

Lineáris regresszió elméleti összefoglaló

Bognár Miklós

Bevezetés az Ökonometriába

0.1 Lineáris algebra, valószínűségi vektorváltozók és mátrixdifferenciálás összefoglaló

A lineáris regresszió megértéséhez elengedhetetlen, hogy tisztában legyünk néhány, lineáris algebrából ismeretes fogalommal és összefüggéssel. Ezen felül nagyon hasznos, ha ismerjük, hogy hogyan kezelendőek a valószínűségi vektorváltozók illetve a mátrixdifferenciálás-kifejezések.

0.1.1 Pszeudoinverzek

Legyen $\mathbf{A} \in \mathbb{R}^{n \times m}$, $n \neq m$ nem négyzetes mátrix. Ha egy $\mathbf{A}x = y$, $x \in \mathbb{R}^{m \times 1}$, $y \in \mathbb{R}^{n \times 1}$ lineáris egyenletrendszer együtthatómátrixaként gondolunk rá, akkor $n \geq m$ vagy $m \geq n$ esetén rendre a *túlhatározottság* vagy *alulhatározottság* esete állna fent, az első esetben általánosságban nem lenne megoldásunk, a második esetben pedig végtelen sok megoldásunk lenne rá. Látszik, hogy az $n \neq m$ esetben nem beszélhetünk \mathbf{A}^{-1} inverzről, helyette egy általánosabb, úgynevezett *pszeudoinverz* kell.

Egy $\mathbf{A} \in \mathbb{R}^{n \times m}$, $n > m$ mátrix *bal oldali pszeudoinverze* (Más néven *Moore-Penrose pszeudoinverz*):

$$\mathbf{A}^\dagger := (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \in \mathbb{R}^{m \times n}$$

Figyeljük meg, hogy ha \mathbf{A}^\dagger -el balról megszorozzuk \mathbf{A} -t, az identitás mátrixot kapjuk, tehát bal oldalról valóban identitásként működik:

$$\mathbf{A}^\dagger \mathbf{A} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{A} = \mathbf{I}$$

Ha jobbról szoroznánk meg:

$$\mathbf{A} \mathbf{A}^\dagger = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

Ez semmi más, mint a *projekció-mátrix* \mathbf{A} oszlopvektorai által kifeszített vektortérre. Ha egy vektor ebben az oszloptérben van, rá persze identitásként hat $\mathbf{A} \mathbf{A}^\dagger$, ha viszont ezen kívül esik, akkor rávetíti az oszloptérre a vektort. Egy túlhatározott $\mathbf{A}x = y$ egyenletrendszert tehát "meg lehet oldani", ha y -t rávetítjük \mathbf{A} oszlopterére, és megoldjuk az $\mathbf{A}x = \tilde{y}$ egyenletrendszert:

$$\tilde{y} = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T y = \mathbf{A} \mathbf{A}^\dagger y$$

$$\mathbf{A}x = \tilde{y} = \mathbf{A} \mathbf{A}^\dagger y$$

$$\mathbf{A}^\dagger \mathbf{A}x = \mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger y$$

$$x = \mathbf{A}^\dagger y$$

Az $n < m$ esetben alulhatározottság áll fenn, itt *jobb oldali pszeudoinverzről* beszélhetünk:

$$\mathbf{A}^\dagger := \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \in \mathbb{R}^{n \times m}$$

Bár ezt nem fogjuk a későbbiekben használni, érdemes lehet megjegyezni, hogy a jobb oldali pszeudoinverzrel való balról szorzás esetén - hasonlóan a bal oldali pszeudoinverzhez - projekciómátrixot kapunk, csak most \mathbf{A} sorvektorai által kifeszített vektortérre.

0.1.2 Valószínűségi vektorváltozók

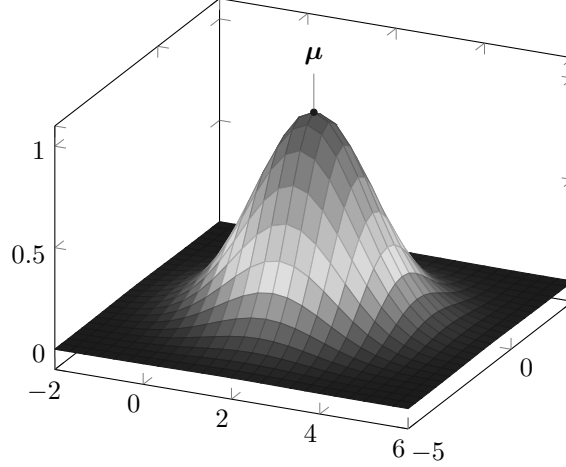
Egy $\boldsymbol{\xi} = [\xi_1, \dots, \xi_n]^T$ vektort *valószínűségi vektorváltozónak* hívunk, ha $\forall i$ -re ξ_i skalárértékű valószínűségi változó. A továbbiakban csak a vektorértékű normális eloszlást követő valószínűségi vektorváltozókkal foglalkozunk, ezek formálisan felírva:

