így is értelmes: ha ismerem egy osztály tanár:diák arányát és a szociális helyzetet, akkor ezek alapján mit mondhatok a teszteredményéről. Ha ezek tényleg hatnak rá, akkor valamit mondhatok, ezt fejezi ki az f -es rész, ε pedig azt, amit nem tudok ezek alapján megmondani, azaz ettől lesz ez csak tipp: mindenképp kell számolnom azzal, hogy a valóság ettől a fent említett okok miatt eltér.

Ez az egyenlet egy **sokasági modell**: azt írja le, hogy a valóság hogyan működik. Pontosan ugyanaz a helyzet, mint bármilyen következtető statisztikai kurzus alapjainál: van a sokaság, amit eloszlásokkal, valószínűségszámítási eszközökkel írunk le, de a tényleges vizsgálatokban mi sem ismerjük. (Tehát nem tudjuk, hogy ezek az eloszlások milyenek.) Ahhoz, hogy megismerjük, veszünk egy mintát, ennek a kezeléséhez már statisztika kell, aminek a feladata épp az lesz, hogy következtessünk a sokaságra.

Most is hasonló a helyzet: mi sem tudhatjuk, hogy milyen *eloszlása* van a teszteredményeknek, csak van egy 420 elemű mintánk rá nézve; és hasonlóan a többi változóval. Most azonban egy pillanatig leszünk valszámos emberek: ne törődjünk azzal a problémával, hogy a sokaságot igazából nem ismerhetjük, játszunk azt, hogy ismerjük (tudjuk mik ezek az eloszlások), és vizsgáljuk meg, hogy ebből mire jutunk! Ugyanúgy mint a következtető statisztikánál, ez nagyon hasznos lesz majd később, a számunkra igazán érdekes – statisztikai – feladat megoldásánál is.

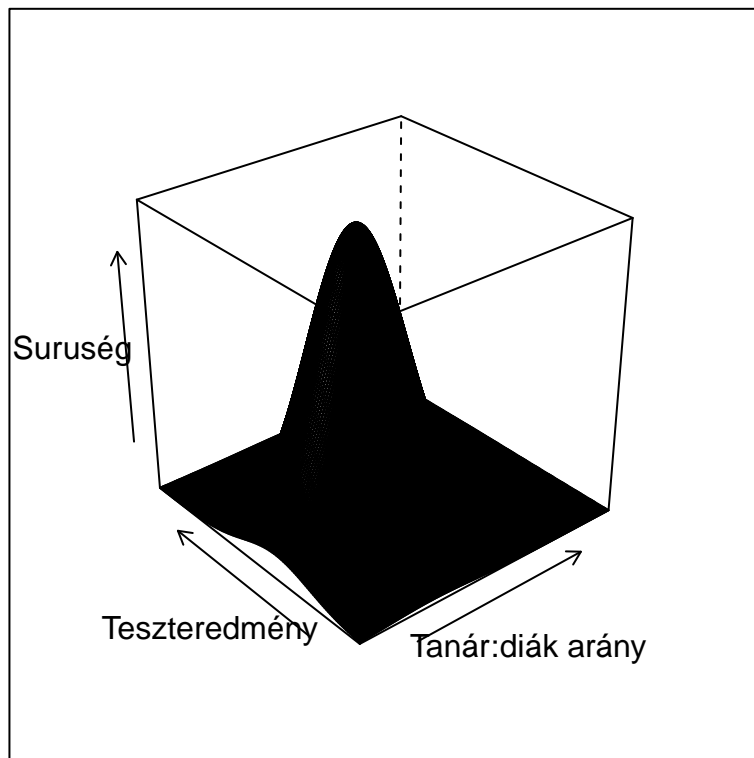
A nem-kísérleti jelleg miatt az az értelmes modell, ha mind az eredményváltozót, mind a magyarázó változókat – és így persze ε -t is – valószínűségi változónak vesszük. (Ezért használtam eddig is nagy betűket!) Bizonyos egyszerűsített tárgyalások úgy tekintik, mintha az X -ek nem valószínűségi változók lennének, hanem rögzített értékek. Ez a kísérletek világában rendben lehet, ahol mi be tudjuk *állítani* az X -ek értékét, de ökonometriában, a társadalmi-gazdasági elemzések világában még közelítő feltevésként is értelmetlen.

1.5 Regresszió a sokaságban

Elsőként tehát le kell írunk a sokaságot: valszámos emberek leszünk, és úgy vesszük mintha ismernénk a sokaságot. Mit jelent ez, mit is ismerünk pontosan? Nem csak Y és X_1, X_2, \dots, X_k eloszlásait (külön-külön), hanem az együttes eloszlásukat is! Ekkor tudunk mindent ezekről (valószínűség-számítási értelemben).

Ezt úgy kell elképzelnünk, mint egy $k + 1$ dimenziós teret: minden pont egy adott magyarázó- és eredményváltozó-kombináció. E fölött értelmezve van egy eloszlás, ami azt mutatja, hogy ha mintát veszünk ebből az eloszlásból, akkor milyen valószínűséggel esünk az adott pont kis környékére.

$k + 1$ dimenziós terekben a legtöbb ember relatíve rosszul tájékozódik, úgyhogy ábrázoljunk egy olyan együttes eloszlást, amikor még átlátható a dolog!



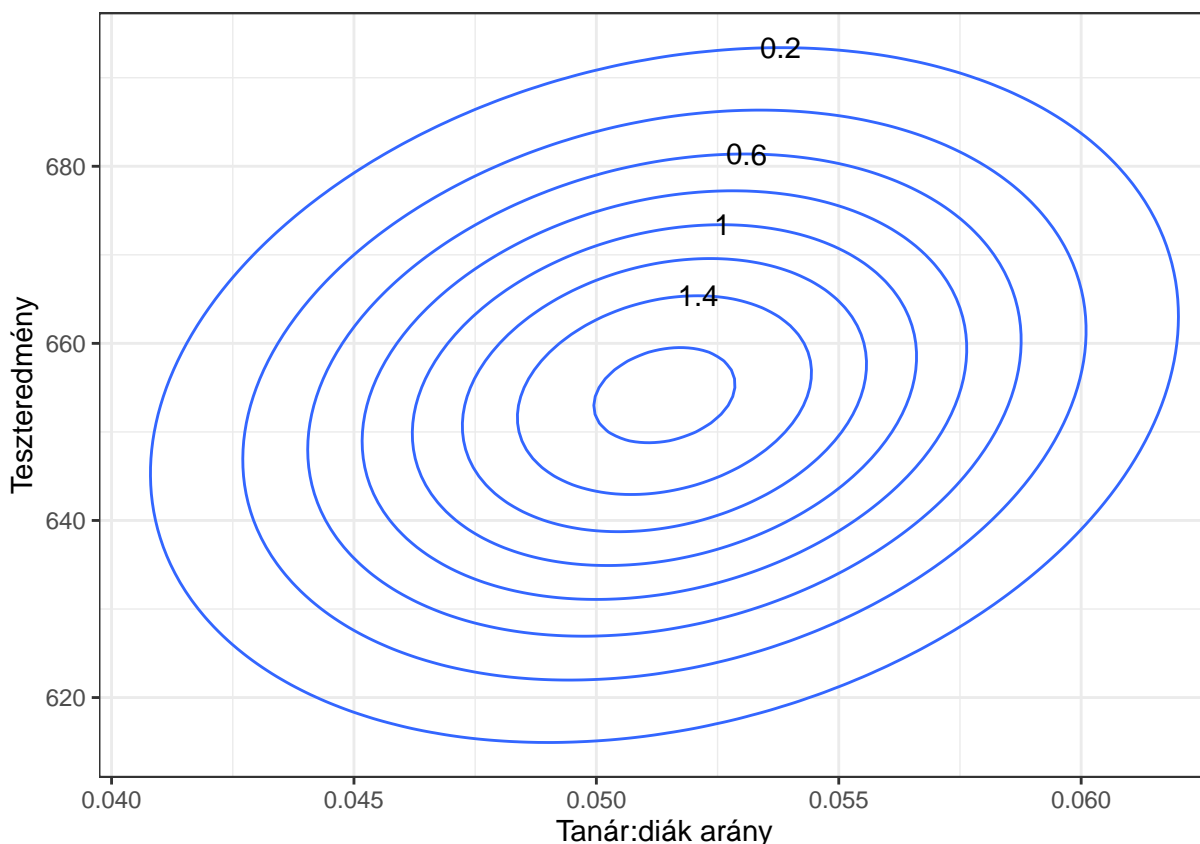
Ez egy kétváltozós eloszlás együttes sűrűségfüggvénye; itt az egyik változó játssza a magyarázó-, a másik az eredményváltozó szerepét. A mintavétel ebből az eloszlásból azt jelenti, hogy kivesszünk egy iskolát (tehát tanár:diák arányt

és teszteredményt egyszerre!); ahol magasan fut a sűrűségfüggvény, arról a környékről gyakran veszünk ki, ahol alacsonyan, ott ritkábban. A hiba mibenléte is jól érthető erről az ábráról: ha kiválasztunk egy adott konkrét X_1 -et, ahhoz csak egyetlen $f(X_1)$ -et adhatunk, mégis Y minden értéket felvehet, tehát lehetetlen, hogy ne hibázzunk. (Csak egyetlen egy pont lesz a végtelen sok közül, ahol nem hibázzunk.) Persze $f(X_1)$ -et majd pont úgy lesz célszerű megválasztani, hogy oda rakjuk, ahol Y gyakran előfordul, hogy a gyakran előforduló esetekben hibázzunk picit, és csak a ritkábbakban nagyobbab – de erről majd kicsit később.

Elárulom, hogy ez az eloszlás többváltozós normális (később ennek majd jelentősége lesz), $\boldsymbol{\mu} = \begin{pmatrix} 654 \\ 0.0514 \end{pmatrix}$ várhatóérték-vektorral és $\mathbf{C} = \begin{pmatrix} 19,1^2 & 0,23 \cdot 19,1 \cdot 0,00515 \\ 0,23 \cdot 19,1 \cdot 0,00515 & 0,00515^2 \end{pmatrix}$ kovariancia-mátrixszal³.

Sajnos ez az ábrázolás nehezen érzékelhető (pláne, ha nem interaktívan nézzük, és nincs módunk forgatni), jobban járunk, ha így rajzoljuk ki:

³Egyszerűen úgy választottam a paramétereket, hogy megfeleljen a kaliforniai példának



Ez ugyanaz mint a fenti sűrűségfüggvény, de “szintvonalakkal” leírva (azaz különböző z magasságokban elmetstettük a sűrűségfüggvényt és a kapott metszeteket ábrázoltuk). Belátható, hogy többváltozós normális esetén ezek mindig ellipszisek⁴. Ezt az ábrázolást szokás “contour plot”-nak nevezni, előnye, hogy – a háromdimenziós érzékeltetéssel szemben – nem érzékeny a nézőpont megválasztására, részek nem takarnak ki másokat stb. (Ám cserében nyilván információvesztéssel jár, ami azzal arányos, hogy milyen sűrűn képezzük a metszeteket.)

Térjünk most vissza az alapkérdésünkre! Úgy vesszük, hogy ez az eloszlás adott, és le akarjuk írni mint $Y = f(\underline{X}) + \varepsilon$; de vajon mi f -re a legjobb választás? Persze egy ilyen kérdést hallva azonnal vissza kell kérdezni: mi a jószág mérőszáma? Hiszen csak ennek ismeretében mondható meg, hogy mi az optimális sokasági regressziófüggvény.

Mivel az ε hibát fejez ki, így azzal valószínűleg kevesen vitatkoznának, hogy az

⁴Úgy, hogy az ellipszis középpontját a várhatóérték-vektor adja meg, a tengelyek a kovariancia-mátrix sajátvektorainak irányába mutatnak, féltengelyeik hossza pedig a kovariancia-mátrix megfelelő sajátértékeivel arányos.

a legjobb f , amely mellett a hiba a legkisebb. Igen ám, de mi az, hogy a hiba a “legkisebb”? Ez nem olyan nyilvánvaló, ennek megértéséhez beszéljünk egy picit a hibáról. A helyzetet a fenti példán úgy kell elképzelnünk, hogy behúzzuk a $f(X_1)$ függvényt; ez olyan mintha rajzolnánk egy görbét az $X_1 - Y$ síkra. Ez után végigmegyünk a sík minden pontján, és megnézzük ott mekkora a hiba: mennyire van távol Y az $f(X_1)$ -től; ez pedig akkora súllyal fog szerepet játszani a hiba eloszlásában, amilyen magasan fut az adott ponton a sűrűségfüggvény. Mindezek természetesen ugyanígy működnek az általános, $k + 1$ dimenziós esetben is.

A hibának tehát egy eloszlása van, így nem egyértelmű, hogy mikor a “legkisebb”. Két dolgot kell tennünk, az egyik választás egyértelmű, de a másik már inkább döntés kérdése. Az első, hogy a hiba helyett annak $\mathbb{E}\varepsilon$ várható értékét tekintjük. Ez jó, mert így a valószínűségi változóból rögtön egy számot kapunk, amire pedig azonnal jobban értjük, hogy mit jelent az, hogy legyen a legkisebb. De igazán azért jó, mert ha összekombináljuk a várható érték fogalmát az előbbi bekezdés végén mondottakkal, akkor látjuk, hogy ez egy nagyon logikus dolgot mond: azt, hogy ott kevésbé számít a hibázás, ahová egyébként is ritkán esünk, és ott számít jobban a hibázás, ami gyakran előfordul!

Azonban még nem végeztük. Ha meggondoljuk, akkor rögtön látjuk, hogy $\mathbb{E}\varepsilon$ még nem lesz jó: a hiba lehet negatív is és pozitív is, de mi⁵ nem mondhatjuk azt, hogy ha egyszer 10-zel fölé lőttünk, egyszer meg 10-zel alá, akkor tökéletesek voltunk. Magyarán: meg kell szabadulni az előjeltől. Itt már van választási lehetőségünk, hogy mit teszünk, most döntünk úgy (és jelen jegyzet túlnyomó többségében ezt adottságnak fogjuk venni), hogy négyzetre emeléssel szabadulunk meg az előjeltől, hiszen a négyzetre emelés függvény tulajdonságai nagyon kellemesek.

Így tehát a megoldandó feladat:

$$\operatorname{argmin}_f \mathbb{E} [Y - f(\underline{X})]^2.$$

Ez első ránézésre nagyon is ijesztően néz ki: optimalizációs feladat – az *összes létező* függvény terében?! Mert azt még érti az ember, hogy x felveszi az összes lehetséges valós számot, és mikor lesz $f(x)$ minimális, na de mi az, hogy valami felveszi az összes létező (k -változós) függvényt...? Hiszen semmi más megközelítés nincs a világon, *akármilyen* k -változós függvény szóba jöhet, semmit nem mondtunk a függvényformáról, összeadhatjuk a változókat, összeszorozhatjuk, hatványozhatjuk, bármilyen műveletet végezhetünk, bármilyen konstans beletelhetünk, és az összes ilyen közül mondjuk meg, hogy ez a kifejezés mikor lesz a legkisebb?!

