

Matematikai statisztika

Informatika alapszak, "A" 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

10. előadás

Regresszióelemzés, lineáris modell – motiváció

- **Gyakorlati szempontból a félév egyik legfontosabb témája!**
- Az 1. órai kérdőíves felmérés alapján mennyire magyarázható jól
 - a hallgatók testmagassága a súlyuk segítségével?
 - a hallgatók testmagassága a súlyuk és a cipőméretük segítségével?
 - a hallgatók statisztika érdemjegye a testmagasságuk segítségével?
- Egy részvény holnapi árfolyamát hogyan jelezzük előre a tegnapi, tegnapelőtti, stb. árfolyamadatok segítségével?
- Egy gazda földvásárlási dilemmája – egy bizonyos földterületen a várható termésátlag mennyire jelezhető előre a földterület fontosabb jellemzői alapján (a talaj kémhatása, a CaCO_3 megjelenési mélysége, a humusztartalom, topográfiai helyzet)?
- Meg lehet-e becsülni annak az esélyét, hogy valaki élete során megbetegszik tüdőrákban? Hogyan modellezzük ezt? Például: megbetegedés esélye \leftarrow dohányzik-e, hány éven át dohányzott élete során, van-e tüdőrákos a közeli rokonságban, van-e egyéb tüdőbetegsége, poros/füstös helyen dolgozik-e?

Regresszióelemzés – bevezetés

Legyenek Y, X_1, \dots, X_p véges szórású valószínűségi változók, amik egy véletlen jelenség egy-egy jellemzői.

A regresszióelemzés célja: a bennünket különösen érdeklő Y valószínűségi változó "minél jobb" közelítése az X_1, \dots, X_p valószínűségi változók segítségével.

Y elnevezései: eredményváltozó, függő változó, endogén változó
 X_i -k elnevezései: magyarázó változók, független változók, exogén változók

Általában megfigyeléseink vannak, amik az $(Y, X_1, \dots, X_p)^T$ valószínűségi vektorváltozó realizációinak tekinthetők:

$(y_i, x_{i,1}, \dots, x_{i,p})^T \quad i = 1, 2, \dots, n$ általában $n \gg p$

Feltehetjük, hogy az y_i megfigyelések rendszerint mérési eredmények, amik sajnos pontatlanok. A mérési hibát ε_i -vel fogjuk jelölni, amiről természetes feltétel, hogy legyen 0 várható értékű és véges σ szórású valószínűségi változó.

Regresszióelemzés: elméleti eredmények

Legyenek Y, X, X_1, \dots, X_p véges szórású valószínűségi változók, c, a, b_1, \dots, b_p valós számok.

Jelölje $\mathbf{X} = (X_1, \dots, X_p)^T$, $\mathbf{b} = (b_1, \dots, b_p)^T$ vektorokat.

	Feladat	Megoldás
a.)	$\min_c E(Y - c)^2$	$\hat{c} = EY$
b.)	$\min_{f: \mathbb{R} \rightarrow \mathbb{R}} E(Y - f(X))^2$ mérhető fv.	$\hat{f}(X) = E(Y X)$
c.)	$\min_{a,b} E(Y - (a + bX))^2$	$\hat{b} = \frac{\text{cov}(X, Y)}{D^2 X}$, $\hat{a} = EY - \hat{b}EX$
d.)	$\min_{f: \mathbb{R}^p \rightarrow \mathbb{R}} E(Y - f(X_1, \dots, X_p))^2$ mérhető fv.	$\hat{f}(X_1, \dots, X_p) = E(Y X_1, \dots, X_p)$
e.)	$\min_{a, b_1, \dots, b_p} E\left(Y - \left(a + \sum_{i=1}^p b_i X_i\right)\right)^2$ [Többváltozós lineáris regresszió]	$\hat{\mathbf{b}} = (\text{cov}(\mathbf{X}, \mathbf{X}))^{-1} \text{cov}(\mathbf{X}, Y)$ $\hat{a} = EY - \sum_{i=1}^p \hat{b}_i EX_i$

$E(Y|X)$: feltételes várható érték

- A modell: $\mathbf{y} = \mathbf{X}\mathbf{b} + \varepsilon$
- $F := \text{Im}X \rightsquigarrow X$ képtere
- $r := \text{rang}(X)$, általában $r \leq p$, teljes rangú esetben $r = p$
- Paraméterbecslés: $\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Projekció az F altérre: $P_F = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
- Becsült értékek: $\hat{\mathbf{y}} := \mathbf{X}\hat{\mathbf{b}}$
- Reziduálisok: $\hat{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}}$
- 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-beclést 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 ($y \sim a + bx$)
 $D^2(\hat{b}) = \frac{\sigma^2}{\sum(x_i - \bar{x})^2}$ és $D^2(\hat{a}) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right]$.
- 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 hátról): egyesével azt vesszük hozzá, amelyikkel a legjobban illeszkedő modellt kapjuk a következő lépésben.

Vége: ha bármelyik, még a modellen kívüli magyarázó változót bevéve, már nem javul a modell illeszkedése.