



# Statisztika

---

Előadások letölthetők a

[http://www.cs.elte.hu/~arato/InfBC\\_stat\\*.pdf](http://www.cs.elte.hu/~arato/InfBC_stat*.pdf)

címről

# Lehetséges hibák

- Elsőfajú hiba:  $H_0$  igaz, de elutasítjuk
- Másodfajú hiba:  $H_0$  hamis, de elfogadjuk

|         |                              | Aktuális helyzet     |                       |
|---------|------------------------------|----------------------|-----------------------|
|         |                              | A nullhipotézis igaz | A nullhipotézis hamis |
| Döntés: | Elfogadjuk a nullhipotézist  | Helyes döntés        | Másodfajú hiba        |
|         | Elutasítjuk a nullhipotézist | Elsőfajú hiba        | Helyes döntés         |



# Alapfogalmak

---

- Emlékeztető:  $\mathbf{X}$  mintatér: a minta lehetséges értékeinek halmaza.
- $\mathbf{X} = \mathbf{X}_e \cup \mathbf{X}_k$
- $\mathbf{X}_k$ : azon lehetséges értékek halmaza, amelyek megfigyelése esetén elutasítjuk a nullhipotézist.
- Gyakran statisztika segítségével határozzuk meg:

$$T(\mathbf{x}) = \begin{cases} 1 & , \mathbf{x} \in \mathbf{X}_k \\ 0 & , \mathbf{x} \notin \mathbf{X}_k \end{cases}$$



# Véletlenített próba

---

- Eddig adott megfigyelés esetén egyértelmű volt a döntésünk:

$$T(\mathbf{x}) = \begin{cases} 1 & , \mathbf{x} \in \mathbf{X}_k \\ 0 & , \mathbf{x} \notin \mathbf{X}_k \end{cases}$$

Véletlenített próba esetén sorolhatunk is:

$$\Psi(\mathbf{x}) = \begin{cases} 1 & , \text{ ha } T(\mathbf{x}) > c \\ \gamma & , \text{ ha } T(\mathbf{x}) = c \\ 0 & , \text{ ha } T(\mathbf{x}) < c \end{cases}$$



## Elsőfajú hiba valószínűsége véletlenített próba esetén

---

$\vartheta \in \Theta_0$ -ra az elsőfajú hiba valószínűsége:

$$P_{\vartheta}(T(\xi) > c) + \gamma P_{\vartheta}(T(\xi) = c) = E_{\vartheta}(\psi(\xi))$$

$\alpha$  a próba terjedelme, ha minden  $\vartheta \in \Theta_0$ -ra

$$E_{\vartheta}(\psi(\xi)) \leq \alpha$$

$\alpha$  a próba szignifikanciaszintje

(másképp: a próba pontos terjedelme),

$$\sup_{\vartheta \in \Theta_0} E_{\vartheta}(\psi(\xi)) = \alpha$$

## Legerősebb próba egyszerű hipotézis esetében

Egyszerű  $H_0$  és  $H_1 : |\Theta_0| = |\Theta_1| = 1$ .

$\psi$  a legerősebb  $\alpha$ -terjedelmű próba, ha:

$$P_{\vartheta_0}(T(\xi) > c) + \gamma P_{\vartheta_0}(T(\xi) = c) = E_{\vartheta_0}(\psi(\xi)) \leq \alpha,$$

továbbá minden más  $\alpha$ -terjedelmű  $\psi'$  próbára, annak másodfajú hibavalószínűsége nagyobb:

$$E_{\vartheta_1}(1 - \psi(\xi)) \leq E_{\vartheta_1}(1 - \psi'(\xi)).$$

# A legerősebb próba

---

- A legegyszerűbb eset:  $H_0$  és  $H_1$  is egyszerű (egyelemű). A valószínűséghányados (vh.) próba:

$$T(\mathbf{x}) = \begin{cases} 1 & \frac{L_1(\mathbf{x})}{L_0(\mathbf{x})} > c \\ \gamma & \frac{L_1(\mathbf{x})}{L_0(\mathbf{x})} = c \\ 0 & \frac{L_1(\mathbf{x})}{L_0(\mathbf{x})} < c \end{cases}$$

- Állítás (Neyman-Pearson lemma): a vh. próba legerősebb a saját terjedelmével. Minden  $0 < \alpha < 1$ -hez létezik ilyen terjedelmű vh. próba. Minden legerősebb próba ilyen alakú.