$$\boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

ahol $\boldsymbol{\mu} \in \mathbb{R}^{n \times 1}$ a várható értékek vektora, $\boldsymbol{\Sigma}$ pedig a *variancia-kovarianca mátrix*. Természetesen $\text{Var}[\boldsymbol{\xi}] = \boldsymbol{\Sigma}$. Természetesen $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ pozitív szemidefinit és szimmetrikus mátrix. Az $n = 1$ esettel analóg módon ξ sűrűségfüggvénye

$$f_{\boldsymbol{\xi}}(\xi_1, \dots, \xi_n) = \frac{e^{-\frac{1}{2}(\boldsymbol{\xi} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\xi} - \boldsymbol{\mu})}}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}}$$

A sűrűségfüggvény $n = 2$ esetben $\boldsymbol{\mu} = [2, -1]^T$ és $\boldsymbol{\Sigma} = \mathbf{I}$ várhatóérték és kovariancia mátrix mellett:



Egy $\mathbf{A} \in \mathbb{R}^{n \times n}$ mátrix mellett a skaláresethez hasonlóan

$$\text{Var}[\mathbf{A}\boldsymbol{\xi}] = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$$

$$\mathbb{E}[\mathbf{A}\boldsymbol{\xi}] = \mathbf{A}\mathbb{E}[\boldsymbol{\xi}]$$

$\boldsymbol{\Sigma}$ kovariancia mátrixot kifejezhetjük várható értékekkel is:

$$\boldsymbol{\Sigma} = \mathbb{E}[(\boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}])(\boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}])^T] = \mathbb{E}[\boldsymbol{\xi}\boldsymbol{\xi}^T] - \mathbb{E}[\boldsymbol{\xi}]\mathbb{E}[\boldsymbol{\xi}^T]$$

$\boldsymbol{\Sigma}$ alakja:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \text{Cov}[\xi_1, \xi_2] & \dots & \text{Cov}[\xi_1, \xi_n] \\ \text{Cov}[\xi_2, \xi_1] & \sigma_2^2 & \dots & \text{Cov}[\xi_2, \xi_n] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[\xi_n, \xi_1] & \text{Cov}[\xi_n, \xi_2] & \dots & \sigma_n^2 \end{bmatrix}$$

ahol $\sigma_1^2, \dots, \sigma_n^2$ rendre ξ_1, \dots, ξ_n varianciái.

0.1.3 Mátrixdifferenciálás nagyon röviden

Legyenek $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{k \times 1}$ vektorok. Ekkor

$$\frac{\partial \mathbf{a}^T \mathbf{b}}{\partial \mathbf{b}} = \frac{\partial \mathbf{b}^T \mathbf{a}}{\partial \mathbf{b}} = \mathbf{a}$$

Ha $\mathbf{A} \in \mathbb{R}^{k \times k}$ mátrix, akkor

$$\frac{\partial \mathbf{b}^T \mathbf{A} \mathbf{b}}{\partial \mathbf{b}} = 2\mathbf{A} \mathbf{b}$$

Ha \mathbf{A} szimmetrikus, akkor ezen felül

$$2\mathbf{A} \mathbf{b} = 2\mathbf{b}^T \mathbf{A}$$

Legyen $\boldsymbol{\beta} \in \mathbb{R}^{k \times 1}$, $\mathbf{A} \in \mathbb{R}^{n \times k}$ és $\mathbf{y} \in \mathbb{R}^{n \times 1}$. Ekkor

$$\frac{\partial 2\boldsymbol{\beta}^T \mathbf{A}^T \mathbf{y}}{\partial \boldsymbol{\beta}} = \frac{\partial 2\boldsymbol{\beta}^T (\mathbf{A}^T \mathbf{y})}{\partial \boldsymbol{\beta}} = 2\mathbf{A}^T \mathbf{y}$$

0.2 A lineáris regresszió és az OLS eljárás

A regresszió kiindulópontja egy \mathcal{X} normális eloszlású sokaság, melynek minden tagja rendelkezik \mathbf{x}_i *featurevektor*-ral, avagy magyarázó változó-vektorral (ezek a *regresszorok*), illetve egy-egy skalár y_i *label*-lel, avagy magyarázott változóval (amiket a regresszorok magyaráznak egy lineáris modell alapján, ezt később

jobban kifejtjük). A sokaságból n darab mintát veszünk (megfigyelést végzünk), a minták iid. normális eloszlásúak, ami persze azt jelenti, hogy minden magyarázó változó-vektor egy vektorértékű normális eloszlású valószínűségi vektorváltozó.

