

Matematikai statisztika

Informatika alapszak, "B" spec.

Zempléni András

Valószínűségelméleti és Statisztika Tanszék
Matematikai Intézet
Természettudományi Kar
Eötvös Loránd Tudományegyetem

Honlap: zempleni.elte.hu

E-mail: andras.zempleni@ttk.elte.hu

Szoba: D 3-310

6. stat. előadás

- A modell: $y = Xb + \varepsilon$
- $F := \text{Im}X \rightsquigarrow X$ képtere
- $r := \text{rang}(X)$, általában $r \leq p$, teljes rangú esetben $r = p$
- Paraméterbecslés: $\hat{b} = (X^T X)^{-1} X^T y$
- Projekció az F altérre: $P_F = X(X^T X)^{-1} X^T$
- Becsült értékek: $\hat{y} := X\hat{b}$
- Reziduálisok: $\hat{\varepsilon} = y - \hat{y}$
- Reziduális négyzetösszeg: $\text{RNÖ} := \text{SSR} \|\hat{\varepsilon}\|^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- Teljes négyzetösszeg: $\text{NÖ} (= \text{SS}) = \sum_{i=1}^n (y_i - \bar{y})^2$
- Determinációs együttható: $R^2 = 1 - \frac{\text{RNÖ}}{\text{NÖ}} = \frac{\text{NÖ} - \text{RNÖ}}{\text{NÖ}} \rightsquigarrow$ az eredményváltozó változékonyságának hány %-át magyarázza a regressziós modell
Értéke 0 és 1 között lehet. Minél nagyobb, annál jobb.

- Korrigált determinációs együttható: $R_{\text{adj}}^2 = 1 - \frac{n-1}{n-p-1} \frac{SSR}{SS}$ \rightsquigarrow egy lehetséges modellválasztási kritérium, minél nagyobb, annál jobb
- Akaike-féle információs kritérium: $AIC = 2(p+1) - 2 \log \hat{L}$, ahol \hat{L} a likelihood-függvény értéke akkor, ha az ML-becslést használjuk (normális eloszlású hibáknál ez megegyezik a legkisebb négyzetes becsléssel)
Ez is egy lehetséges modellválasztási kritérium, minél kisebb, annál jobb.

- t -próba az egyes együtthatókra (feltételezzük a hibák normális eloszlását): $H_0 : b = 0$, $H_1 : b \neq 0$
- A próbastatisztika: $t = \frac{\hat{b}}{D(\hat{b})}$, ez $n - 1$ szabadságfokú t -eloszlású, ha igaz a H_0 .
- Az egy magyarázó változós esetben ($y \sim a + bx$) $D^2(b) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$ és $D^2(a) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]$.
- A reziduálisok: $r_i = y_i - \hat{y}_i$, ebből becsülhető a hiba szórásnégyzet: $\hat{\sigma}^2 = \sum_{i=1}^n r_i^2 / (n - p)$
- A szokott módon konfidencia intervallum is konstruálható a becslésekhez

A regressziós modell "felépítése"

Ha p magyarázó változónk van, akkor 2^p modell közül kell a legjobbat megkeresni. Több módszer közül lehet választani:

- Nagyról kicsire (hátról előre): először az összes magyarázó változót be vesszük, majd egyenként a legkevésbé szignifikánsat kivesszük egészen addig, míg mindegyik szignifikáns lesz
- Kicsiről nagyra (előlről hátra): egyesével azt vesszük hozzá, amelyekkel a legjobban illeszkedő modellt kapjuk a következő lépésben. Vége: ha bármelyik, még a modellen kívüli magyarázó változót bevéve, már nem javul a modell illeszkedése.