# Próbák a normális eloszlás várható értékére: t próba.

---

- $H_0: m=m_0$  ,  $H_1: m \neq m_0$  . Ha nem ismert a szórás (t-próba):

$$t = \sqrt{n} \frac{\bar{X} - m_0}{\hat{\sigma}}$$

- ahol  $\hat{\sigma} = \sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 / (n-1)}$
- Kritikus tartomány:  $|t| > t_{1-\alpha/2, n-1}$ . ( $H_0$  esetén a próbatatisztika  $n-1$  szabadságfokú, t-eloszlású.)
- Ha egyoldali az ellenhipotézis, akkor a kritikus tartomány  $t > t_{1-\alpha, n-1}$  ( $m > m_0$ ), illetve  $t < -t_{1-\alpha, n-1}$  alakú ( $m < m_0$ ). Ezek is legerősebb próbák!





## Megjegyzések

---


- A kétoldali esetre kapott próba nem a legerősebb (ilyenkor nincs is ilyen).
- Ha a minta elemszáma nagy, a t-próba helyett az u-próba is használható (ekkor még a normális eloszlásúságra sincs szükség a centrális határeloszlás tétel miatt).



# Kétoldali próbák és konfidencia intervallumok

---

- A normális eloszlásnál a várható értékre vonatkozó  $\alpha$  terjedelmű próbánál láttuk, hogy a  $H_0: m=m_0$  hipotézist a  $H_1: m \neq m_0$  hipotézissel szemben pontosan akkor fogadjuk el, ha  $m_0$  benne van az  $1 - \alpha$  megbízhatóságú konfidencia intervallumban.



# Kétmintás eset: párosított megfigyelések

---

- Példa: Van-e különbség Budapest és Cegléd napi átlaghőmérséklete között?  
 $H_0: m_1 = m_2$  a nullhipotézis.
- Ha ugyanazon napokról van megfigyelésünk mindkét helyen: nem függetlenek a minták. Ekkor a párok tagjai közötti különbséget vizsgálva, az előző egymintás esetre vezethető vissza a feladat.  $H_0^*: m = 0$ ,  $H_1^*: m \neq 0$  az új hipotézisek.

## Kétmintás eset: független minták

---

Ha ismert a szórás: ( $\bar{X}$  n elemű,  $\sigma_1$  szórású,  $\bar{Y}$  m elemű,  $\sigma_2$  szórású), alkalmazható a kétmintás u-próba

$$u = \frac{\bar{X} - \bar{Y}}{\sqrt{\sigma_1^2 / n + \sigma_2^2 / m}}$$

Kritikus tartomány: mint az egymintás esetben  
Ha ismeretlenek, de azonosak a szórások:

$$t_{n+m-2} = \sqrt{\frac{nm(n+m-2)}{n+m}} \frac{\bar{X} - \bar{Y}}{\sqrt{\sum (X_i - \bar{X})^2 + \sum (Y_i - \bar{Y})^2}}$$



# A szórás vizsgálata kétmintás esetben: F-próba

---

- $H_0: \sigma_1 = \sigma_2$
- Két független,  $n$ , illetve  $m$  elemű normális eloszlású minta alapján a próbastatisztika:  
(a korrigált tapasztalati szórásnégyzetek hányadosa) 
$$F = \max\left(\frac{s_1^2}{s_2^2}, \frac{s_2^2}{s_1^2}\right)$$
- Kritikus érték: az  $n-1, m-1$  szabadságfokú F eloszlás  $1-\alpha/2$  kvantilise ( $n$  a számlálóbeli,  $m$  pedig a nevezőbeli minta elemszáma).