Az érdekes az, hogy bármilyen abszurdan is néz ki, a dolognak van megoldása! Ráadásul a végeredmény nem is túl bonyolult: f legjobb megválasztása adott pontban Y feltételes várható értéke lesz a kérdéses pontban:

⁵A statisztikusokról szóló viccekkel szemben.

$$f_{\text{opt}}(\mathbf{x}) = \mathbb{E}(Y \mid \underline{X} = \mathbf{x}).$$

Bizonyítsuk is be ezt! Legyen f_{opt} a feltételes várható érték, f pedig egy tetszőleges k -változós függvényt. Alakítsuk át a kritériumfüggvényt:

$$\begin{aligned} \mathbb{E}[Y - f(\underline{X})]^2 &= \mathbb{E}[Y - f_{\text{opt}}(\underline{X}) + f_{\text{opt}}(\underline{X}) - f(\underline{X})]^2 = \\ &= \mathbb{E}[Y - f_{\text{opt}}(\underline{X})]^2 + \mathbb{E}\{[Y - f_{\text{opt}}(\underline{X})][f_{\text{opt}}(\underline{X}) - f(\underline{X})]\} + \\ &+ \mathbb{E}[f_{\text{opt}}(\underline{X}) - f(\underline{X})]^2. \end{aligned}$$

A középső tag szerencsére nulla, ezt toronyszabállyal láthatjuk be:

$$\begin{aligned} \mathbb{E}\{[Y - f_{\text{opt}}(\underline{X})][f_{\text{opt}}(\underline{X}) - f(\underline{X})]\} &= \\ = \mathbb{E}\{\mathbb{E}\{[Y - f_{\text{opt}}(\underline{X})][f_{\text{opt}}(\underline{X}) - f(\underline{X})] \mid \underline{X}\}\} &= \\ = \mathbb{E}\{[f_{\text{opt}}(\underline{X}) - f_{\text{opt}}(\underline{X})]\mathbb{E}[f_{\text{opt}}(\underline{X}) - f(\underline{X}) \mid \underline{X}]\} &= 0, \end{aligned}$$

így azt kaptuk, hogy

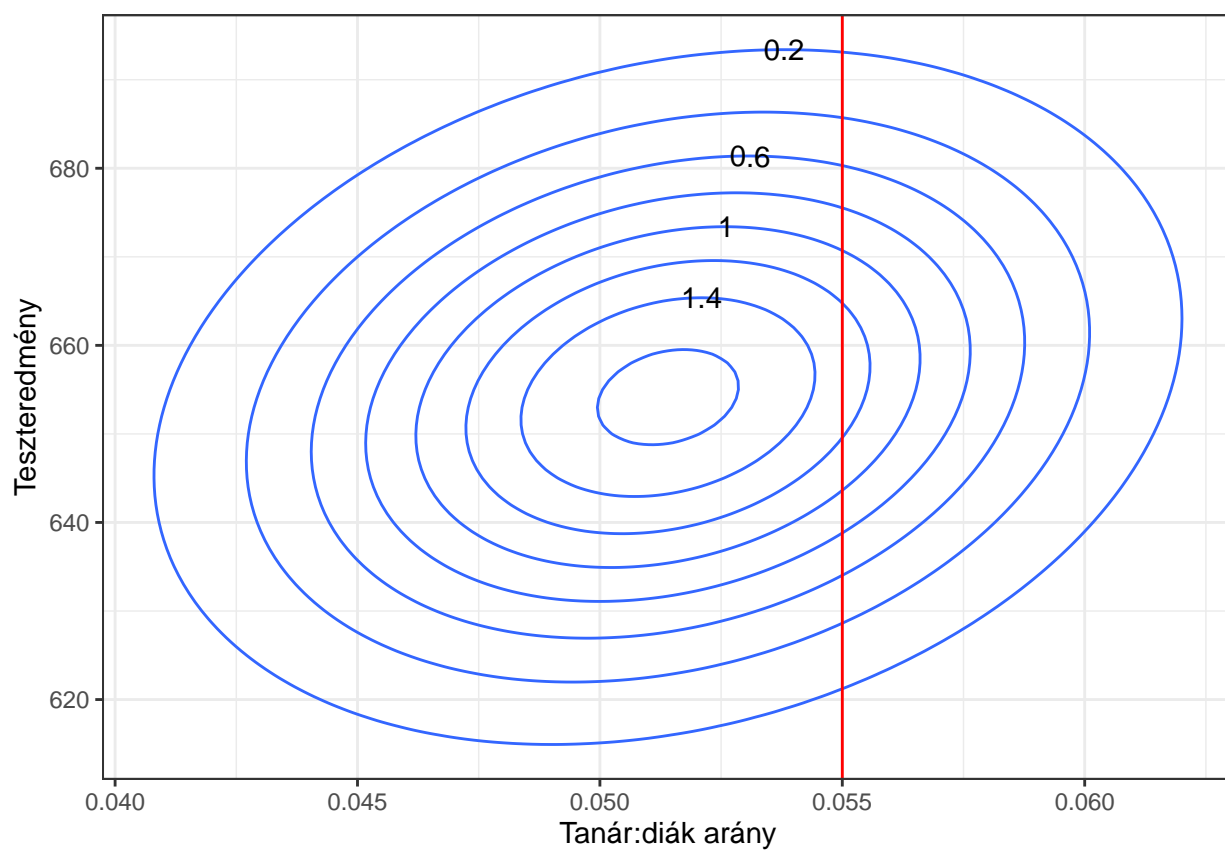
$$\mathbb{E}[Y - f(\underline{X})]^2 = \mathbb{E}[Y - f_{\text{opt}}(\underline{X})]^2 + \mathbb{E}[f_{\text{opt}}(\underline{X}) - f(\underline{X})]^2,$$

amiből már csakugyan látható, hogy f_{opt} a legjobb választás, hiszen az első tagra nincsen ráhatásunk (mi ugye f -et állítjuk), a második tag pedig egy négyzet várható értéke, így 0-nál kisebb nem lehet, de az csakugyan elérhető, ha f -nek f_{opt} -ot választjuk.

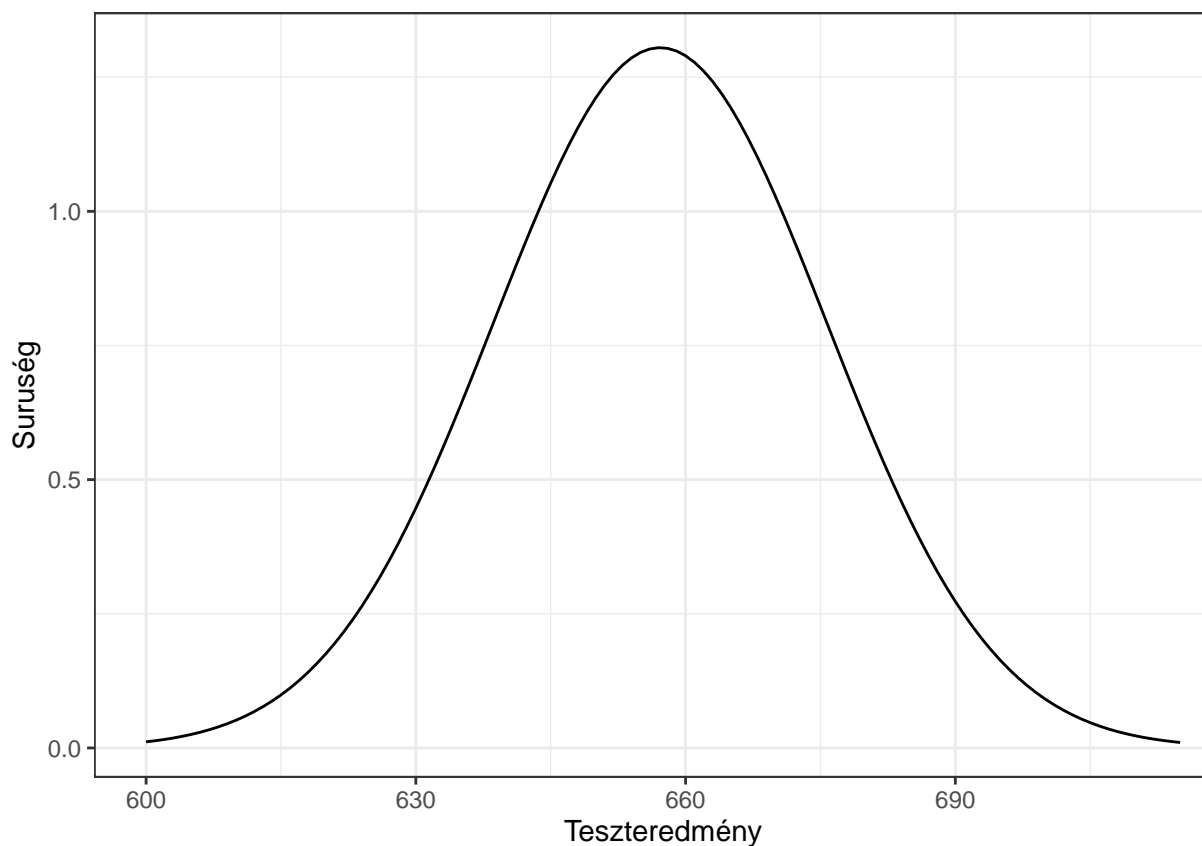
Látható tehát, hogy ez az eredmény *teljesen univerzális*, semmit nem tételeztünk fel f -ről!

Talán nem felesleges feleleveníteni ezen a ponton a feltételes várható érték fogalmát.

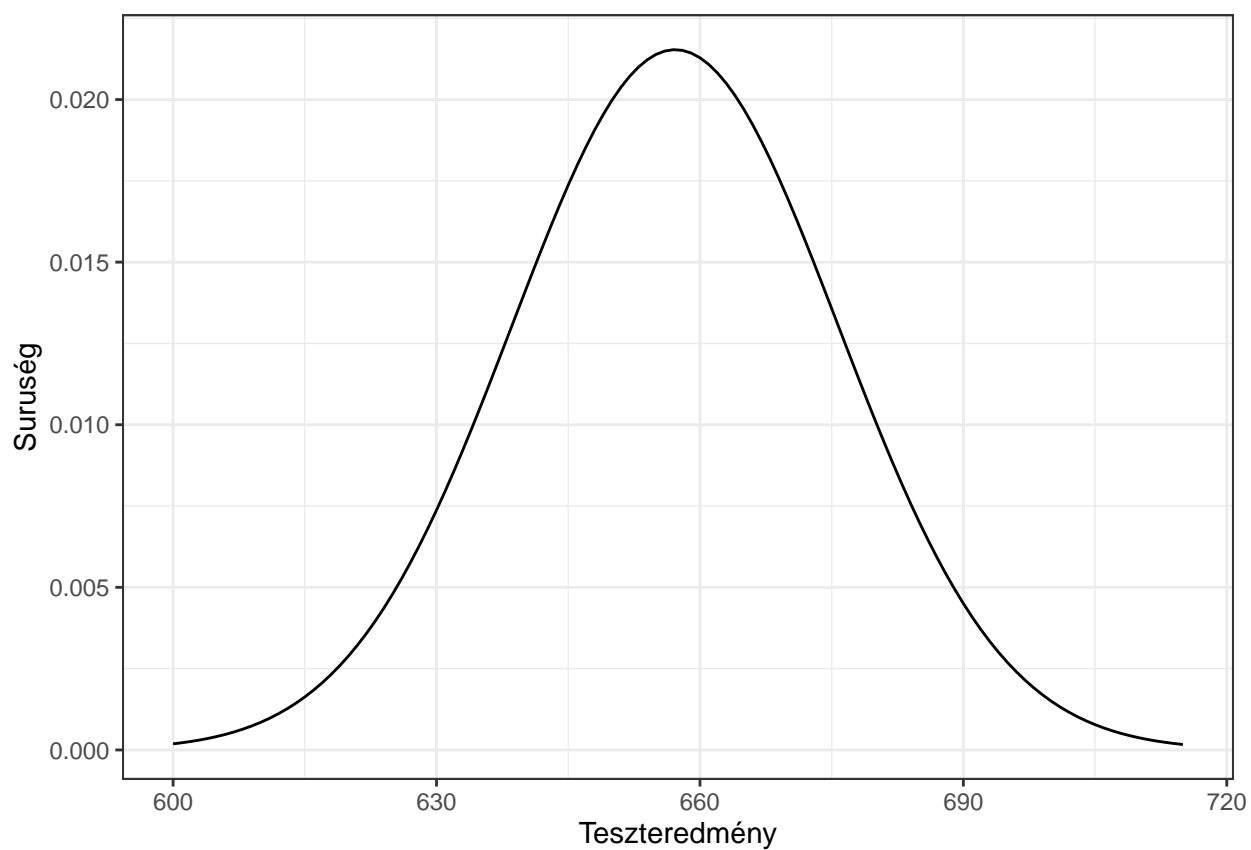
A kiindulópont a feltételes eloszlás, amit úgy kapunk, hogy fogjuk az együttes eloszlást, és egy adott ponton (ami a feltétel) átmetszük. Mondjuk legyen a feltétel az, hogy $X_1 = 0,055$:



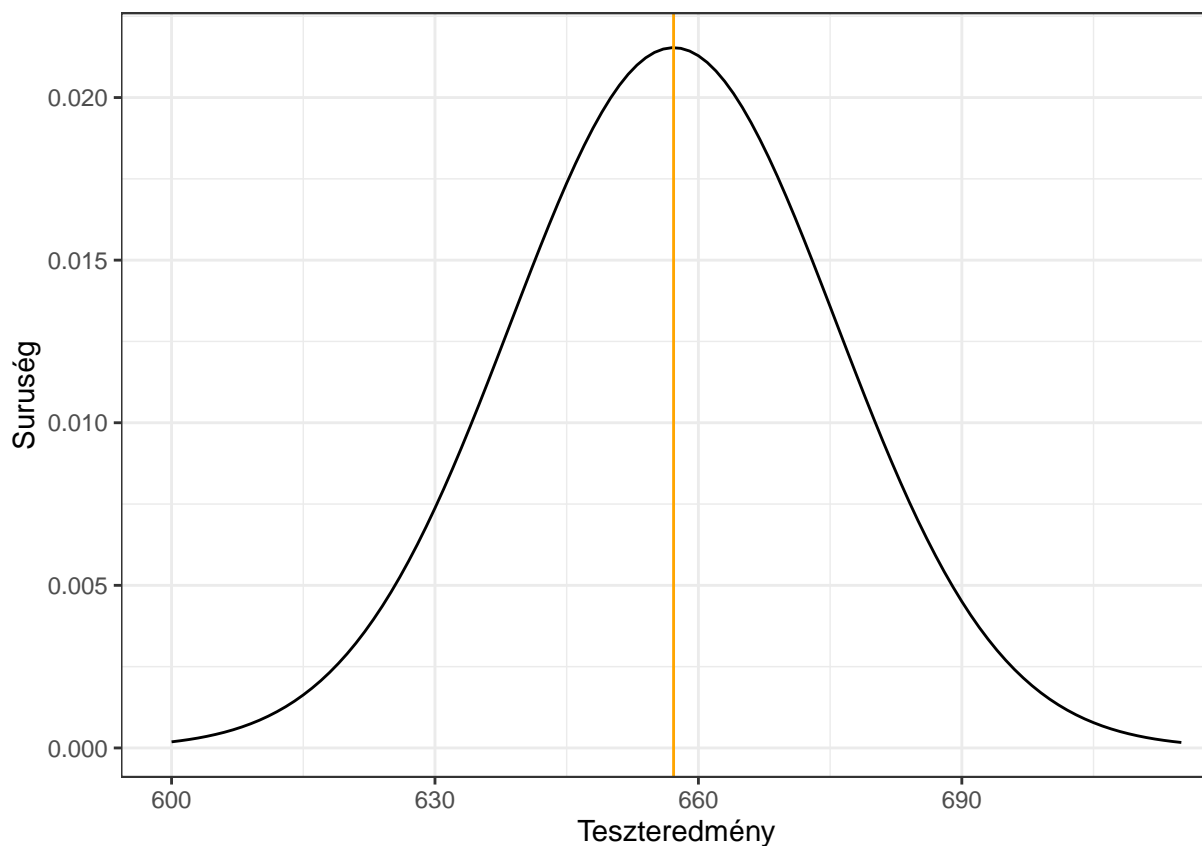
Az együttes sűrűségfüggvény, ne feledjük, egy hegy (aminek a szintvonalait mutatja az ábra), tehát arról van szó, hogy fogunk egy nagy kést, és a piros vonal mentén végigvágjuk a hegyet. Így ezt kapjuk:



Vigyázat, ez még nem sűrűségfüggvény, hiszen nem 1 a görbe alatti területe! De már majdnem megvagyunk, nincs más feladatunk, mint átnormálni (elosztani alkalmas konstanssal), hogy 1 legyen a görbe alatti terület, ez az alkalmas konstans persze a jelenlegi görbe alatti területe lesz, ami nem más, mint a vetületi eloszlás értéke a feltétel pontjában. Elvégezve ezt kapjuk a feltételes eloszlást:

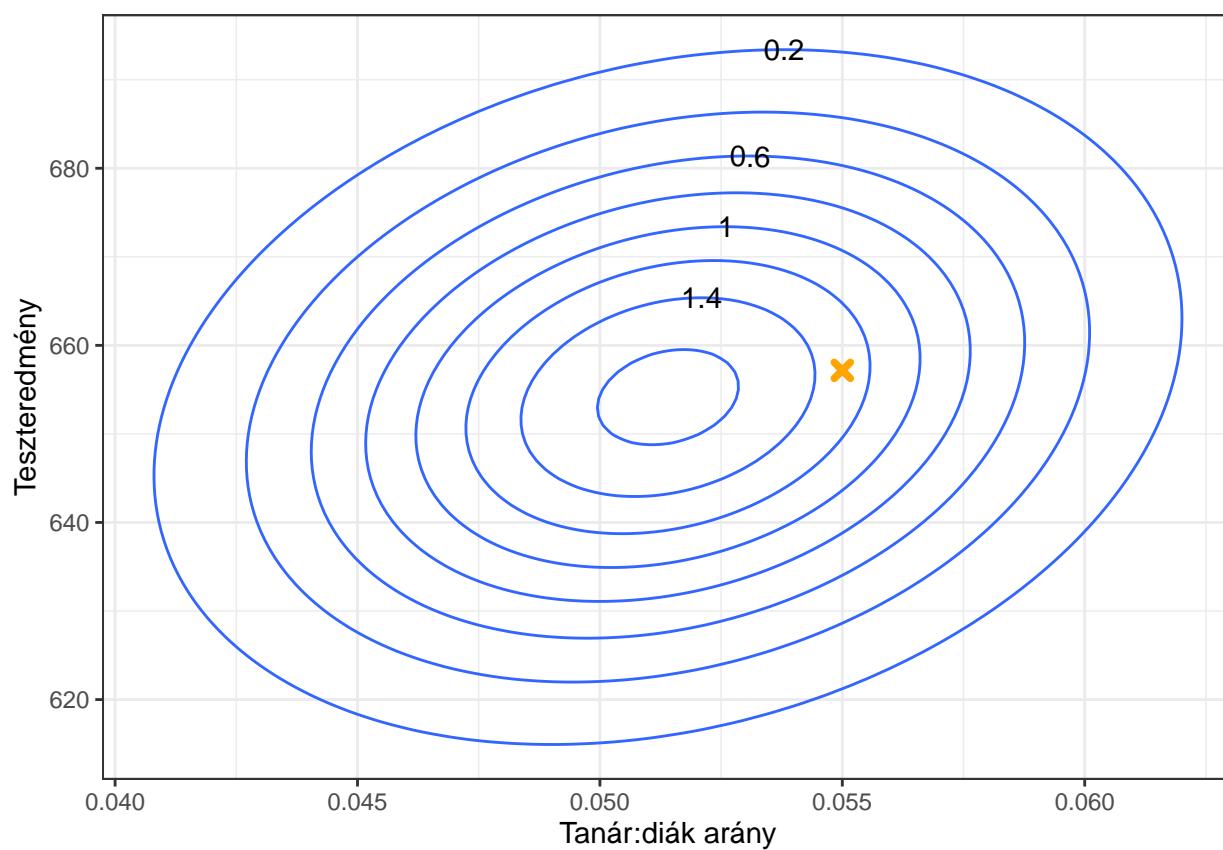


A feltételes várható érték nem más, mint a feltételes eloszlás várható érték – tehát ennek a fenti függvénynek a várható érték. Bejelölve rajta:



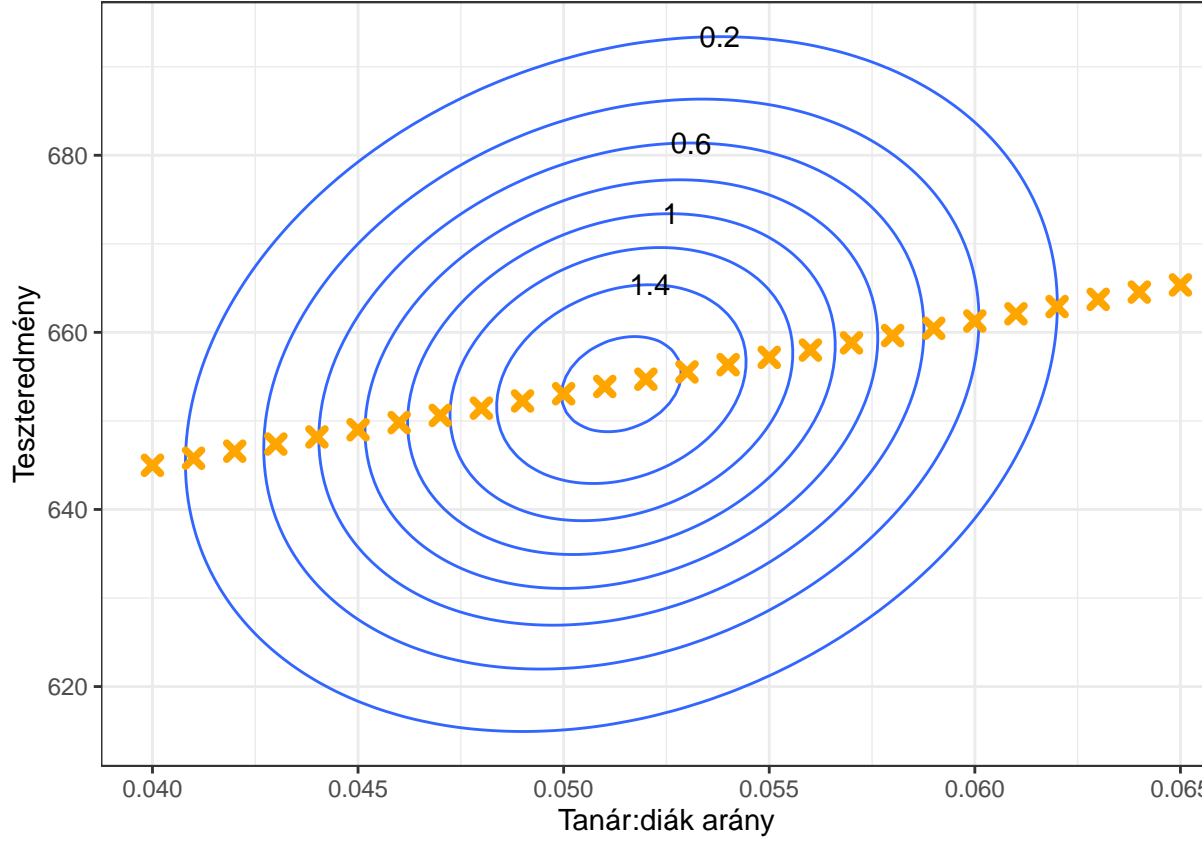
A feltételes várható érték tehát nagyjából 657, és ne feledjük, ez ahhoz a feltételhez tartozik, hogy a tanár:diák arány értéke 0,055. Ahogy az előbb megállapítottuk: ha valaki azt kérdezi, hogy ekkora tanár:diák arány mellett mi a legjobb tippünk a teszteredményre, akkor válaszoljunk 657-et! Ezzel is hibázhatunk persze, de így is ekkor járunk a legjobban (elfogadva persze, hogy négyzetes hibázást minimalizálunk).

Jelöljük is be ezt az értéket az eredeti ábránkon:



Így ni: ha 0,055-ben kérdeznék meg minket, akkor ez a legjobb tippünk.

De az ember itt már vérszemet kap: vajon mi történik, ha kiszámoljuk az összes többi pontban is, hogy mi a legjobb tippünk, tehát a feltételes várhatóértéket?!
Íme:



Nem lehet nem észrevenni: ezek mind egy egyenesre⁶ illeszkednek! A dolog természetesen nem véletlen, és azért van így, mert az eloszlás többváltozós normális volt. Ez esetben az optimális sokasági regressziófüggvény csakugyan mindig lineáris, ez tételként kimondható⁷: ha Y és \underline{X} együttes eloszlása normális⁸, akkor

$$\mathbb{E}(Y \mid \underline{X}) = \mathbb{E}Y + \mathbf{C}_{Y\underline{X}} \mathbf{C}_{\underline{X}\underline{X}}^{-1} (\underline{X} - \mathbb{E}\underline{X}).$$

Egy pillanatra álljunk meg. Eddig feltételes eloszlást csak úgy írtunk, hogy a feltétel az egy konkrét érték (szám vagy vektor) volt: $\mathbb{E}(Y \mid \underline{X} = \mathbf{x})$. De itt

⁶Érdemes megfigyelni (ez kétváltozós esetben jó szemmértékkel még érzékelhető vizuálisan is), hogy a regressziófüggvény *nem* az ellipszisek nagytengelye – tehát a korrelációs mátrix megfelelő sajátvektora – irányába mutat! Hanem az ellipszis “vízszintesen szélső” pontjain megy át.)

⁷A bizonyítást itt elhagyom, lásd például: Bolla-Kráml: Statisztikai következtetések elmélete. Typotex, 2005. 207-208. oldal.

⁸Jelölje $\mathbf{C}_{\underline{X}\underline{X}}$ az X -ek szokásos kovarianciamátrixát, $\mathbf{C}_{Y\underline{X}}$ pedig azt az oszlopvektort, amely sorban az összes X kovariációját tartalmazza Y -nal.

valami más szerepel! A magyarázathoz elevenítsünk fel egy valság definíciót: a $\mathbb{E}(Y \mid \underline{X} = \mathbf{x})$ egy h transzformációt definiál (hiszen adott \mathbf{x} -hez hozzárendel egy valós számot), és $\mathbb{E}(Y \mid \underline{X})$ alatt $h(\underline{X})$ -et értjük. Tehát van értelme az $\mathbb{E}(Y \mid \underline{X})$ objektumnak is, és ez egy valószínűségi változó lesz. Számunkra ebből annyi fontos, hogy ha $\mathbb{E}(Y \mid \underline{X})$ -t látunk, azt értsük úgy, mint valami, ami *minden* \mathbf{x} esetén működik, bármikor beírható, hogy így $\mathbb{E}(Y \mid \underline{X} = \mathbf{x})$ legyen belőle. Természetesen ez fontos, hogy ha egy egyenletben szerepel, akkor ezt az átírást mindenhol megtegyük, pl. írhatjuk, hogy $\mathbb{E}(Y \mid \underline{X} = \mathbf{x}) = \mathbb{E}Y + \mathbf{C}_{Y\underline{X}} \mathbf{C}_{\underline{X}\underline{X}}^{-1} (\underline{X} - \mathbf{x})$ (hiszen \mathbf{x} várható értéke saját maga). Ez a jelölés tehát egyfajta általánosítás.

Egy tulajdonságát már ennyi alapján is rögtön láthatjuk a sokasági regresszió-függvénynek: hogy átmegy a várhatóértékek pontján. Hiszen ha a magyarázó változók értéke épp a várhatóértékük, akkor a második tag kiesik, és azt kapjuk, hogy a regressziófüggvény pont az eredményváltozó várható értékét veszi fel.

Visszatérve, ha bevezetjük a

$$\beta_0 = \mathbb{E}Y - \mathbf{C}_{Y\underline{X}} \mathbf{C}_{\underline{X}\underline{X}}^{-1} \mathbb{E}\underline{X}$$

és a

$$(\beta_1 \quad \beta_2 \quad \cdots \quad \beta_k)^T = \mathbf{C}_{Y\underline{X}} \mathbf{C}_{\underline{X}\underline{X}}^{-1} \underline{X}$$

jelöléseket, akkor írhatjuk, hogy

$$\mathbb{E}(Y \mid \underline{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k.$$

Ebből talán még világosabban látszik a korábbi állítás: hogy többváltozós normális eloszlásnál speciálisan a regressziófüggvény lineáris lesz.

Érdemes megnézni, hogy a még áttekinthető kétváltozós ($Y, X_1 = X$) esetben ez az általános eredmény mire specializálódik: ekkor azt kapjuk, hogy $\mathbb{E}(Y \mid X) = \mathbb{E}Y + \frac{\text{cov}(X, Y)}{\mathbb{D}^2 X} \cdot (X - \mathbb{E}X)$. Két dolgot vegyünk észre:

- Korreláció megjelenése:

$$\mathbb{E}(Y \mid X) = \mathbb{E}Y + \frac{\text{cov}(X, Y)}{\mathbb{D}^2 X} (X - \mathbb{E}X) = \mathbb{E}Y + \frac{\mathbb{D}Y}{\mathbb{D}X} \cdot \text{corr}(X, Y) \cdot (X - \mathbb{E}X).$$

- A linearitás megjelenése itt:

$$\mathbb{E}(Y \mid X) = \mathbb{E}Y + \frac{\text{cov}(X, Y)}{\mathbb{D}^2 X} (X - \mathbb{E}X) = \left(\mathbb{E}Y - \frac{\text{cov}(X, Y)}{\mathbb{D}^2 X} \cdot \mathbb{E}X \right) + X \cdot \frac{\text{cov}(X, Y)}{\mathbb{D}^2 X},$$

azaz $\mathbb{E}(Y | X) = \beta_0 + \beta_1 X$, ha $\beta_0 = \left(\mathbb{E}Y - \frac{\text{cov}(X,Y)}{\mathbb{D}^2 X} \cdot \mathbb{E}X \right)$ és $\beta_1 = \frac{\text{cov}(X,Y)}{\mathbb{D}^2 X}$.

És most egy borzasztó fontos dolog következik. Rakjuk össze a puzzle darabjait: egyfelől tudjuk, hogy $Y = f(X_1, X_2, \dots, X_k) + \varepsilon$, másrészt most már azt is megállapítottuk, hogy az itt szereplő f legjobb értéke $\mathbb{E}(Y | \underline{X})$, és ez mindig⁹ igaz. Tehát azt kaptuk, hogy:

$$Y = \mathbb{E}(Y | \underline{X}) + \varepsilon$$

Ez a dekompozíciót szokás a regresszió “hibaalakjának” (error form) nevezni. Lényegében arról van szó, hogy szétbontjuk az eredményváltozó alakulását egy “magyarázóváltozókkal elérhető legjobb becslés” (már láttuk: a feltételes várhatóérték) és egy “maradék hiba” részre (ami marad). A regresszióanalízis a *feltételes* eloszlásra koncentráll! Ezért elvileg olyasmit kéne írunk, hogy “ $(Y | X) = \mathbb{E}(Y | \underline{X}) + \varepsilon$ ”, de ezt nem tesszük (az $(Y | \underline{X})$ objektumot nem szokás definiálni), ehelyett a bal oldalra simán Y -t írunk (de ne feledjük, hogy ez *feltételes*).

Nagyon fontos látni, hogy a regresszió *mindig* felírható így! És ráadásul ez *biztosan* optimális felírás. A mi választásunk az lesz, hogy majd $\mathbb{E}(Y | \underline{X})$ helyébe mit írunk be, például *ha* tudjuk hogy minden változó együttes eloszlása többváltozós normális, akkor azt, hogy $\mathbb{E}(Y | \underline{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$, és így azt kapjuk, hogy

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon,$$

de vigyázat, ez már – szemben az előző formulával – *nem* univerzális, csak normalitás esetén érvényes.