A megfigyelt magyarázó változó-vektorokat soronként egymásra rakva felépítünk egy úgynevezett *design mátrixot*, melyet mostantól \mathbf{X} -el jelölünk. Minden \mathbf{x}_i magyarázó változó-vektor első eleme konstans 1, ez tölti be az intercept, avagy kétdimenziós esetben az y-tengellyel való metszéspont szerepét. n darab megfigyelés és p elemszámú magyarázó változó-vektorral \mathbf{X} alakja a következő:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,p-1} \\ 1 & x_{2,1} & \dots & x_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,p-1} \end{bmatrix}_{n \times p}$$

A megfigyelt magyarázott változókat szintén sorokba tömörítjük, így mivel mindegyik skalár, egy vektort kapunk:

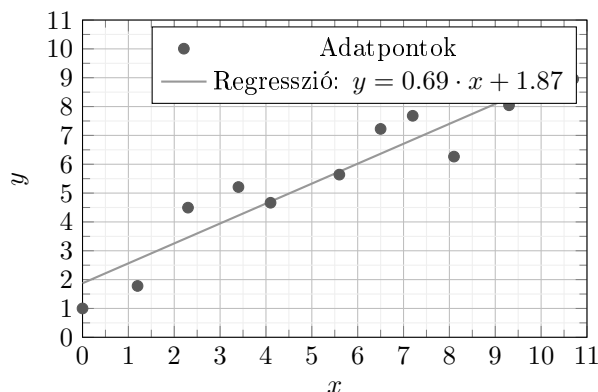
$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

A lineáris regresszió feladata, hogy a lineáris

$$\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{y}$$

modell mellett megtalálja azt a $\boldsymbol{\beta} \in \mathbb{R}^{p \times 1}$ együttható-vektort, amelyre az $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$ egyenletrendszer "legjobban teljesül". $\boldsymbol{\epsilon}$ az úgynevezett *hibavektor*. Az OLS, avagy *Ordinary Least Squares* becslési eljárást pontosan ezt a $\boldsymbol{\beta}$ paramétervektort becsüli, a paraméterbecslést kétféleképpen is levezetjük. Az OLS elnevezés valójában az analitikus levezetésből nyer értelmet a legkönnyebben, de először a - szerintem intuitívabb - projekciós módszert nézzük meg.

A lineáris regresszió egy darab regresszor (magyarázó változó) esetén az alábbi ábrával szemléltethető:



Itt $\hat{\boldsymbol{\beta}}$ paraméterbecslés vektor alakja

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} 1,87 \\ 0,69 \end{bmatrix}$$

Azt, hogy hogyan kaptuk meg $\hat{\boldsymbol{\beta}}$ paraméterbecslést, a következő fejezetek tárgyalják részletesen. Ezen kívül külön foglalkozunk majd a fenti egyváltozós regresszióval is (a $p = 2$ -es eset).

0.2.1 Az OLS-becslés geometriai értelmezése

Szinte mindig $n > p$, így az egyenletrendszer *túlhatározott*, és nagyon specifikus esetektől eltekintve nem létezik egzakt megoldása. Az első fejezetben azonban láttuk, hogy a bal oldali pszeudo inverz pontosan ezt a problémát orvosolja. A jelölési konvenció a megoldásból nyert *paraméter-becslésre* $\hat{\beta}$, ami a mintavétel véletlenszerűségéből adódóan maga is vektorértékű valószínűségi változó ($\hat{\beta}$ pontos eloszlásáról a későbbiekben lesz szó):

$$\hat{\beta} = \mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Ebben az esetben \mathbf{y} -t az \mathbf{X} design mátrix oszlopterére vetítettük. Legyen \mathbf{e} a valós y_i -k és a $\mathbf{X}\hat{\beta} = \hat{\mathbf{y}}$ modellbecslés által prediktált \hat{y}_i -k közti eltérések vektora (sokszor \mathbf{e} -t $\hat{\mathbf{e}}$ -ként is jelölik):

