

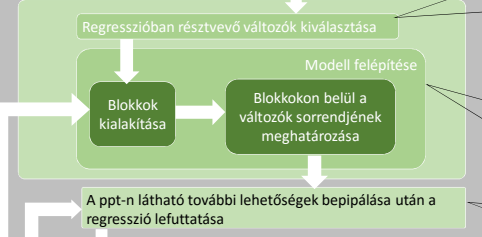
# Kezdőpont

Az adatokat felkészítése, adattisztítás, adatokkal való általános megismerkedés (leíró statisztikák és grafikonok) megtörtént

A flowchart a regresszió ppt-vel együtt használandó, itt csak a lépéseket sorolom fel rövid magyarázattal, a részletek, definíciók a ppt-ben olvashatóak. A karikában elhelyezett számok a ppt-ben felsorolt feltételeket jelölik.

Az adatfeldolgozás lépéseit nem részletezem, megtalálható az adatfeldolgozas\_flowchart-on.

## A regressziós modell



1 A kimutatható hatások függnek a prediktor változók számától és az elemszámtól: minél több prediktor változónk van, annál nagyobb elemszám szükséges a hatások kimutatásához. Ezért nem érdemes az adatbázisunkban lévő minden változót a regresszióba beleerőltetni, csak azokat, melyekre hipotézisünk van. Érdemes lehet az egymáshoz közeli fogalmakat mérő változókat akár előzetesen összevonni főkomponens elemzéssel.

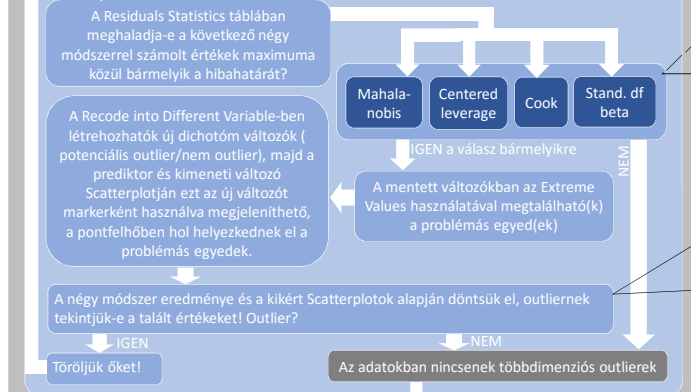
2 Független változónak csak skála típusú változó választható, prediktor változónak skála típusú vagy dichotóm. Lehetséges nominális változót is prediktorként használni, ha dummy változóra bontjuk.  
3 Triviális, hogy nulla varianciájú változót ne tegyünk a modellbe.

Előzetes tudás alapján alkossunk a prediktorközből blokkokat, ezáltal megadva, mely változók kerüljenek előbb vagy később a modellbe. Egy blokkon belül a változók modellbe lépését beállíthatjuk ENTER-re (az egy blokkon belüli változók egyszerre lépnek a modellbe), vagy választhatunk több STEPWISE módszer (a változók modellbe lépési sorrendjét matematikailag határozzuk meg) közül. A pszichológiában ENTER-t szoktuk választani, mert így a prediktorkor összehasonlíthatóak lesznek. A stepwise kutatásra kevésbé alkalmas módszer, mert a változók modellbekerüléséről egymás után döntünk, így előfordulhat, hogy egy változó azért nem kerül a modellbe, mert amit ő meg tudna a varianciából magyarázni, az a többi változó már megmagyarázta.

Nem részletezem, mert a modell felépítésén kívül nincs más beállítási lehetőségünk, a további lehetőségeket a regressziós modellen nem változtatnak, csak azt állítja be, milyen plusz táblázatokat (például a feltételek ellenőrzéséhez tartozó táblákat) kérünk. Ez a ppt-ből követhető.