Lehetne $\mathbb{E}(Y | \underline{X})$ helyébe más is beírni, de mindaddig, amíg jellegre ilyen megoldást használunk, tehát megadjuk a függvény formáját, csak egy vagy több paraméter az, ami meghatározza a konkrét függvényt, szokás **paraméteres regresszióról** beszélni. Ez nem kötelező, lehetne az $\mathbb{E}(Y | \underline{X})$ anélkül próbálni közelíteni, hogy bármilyen *konkrét* függvényforma mellett elköteleződne, ekkor beszélünk nem-paraméteres regresszióról. Ilyenekkel most nem foglalkozunk, csak az érzékeltetés kedvéért egy lehetőség: TODO

Tegyük még egy megállapítást, ami most nem tűnik nagyon izgalmasnak, de a későbbiekben rettentő fontos lesz. Annyit kell tudnunk hozzá, hogy a feltételes várható érték is lineáris, valamint, hogyha valaminek kétszer vesszük a várható értékét, az ugyanaz mintha egyszer vennénk, és ez nem csak a szokásos várható értékre, hanem a feltételes várható értékre is igaz:

⁹Ha szigorúak akarunk lenni, akkor azért annyit hozzá kell tennünk, hogy *ha* létezik egyáltalán a feltételes várható érték. Vannak eloszlások, amiknek egyszerűen nem létezik várható érték, úgyhogy ez elvileg nem mindegy, de mi most ilyen helyzetekkel nem fogunk foglalkozni.

$$\mathbb{E}(\varepsilon \mid \underline{X}) = \mathbb{E}(Y - \mathbb{E}(Y \mid \underline{X}) \mid \underline{X}) = \mathbb{E}(Y \mid \underline{X}) - \mathbb{E}[\mathbb{E}(Y \mid \underline{X}) \mid \underline{X}] = \mathbb{E}(Y \mid \underline{X}) - \mathbb{E}(Y \mid \underline{X}) = 0.$$

Azaz azt kapjuk, hogy $\mathbb{E}(\varepsilon \mid \underline{X}) = 0$. Nagyon fontos, hogy értsük, hogy most mit mondunk: *ha* (!) tényleg – valóságban helyes – $\mathbb{E}(Y \mid \underline{X})$ -t használjuk, *akkor* $\mathbb{E}(\varepsilon \mid \underline{X}) = 0$ *kell* legyen. Ez azért lesz izgalmas, mert $\mathbb{E}(Y \mid \underline{X})$ -t mi sem tudhatjuk biztosan, majd be kell valamit írunk a helyébe (például azt, hogy $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$). Ez a megállapítás tehát azt mondja, hogy *ha* beletrafálunk a dologba, *akkor* $\mathbb{E}(\varepsilon \mid \underline{X}) = 0$ *kell* legyen. *De* ha nem (például ezt írjuk be, csak épp közben nem normális az eloszlás), akkor már ez egyáltalán nem biztos, hogy igaz lesz!

Összefoglalva, ott tartunk, hogy *ha* az eloszlás normális akkor $\mathbb{E}(Y \mid \underline{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$ és így $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$. A bökkendő persze ott van, hogy azt mi magunk sem tudhatjuk, hogy milyen a változóink eloszlása. És itt jön egy fontos döntés. Munkánkat úgy fogjuk megkezdeni, hogy azt mondjuk *akármilyen* is az eloszlás, mi *mindenképp* az $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$ modellt használjuk! Ezt fogjuk **lineáris regresszió**nak nevezni. Még egyszer: ez egy helyes döntés, ha normális a változóink eloszlása, de különben nem. Ha nem ismerjük a változóink eloszlását, akkor ez többé már nem egy matematikailag levezethető szükségszerűség, hanem egy *választás* a részünkről. De több ok szól e választás mellett:

- Többváltozós normalitásnál egzaktan ez a helyzet
- Más esetekben ugyan nem, de cserében nagyon kellemesek a tulajdonságai, különösen ami az interpretációt illeti
- Az is elmondható, hogy – a Taylor-sorfejtés logikáját követve – bármi más is a jó függvényforma, legalábbis *lokálisan* ez is jó közelítés kell legyen
- Bár első ránézésre ez vegytisztán lineáris, valójában majd látni fogjuk, hogy egy sor nemlineáris modell is *visszavezethető* erre a modellre
- És végezetül egy lényeges szempont: lesznek majd eszközeink arra, hogy észrevegyük, hogy rossz volt ez a választás, és megpróbáljuk kijavítani

De újfent nagyon fontos hangsúlyozni, hogy a valós munka során, ahol nem tudjuk mik az eloszlások, azt, hogy az $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$ fennáll, már nem kezelhetjük matematikai szükségszerűségnek, hanem mint feltételt fel kell tennünk!

1.6 A szálak összeérnek

Ezen a ponton összeérnek a szálak. Vegyük csak újra az előbbi alakot (és feltételezzük, hogy a szükséges feltevések teljesültek):

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon.$$

2. fejezet

Regresszió a mintában: következtetés

Pár fogalmat talán érdemes feleleveníteni következtető statisztikából. Az alap-probléma: a halmaz amire a kérdésünk irányul, a **sokaság** sajnos azonban ennek minden elemét nem tudjuk megfigyelni (azaz lemérni), csak egy részét, a kisebb részhalmaz a **minta**. Ez előfordulhat akkor, ha a sokaság TODO

2.1 A hagyományos legkisebb négyzetek (OLS) elve

Ilyen becslési elv a hagyományos legkisebb négyzetek (ordinary least squares, OLS) elve. Mint általános becslési el, nem kell hozzá semmilyen regresszió, a legközségesebb következtető statisztikai példán is elmondható. Példaként vegyük az egyik legelemibb kérdést: sokasági várható érték becslése normalitás esetén, tehát a sokaság eloszlása normális (az egyszerűség kedvéért legyen a szórás is ismert, tehát azt nem kell becsülnünk). Ami fontos: bár egy alap következtető statisztika kurzuson nem szokták mondani, de lényegében itt is az a helyzet, hogy egy *modellt* feltételezünk a sokaságra, jelesül: $Y \sim \mathcal{N}(\mu, \sigma_0^2)$, amit nem mellesleg úgy is írhatnánk, hogy $Y = \mu + \varepsilon$, ahol $\varepsilon \sim \mathcal{N}(0, \sigma_0^2)$. Most μ megbecslése céljából veszünk egy n elemű fae (független, azonos eloszlású) mintát a sokaságból; ekkor feltevésünk szerint $Y_i = \mu + \varepsilon_i$ lesz az i -edik mintaelem. (A feltevésünk igazából azt jelentette, hogy az ε_i változók függetlenek és azonos eloszlásúak.) Figyeljünk oda a kis és nagybetűkre! A nagy betű valószínűségi változó, valami aminek eloszlása van, sokasági dolog. Kisbetű egy konkrét szám, nem valószínűségi, nincsen eloszlása, mintabeli dolog. Most valaki megkérdezhetné, hogy oké, azt értem, hogy Y miért nagy betű, de az Y_i miért az? Hiszen azt mondtuk, hogy az az egyik mintaelem...! Talán a legjobban úgy lehet ezt elképzelni, hogy a véletlen mintavétel az, hogy megkeverjük az urnát, hogy kihúzzunk belőle

egy golyót. Megáll a keverés, nyúlunk bele az urnába, hogy húzzunk: ekkor számunkra az még egy véletlen dolog, hogy mi lesz az elsőként húzott elem, annak eloszlása van (fae mintavétel esetén – tehát ha a golyókat mindig visszadobjuk, és az urnát mindig jól átkeverjük – ugyanaz, mint a sokaság, tehát mint az egész urna eloszlása). Ekkor ez még Y_1 számunkra. Ekkor kihúzzuk a golyót, és meglátjuk a konkrét értéket: ez lesz y_1 . Kicsit matematikusabban szólva: kaptunk egy realizációt Y_1 -ből, ez lesz az y_1 .

A másik ami fontos: a modellből következik egy *becsült érték* minden mintabeli elemhez, jelen esetben, ha m egy feltételezett érték az ismeretlen sokasági várható értékre, akkor

$$\hat{y}_i = m.$$

(Itt persze elvileg beszélni kellene arról, hogy még ha tudjuk is, hogy a sokasági várható érték m , miért pont az lesz a becslésünk is minden mintaelemre. Fogadjuk el intuitíve, egyébként olyan – négyzetes hiba minimalizálására hivatkozó – érvelést használhatnánk mint az előző fejezetben, úgy, hogy az egyetlen magyarázó változónk az $X_0 = 1$.)

Egy kitérő megjegyzés: ha jobban megnézzük a fentieket, akkor láthatjuk, hogy az OLS-elv alkalmazásához igazából nem is kell, hogy a sokasági eloszlást ismerjük, csak annyi a fontos, hogy legyen egy modellünk, és belőle tudjunk becsült értékeket származtatni a ténylegesen is ismert megfigyelésekhez.

És akkor jöhet az OLS-elv! Egy mondatban összefoglalva: az ismeretlen sokasági paraméterre az a becsült érték, amely mellett a tényleges mintabeli értékek, és az adott paraméter melletti, modellből származó becsült értékek a legközelebb vannak egymáshoz, amit úgy fordítunk le, hogy a köztük lévő eltérések négyzetének összege a legkisebb! A megoldandó – optimalizációs jellegű – feladat tehát matematikailag:

$$\hat{\mu} = \operatorname{argmin}_m \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \operatorname{argmin}_m \sum_{i=1}^n (y_i - m)^2.$$

És ennek megoldása természetesen $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$ ebben a példában.

Egyetlen kiegészítést kell tenni a fentiekhez. Megkaptuk ugyan a becslőt, csak-hogy az \bar{y} egyetlen konkrét szám. (Hát persze, mert egy konkrét mintához, a $\{y_1, y_2, \dots, y_n\}$ mintához – kisbetűk! – tartozik.) Minket azonban alapvető fontossággal fog érdekelni a becsülő **mintavételi eloszlása**, tehát, hogy ha újra meg újra mintát veszünk ugyanabból a sokaságból, és mindegyik mintából kiszámoljuk a becsülőfüggvény értékét (jelen esetben a mintaátlagot), akkor annak mi lesz az eloszlása. A becsülőfüggvényünk az igazából egy *transzformáció* a mintaelemekkel (“add össze őket és oszd el a mintanagysággal”), de ha egyszer ez a transzformáció megvan, azt nyugodtan ráereszthetjük valószínűségi változókra

is, nem csak számokra! Ami magyarul azt fogja jelenteni, hogy felírjuk ugyanazt – csak épp kisbetűk helyett nagybetűkkel. Jelen példában a becslőfüggvényünk: $\frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}$, és íme, ennek már nagyon is eloszlása van, hiszen egy valószínűségi változó maga is – ez az eloszlás lesz a mintavételi eloszlás. Megvizsgálhatóak a tulajdonságai, megnézhetjük, hogy a várható értéke egyezik-e a sokasági paraméterrel (torzítatlanság), hogy mekkora a szórása (hatásosság), hogy hogyan viselkedik, ha n egyre nagyobb (konzisztencia) és így tovább.

2.2 Lineáris regresszió becslése OLS-elven

Most vegyük elő a lineáris regressziókat! (Ahol ezt közszemérem-sértés veszélye nélkül megtehetjük.) Azt látjuk, hogy ott eddig a sokaságról beszéltünk, feltettünk egy modellt (*ugyanúgy mint az előbbi példában*), jó, lehet, hogy egy kicsit bonyolultabbat, de akkor is, ugyanúgy egy sokaságra vonatkozó modell, amiből, megint csak pontosan ugyanúgy mint az előbbi példában, tudunk egy becslést előállítani minden mintaelemhez. Ez lehetővé teszi, hogy az ismeretlen paramétereket OLS-elven megbecsüljük!

Lássuk a részleteket. A változóink az $(Y, X_1, X_2, \dots, X_k)$, ezekre vegyünk egy n elemű mintát; az i -edik mintaelemet jelölje $(Y_i, X_{i1}, X_{i2}, \dots, X_{ik})$. Természetesen a modellünk ezekre is igaz lesz, tehát írhatjuk, hogy

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \varepsilon_i.$$

Ez minden i -re teljesül, tehát ha nagyon elszántak vagyunk, akkor n ilyen egyenletet írhatnánk fel:

$$\begin{aligned} Y_1 &= \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \dots + \beta_k X_{1k} + \varepsilon_1 \\ Y_2 &= \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \dots + \beta_k X_{2k} + \varepsilon_2 \\ &\dots \\ Y_n &= \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \dots + \beta_k X_{nk} + \varepsilon_n \end{aligned}$$

Az i -edik mintaelem realizációja az $(y_i, x_{i1}, x_{i2}, \dots, x_{ik})$. (A minta egyelőre legyen fae – hogy ez mennyire jó feltevés, arról később még fogunk beszélni.)

Ha b_0, b_1, \dots, b_k -val jelöljük a feltételezett sokasági paramétereket, akkor a becslés

$$\hat{y}_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_k x_{ik}$$

lesz az i -edik mintaelemre. (Itt szerencsére nincs mit gondolkozni, hiszen azt az előző fejezetben részletesen levezettük, hogy ez lesz a legjobb becslés adott \mathbf{x} mellett.)

Most hogy megvannak a becsült értékek (\widehat{y}_i) és a tényleges értékek (y_i), betű szerint ugyanazt az optimalizációs feladatot kell felírunk, mint az előbb, csak \widehat{y}_i lesz kicsit hosszabb, ha kifejtjük:

$$\begin{aligned} (\widehat{\beta}_0, \widehat{\beta}_1, \widehat{\beta}_2, \dots, \widehat{\beta}_k) &= \operatorname{argmin}_{b_0, b_1, b_2, \dots, b_k} \sum_{i=1}^n (y_i - \widehat{y}_i)^2 = \\ &= \operatorname{argmin}_{b_0, b_1, b_2, \dots, b_k} \sum_{i=1}^n [y_i - (b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_k x_{ik})]^2 \end{aligned}$$

Annyi bonyolódottság van, hogy itt most *több* paramétert kell becsülni, de ez csak a kivitelezést nehezíti, elvileg teljesen ugyanaz a feladat.

Össze ne keverjük β_i -t, b_i -t és $\widehat{\beta}_i$ -t! β_i a kérdéses sokasági paraméter valódi, tényleges értéke, egy adott, konkrét szám (csak mi nem tudjuk mennyi), b_i egy általunk feltételezett érték rá, mi állítjuk be, választhatunk nagy számot is, kis számot is, tetszés szerint, a fenti optimalizációban végig fogunk vele futni az összes lehetséges értékén, $\widehat{\beta}_i$ pedig a megoldásként kapott *legjobb tippünk* β_i -re, de ettől még csak tipp, azaz eloszlása lesz, hiszen a mintától is függeni fog, mintáról mintára ingadozni fog (miközben a valódi érték ugyebár állandó – ez lesz a mintavételi hiba forrása).

Ezt az optimalizációs problémát kell tehát most megoldanunk. Ezt megtehetnénk a fenti formában is, de célszerűbb, ha már most áttérünk a vektoros/mátrixos jelölésekre. Ez eleinte kicsit kényelmetlennek tűnhet, de a magasabb absztrakciós szint később ki fog fizetődni: lehet, hogy most kicsit nehezebben indulunk, de cserében a bonyolultabb problémák sem lesznek sokkal nehezebbek.

Fogjunk tehát össze mindent értelemszerű vektorokba és mátrixokba! A jelölésrendszer teljes bemutatása végett felírom a mintavétel előtti – valószínűségi változós – és a realizálódott értékes alakokat is¹. Az eredményváltozók:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix}.$$

A magyarázó változókat nyilván mátrixba kell összefogni, de itt egy kis cselre lesz szükségünk: hozzáveszünk az elejéhez egy csupa 1 oszlopot. (Az így kapott mátrixot a regresszió **design mátrixának** szokás nevezni.) Íme:

¹A jelölésrendszer sajnos nem tökéletesen konzisztens, hiszen \mathbf{X} nagybetű, és mégis kisbetűs dolgokat fog össze. Nem akartam szakítani a lineáris algebra hagyományával, hogy a mátrixot nagybetű jelöli, bár ez tényleg keveredik a valószínűségszámítás nagybetűjével. Abból azonban, hogy vastagítás vagy aláhúzás van, mindenképpen világos lesz, hogy valószínűségi változóról vagy realizálódott értékről van szó, még ha a kis és nagy betű nem is segít.