$$\mathbf{e} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

0.2.2 Az OLS-becslés mint szélsőérték-feladat

$\hat{\beta}$ paraméterbecslés-vektort megkaphatjuk úgy is, ha tekintjük az alábbi minimalizálási feladatot:

$$\mathbf{e}^T \mathbf{e} \rightarrow \min_{\hat{\beta}}$$

azaz minimalizáljuk a becsült \hat{y}_i és tényleges y_i magyarázott változók közötti négyzetösszeget. $\mathbf{e}^T \mathbf{e}$ -t RSS, azaz *sum of squared residuals* néven is emlegetik. Írjuk ki a hiba-négyzetösszeg teljes alakját:

$$\mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{y}^T \mathbf{y} - \hat{\beta}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\hat{\beta} + \hat{\beta}^T \mathbf{X}^T \mathbf{X}\hat{\beta} = \mathbf{y}^T \mathbf{y} - 2\hat{\beta}^T \mathbf{X}^T \mathbf{y} + \hat{\beta}^T \mathbf{X}^T \mathbf{X}\hat{\beta}$$

Itt felhasználtuk, hogy a transzponálás "megfordítja a szorzatot", illetve hogy skalár transzponáltja önmaga, így $\mathbf{y}^T \mathbf{X}\hat{\beta} = (\mathbf{y}^T \mathbf{X}\hat{\beta})^T = \hat{\beta}^T \mathbf{X}^T \mathbf{y}$. A minimalizációhoz vennünk kell a kifejezés $\hat{\beta}$ szerinti deriváltját, majd 0-val egyenlővé tenni:

$$\frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \hat{\beta}} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\hat{\beta} = 0$$

Ebből megkapjuk az úgynevezett *normálegyenletet*:

$$(\mathbf{X}^T \mathbf{X})\hat{\beta} = \mathbf{X}^T \mathbf{y}$$

$(\mathbf{X}^T \mathbf{X})$ szimmetrikus, és ha feltesszük, hogy létezik inverze, akkor balról beszorozva mindét oldalt:

$$(\mathbf{X}^T \mathbf{X})^{-1}(\mathbf{X}^T \mathbf{X})\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Látható, hogy pontosan ugyanaz jött ki, mint a pszeudo inverzes levezetésben. Míg ez utóbbi pusztán analitikus úton jutott el $\hat{\beta}$ -hoz, a pszeudo inverzes módszert geometrikus úton is el lehet képzelni.

0.3 Az OLS-becslés tulajdonságai

Vegyük az OLS paraméterbecslés normálegyenletét, és figyeljük meg, hogy $\mathbf{X}^T \mathbf{e} = \mathbf{0}$:

$$(\mathbf{X}^T \mathbf{X})\hat{\beta} = \mathbf{X}^T \mathbf{y}$$

A modellből adódóan $\mathbf{y} = \mathbf{X}\hat{\beta} + \mathbf{e}$ behelyettesítéssel:

$$(\mathbf{X}^T \mathbf{X})\hat{\beta} = \mathbf{X}^T (\mathbf{X}\hat{\beta} + \mathbf{e})$$

$$(\mathbf{X}^T \mathbf{X})\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})\boldsymbol{\beta} + \mathbf{X}^T \mathbf{e}$$

$$\mathbf{X}^T \mathbf{e} = \mathbf{0}$$

valóban. Ez azt jelenti, hogy *minden magyarázó változó (regresszor) korrelálatlan a hibával*, pontosabban megfogalmazva *a regresszorok és a hibák mintakorrelációja zérus*. Mivel \mathbf{X} mátrix első oszlopa konstans 1-eket tartalmaz, így $\hat{\beta}_0$ maga az intercept lesz, és emiatt

$$\sum_{i=1}^n e_i = 0$$

azaz a hibák összege 0. Ha leosztunk n -nel:

$$\frac{1}{n} \sum_{i=1}^n e_i = \bar{e}$$

azaz a hibatagok (*rezidiumok*) mintaátlagja - ami persze torzítatlan becslése a várható értéknek - 0, tehát $\mathbb{E}[e] = \mathbf{0}$.