## Feltételek ellenőrzése az output táblázatai alapján

### Többdimenziós outlierek kiszűrése



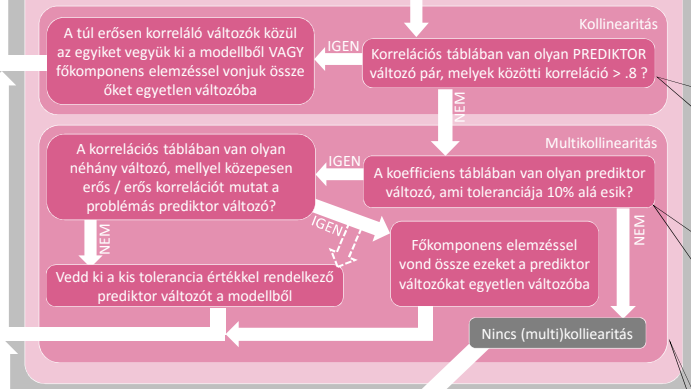
A regresszió feltételeinek teljesülése hat egymásra, ezért teljesülésüket egyszerre kell figyelni. Itt egy viszonylag optimális sorrendben követik egymást az ellenőrzés lépései. Ha valamelyik feltétel sérülése miatt kivesszünk egy egyedet a mintából vagy egy változót a vizsgálatból, akkor az egész elemzést újra kell futtatni, és minden feltételt újra kell ellenőrizni.

1 A többdimenziós outlierek olyan értékek, melyek egy dimenzióban nem szélsőségesek, csak kettő vagy több dimenzióban való együtt előfordulásuk az (pl. 190cm és 50kg). Szűrőnkre négy módszer is van. Mindegyik módszernek eltérő kritériumszintje van, mely fölött outlierek itélik az egyedet. Ezeket a vizsgálatonként változó határértékeket ki kell a ppt-ben olvasható képletek alapján számolni. Mind a négy módszer létrehoz egy adatbázisunkban egy új változót (tehát létrejön egy Mahalanobis, Cook, stb. változó).  
2 A Residual Statistics tábla az új változók maximumát mutatja (kivéve a Stand. Df beta, mely maximumát nekünk kell kikérni). Tudjuk, hogy a mintában egy adott módszer alapján nincs outlier, ha a táblázatban a maximum a módszerhez tartozó határérték alatt van. Ha a maximum meghaladja a határértéket, akkor legalább egy többdimenziós outlier jelez a módszer. Hogy a mintában hol vannak ezek a többdimenziós outlierek, azt az Extreme Values listájából tudhatjuk könnyen meg, amit az Analyze / Descriptive Statistics / Explore-on belül az Opciókban az Outliers beépítésével kérhetünk ki.

A nehézséget az jelenti, hogy a négy módszer eredménye ritkán fed át egymással, és nekünk kell eldönteni, melyik eredménynek higgyünk, mikor mérlegeljük, kiket tekintünk outlierek. Érdemes mind a négy módszert csupán indikátorként tekinteni, és együttes értelmezésükkel, valamint grafikonok elemzésével dönteni az outlierekről. Grafikon a következőképpen lehet kialakítani: Tegyük fel, hogy a Cook távolság alapján van a mintában öt outlier, a Cook távolság határértéke 1. A Recode into Different Variable-ben létrehozható új dichotóm változó, aminek az értéke 0, ha a Cook távolság < 1, ha a Cook távolság > 1, nevezük az új változót Cook indikátornak. Ezt követően kérjük ki Scatterplotokat, melyek a prediktor változó és kimeneti változó kapcsolatát ábrázolják, és Markernek állítsuk be az új változót, ekkor a pontfelhőben ezek az értékek más színnel lesznek megjelenítve, és így láthatjuk, hogy a pontfelhőben hol helyezkednek el a potenciális többdimenziós outlierek.

A négy módszer eredménye és a kikért Scatterplotok alapján döntünk el, outlierek tekintjük-e a talált értékeket, ha igen, töröljük őket. **Ha törölünk, az egész elemzést újra kell futtatni!**