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix}, \underline{\underline{X}} = \begin{pmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1k} \\ 1 & X_{21} & X_{22} & \cdots & X_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{nk} \end{pmatrix}.$$

Ez a csupa 1 oszlop azért lesz célszerű, mert ha a regressziós koefficienseket egy

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

vektorba, a hibatagokat pedig egy

$$\boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_1 \\ \vdots \\ \varepsilon_k \end{pmatrix}$$

vektorba fogjuk össze, akkor a korábbi, n darab egyenletből álló, igencsak terjengős felírás helyett nemes egyszerűséggel ezt írhatjuk:

$$\underline{Y} = \underline{\underline{X}}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

És ennyi, pontosan ugyanaz van leírva!

Látható tehát, hogy a csupa 1 oszlop azért kellett, hogy a vektorral való rászor-zásnál az legyen a β_0 szorzója, így az egyenletben tényleg egyszerűen β_0 fog megjelenni.

Menjünk most vissza az OLS optimalizációs problémájára! Ezekkel a jelölésekkel a kezünkben ugyanis azt is sokkal egyszerűbben felírhatjuk:

$$\underset{\mathbf{b}}{\operatorname{argmin}} \hat{\mathbf{y}}^T \hat{\mathbf{y}} = \underset{\mathbf{b}}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b}),$$

hiszen számok négyzetösszegét megkapjuk, ha összefogjuk őket egy vektorba, és vesszük ezen vektor saját transzponáltjával vett szorzatát. ($\hat{\mathbf{y}}$ és \mathbf{b} az értelemszerű vektorok, \hat{y}_i -ket és b_i -ket fogják össze.)

Az $(\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b})$ hibanégyzetösszeget *ESS*-sel (error sum of squares) is fogjuk jelölni².

²Sajnos néhány irodalom az általunk használt *ESS*-re inkább az *RSS*-t (residual sum of squares) rövidítést használja, ami a jelölési zavarok legszerencsétlenebb típusa, ugyanis az *RSS*-t majd később mi is fogjuk használni, csak épp másra. Éppen ezért, ha ilyenekről olvasunk, mindig tisztázni kell, hogy a könyv vagy program írói mit értenek alatta.

És akkor essünk neki: oldjuk meg ezt az optimalizációt! Először alakítsuk át a célfüggvényt, bontsuk fel a zárójeleket:

$$\operatorname{argmin}_{\mathbf{b}} (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b}) = \operatorname{argmin}_{\mathbf{b}} [\mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}].$$

Itt egyszerű algebrai átalakításokat végzünk (és a definíciókat használjuk), hiszen a zárójeleket felbontani, műveleteket elvégezni, mátrixokkal/vektorokkal is hasonlóan kell mint valós számokkal. (A transzponálás tagonként elvégezhető, azaz $(\mathbf{a} - \mathbf{b})^T = \mathbf{a}^T - \mathbf{b}^T$.) Egyedül annyit kell észrevenni, hogy a $\mathbf{y}^T \mathbf{X} \mathbf{b}$ egy egyszerű valós szám, ezért megegyezik a saját transzponáltjával, $\mathbf{b}^T \mathbf{X}^T \mathbf{y}$ -nal. Ezért írhattunk $-(\mathbf{X} \mathbf{b})^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{b}$ helyett egyszerűen – például $-2\mathbf{b}^T \mathbf{X}^T \mathbf{y}$ -t. (Itt mindenhol felhasználtuk, hogy a transzponálás megfordítja a szorzás sorrendjét: $(\mathbf{A} \mathbf{B})^T = \mathbf{A}^T \mathbf{B}^T$.)

Most jön a minimum megkeresése. Az ember rávágja, hogy deriválni kell, de itt ez picit zűrösebb, hiszen a függvényünk többváltozós (ráadásul az is határozatlan, hogy pontosan hányváltozós). Itt jelentkezik igazán a mátrixos jelölésrendszer előnye. A $\mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}$ lényegében egy “másodfokú kifejezés” többváltozós értelemben (az $ax^2 + bx + c$ többváltozós megfelelője), és ami igazán szép: pont ahogy az $ax^2 + bx + c$ lederiválható a változója (x) szerint (eredmény $2ax + b$), ugyanúgy ez is lederiválható a változója (azaz \mathbf{b}) szerint... és az eredmény az egyváltozóssal teljesen analóg lesz, ahogy fent is látható! (Ez persze bizonyítást igényel! – lásd többváltozós analízisből.) Bár ezzel átléptünk egyváltozóról többváltozóra, a többváltozós analízisbeli eredmények biztosítanak róla, hogy formálisan ugyanúgy végezhető el a deriválás. (Ezt írja le röviden a “vektor szerinti deriválás” jelölése. Egy \mathbf{b} vektor szerinti derivált alatt azt a vektort értjük, melyet úgy kapunk, hogy a deriválandó kifejezést lederiváljuk \mathbf{b} egyes b_i komponensei szerint – ez ugye egyszerű skalár szerinti deriválás, ami már definiált! –, majd ez eredményeket összefoglaljuk egy vektorba. Látható tehát, hogy a vektor szerinti derivált egy ugyanolyan dimenziós vektor, mint ami szerint deriváltunk.) Ami igazán erőteljes ebben az eredményben, az nem is egyszerűen az, hogy “több” változónk van, hanem, hogy nem is kell tudnunk, hogy mennyi – mégis, általában is működik! Az eredmény tehát:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{b}} [\mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}] &= \\ &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{b} = 0 \Rightarrow \widehat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \end{aligned}$$

ha $\mathbf{X}^T \mathbf{X}$ nem szinguláris.

Azt, hogy a megtalált stacionaritási pont tényleg minimumhely, úgy ellenőrizhetjük, hogy megvizsgáljuk a Hesse-mátrixot a pontban. A mátrixos jelölésrendszerben ennek az előállítása is egyszerű, még egyszer deriválni kell a függvényt a változó(vektor) szerint:

$$\frac{\partial^2}{\partial \mathbf{b}^2} [\mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}] = \frac{\partial}{\partial \mathbf{b}} [-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{b}] = 2\mathbf{X}^T \mathbf{X}.$$

Az ismert tétel szerint a függvénynek akkor van egy pontban ténylegesen is (lokális, de a konvexitás miatt egyben globális) minimuma, ha ott a Hesse-mátrix pozitív definit. Esetünkben ez minden pontban teljesül. A $\mathbf{X}^T \mathbf{X}$ ugyanis pozitív szemidefinit (ez egy skalárszorzat-mátrix, más néven Gram-mátrix, amelyek mindig pozitív szemidefinitek), a kérdés tehát csak a határozott definitiség. Belátható azonban, hogy ennek feltétele, hogy $\mathbf{X}^T \mathbf{X}$ ne legyen szinguláris – azaz itt is ugyanahhoz a feltételhez értünk! Megjegyezzük, hogy ez pontosan akkor valósul meg, ha az \mathbf{X} teljes oszloprangú. (Erre a kérdésre a modellfeltevések tárgyalásakor még visszatérünk.)

Végül egy számítástechnikai megjegyzés: az együtthatók számításánál a fenti formula direkt követése általában nem a legjobb út, különösen ha sok megfigyelési egység és/vagy változó van. Ekkor nagyméretű mátrixot kéne invertálni, amit numerikus okokból (kerekítési hibák, numerikus instabilitás stb.) általában nem szeretünk. Ehelyett, a különféle programok igyekeznek a direkt mátrixinverziót elkerülni, tipikusan az \mathbf{X} valamilyen célszerű mátrix dekompozíciójával (QR-dekompozíció, Cholesky-dekompozíció). Extrém esetekben még az is elképzelhető, hogy az egzakt, zárt alakú megoldás előállítás helyett valamilyen iteratív optimalizálási algoritmus (gradiens módszer, Newton–Raphson-módszer) alkalmazása a gyakorlatban járható út, annak ellenére is, hogy elvileg van zárt alakban megoldása.

A kapott eredmény nem más, mintha \mathbf{X} Moore–Penrose pszeudoinverzével szoroznánk \mathbf{y} -t.

TODO

Végezzük el a fenti műveleteket közvetlenül lekódolva R-ben a már látott kaliforniai iskolás példára, ha a pontszámot a tanár:diák arányt a pontszámmal és a jövedelemmel regresszáljuk:

```
y <- CASchools$score
X <- cbind( 1, CASchools$tsratio, CASchools$income )
solve( t(X)%*%X )%*%t(X)%*%y
```

```
##      [,1]
## [1,] 614.0
## [2,] 233.4
## [3,]  1.8
```

Egy mátrixot a `t` függvénnyel transzponálhatunk és a `solve` függvénnyel invertálhatunk, a `cbind` pedig vektorokat, mint oszlopvektorokat fűz egybe mátrixszá. (Valaki megkérdezheti, hogy akkor az 1 miért működik, hiszen az nem vektor: ez az R egyik jellemző – kétélű fegyverként viselkedő – tulajdonsága: megengedi

a trehányyságot, ugyanis érzékeli, hogy mi a helyzet, és automatikusan egymás alá rakja annyiszor, mint amilyen hosszúak a többi vektorok.)

Természetesen az R tartalmaz beépített parancsot regressziók becslésére:

```
lm( score ~ tsratio + income, data = CASchools)
```

```
##  
## Call:  
## lm(formula = score ~ tsratio + income, data = CASchools)  
##  
## Coefficients:  
## (Intercept)      tsratio      income  
##      613.98      233.41      1.84
```

Az `lm` a lineáris modell rövidítése. Első argumentumban a regressziós egyenletet kell megadnunk, mint egy R formula (tehát `~` felel meg az egyenlőségjelnek, bal oldalán az eredményváltozó, jobb oldalán a magyarázó változók felsorolása, `+` jellel elválasztva.) Az R konstans alapbeállításként rak a modellbe, azt kell külön kérnünk ha nem szeretnénk (egy `-1` hozzáfűzésével az utolsó magyarázó változó után). A `data` argumentum tartalma a szokásos: ha használjuk, akkor a formulában elég a változónevet leírni, nem kell jelölni, hogy melyik adatkeretre vonatkoznak, mert az R úgy érti, hogy mind a `data` argumentumban megadottra értendő.

3. fejezet

Kategoriális magyarázó változók

3.1 Regresszió csak minőségi változóval (ANO-VA)

3.1.1 Minőségi változók a regresszióban

- A kérdés, ami mostani kutatásainkat motiválja: hogyan szerepeltethetünk egy *minőségi* (nominális vagy ordinális, szokás kategoriális változónak is nevezni) tulajdonságot, pl. férfi–nő, egészséges–beteg, alapfokú–középfokú–felsőfokú végzettségű stb. egy regressziós modellben
- A regresszió csak számszerű adatokat tud felhasználni \rightarrow valahogy *kódolni* kell a kategoriális tulajdonság lehetséges értékeit (kimeneteit, csoportjait)
- Eddig csak mennyiségi tulajdonságokkal foglalkoztunk, aminek kódolása triviális volt: a naturáliában kifejezett értékével (m^2 , eFt stb.)
- Pl. férfi = 0, nő = 1 elég kézenfekvő, de mi van az iskolai végzettséggel?
- Az alap = 0, közép = 1, felső = 2 belekódolja az adatokba, hogy a felső és a közép közti különbség *kénytelen* ugyanakkora lenni, mint a közép és alap közötti (ha felső = 3, akkor kétszer akkora stb.)
- De mi semmi ilyet nem akarunk, hiszen azt szeretnénk, hogy ezt az adatok mondják meg!

3.1.2 Dummy változó fogalma

- A kódolást megvalósíthatjuk olyan változóval vagy változókkal, melyek *csak* 0 vagy 1 értéket vehetnek fel
- Az ilyen változókat nevezzük dummy (bináris vagy indikátor) változónak
- Ha két kimenet van, akkor a kódolás teljesen kézenfekvő: egy dummy

változóra van szükségünk, mely (például) 0 értéket vesz fel férfira, 1-et nőre

- Bonyolultabb a helyzet, ha több kimenet van
- Triviális kódolás:

| | D_A | D_B | D_C |
|---|-------|-------|-------|
| A | 1 | 0 | 0 |
| B | 0 | 1 | 0 |
| C | 0 | 0 | 1 |

3.1.3 Kódolás

Ezen “kódolási tábla” alapján a kódolás (pl. X_1 : jövedelem, X_2 : iskolai végzettség, X_3 : életkor):

$$\begin{bmatrix} & X_1 & X_2 & X_3 \\ \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} & \begin{bmatrix} 213 \\ 311 \\ \vdots \\ 128 \end{bmatrix} & \begin{bmatrix} B \\ C \\ \vdots \\ B \end{bmatrix} & \begin{bmatrix} 32 \\ 41 \\ \vdots \\ 18 \end{bmatrix} \end{bmatrix} \rightarrow \begin{bmatrix} & X_1 & D_A & D_B & D_C & X_3 \\ \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} & \begin{bmatrix} 213 \\ 311 \\ \vdots \\ 128 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} & \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} & \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} & \begin{bmatrix} 32 \\ 41 \\ \vdots \\ 18 \end{bmatrix} \end{bmatrix}$$

Itt már minden tisztán numerikus, működhet a regresszió!

3.1.4 Referencia-kódolás

- ... ám vegyük észre, hogy 3 csoporthoz *nem* kell 3 dummy változó, kódolható 2-vel is!
- Általában k kimenet kódolása megoldható $k - 1$ dummy változóval az ún. referencia-kódolás logikájával
- Itt kiválasztunk egy kimenetet, aminél mind a $k - 1$ darab dummy változó 0 értéket vesz fel (kontrollcsoport vagy referenciacsoport), és a többi $k - 1$ csoportot az jelzi, hogy a $k - 1$ dummy változó közül *melyik* vesz fel 1 értéket (mindig csak 1!)
- Például (3 kimenetre):

| | R_A | R_B |
|---|-------|-------|
| A | 1 | 0 |
| B | 0 | 1 |
| C | 0 | 0 |

- Itt C a referenciacsoport, R_A és R_B a két szükséges (ugye $k = 3$!) magyarázó változó

- Vegyük észre, hogy $R_A \equiv D_A$ és $R_B \equiv D_B$ (tehát a két kódoláshoz pontosan ugyanazon dummykra van szükség, csak a referencia-kódolásnál eldobjuk az egyiket – ez lesz a kontrollcsoport)

3.1.5 Dummy változó csapda

- Ha van konstans a modellben, akkor *tilos* is k csoporthoz k dummyt használni a kódoláshoz
- Ellenkező esetben egzakt multikollinearitás jön létre (gondoljuk végig, hogy a dummy változókhoz mi tartozik a design mátrixban, ld. előbb!); ez az ún. *dummy változó csapda*
- Ha k csoportot mégis k dummyval kódolunk (“triviális kódolás”), akkor viszont nem szerepeltethetünk konstans

3.1.6 Triviális kódolás konstans nélkül

- A két kódolási mód (k darab dummy, nincs konstans és $k - 1$ darab dummy, van konstans) jól szemléltethető egy csak a nominális tulajdonsággal magyarázó regresszióval
- k darab dummy, nincs konstans:

| | D_A | D_B | D_C |
|---|-------|-------|-------|
| A | 1 | 0 | 0 |
| B | 0 | 1 | 0 |
| C | 0 | 0 | 1 |

$$Y = \beta_A D_A + \beta_B D_B + \beta_C D_C + \varepsilon$$

- Együtthatók értelmezése: ha az A csoportban vagyunk, akkor a fenti egyenlet $Y = \beta_A + \varepsilon$ lesz $\Rightarrow \beta_A$ az A csoport csoportátlaga (legkisebb négyzetes elv!); hasonlóan a többi