Egy másik, ugyancsak fontos tulajdonság a predikciós formulából következik:

$$\hat{\mathbf{y}}^T \mathbf{e} = (\mathbf{X}\hat{\boldsymbol{\beta}})^T \mathbf{e} = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{e} = 0$$

azaz *a becsült \hat{y}_i -ok korrelálatlanok a rezidiumokkal*. Így azt is beláthatjuk, hogy *a modell által prediktált és a tényleges magyarázott változók mintaátlagai megegyeznek*:

$$\bar{\mathbf{y}} = \bar{\hat{\mathbf{y}}}$$

Felmerülhet a kérdés, hogy mindig létezik-e $(\mathbf{X}^T \mathbf{X})^{-1}$. Abban az esetben, ha \mathbf{X} oszloprangja kisebb, mint p , tehát *tökéletes multikollinearitás* áll fenn, akkor \mathbf{X} szinguláris értékei között lesz 0, így $\mathbf{X}^T \mathbf{X}$ sajátértékei között is, azaz $\mathbf{X}^T \mathbf{X}$ nem lesz invertálható. Ezentúl tehát feltételezzük, hogy nem áll fenn tökéletes multikollinearitás.

0.3.1 A Gauss-Markov feltételezések

A Gauss-Markov feltételezések biztosítják, hogy az OLS eljárással kapott $\hat{\boldsymbol{\beta}}$ paraméterbecslésünk *BLUE*, azaz *Best Linear Unbiased Estimator* lesz. Ez azt jelenti, hogy nem fogunk tudni találni olyan - nem az OLS eljárással kapott - paraméterbecslést $\boldsymbol{\beta}$ -ra, ami lineáris, torzítatlan, és kisebb mintavarianciával rendelkezik, mint $\hat{\boldsymbol{\beta}}$.

Formálisan kimondva az első Gauss-Markov feltétel a már látott modellegyenlet:

$$\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{y}$$

A második Gauss-Markov feltétel szerint \mathbf{X} oszloprangja megegyezik oszlopainak számával, az oszlopok mind lineárisan függetlenek, azaz nincs zérus szinguláris értéke. Ezt $(\mathbf{X}^T \mathbf{X})^{-1}$ létezésénél már feltételeztük, formálisan ez is egyike a feltételeknek.

A harmadik feltétel szerint

$$\mathbb{E}[\boldsymbol{\epsilon} \mid \mathbf{X}] = \mathbf{0}$$

$$\mathbb{E} \begin{bmatrix} \epsilon_1 \mid \mathbf{X} \\ \epsilon_2 \mid \mathbf{X} \\ \vdots \\ \epsilon_n \mid \mathbf{X} \end{bmatrix} = \begin{bmatrix} \mathbb{E}[\epsilon_1] \\ \mathbb{E}[\epsilon_2] \\ \vdots \\ \mathbb{E}[\epsilon_n] \end{bmatrix} = \mathbf{0}$$

Ez azt jelenti, hogy a modell szerinti hibatag várható értékét nem befolyásolja egyik magyarázó változó sem. Ebből következőleg

$$\mathbb{E}[\mathbf{y} \mid \mathbf{X}] = \mathbb{E}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \mid \mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$$

A negyedik feltétel a hibák kovariancia mátrixára vonatkozik, mégpedig

$$\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T \mid \mathbf{X}] = \sigma^2 \mathbf{I}$$

A hibatagok *homoszkedasztikusak és korrelálatlanok*, azaz azonosan σ^2 varianciájúak és $\forall i \neq j : Cov[\epsilon_i, \epsilon_j] = 0$. Ha kiírjuk $\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T$ mátrixformáját:

$$\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T \mid \mathbf{X}] = \mathbb{E} \begin{bmatrix} \epsilon_1^2 \mid \mathbf{X} & \epsilon_1\epsilon_2 \mid \mathbf{X} & \dots & \epsilon_1\epsilon_n \mid \mathbf{X} \\ \epsilon_2^1 \mid \mathbf{X} & \epsilon_2\epsilon_2 \mid \mathbf{X} & \dots & \epsilon_2\epsilon_n \mid \mathbf{X} \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_n^1 \mid \mathbf{X} & \epsilon_n\epsilon_2 \mid \mathbf{X} & \dots & \epsilon_n^2 \mid \mathbf{X} \end{bmatrix}$$

és persze $\forall i : \mathbb{E}[\epsilon_i \mid \mathbf{X}] = 0$ miatt a fenti mátrix diagonálisában ϵ_i -k varianciái, a többi helyen pedig a kovarianciák, amik a feltétel szerint 0-k, így $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T \mid \mathbf{X}]$ kovarianciamátrix valóban diagonális, a homoszkedaszticitás feltétele mellett pedig minden diagonális elem σ^2 . Mostantól a hibatagok varianciáját $\boldsymbol{\Sigma}$ fogja jelölni, $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$.