# Kétmintás t-próba ismét

---

- Alkalmazható, ha az F-próba elfogadja a szórások azonosságát.
- Ha nem, akkor Welch-próba:

$$t' = \frac{\bar{X} - \bar{Y}}{\sqrt{s_1^2 / n_1 + s_2^2 / n_2}}$$

- $H_0$  esetén közelítőleg t eloszlású f szabadságfokkal, ahol

$$\frac{1}{f} = \frac{c^2}{n-1} + \frac{(1-c)^2}{m-1} \quad c = \frac{s_1^2 / n}{s_1^2 / n + s_2^2 / m}$$

## $\chi$ -négyzet próba

---

- $H_0$  hipotézis: az  $A_1, A_2, \dots, A_r$  teljes eseményrendszerre teljesül  $P(A_1)=p_1, P(A_2)=p_2, \dots, P(A_r)=p_r$
- A tesztstatisztika:

$$\sum_{i=1}^r \frac{(v_i - np_i)^2}{np_i}$$

ami aszimptotikusan  $r-1$  szabadságfokú  $\chi$ -négyzet eloszlású, ha igaz a nullhipotézis.

- Kritikus tartomány: ha a statisztika értéke nagyobb, mint az  $r-1$  szabadságfokú  $\chi$ -négyzet eloszlás  $1-\alpha$  kvantilise, elutasítjuk a nullhipotézist.

# $\chi$ -négyzet próba (folytatás)

---

- Miért is ez a határeloszlás?

$r = 2$ ,  $H_0 : P(A) = p$ ,  $\nu : A$  gyakorisága  $n$  kísérletből

$$\chi^2 = \frac{(\nu - np)^2}{np} + \frac{((n - \nu) - n(1 - p))^2}{n(1 - p)} = \frac{(\nu - np)^2}{np} + \frac{(\nu - np)^2}{n(1 - p)} = \frac{(\nu - np)^2}{np(1 - p)}$$

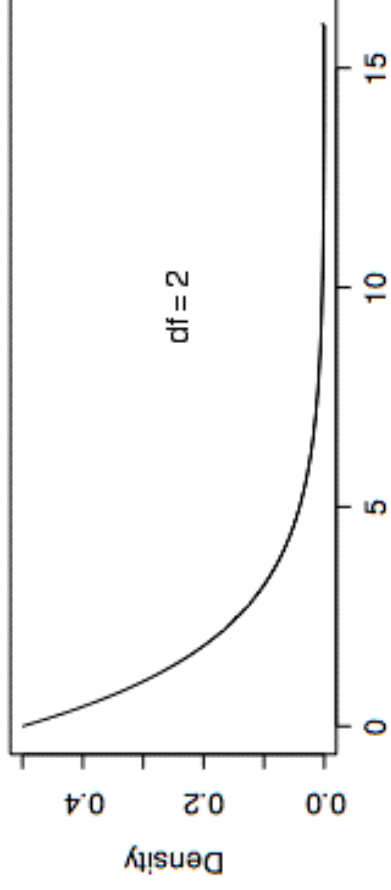
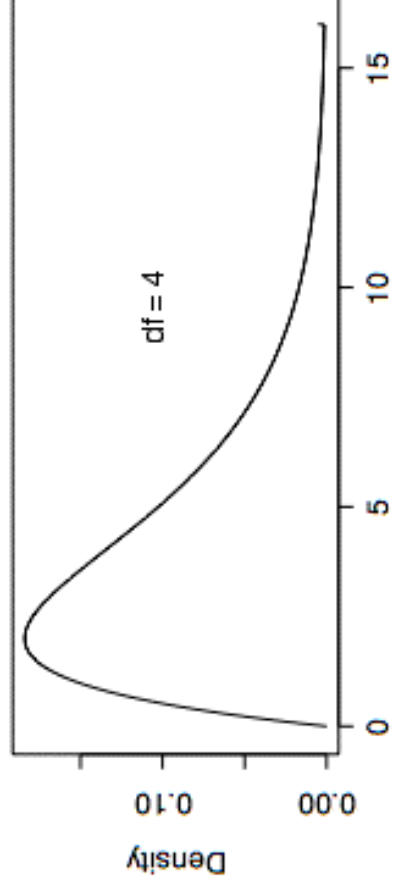
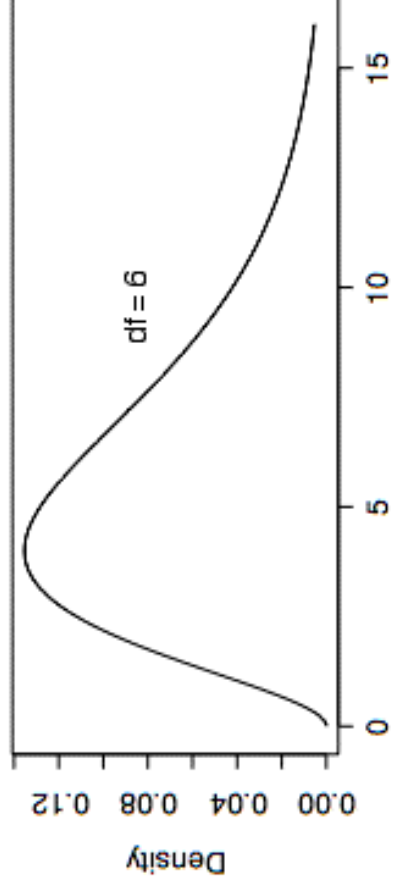
$\xi_i = 1$ , ha az  $i$ . kísérletnél  $A$  bekövetkezik, 0 különben