3.1.7 Referencia-kódolás konstanssal

- $k - 1$ darab dummy, van konstans:

| | D_A | D_B |
|---|-------|-------|
| A | 1 | 0 |
| B | 0 | 1 |
| C | 0 | 0 |

$$Y = \beta^* + \beta_A^* D_A + \beta_B^* D_B + \varepsilon$$

3.1.8 Együtthatók értelmezése referencia-kódolásnál

- Értelmezésnél egy dolgot tartsunk mindig szem előtt: ugyanarra a csoportra ugyanannak az értéknek kell kijönnie, akárhog kódolunk!
- Így $\beta_C = \beta^*$
- Továbbá (a B csoport példáján):

$$\beta_B = \beta^* + \beta_B^* = \beta_C + \beta_B^* \Rightarrow \beta_B^* = \beta_B - \beta_C$$

- Tehát az együtthatók az *eltéréseket* jelentik a referenciacsoporttól (ami pedig a konstansba kerül)
- Vegyük észre, hogy a változónkénti szignifikanciák eltérhetnek (mert másra fognak vonatkozni), de az előrejelzése – és így a modellminősítő mutatók – nem

3.1.9 Fontos hipotézisvizsgálatok

- Egyrészt: szignifikáns-e egy adott csoport átlagának eltérése a referenciacsoport átlagától
- Ez itt nem más, mint β_A^* vagy β_B^* relevanciája
- Egyszerűen t -próbával ellenőrizhető!
- Másrészt: van-e egyáltalán bármilyen csoportok közötti eltérés:

$$H_0 : \beta_A^* = \beta_B^* = \dots = 0$$

$$H_1 : \exists j : \beta_j^* \neq 0$$

- Több csoport átlaga eltér-e? De hát az az ANOVA!
- Az egyezés nem pusztán formai, teljes tartalmi egyezés van (ez nem csak hasonló, hanem ugyanaz: az ANOVA elmondása regressziós “keretben”)

3.1.10 Egynél több kategoriális magyarázó változó

- Ha egynél több kategoriális magyarázó változó van, akkor nem kódolható mindegyik triviálisan, ilyenkor már a konstans eltávolítása sem segít
- (Nem az lesz a baj, hogy valamelyik összege a konstans, hanem, hogy a kettő összege ugyanaz – ez elvileg is megoldhatatlan)
- Referencia-kódolás minden további nélkül használható
- A kétszemponos ANOVA megfelelője regressziós keretben!
- Természetesen feltételezhető interakció is, ez esetben a dummy-kat az összes lehetséges kombinációban szorozni kell

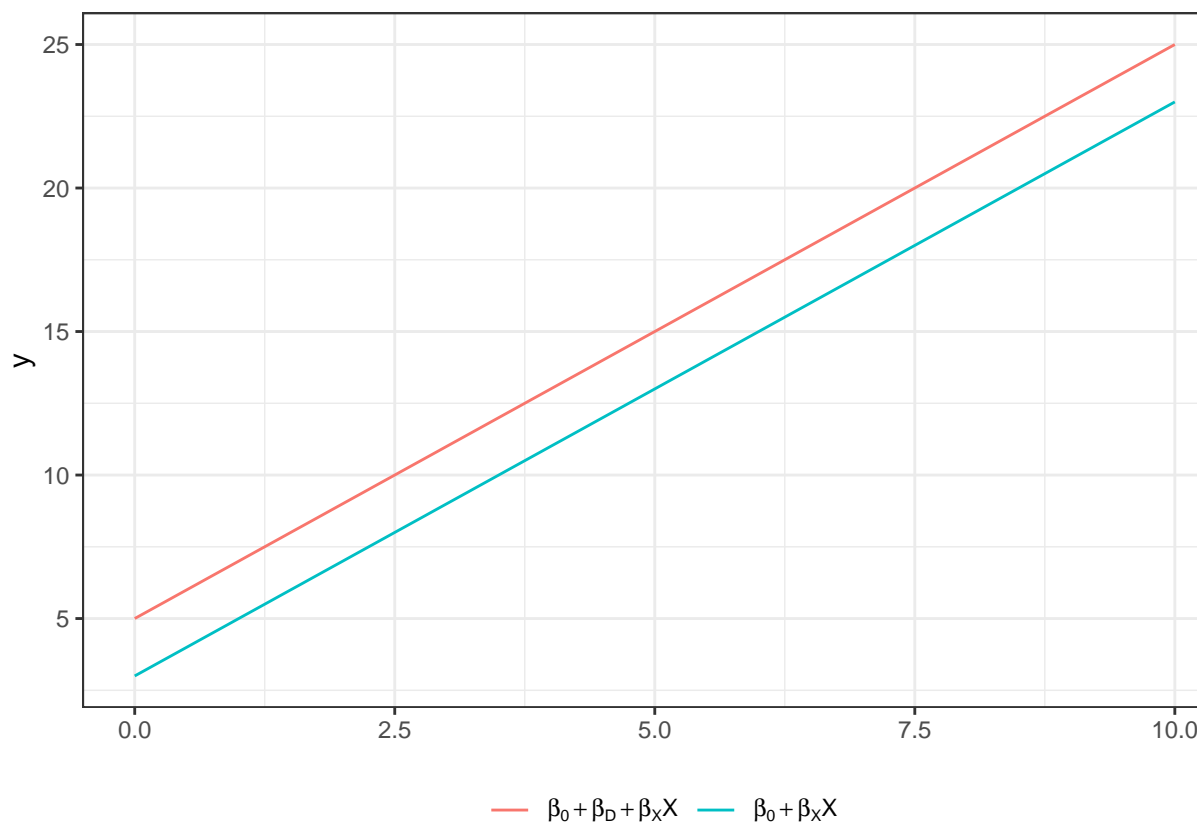
3.2 Regresszió minőségi és mennyiségi magyarázó változóval (ANCOVA)

3.2.1 Dummyzás folytonos magyarázó változó jelenléte mellett

- Amit eddig csináltunk az lényegében az volt, amit *konstans dummyzásának* nevezhetünk: csoportonként eltérő (de konstans) értékkel becsültük az eredményváltozót
- Mi van, ha bevonunk egy folytonos magyarázó változót?
- Azaz ekkor már nem egy konstans becsülünk az egyes csoportokra, hanem egy egyenest (a folytonos magyarázó változó függvényében)
- Dummyzással (tehát a csoporttagság szerint) eltéríthetjük az egyenesek tengelymetszetét és meredekségét is!
- Lehet csoportonként különböző
 - +1 egység magyarázó változó hatása
 - a 0 magyarázó változóhoz tartozó eredményváltozó
- E feladat neve: ANCOVA

3.2.2 Eltérő tengelymetszet

Ha csak a tengelymetszetet térítjük el (+1 egység magyarázó változó hatása ugyanaz minden csoportban, de nem ugyanannyi a 0 magyarázó változóhoz tartozó eredményváltozó):

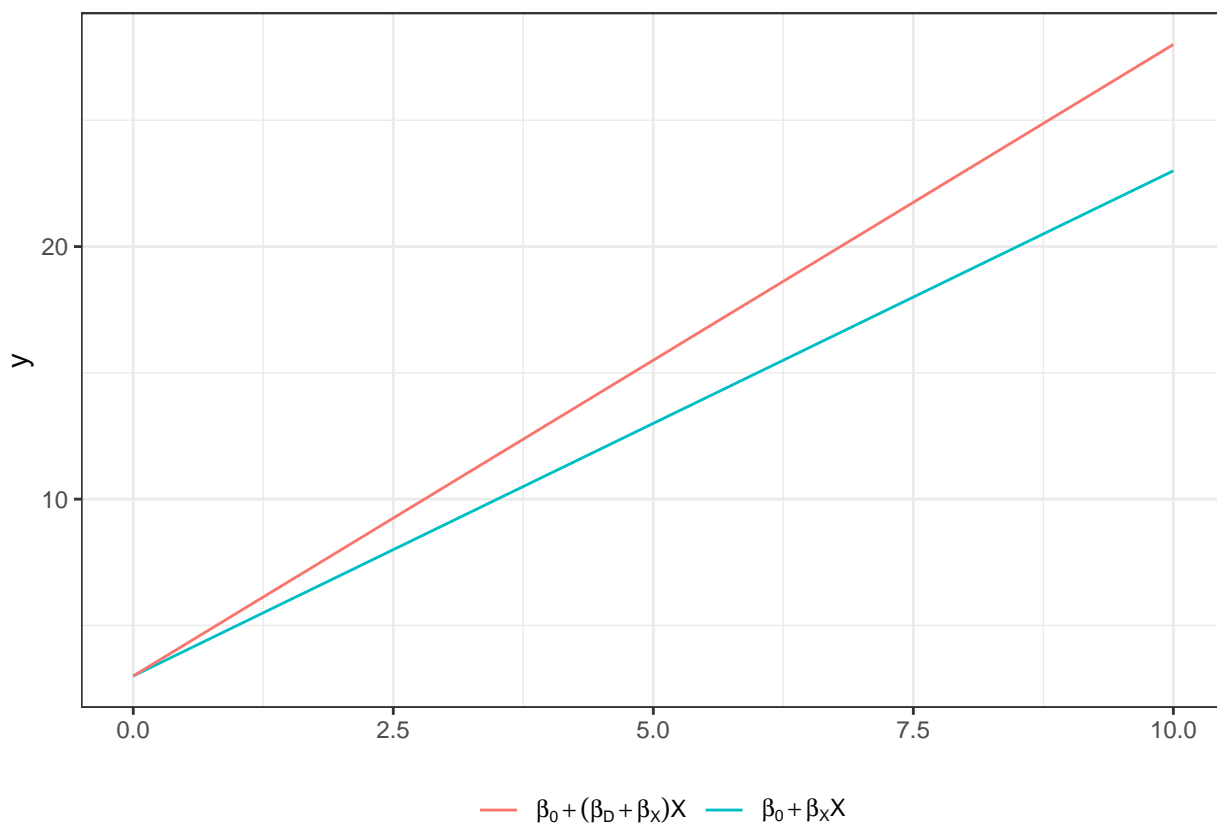


Algebrailag:

$$Y = \beta_0 + \beta_D D + \beta_X X + \varepsilon$$

3.2.3 Eltérő meredekség

Ha csak a meredekséget térítjük el (0 magyarázó változóhoz ugyanaz az eredményváltozó tartozik, de +1 egység magyarázó változó hatása csoportonként eltérő):



Algebrailag:

$$Y = \beta_0 + (\beta_X + \beta_D D) X + \varepsilon$$

3.2.4 Eltérő tengelymetszet és meredekség

- Akár a tengelymetszet és a meredekség is lehet különböző
- Ahogy előbb láttuk, csak a módszereket kell kombinálni: a konstans és a meredekséget is megdummyzzuk:

$$Y = \beta_1 + \beta_2 X + \varepsilon,$$

de úgy, hogy $\beta_1 = \alpha + \alpha_A D_A + \alpha_B D_B$ és $\beta_2 = \gamma + \gamma_A D_A + \gamma_B D_B$

- Nagyon fontos észrevenni, hogy a meredekség dummyzása a dummy és a mennyiségi változó közti interakcióra vezet:

$$Y = \alpha + \alpha_A D_A + \alpha_B D_B + \gamma X + \gamma_A (D_A X) + \gamma_B (D_B X) + \varepsilon$$

- Logikus is: az egyik változó (folytonos) hatása eltér a szerint, hogy a másik változónak (kategoriális) mi a szintje: különböző meredekségek
- Avagy fordítva elmondva (egyenértékűen, hiszen az interakció ugye szimmetrikus): az egyik változó (kategoriális) hatása eltér aszerint, hogy a másik változónak (folytonos) mi a szintje: az egyenesek közti különbség függ attól, hogy hol nézzük

3.2.5 Eltérő tengelymetszet és meredekség

- De hát ez megoldható a minta szétszedésével is!
- A két módszer – természetesen – ugyanarra az eredményre vezet
- A dummyzás mégis jobb a minta szétszedésénél; vajon miért? Mert messze-menően több lehetőségünk van a dummyzott (egybenlévő) modellel \rightarrow gazdaságilag releváns hipotézisek vizsgálhatóak egyszerűen (ld. mindjárt)

3.2.6 Hipotézisvizsgálat a dummyzott modellben

- Pl.: van-e egyáltalán bármilyen eltérés a csoportok között? (Értsd: eltér-e a becsült egyenes (bármilyen szempontból) a csoportok között, vagy mindegyikben teljesen ugyanaz?)
- Ez az ún. *strukturális törés*, hipotézispárja: $H_0 : \alpha_A = \alpha_B = \gamma_A = \gamma_B = 0$, H_1 : valamelyik ezek közül nem nulla, tehát van strukturális törés
- És most jön a szép rész: ha a fenti modellt megbecsültük (sima OLS-sel), akkor ez a hipotézis egyszerűen egy közösleges Wald- (vagy hasonló) próbát jelent!
- Hasonlóképp: nem lehet, hogy csak a tengelymetszetek eltérőek? \rightarrow ez az ún. *párhuzamos ráták* hipotézise, $H_0 : \gamma_A = \gamma_B = 0$; szintén Wald-tesztel elintézhető
- Minden hasonló, gazdasági kérdés *lefordítható* ökonometriailag, például változó vagy változók relevanciájának tesztelésére

3.2.7 Kontraszt-kódolás

- Kontraszt-kódolás: trükkös kódolás úgy kitalálva, hogy a dummy-k együtthatója ne a referencia-csoportéhoz, hanem az átlaghoz képesti eltérést jelentse
- A megoldás:

| | C_A | C_B |
|---|-------|-------|
| A | 1 | 0 |
| B | 0 | 1 |
| C | -1 | -1 |

- (A dummy változó nem 0 és 1 értéket vehet csak fel)
- Miért fog ez működni?

Mert:

$$\beta_0 + \beta_{C_A} + 0 = \bar{y}_A \quad (3.1)$$

$$\beta_0 + 0 + \beta_{C_B} = \bar{y}_B \quad (3.2)$$

$$\beta_0 - \beta_{C_A} - \beta_{C_B} = \bar{y}_C \quad (3.3)$$

És így:

- (1)+(2)+(3) $\Rightarrow 3\beta_0 = \bar{y}_A + \bar{y}_B + \bar{y}_C \Rightarrow \beta_0$ tényleg a főátlag (ha azonosak a csoportok elemszámai! különben ún. súlyozott kontraszt kellene, ahol a dummy változók már nem is feltétlenül egész értékeket vennének fel)
- (2)+(3) $\Rightarrow 2\beta_0 - \beta_{C_A} = \bar{y}_B + \bar{y}_C \Rightarrow \beta_{C_A} = 2\beta_0 - (\bar{y}_B + \bar{y}_C) = 2\beta_0 - (3\beta_0 - \bar{y}_A) \Rightarrow \beta_{C_A} = \bar{y}_A - \beta_0 \Rightarrow$ tényleg az átlagtól való eltérés (és hasonlóan a másik)

3.2.8 Egy terminológiai megjegyzés

- Az angol irodalomban az általunk kontrasztkódolásnak nevezett módszert nagyon gyakran “effect coding”-nak nevezik...
- ...a kontraszt pedig az, amikor a csoportok tetszőleges – általunk meghatározott – lineáris kombinációját teszteljük

4. fejezet

Nemlineáris modellek

4.1 Elöljáróban: a marginális hatás általánosabb értelmezése

4.1.1 A marginális hatás fogalma

- Marginális hatás: a magyarázó változó kis növelésének hatására mekkora az eredményváltozó *egységnyi magyarázóváltozó-növelésre jutó* változása
- Tipikus egyszerűsítés: a magyarázó változó egységnyi növelésének hatására mennyit változik az eredményváltozó
- (Hiszen a kettő ugyanaz, ha a változó hatása lineáris)
- Idáig az i -edik magyarázó változó ilyen módon értelmezett marginális hatása és a β_i számértéke gyakorlatilag szinonima volt

4.1.2 A marginális hatás precízebben

- Definíció alapján a marginális hatás: $\frac{\Delta Y}{\Delta X_j}$, ha ΔX_j kicsiny
- Ugye egyetemen vagyunk \rightarrow a marginális hatás $\frac{\partial Y}{\partial X_j}$
- A többváltozós lineáris regresszió eddigi (sokasági) modelljében $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon$, ezért

$$\begin{aligned}\frac{\partial Y}{\partial X_j} &= \frac{\partial}{\partial X_j} [\beta_0 + \beta_1 X_1 + \dots + \\ &\quad + \dots + \beta_{j-1} X_{j-1} + \beta_j X_j + \beta_{j+1} X_{j+1} + \dots + \beta_k X_k + \varepsilon] = \\ &= \beta_j\end{aligned}$$

- ... hát ezért tekinthettük eddig a marginális hatást és a becsült regressziós koefficiens szinonimának!