Az utolsó feltétel szerint a hibatagok normális eloszlást követnek:

$$\boldsymbol{\epsilon} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

Kijelenthetjük tehát, hogy y_i -k varianciáját nem csak \mathbf{x}_i -ek magyarázzák, hanem σ^2 *magyarázatlan variancia* is. Úgy is megfogalmazhatjuk, hogy a modell szerint minden \mathbf{y} magyarázott változó-vektor regresszorok szerinti feltételes eloszlása

$$\mathbf{y} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma})$$

Lássuk be, hogy a feltételek teljesülése mellett $\hat{\boldsymbol{\beta}}$ valóban torzítatlan becslést ad $\boldsymbol{\beta}$ -ra! Láttuk, hogy $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, és a modell szerinti $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ behelyettesítéssel

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})$$

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon},$$

mindkét oldalon véve a várható értéket:

$$\mathbb{E}[\hat{\boldsymbol{\beta}}] = \mathbb{E}[\boldsymbol{\beta}] + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] = \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbb{E}[\mathbf{X}^T \boldsymbol{\epsilon}]$$

Mivel a Gauss-Markov feltételek egyike, hogy $\mathbb{E}[\mathbf{X}^T \boldsymbol{\epsilon}] = \mathbf{0}$, így

$$\mathbb{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$$

ezzel készen is vagyunk.

0.3.2 $\hat{\boldsymbol{\beta}}$ varianciája

A hibavektor variancia-kovariancia mátrixához hasonlóan képezhetjük $\hat{\boldsymbol{\beta}}$ valószínűségi vektorváltozó variancia-kovariancia mátrixát:

$$Var[\hat{\boldsymbol{\beta}}] = \mathbb{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T]$$

Láttuk, hogy

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \implies \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \\ \mathbb{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T] &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon})^T] \end{aligned}$$

A transzponálás "szorzatmegfordító" tulajdonságából következően, illetve $\mathbf{X}^T \mathbf{X}$ szimmetrikus voltából

$$\text{Var}[\hat{\beta}] = \mathbb{E} \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \right]$$

$$\text{Var}[\hat{\beta}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

Itt válik igazán fontossá, hogy $\mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]$ variancia-kovariancia mátrix alakja $\sigma^2 \mathbf{I}$, így σ^2 kiemelhető a mátrixszorzások elé, az identitást pedig triviálisan nem szükséges kiírni:

$$\text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

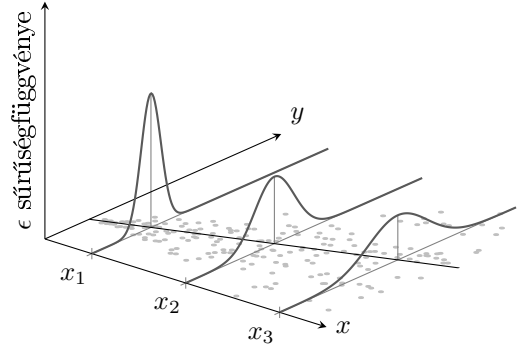
A mátrixszorzás asszociativitásából pedig a

$$\text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

végleges alakot kapjuk. Ugyanez megkapható az első fejezetben bemutatott $\text{Var}[\mathbf{A}\boldsymbol{\xi}] = \mathbf{A} \text{Var}[\boldsymbol{\xi}] \mathbf{A}^T$ transzformált variancia képlettel is, $\boldsymbol{\xi}$ helyett \mathbf{y} , \mathbf{A} helyett pedig $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ transzformáció mátrixsal (már ha \mathbf{X} -eket fixnek tekintjük). A várható értékes felírásból látszik, hogy persze $\text{Var}[\hat{\beta}]$ alakja

$$\mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T] = \begin{bmatrix} \text{Var}[\hat{\beta}_1] & \text{Cov}[\hat{\beta}_1, \hat{\beta}_2] & \dots & \text{Cov}[\hat{\beta}_1, \hat{\beta}_p] \\ \text{Cov}[\hat{\beta}_2, \hat{\beta}_1] & \text{Var}[\hat{\beta}_2] & \dots & \text{Cov}[\hat{\beta}_2, \hat{\beta}_p] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[\hat{\beta}_p, \hat{\beta}_1] & \text{Cov}[\hat{\beta}_p, \hat{\beta}_2] & \dots & \text{Var}[\hat{\beta}_p] \end{bmatrix}$$

Ha n elég nagy, akkor $\hat{\beta}$ eloszlása *megközelítőleg normális lesz*. Csupán érdekesség, de el lehet képzelni, hogy heteroszkedaszticitás ($\exists i, j : \sigma_i^2 \neq \sigma_j^2$) és $p = 2$ mellett a modell az alábbi ábrával szemléltethető:



A $\hat{\beta}$ varianciája formulában szereplő σ^2 hibavariancia maga is becslésre szorul, ennek $\hat{\sigma}^2$ torzítatlan becslése a tényleges \mathbf{e} hibatagokkal számolható:

$$\hat{\sigma}^2 = \frac{\mathbf{e}^T \mathbf{e}}{n - p}$$

ezt azonban nem fogjuk belátni.