- Gyakran kell valószínűséget becsülnünk/osztályoznunk (két osztályba)
 - Betegség kialakulásának valószínűsége
 - Hitelező csődbemenetelének valószínűsége
 - A vizsga sikeres letételének valószínűsége
- Itt a hagyományos lineáris modell nem célravezető (könnyen adódnak negatív vagy 1-nél nagyobb értékek)
- A leggyakrabban használt, logisztikus függvény:

$$P(Y = 1) = \frac{1}{1 + e^{-(b_0 + b_1x_1 + \dots + b_px_p)}}$$

ahol b_i a becsülendő együtthatók

- Az ún. odds-hányados: $P(Y = 1)/P(Y = 0) = e^{b_0 + b_1x_1 + \dots + b_px_p}$
- Paraméterbecslés: pl. maximum likelihood módszerrel (numerikus módszerekkel lehet megkapni)

- Az extrapoláció (a megfigyelt adatok tartományán kívülre történő előrejelzés) egyáltalán nem megbízható!
- Ha nem találtunk jól közelítő egyszerű függvényt, alkalmazhatunk nemparaméteres közelítést is a feltételes várható értékre (ez egyáltalán nem használható extrapolációra):

$$E(\widehat{Y|X=x}) = \frac{\sum_{i=1}^n Y_i k((x - X_i)/h_n)}{\sum_{i=1}^n k((x - X_i)/h_n)}$$

ahol k a magfüggvény, h_n az ablakszélesség

- Elnevezés: Nadarajah-Watson módszer
- Az ablakszélesség lényeges (nem könnyű a jó megválasztása)
 - Ha túl kicsi, az egyedi megfigyelések zajosságát követi le a közelítés
 - Ha túl nagy, túlságosan sima eredményt kapunk

- A lineáris modell egyik legfontosabb alkalmazása (faktorok különböző szintjeinek van-e hatása?)
- Motivációs példák:
 - Hatással van-e egy vállalatnál a (bruttó) fizetésekre az, hogy valaki nő-e, avagy férfi?
 - Különböző vetőmagokra megnézték a termésátlagot egy nagyobb földterület különböző részein. Vajon hatással van-e a vetőmag fajtája a termésátlagra?
 - Hatással van-e a valszám gyakorlati összpontszámra, hogy a hallgatónak ki a gyakorlatvezetője?
- a megfigyelések y_{ij} az i -edik "szinten" mért j -edik érték

- szórásfelbontás: Teljes négyzetösszeg:
$$SST = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\bullet\bullet})^2,$$

külső (csoportok közötti) négyzetösszeg:
$$SSK = \sum_{j=1}^p n_j (\bar{y}_{j\bullet} - \bar{y}_{\bullet\bullet})^2,$$

belső (csoporton belüli) négyzetösszeg:
$$SSB = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{j\bullet})^2$$

Szórásnégyzet-hányados

$$H^2 = 1 - \frac{SSB}{SST} = \frac{SSK}{SST}$$

Megjegyzés: ez nem más, mint a regressziónál az R^2

Tulajdonságai:

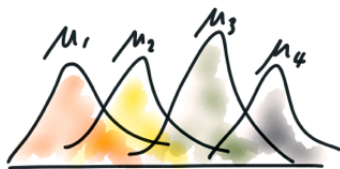
- $H^2 = 0$ esetén a két ismerv között nincs a szórásfelbontással kimutatható kapcsolat, DE (!!) ezzel nem bizonyítottuk be, hogy függetlenek egymástól (analógia: korrelálatlanságból nem következik a függetlenség)
- $H^2 = 1$ esetén a két ismerv között függvényszerű kapcsolat van
- $0 < H^2 < 1$ esetén a két ismerv között sztochasztikus kapcsolat van
- erős a kapcsolat, ha H^2 közel van 1-hez és gyenge a kapcsolat, ha 0-hoz

Szórásелеmzés (ANOVA)

- Elnevezései: szórásелеmzés = variancia-analízis = ANOVA (analysis of variance)
- A szórásелеmzési feladat fő kérdése: hatással van-e az eredményváltozó értékére, hogy a faktor melyik szintjén vagyunk? Jelölje b_i az i -edik szinten a várható értéket

$$H_0 : b_1 = b_2 = \dots = b_p$$

$$H_1 : \text{nem igaz } H_0$$



ANOVA

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 ?$$

ANOVA táblázat:

Szóródás forrása	Szabadságfok	Négyzetösszegek	Tapasztalati szórásnégyzetek	
Külső	$p - 1$	SSK	$MSK = \frac{SSK}{p-1}$	$F = \frac{MSK}{MSB} = \frac{SSK}{SSB} \cdot \frac{p-1}{n-p}$
Belső	$n - p$	SSB	$MSB = \frac{SSB}{n-p}$	
Teljes	$n - 1$	SST	–	