### A (multi)kollinearitás

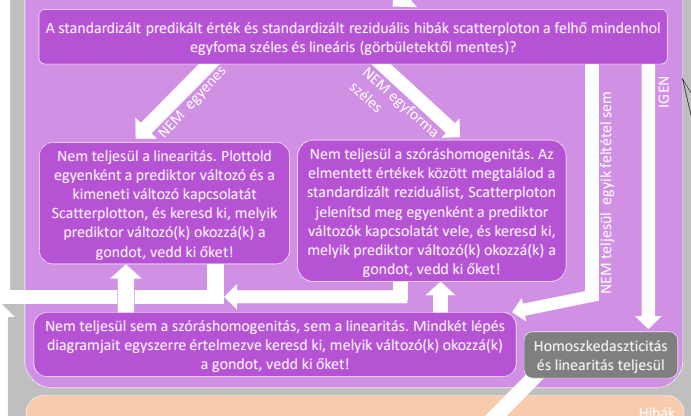


3 Következésként a kollinearitás és multikollinearitás hiányát érdemes ellenőrizni, mert a kettő feltétel még nem (kollinearitás) vagy csak kis mértékben (multikollinearitás) a kimeneti változó függvénye, jóval inkább a prediktor változók egymással való kapcsolatához köthető. Mind a kettő azért probléma, mert a prediktor változók hatásának értelmezését teszik nehézé.

Kollinearitásról két változó közötti nagyon erős ( $r > .8$ ) korreláció esetén beszélünk, ilyenkor a két változó „szinte ugyanazt méri”. Probléma, mert a kimeneti változó varianciájának is nagyrészt ugyanazon részét fogják megragadni, a variancia ugyanazon részéért versengenek. Ha enter módszerrel dolgozunk, kollinearitás esetén minden prediktor változó különböző hatását alul fogja becsülni, stepwise módszer esetén pedig az a változó, mely később kerül sorra, más nem tud jelentősen hozzáadni a megmagyarázott varianciához, ezért végül be sem kerül a modellbe, amit könnyű tévesen gyanú magyarázóerőként értelmezni. Kollinearitás esetén a két változó közül az egyiket ki kell hagyni az elemzésből vagy főkomponens elemzéssel egyetlen közös változóvá alakítani őket. **Mindektő esetén az egész elemzést újra kell futtatni!**

Multikollinearitás akkor fordul elő, ha több prediktor változó közepesen / erősen korrelál egymással. Ilyenkor előfordulhat, hogy lesz olyan prediktor változó, melynek varianciáját a többi változó szinte teljesen megmagyarázza. Ilyenkor kicsi lesz a változó független varianciája (az a varianciarész, amit csak ő magyaráz a kimeneti változóból), mert az általa megmagyarázott variancia nagy részét a többi változó is lefedti. A probléma az, hogy az ilyen változó magyarázó erejét alulbecsüljük, és a hozzá tartozó regressziós együtthatók sokszor teljesen értelmezhetetlenné válnak. Ellenőrizni a Koefficiens tábla tolerancia értékével lehet, ami százalékban adja meg, hogy a prediktor változó független varianciája mekkora része a megmagyarázott varianciának. A tolerancia 10% alatt jelez problémát. Ekkor a problémás prediktor változót ki kell venni a modellből vagy ha a korrelációs tábla alapján kiválasztható több erősen korreláló változó, akkor ezeket érdemes lehet főkomponens elemzéssel összevonni. **Mindektő esetén az egész elemzést újra kell futtatni!**

### Homoskedaszticitás és linearitás



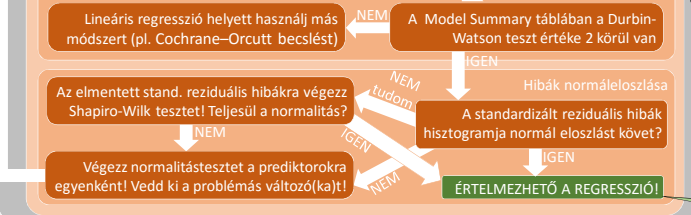
4 Bár a feltétel teljesüléséért inkább a vizsgálat megtervezésekor kell biztosítani, itt említem meg a külső változó problémáját, mely egy olyan tulajdonság, melyet a vizsgálatban nem mérünk, de több prediktor változóra is hatással van (például a kor változó gyerekeknel a súly és magasság prediktora), ezáltal a prediktor változók között (multi)kollinearitást okozhat. Végezzünk a probléma ellen jól átgondolt vizsgálati elrendezéssel lehet.

A homoskedaszticitás és linearitás feltételét egyszerre tárgyalom, mert ugyanazon a pontdiagramon tudjuk a két feltétel teljesülését ellenőrizni.