0.4 A $p = 2$ -es egyszerű modell

Nézzük meg, hogy eddig látott paraméterbecslés és becslés-variancia hogy néz ki a legegyszerűbb, egy darab konstans interceptet és egy darab magyarázó változót tartalmazó OLS-el becsült modellben. A modell egyenlete minden $i = 1 \dots n$ megfigyelésre

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Az \mathbf{X} design mátrixunk most

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \in \mathbb{R}^{n \times 2}$$

lesz, $\hat{\beta}$ paraméterbecslés pedig

$$\hat{\beta} = \left(\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix}$$

A 2×2 -es mátrixok invertálása könnyen megy:

$$\begin{aligned} \hat{\beta} &= \frac{1}{n \sum_i x_i^2 - (\sum_i x_i)^2} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{bmatrix} \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix} = \frac{1}{n \sum_i x_i^2 - (\sum_i x_i)^2} \begin{bmatrix} \sum_i x_i^2 \sum_i y_i - \sum_i x_i \sum_i x_i y_i \\ -\sum_i x_i \sum_i y_i + n \sum_i x_i y_i \end{bmatrix} = \\ &= \begin{bmatrix} \frac{n(\frac{1}{n} \sum_i x_i^2) \cdot n(\frac{1}{n} \sum_i y_i) - n(\frac{1}{n} \sum_i x_i) \cdot n(\frac{1}{n} \sum_i x_i y_i)}{n^2(\frac{1}{n} \sum_i x_i^2) - n^2(\frac{1}{n} \sum_i x_i)^2} \\ \frac{n^2 \frac{1}{n} \sum_i x_i y_i - n(\frac{1}{n} \sum_i x_i) \cdot n(\frac{1}{n} \sum_i y_i)}{n^2(\frac{1}{n} \sum_i x_i^2) - n^2(\frac{1}{n} \sum_i x_i)^2} \end{bmatrix} \end{aligned}$$

Az n elemű mintából képzett *mintaátlag* semmi más, mint $\frac{1}{n} \sum_i x_i$ illetve $\frac{1}{n} \sum_i y_i$, a kovariancia x és y között pedig $\mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$, n elemű - a várható értéket torzítatlanul becsülő - mintaátlagokkal ez persze semmi más, mint az *empirikus kovariancia* $\text{empcov}[x, y] = \frac{1}{n} \sum_i x_i y_i - (\frac{1}{n} \sum_i x_i)(\frac{1}{n} \sum_i y_i)$. x varianciája $\mathbb{E}[x^2] - \mathbb{E}^2[x]$ -ként áll elő, $\mathbb{E}[x^2]$ empirikus becslése pedig $\frac{1}{n} \sum_i x_i^2$. A vektor mindkét elemében n^2 -el leosztva látható, hogy a nevezőkben pontosan x mintából számolt varianciája (*empvar*) van, míg a vektor második elemének számlálója pontosan x és y mintából számolt kovarianciája. A vektor első elemének számlálójában $\overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy}$ áll. Jelölje mostantól a mintából számolt varianciát és kovarianciát \widehat{Var} és \widehat{Cov} , ezzel a paraméterbecslés alakja

$$\hat{\beta} = \begin{bmatrix} \frac{\overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy}}{\widehat{Var}[x]} \\ \frac{\widehat{Cov}[x, y]}{\widehat{Var}[x]} \end{bmatrix}$$

Azt kaptuk tehát, hogy a legegyszerűbb egyváltozós regresszió becsült paraméterei

$$\hat{\beta}_0 = \frac{\overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy}}{\widehat{Var}[x]}$$

$$\hat{\beta}_1 = \frac{\widehat{Cov}[x, y]}{\widehat{Var}[x]}$$

Sokszor a mintaszámmal normálatlan empirikus kovarianciát és varianciát S_{xy} és S_{xx} jelöléssel látják el:

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

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2,$$

ezekkel felírva β_1 becslését:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

β_0 becslésének alakja β_1 ismeretében is kiszámolható, és sokszor ez a módszer sokkal kényelmesebb (már ha ismerjük $\hat{\beta}_1$ értékét):

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Ez nem csak intuitívan értelmezhető ("Az átlagos y semmi más, mint az y -tengellyel való metszéspont és $\hat{\beta}_1 \bar{x}$ összege"), hanem formálisan is levezethető a modell egyenletéből (meg abból, hogy beláttuk, hogy a paraméterbecslés torzítatlan a feltevéseink mellett, illetve hogy a hibatagok várható értéke 0):

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$\mathbb{E}[y] = \beta_0 + \beta_1 \mathbb{E}[x]$$

$$\beta_0 = \mathbb{E}[y] - \beta_1 \mathbb{E}[x]$$

A várhatóérték-operátor helyett persze a mintaátlagokkal dolgozva:

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

valóban.