$$\nu = \sum_{i=1}^n \xi_i, \quad E\xi_i = p, \quad D^2\xi_i = p(1 - p),$$

$$\chi^2 = \left( \frac{\sum_{i=1}^n \xi_i - nE\xi_1}{\sqrt{nD\xi_1}} \right)^2 \xrightarrow[n \rightarrow \infty, \text{eloszlásban}]{} \chi_1^2$$



Chi Square



# Példa (kockadobás)

- 36 kockadobás eredménye

| Szám | Megfigyelt | $np_i$ | $\frac{(v_i - np_i)^2}{np_i}$ |
|------|------------|--------|-------------------------------|
| 1    | 8          | 6      | 0.667                         |
| 2    | 5          | 6      | 0.167                         |
| 3    | 9          | 6      | 1.500                         |
| 4    | 2          | 6      | 2.667                         |
| 5    | 7          | 6      | 0.167                         |
| 6    | 5          | 6      | 0.167                         |


$$n=36, r=6$$

---

$$\sum_{i=1}^6 \frac{(v_i - np_i)^2}{np_i} \sim \chi_5^2$$

$$\sum_{i=1}^6 \frac{(v_i - np_i)^2}{np_i} = 5.333$$

$$P(\chi_5^2 > 5.333) = 0.377 \Rightarrow$$

Nem tudjuk a szabályosság hipotézisét elutasítani!

## Példa (számítógépek népszerűsége)

---

- 100 amerikai diák

| Számí-<br>tógép | Megfigyelt | $np_i$  | $\frac{(v_i - np_i)^2}{np_i}$ |
|-----------------|------------|---------|-------------------------------|
| IBM             | 47         | 33.3333 | 5.604                         |
| Macintosh       | 36         | 33.3333 | 0.213                         |
| Egyéb           | 17         | 33.3333 | 8.003                         |


$$n=100, r=3$$

---

$$\sum_{i=1}^3 \frac{(v_i - np_i)^2}{np_i} \sim \chi^2$$

$$\sum_{i=1}^3 \frac{(v_i - np_i)^2}{np_i} = 13.820$$

$$P(\chi^2 > 5.99) = 0.05 \Rightarrow$$

Elutasítjuk az egyforma kedveltség hipotézisét!