5 A homoskedaszticitás tulajdonképpen a szóráshomogenitás feltételét jelenti, de két vagy több prediktor változó esetén nehezebb átlátni, mit is jelent. A regresszió feltétele, hogy a predikált értékekhez tartozó reziduális hiba mindenhol homogén legyen: ha végighaladunk a predikált értékeken (tehát a regressziós egyenesen vagy síkon vagy többdimenziós téren), akkor az egyenes vagy sík körül a pontfelhő mindenhol egyforma vastag, a predikált értékek és mért értékek távolsága mindenhol hasonló. Ha nem teljesül a homoskedaszticitás, az azt fogja eredményezni, hogy a regresszió által tett becsülésem mindenhol lesz egyforma pontos: ahol kicsi a hiba, ott pontos a becsülés, ahol nagy az eltérés a predikált és mért értékek között, ott pontatlan a becsülés.

6 A linearitás szintén a predikált értékek és reziduálisok pontdiagramján ellenőrizhetjük. Ezen a diagramon azt várjuk, hogy a pontfelhő vízszintesen, egyenes vonalban helyezkedjen el. Kikértünk egy másik pontdiagramot is, mely a predikált érték és a kimeneti érték kapcsolatát szemlélteti, itt triviális módon pozitív korrelációra jellemző egyenes, növekedő pontfelhőt várunk. A linearitás feltételét fontos ellenőrizni, mert ha a prediktor változók nem lineáris kapcsolatban vannak a kimeneti változóval, akkor a lineáris regresszióval alul fogjuk a magyarázóerőt becsülni. Ha a két feltétel bármelyike nem teljesül, a dolgunk ugyanaz: menjünk végig az összes prediktor változón, nézzük meg pontdiagrammal a kimeneti változóval vagy a reziduálisval való kapcsolatát, és nézzük meg, melyik prediktor változónál sérül a kérdéses feltétel! A problémás prediktor változót vegyük ki a modellből, majd futtassuk újra az elemzést!

### Hibák



A hibák tulajdonságainak ellenőrzése a végére került, mert e feltétel teljesülését szinte az összes korábban felsorolt feltétel befolyásolja. Eddig is használtuk, de még nem definiáltuk: a reziduális hiba a kimeneti változó varianciájának azon része, melyet a regressziós modellünkkel nem tudunk megmagyarázni, a kimeneti változóban olyan változatosság, melyet a prediktor változók nem tudnak bejósolni. Ez már elméletileg csak zaj, egyéni különbségekből, véletlenekből adódó változatosság, tehát nem szabad semmilyen szisztematikuságot, mintzatot tartalmaznia.

7 A hibák függetlensége alapján a hibák nem korrelálhatnak, az egyik értékhez tartozó hiba nem befolyásolhatja a másikat. Pszichológiában a hibák korrelálatlansága ritkán jelent problémát, mivel nem szoktunk idősorokat elemezni. Más a helyzet a gazdaságban, ahol például ha egy részvény változását szeretnénk bejósolni más részvényekből, hiteken át gyűjtöm az adatokat, evidens, hogy a részvények egyik napi értéke nem lesz független a következő napitól.

8 A hibák normáloszlásának feltétele szintén abból adódik, hogy az feltételezzük, hogy a regresszió végén megmaradó megmagyarázatlan variancia már csak zaj, véletlenszerű, tehát normál eloszlást követ.

Ha ide eljutottál, a lineáris regresszió feltételei teljesülnek. Most jön az igazi munka.

### Lineáris regresszió feltételei

- 1) Változók típusa
- 2) Nem-nulla variancia
- 3) Nincs (multi)kollinearitás
- 4) Nincs külső változó
- 5) Homoskedaszticitás
- 6) Linearitás
- 7) Független hibák
- 8) Hiba normális eloszlása
- 9) Elemszám
- 10) Nincsenek (több)dimenziós outlierek

Ne feledjétek, ez a flowchart csak útmutatóként használható, a statisztikai elemzés mindig értelmező döntések sorozatából áll, így az itt felvázolt lépések mechanikus követése nem feltétlenül vezet jó megoldáshoz.

Várhelyi Klára – stathelp.hu