0.4.1 $\hat{\beta}$ varianciája és az R^2 mutató

Láttuk, hogy a paraméterbecslés varianciája

$$\text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

A már levezetett $p = 2$ -es design mátrixsal dolgozva:

$$\text{Var}[\hat{\beta}] = \sigma^2 \left(\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \right)^{-1} = \sigma^2 \frac{1}{n \sum_i x_i^2 - (\sum_i x_i)^2} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{bmatrix}$$

Használjuk ki az empirikus variancia képletét:

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

Innen könnyen látszik, hogy

$$\text{Var}[\hat{\beta}_0] = \sigma^2 \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Var}[\hat{\beta}_1] = \sigma^2 \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Kimondhatjuk tehát, hogy ahogy σ^2 nő, úgy nő a paraméterbecslésünk varianciája, avagy *bizonytalansága* is. Hasonlítsuk össze az általános esetben kapott $\hat{\beta}$ variancia képletét β_1 varianciáéval:

$$\sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

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

A 2×2 -es mátrixszorzást elvégezve tényleg azt kaptuk, hogy az egyváltozós regresszió esetén S_{xx} semmi más, mint az $\mathbf{X}^T \mathbf{X}$ centralizálatlan regresszor-kovariancia mátrix.

Nagyon fontos - és ezért itt is kihangsúlyozandó - hogy y varianciája kettő forrásból jön: a regresszorok varianciájából és a regresszorok által nem magyarázott hibavarianciából. Írjuk ezt az összefüggést fel a mi esetünkben a modellegyenlet segítségével (persze a regresszorok és a hibák korrelálatlansága mellett):

$$\text{Var}[\mathbf{y}] = \beta_1^2 \text{Var}[\mathbf{x}] + \text{Var}[\epsilon]$$

Itt kihasználtuk, hogy a modell szerint β_0 konstans, így zérus varianciája van. $Var[\epsilon]$ hibavariancia az a része y varianciájának, amit nem magyaráznak a regresszorok. Ha $Var[\epsilon]$ kicsi, ez annyit jelent, hogy a becsült \hat{y} -ok és a tényleges y -ok közel vannak egymáshoz, azaz a regresszióval nagyon jól becsülhetjük a valódi y értékeket.

Legyen

$$R^2 := \frac{\beta_1^2 Var[\mathbf{x}]}{Var[\mathbf{y}]}$$

az arány, amiben a regresszorok varianciája magyarázza a magyarázott változó teljes varianciáját. R^2 0 és 1 közötti szám, minél közelebb van 1-hez, annál jobban becsülhető y a regresszorokkal. β_1 becslését beírva adódik:

$$R^2 = \frac{|Cov[\mathbf{x}, \mathbf{y}]|^2}{Var[\mathbf{x}]Var[\mathbf{y}]}$$

R^2 a regresszió "erősségét" mutatja, így a normálatlan empirikus kovarianciákkal és varianciákkal (S_{xy} , S_{xx} , S_{yy}):

$$R^2 = \frac{S_{xy}^2}{S_{xx}S_{yy}}$$

Itt persze $S_{yy} = \sum_i (y_i - \bar{y})^2$ Vezessük be az alábbi jelöléseket:

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SSE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SSR = \sum_{i=1}^n e_i^2$$

SST a *Sum of Squares Total*, SSE a *Sum of Squares Explained*, SSR pedig a *Sum of Squares Residual*. Az előbbi varianciafelbontásból könnyen látszik, hogy mivel SSE a regresszorok által magyarázott variancia, SSR pedig a magyarázatlan variancia:

$$SST = SSE + SSR$$

R^2 -et az előbbihez hasonlóan, csak most az új jelölésekkel felírva:

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

(Az irodalomban néha - zavaró módon - Az SSE a hibák négyzetösszegét jelenti, mint Sum of Squares Error, és az SSR jelenti a magyarázott varianciát, mint Sum of Squares Regression.)