## $\chi$ -négyzet próba illeszkedésvizsgálatra

---

- Illeszkedésvizsgálat:

$H_0 : \xi_1, \dots, \xi_n \text{ } F \text{ eloszlásfüggvényűek}$

- Visszavezetjük az előző esetre

$$A_i = \{\xi \in C_i\}, i = 1, 2, \dots, r, \bigcup_i C_i = \mathbf{R}$$

Diszkrét esetben gyakran:  $A_i = \{\xi = x_i\}, i = 1, 2, \dots, r$



## Példa

- Mi lehet egy vezető által okozott károk számának eloszlása?
- Poisson eloszlású-e?

| Kár-szám       | 0      | 1     | 2    | 3   | 4  | 5 | 6 | 7 | >7 | Össze-sen |
|----------------|--------|-------|------|-----|----|---|---|---|----|-----------|
| Veze-tők száma | 129524 | 16267 | 1966 | 211 | 31 | 5 | 1 | 1 | 0  | 148006    |

## Becsléses $\chi$ -négyzet próba

---

- $H_0$  hipotézis: az  $A_1, A_2, \dots, A_r$  teljes eseményrendszerre teljesül:

$$P(A_i) = p_i(\vartheta_1, \dots, \vartheta_s), i = 1, 2, \dots, r$$

$\vartheta_1, \dots, \vartheta_s$  ismeretlen paraméterek.

A tesztstatisztika:

$$\chi^2 = \sum_{i=1}^r \frac{(v_i - n\hat{p}_i)^2}{n\hat{p}_i} \xrightarrow{n \rightarrow \infty} \chi_{r-s-1}^2,$$

ahol

$$\hat{p}_i = p_i(\hat{\vartheta}_1, \dots, \hat{\vartheta}_s).$$



## Példa (folyt.)

| Kár-<br>szám               | 0       | 1      | 2     | 3   | 4   | 5    | 6     | 7     | >7    | Össze-<br>sen |
|----------------------------|---------|--------|-------|-----|-----|------|-------|-------|-------|---------------|
| Veze-<br>tők<br>száma      | 129524  | 16267  | 1966  | 211 | 31  | 5    | 1     | 1     | 0     | 148006        |
| $np_i$<br><i>Poisson</i>   | 128 433 | 18 218 | 1 292 | 61  | 2,2 | 0,06 | 0,001 | 3E-05 | 5E-07 |               |
| $Np_i$<br><i>Neg. bin.</i> | 129 541 | 16 237 | 1 962 | 234 | 28  | 3,3  | 0,39  | 0,05  | 0,006 |               |


$$n=148006, r=5$$

$$A_i = \{\xi = i\}, i = 0, 1, 2, 3$$

$$~~A_4 = \{\xi \geq 4\}~~$$

---

Poisson eset:

$$\hat{\lambda}=0.709$$

$$\sum_{i=0}^4 \frac{(v_i - n\hat{p}_i)^2}{n\hat{p}_i} \sim \chi_{5-1-1}^2$$

$$\sum_{i=0}^4 \frac{(v_i - n\hat{p}_i)^2}{n\hat{p}_i} > 200$$

$$P(\chi_3^2 > 17.7) = 0.05\% \Rightarrow$$

Elutasítjuk Poisson eloszlás hipotézisét!



## Az illeszkedésvizsgálat alkalmazása folytonos eloszlásokra

---

- A teljes eseményrendszer a számegyenes felosztása révén jön létre.
- Ügyeljünk arra, hogy minden intervallum közel azonos valószínűségű legyen.
- Ha paraméterbecslés szükséges, ML módszer alkalmazható.

# $\chi$ -négyzet próba homogenitásvizsgálatra

---

- Homogenitásvizsgálat:

$H_0 : \xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_m$  ugyanolyan eloszlásúak

- Hasonlóan járunk el, mint korábban

$$\bigcup_{i=1}^r C_i = \mathbf{R}$$

$$v_i = \left| \{j : \xi_j \in C_i\} \right|, \mu_i = \left| \{j : \eta_j \in C_i\} \right|, i = 1, 2, \dots, r,$$

A tesztstatisztika:

$$\chi^2 = nm \sum_{i=1}^r \frac{\left( \frac{v_i}{n} - \frac{\mu_i}{m} \right)^2}{\frac{v_i}{n} + \frac{\mu_i}{m}} \xrightarrow{n, m \rightarrow \infty} \chi_{r-1}^2$$