4.2 A linearitás feloldása

4.2.1 Emlékeztetőül: a linearitás következményei

- A linearitás két dolgot vont maga után:
 - Mindegy, hogy honnan indulva növelem a változót egy egységgel (a változó hatása lineáris)
 - Mindegy, hogy a többi változó milyen szinten van rögzítve (additivitás)
- E kettőt fogjuk most feloldani

4.2.2 Az additivitás feloldása: az interakció

- Eddigi modellünkben a marginális hatások a többi változó szintjétől függetlenül állandóak voltak
- Például: 1 Ft pluszjövedelem taglétszámtól függetlenül azonos többletkiadást jelent...?
- Ha nem, akkor azt mondjuk, hogy a két változó között *interakció* van: az egyik marginális hatásának *nagyságát* befolyásolja a másik *szintje*
- A kapcsolat tehát marginális hatás és szint között van (nem marginális hatás és marginális hatás vagy szint és szint között!)
- Kézenfekvő indulás: az egyik változó szintje *lineárisan* hasson a másik marginális hatására; sokaságban felírva:

$$(\beta_J + \beta_{JT}\text{Tag}) \text{Jov},$$

ahol β_{JT} az interakció hatását kifejező (lineáris) együttható - Helyezzük ezt be a (sokasági) regresszióba:

$$Y = \beta_0 + (\beta_J + \beta_{JT}\text{Tag}) \text{Jov} + \beta_T \text{Tag} + \varepsilon,$$

azonban felbontva a zárójelet:

$$\begin{aligned} Y &= \beta_0 + \beta_J \text{Jov} + \beta_{JT} \text{Tag} \cdot \text{Jov} + \beta_T \text{Tag} + \varepsilon = \\ &= \beta_0 + \beta_J \text{Jov} + (\beta_T + \beta_{JT} \text{Jov}) \text{Tag} + \varepsilon \end{aligned}$$

- Tehát az interakció *szükségképp*, automatikusan “szimmetrikus”: ha az egyik változó szintje hat a másik marginális hatására akkor szükségképp fordítva is: a másik szintje is hatni fog az előbbi marginális hatására
- Azaz “egyszerre” lesz igaz, hogy $(\beta_J + \beta_{JT}\text{Tag}) \text{Jov}$ és $(\beta_T + \beta_{JT}\text{Jov}) \text{Tag}$: attól függően, hogy milyen szempontból nézzük (melyik marginális hatását vizsgáljuk, ezt még ld. később is)

- A regresszióban így elég egyszerűen ennyit írni:

$$\beta_T \text{Tag} + \beta_J \text{Jov} + \beta_{JT} (\text{Jov} \cdot \text{Tag}).$$

- ... mindkét – másik szintjétől függő – marginális hatás ebből kiadódik, függően attól, hogy hogyan bontjuk fel a zárójelet (melyik változót vizsgáljuk)
- Ez a marginális hatás pontosabb értelmezése mellett még szebben látható

4.2.3 A marginális hatás interakciók esetén

- Ha interakció van, például a l -edik és az m -edik tag között, akkor az l -edik marginális hatása:

$$\begin{aligned} \frac{\partial Y}{\partial X_l} &= \frac{\partial}{\partial X_l} [\beta_0 + \beta_1 X_1 + \dots + \\ &\quad + \dots + \beta_l X_l + \dots + \beta_m X_m + \dots + \beta_k X_k + \beta_{lm} X_l X_m + \varepsilon] = \\ &= \beta_l + \beta_{lm} X_m \end{aligned}$$

- Így precíz az előbbi állításunk arról, hogy ha az egyik szerint vizsgáljuk a marginális hatást, akkor az a másik szintjétől fog függeni (gondoljuk hozzá a másik szerinti deriválást is!)

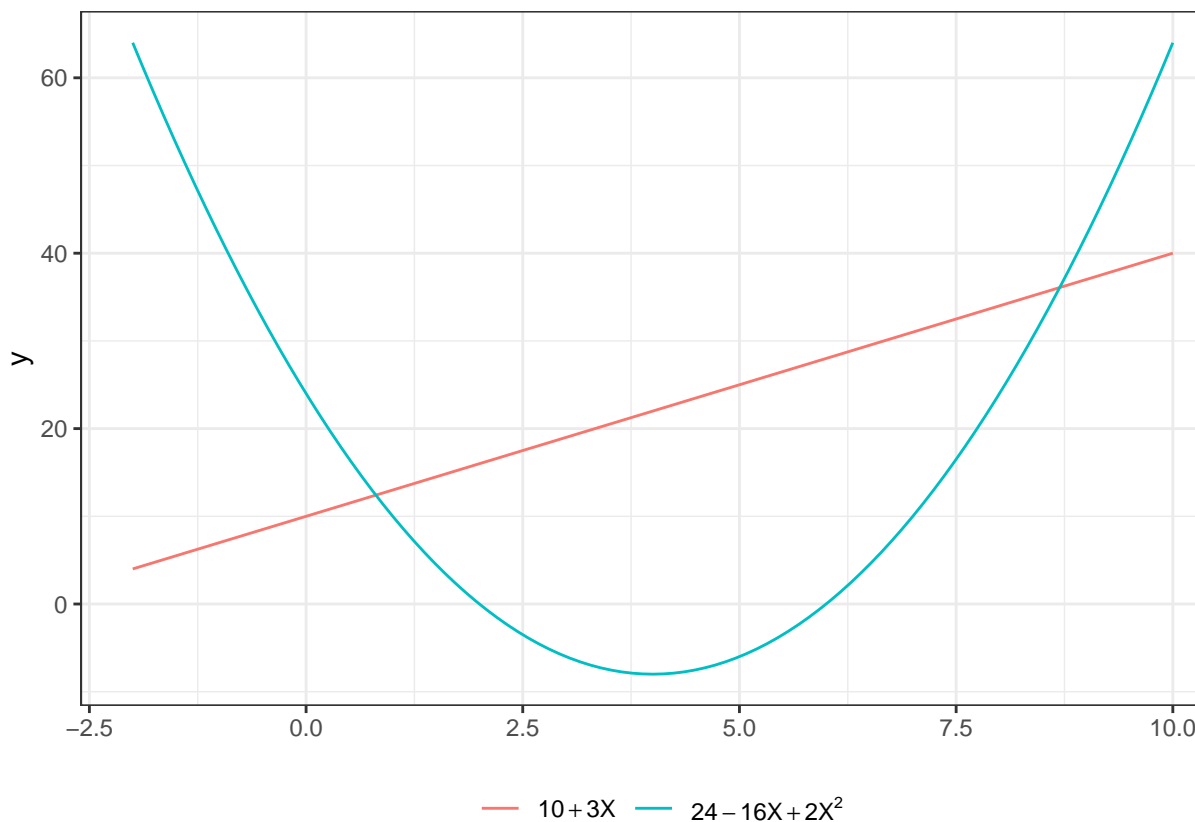
4.2.4 A változónkénti linearitás feloldása: egy motiváló példa

- Már volt: mit jelent az, ha megsértjük a “marginális hatás nem függ attól, hogy a többi magyarázó változót milyen szinten rögzítjük” következményét a linearitásnak
- És ha a “marginális hatás nem függ attól, hogy milyen szintről indulva növeljük a változót” következményt szeretnénk feloldani?
- Például: 1 évvel idősebb életkor kiinduló életkortól függetlenül azonos kiadásváltozást jelent...?
- Használjunk a lineáris függvényforma helyett mást, például négyzeteset (parabolát):

$$\frac{\partial}{\partial X_j} [\dots + \beta_j X_j + \beta_{jj} X_j^2 + \dots] = \beta_j + 2\beta_{jj} X_j$$

4.2.5 A változónkénti linearitás feloldása: egy motiváló példa

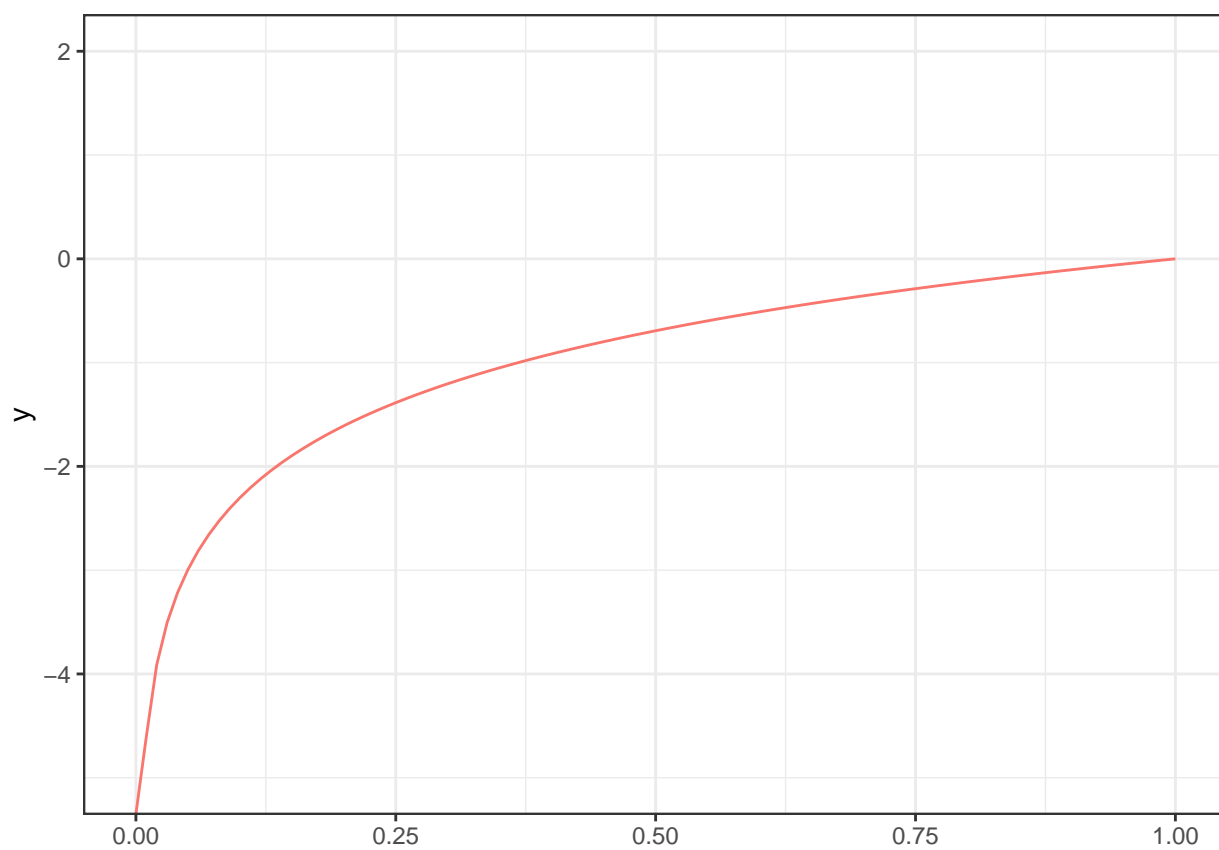
Szemléletesen az egy magyarázó változós esetben:

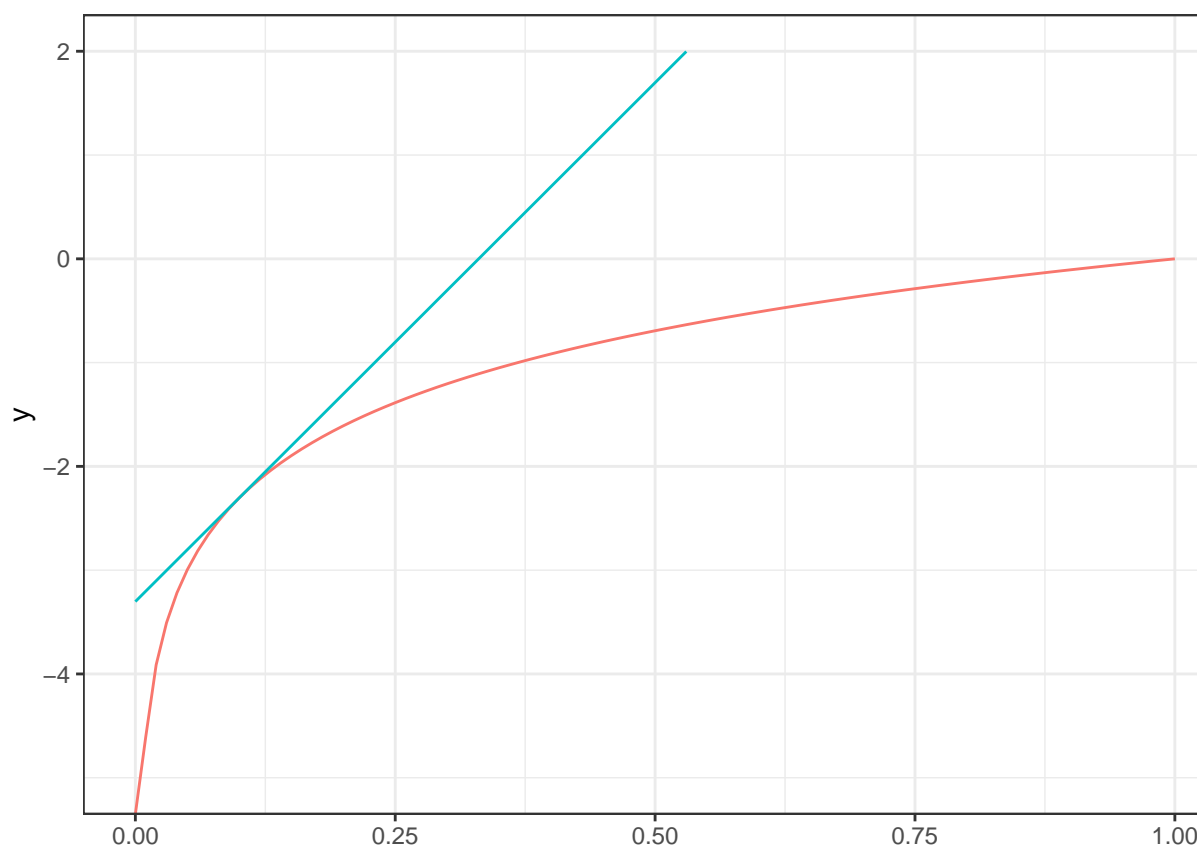


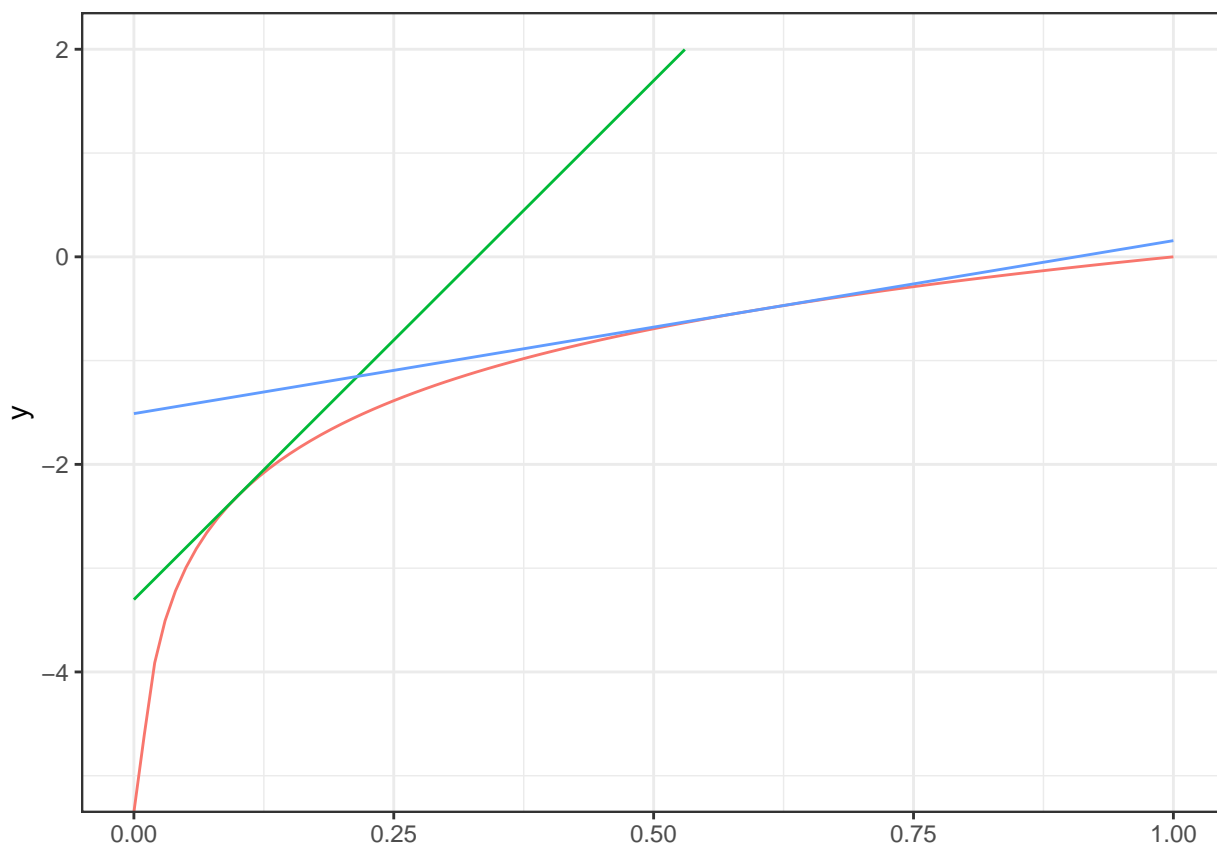
Szélsőérték hely nyilvánvaló (első derivált előjelet vált): $\beta_j + 2\beta_{jj}X_j = 0 \Rightarrow X_j = -\frac{\beta_j}{2\beta_{jj}}$

4.2.6 Linearitás, mint közelítés

- Az élet általában nemlineáris
- Miért használunk mégis lineáris modelleket: mert sokszor nem térnek el (nagyon) a valóságtól, de mégis sokkal könnyebben kezelhetőek matematikailag
- Ez tehát az esetek többségében egy közelítés
- Mint ilyen: vizsgálni kell az érvényességi határokat
- “Munkaponti linearizálás”







4.2.7 Érvényességi határok

- Az érvényességi határokat az eddig látott modellekben is érdemes végiggondolni
- Azonnal kézenfekvő példa: a konstans (nagyon sok esetben)
- De sok meredekségnél is megragadható ez (fogyasztási függvény példája)
- Ez is egyfajta munkaponti linearizálás

4.2.8 Nemlinearitás fajtái

- Az $\beta_1 + \beta_2 X + \beta_3 X^2$ egy nemlineáris kifejezés (matematikailag)
- De figyelem: ennek ellenére minden további nélkül, tökéletesen kezelhető pusztán az eddig látott (lineáris!) eszköztárral, hiszen az OLS-nek mindegy, hogy a második magyarázó változó értékei történetesen épp az első négyzetei
- (Egészen addig nincs baj, amíg a kapcsolat nem lineáris)
- Nem úgy mint a $\beta_1 X^{\beta_2} \rightarrow$ ez nem becsülhető OLS-sel
- A megkülönböztetés végett az első esetet változójában, a másodikat paraméterében nemlineáris modellnek nevezzük

- Mi van “nemlinearitást okozó pozícióban”

4.2.9 Változójában nemlineáris modell

- Jellemző: továbbra is fennáll a “változók konstansokkal szorozva majd összeadva” (tehát: lineáris kombinációs) struktúra
- De elképzelhető, hogy egy változó egy “eredeti” változó transzformáltja
- Itt szükségképp nemlineáris transzformációról beszélünk!
- Vegyük észre, hogy az “eredeti” és “transzformált” közti megkülönböztetés teljesen mesterséges (csak mi tudjuk, hogy mi volt az adatbázisban bemenő adatként), az OLS-nek mindegy
- Ide tartozik a kvadratikus hatás, általában az X^a magyarázó változók, a $\log_a X$, az a^X stb., ahol a konstans
- Az előzőek miatt a becslés ugyanaz, egyedül az interpretálás igényel további tárgyalást

4.2.10 Paramétereiben nemlineáris modell

- Megsérti a lineáris kombináció struktúráját: paraméter nem csak szorzóként szerepel a regresszióban
- Például X^β , $\log_\beta X$ stb.
- Ez már nem becsülhető OLS-sel: az eredményváltozó nem állítható elő mátrixműveletekkel
- Más módszert fogunk használni

4.2.11 Interakció és kvadratikus hatás revisited

- Az előzőek fényében nyilvánvaló: a kvadratikus hatás egyfajta (igen egyszerű) változójában nemlineáris modell
- Az interakció szintén változóbeli nemlinearitás, de nem annyira kézenfekvő módon (mindenképp indokolt a külön tárgyalása)

4.2.12 Nemlinearitás kezelése: NLS

- Vegyük észre, hogy a $\min_\beta ESS$ célfüggvény akkor is tartható, ha nemlineáris modellt specifikálunk!
- (Csak az ESS számításához szükséges \hat{Y} -ok másképp jönnek ki, de ez a fenti optimalizáció szempontjából *teljesen mindegy*)
- Oldjuk meg ezt az optimalizációs feladatot!
- Ez a nem-lineáris legkisebb négyzetek (NLS, non-linear least squares) módszere
- Sajnos mondani könnyebb, mint a gyakorlatban kivitelezni; szemben a lineáris specifikációval, a kritériumfelület nem kvadratikus, emiatt nincs egyetlen művelettel megtalálható optimum
- Van-e egyáltalán egyértelmű (globális) optimum? Mi van, ha több lokális optimum létezik?

4.3. NÉHÁNY NEVEZETES, PARAMÉTERÉBEN NEMLINEÁRIS MODELL 61

- Ettől el is tekintve, a konkrét optimalizáció számos gyakorlati problémát vehet fel, mivel valamilyen iteratív algoritmus kell
- Több lehetőség van, különféle előnyökkel és hátrányokkal (Gauss–Newton keresés, Levenberg–Marquardt algoritmus, konjugált gradiens keresés stb.), de mind rengeteg numerikus kérdést vet fel:
 - Meg tudjuk találni az optimumot? Biztosan? (Lehet-e baj a konvergenciával? Mi legyen a konvergencia-kritérium?)
 - Mennyi idő alatt találjuk meg?
 - Milyen kezdőértékből induljunk? (Milyen a módszer numerikus stabilitása?)
 - stb. stb. stb.

4.2.13 Nemlinearitás kezelése: algebrai linearizáció

- Mi a fenti (mindig alkalmazható) módszerrel szemben egy másik (könnyebb, de nem mindig alkalmazható) módszert fogunk vizsgálni: algebrai linearizálás
- Alkalmas transzformációval a nemlineáris problémát lineárisra alakítjuk, azt OLS-sel megoldjuk, majd a kapott eredményeket visszatranszformáljuk az eredeti transzformáció inverzével
- Például: $Y = \beta_1 X^{\beta_2} \varepsilon$ paramétereiben nemlineáris. . .
- . . . de mindkét oldal logaritmusát véve $\log Y = \log \beta_1 + \beta_2 \log X + \varepsilon'$ már az!
- Adatbázis logaritmálása, eredmények visszahatványozása
- (Amint mondtuk, nem mindig alkalmazható, de azért nagyon sok, gyakorlatilag fontos esetben igen)
- Természetesen itt is eltérő, specifikus értelmezések jelenthetnek meg

4.3 Néhány nevezetes, paraméterében nemlineáris modell

4.3.1 Log-log modell

- Például a Cobb–Douglas termelési modell:

$$Y = \beta_1 L^{\beta_L} K^{\beta_K} \varepsilon,$$

ahol Y a kibocsátás, L a munka, K a tőke (ill. általában a termelési tényezők) felhasználása - Elaszticitása:

$$\text{El}_L(L, K) = \frac{\frac{dY}{Y}}{\frac{dL}{L}} = \frac{dY}{dL} \frac{L}{Y} = \beta_1 \beta_L L^{\beta_L-1} K^{\beta_K} \frac{L}{\beta_1 L^{\beta_L} K^{\beta_K}} = \beta_L$$

- Ezért nevezik *konstans elaszticitású* modellnek is
- Kezelése linearizálással: mindkét oldalt logaritmáljuk

$$\log Y = \log \beta_1 + \beta_L \log L + \beta_K \log K + \varepsilon'$$

- Minden változót (eredmény és összes magyarázó is) logaritmálni kell
- Innen a modell neve
- Csak a konstans lesz logaritmálva, a többi koefficiens a transzformáció ellenére (ill. épp azért...) közvetlenül kapjuk
- Volumenhozadék (skáláhozadék): $\beta_K + \beta_L$ viszonya 1-hez

4.3.2 Log-lin modell

- Például a jövedelem alakulása:

$$Y = e^{\beta_1 + \beta_2 X + \varepsilon}$$

- Linearizálás ismét mindkét oldal logaritmálásával:

$$\log Y = \beta_1 + \beta_2 X + \varepsilon$$

- Elnevezés logikája így már látható: az eredményváltozó logaritmálva, de a magyarázó változók maradnak szintben
- Növekedési ráta: $e^{\beta_1 + \beta_2(X+1)+u} = Y e^{\beta_2}$, pillanatnyi növekedési ütem: $\beta_2 = \frac{d \log Y}{dX} = \frac{1}{Y} \frac{dY}{dX}$
- Elaszticitás: $\text{El}_X(X) = \frac{dY}{dX} \frac{X}{Y} = \beta_2 X$, tehát csak X -től függ

4.3.3 Lin-log modell – kakukktojás!

- Az előzőek alapján már világos a jelentése (pl. terület és kínálati ár összefüggése):

$$Y = \beta_1 + \beta_2 \log X + \varepsilon$$

- Miért kakukktojás?
- β_2 értelmezése:

$$\frac{dY}{dX} = \frac{\beta_2}{X} \Rightarrow \beta_2 = \frac{dY}{dX/X}$$

- Elaszticitás:

$$\text{El}_X(X) = \frac{\beta_2}{X} \frac{X}{Y} = \frac{\beta_2}{Y},$$

tehát csak Y -től függ (közvetlenül)

A “közvetlenül” itt arra utal, hogy természetesen az Y helyébe beírható lenne a – csak X -eket tartalmazó – felírása.

4.3.4 Reciprok modell – kakukktójas

- Például keresleti modell:

$$Y = \beta_1 + \frac{\beta_2}{X} + \varepsilon$$

- Miért kakukktójas?
- Aszimptotikusan: $\lim_{X \rightarrow \infty} \mathbb{E}(Y | X) = \beta_1$
- Határkiadás: $\frac{dY}{dX} = -\frac{\beta_2}{X^2}$
- Elaszticitás:

$$\text{El}_X(X) = -\frac{\beta_2}{X^2} \frac{X}{Y} = \frac{\beta_2}{XY}$$

- Paraméterek értelmezése, β_2 előjelének jelentősége az “aszimptotikus” viselkedés szempontjából: az élvezeti cikkek példája

4.3.5 Egy komplex példa vegyes modellre

- Tekintsük a következő kiadási modellt:

$$Y = \beta_1 e^{\beta_N N} J^{\beta_J + \beta_{JN} N} \varepsilon,$$

ahol Y a kiadás, N a nem (dummy, 1 ha férfi), J a jövedelem - Alapvetően log-log jellegű (bár a β_N felfogható exponenciálisként is), ráadásul még dummyzva is, hogy minden eltérő legyen nemenként (strukturális törés) - Ez utóbbi miatt mondhatjuk egyszerűen, hogy log-log jellegű - Linearizálás:

$$\log Y = \log \beta_1 + \beta_N N + \beta_J J + \beta_{JN} (J \cdot N) + \varepsilon'$$

- Paraméterértelmezések:
 - β_1 és $\beta_N \rightarrow$ autonóm fogyasztás *mindkét* nemre teljesen külön (exponenciális jelleg ebben)
 - * A jövedelemnél célszerű az elaszticitást megragadni (a log-log jelleg miatt), hiszen: $\text{El}(Y, J) = \beta_J + \beta_{JN} N$, azaz az elaszticitás nemtől függő (avagy: mindkét nemre teljesen külön elaszticitás, referencia-kódolás jelleggel)

4.4 Specifikációs tesztek

4.4.1 A specifikációs tesztek

- Itt már nagyon erősen felmerül a kérdés: hogyan dönthetünk a különféle függvényformák között?
- Ld. a termelési függvény példáját \rightarrow megadható lineárisan és Cobb-Douglas jelleggel (eredmény nagyon nem mindegy)

- Hogyan lehet *analitikusan* dönteni?
- Az előző példára: BM-teszt, PE-teszt stb.
- Általánosságban (nem csak log/lin kérdésekre, mint az előzőek): ún. specifikációs tesztek

4.4.2 Egy egyszerű specifikációs teszt

- Egy egyszerű ötlet: adjuk hozzá a magyarázó változókhoz a magyarázó változók valamilyen nemlineáris transzformáltját (tipikusan négyzeteiket vagy logaritmusait)...
- ...és nézzük meg, hogy *együttesen* szignifikánsak-e!
- Ha igen, az specifikációs hibára utal
- (Tehát figyelem, ezt alapvetően nem arra használjuk, hogy arra következtessünk, hogy egy adott konkrét változó négyzetét vagy logaritmusát hozzá kell-e adni a függvényformához, hanem összességében nézzük őket, specifikációs tesztként)
- Hátrány: sok szabadságfokot használ el, és csak elég speciális alakú nemlinearitásokkal tesztel

4.4.3 Ramsey RESET-je

- A modellspecifikáció *általános* tesztje; emiatt előnye: nem egy adott specifikációs kérdésre keres választ, hanem általában vizsgálja, hogy a specifikáció jó-e; hátránya, hogy ha nemleges választ ad, nem derül ki, hogy pontosan mi a specifikáció baja
- Trükk: új regressziót becsül, melynek eredményváltozója ugyanaz, de a magyarázó változókhoz hozzáadja az eredeti regresszió becsült eredményváltozójának magasabb hatványait (\hat{Y}^3 -ig néha \hat{Y}^4 -ig is):

$$Y = \beta'_0 + \beta'_1 X_1 + \beta'_2 X_2 + \dots + \beta'_k X_k + \gamma_1 \hat{Y}^2 + \gamma_2 \hat{Y}^3 + \varepsilon'$$

- Amit tesztelni kell az előzőhöz hasonlóan: $H_0 : \gamma_1 = \gamma_2 = 0$ (F -próba és LM-próba is van, heteroszkedaszticitásra is robusztussá tehető)
- Ilyen módon takarékos a szabadsági fokokkal, csak 2-3-at használ el (a fenti trükkel “összesűríti” a magyarázó változókat), ráadásul általánosabb alakú nemlinearitások is beleférnek ebbe, mint a csak négyzetekkel/logaritmusokkal tesztelés
- Specifikációs teszt, tehát kihagyott változó detektálására általában nem alkalmas