# Ki tanul jobban?

---

2009. január 5-ei vizsga

| Jegy     | Férfi | Nő  | Összesen |
|----------|-------|-----|----------|
| 1        | 47    | 4   | 51       |
| 2        | 11    | 1   | 12       |
| 3        | 11    | 2   | 13       |
| 4        | 9     | 2   | 11       |
| 5        | 8     | 2   | 10       |
| Összesen | 86    | 11  | 97       |
| Átlag    | 2,1   | 2,7 | 2,1      |

$$C_1 = \{1;2\}, C_1 = \{3;4;5\}$$

$$\overline{v_i = \left| \left\{ j : \xi_j \in C_i \right\} \right|}, \mu_i = \left| \left\{ j : \eta_j \in C_i \right\} \right|, i = 1, 2,$$

$$v_1 = 58, v_2 = 28, \mu_1 = 5, \mu_2 = 6, n = 86, m = 11$$

A tesztstatisztika:

$$\chi^2 = 86 \cdot 11 \left( \frac{\left( \frac{58}{86} - \frac{5}{11} \right)^2}{58 + 5} + \frac{\left( \frac{28}{86} - \frac{6}{11} \right)^2}{28 + 6} \right) = 2.071$$

$$P(\chi_1^2 > 2.71) = 10\% \Rightarrow$$

Nem tudjuk elutasítani az egyforma képesség hipotézisét!

## $\chi$ -négyzet próba függetlenségvizsgálatra

---

- $H_0$  hipotézis: az  $A_1, A_2, \dots, A_r$  és  $B_1, B_2, \dots, B_s$  teljes eseményrendszerekre teljesül a függetlenség.

$$\sum_{i,j} \frac{(v_{ij} - np_{ij})^2}{np_{ij}}$$

- Kritikus tartomány: ha a statisztika értéke nagyobb, mint az  $rs-1$  szabadságfokú  $\chi$ -négyzet eloszlás  $1-\alpha$  kvantilise, elutasítjuk a nullhipotézist.



## Becsléses eset

---

- Általában, ha az illesztendő eloszlást nem ismerjük – csak a családját – becsljük a paramétereit. Ekkor a próbatisztika szabadságfoka annnyival csökken, ahány paramétert becsltünk.
- Függetlenségvizsgálatnál általában nem ismerjük a teljes eseményrendszer tagjainak valószínűségét, így  $r-1+s-1$  valószínűséget kell becslnünk. A szabadságfok ekkor tehát  $rs-1-r-s+2=(r-1)(s-1)$ .





$v_{ij} : A_i B_j$  gyakorisága

$v_{i\bullet} : A_i$  gyakorisága

---

$v_{\bullet j} : B_j$  gyakorisága

A tesztstatisztika

$$n \sum_{i,j} \frac{\left( v_{ij} - \frac{v_{i\bullet} v_{\bullet j}}{n} \right)^2}{v_{i\bullet} v_{\bullet j}} \xrightarrow{n \rightarrow \infty} \chi^2_{(r-1)(s-1)}$$

$r = s = 1$  esetben

$$n \frac{(v_{11} v_{22} - v_{12} v_{21})^2}{v_{1\bullet} v_{2\bullet} v_{\bullet 1} v_{\bullet 2}} \xrightarrow{n \rightarrow \infty} \chi^2_1$$



# Példa

---

- [http://onlinestatbook.com/case\\_studies/diet.html](http://onlinestatbook.com/case_studies/diet.html)
- <http://onlinestatbook.com/chapter14/contingency.html>



## $Y$ közelítése $X$ függvényével

---

- Gyakori eset, hogy nem ismerjük a számunkra érdekes mennyiség ( $Y$ ) pontos értékét (pl. holnapi részvényárfolyam, vízállás, időjárás). Van viszont információnk hozzá kapcsolódó mennyiségről ( $X$ , mai értékek).
- Feladat: olyan  $f_0$  megtalálása, amelyre  $f_0(X)$  a lehető legjobb közelítése  $Y$ -nak.
- Matematikailag:  $f_0$  a megoldása a szélsőérték-problémának (legkisebb négyzetes becslés).  $\min_f E(Y - f(X))^2$

# Valószínűségszámításból tanultak

---

$E(Y - a)^2$  minimumhelye:  $EY$

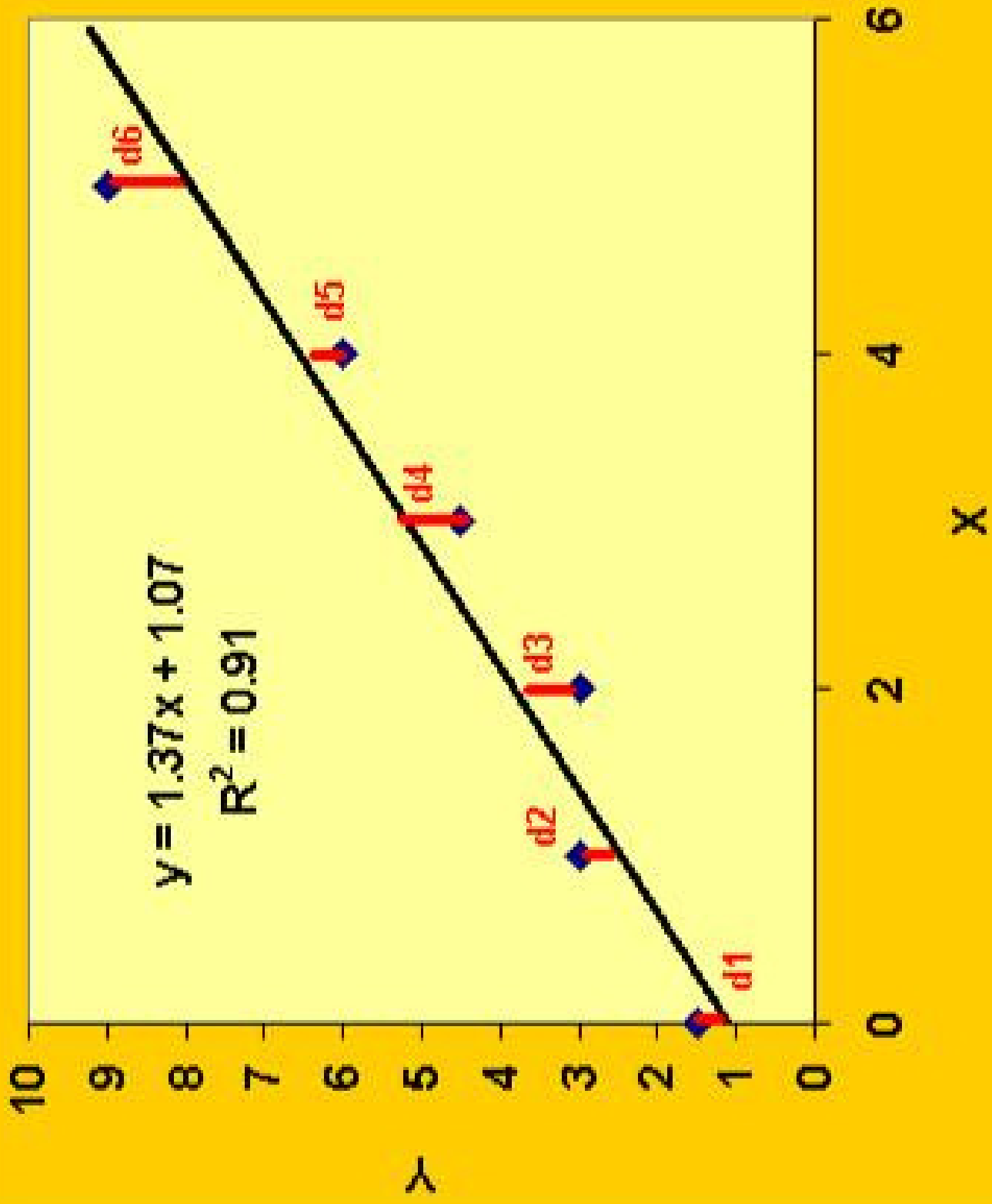
$E(Y - f(X))^2$  minimumhelye:  $f_0(x) = E(Y | X = x)$

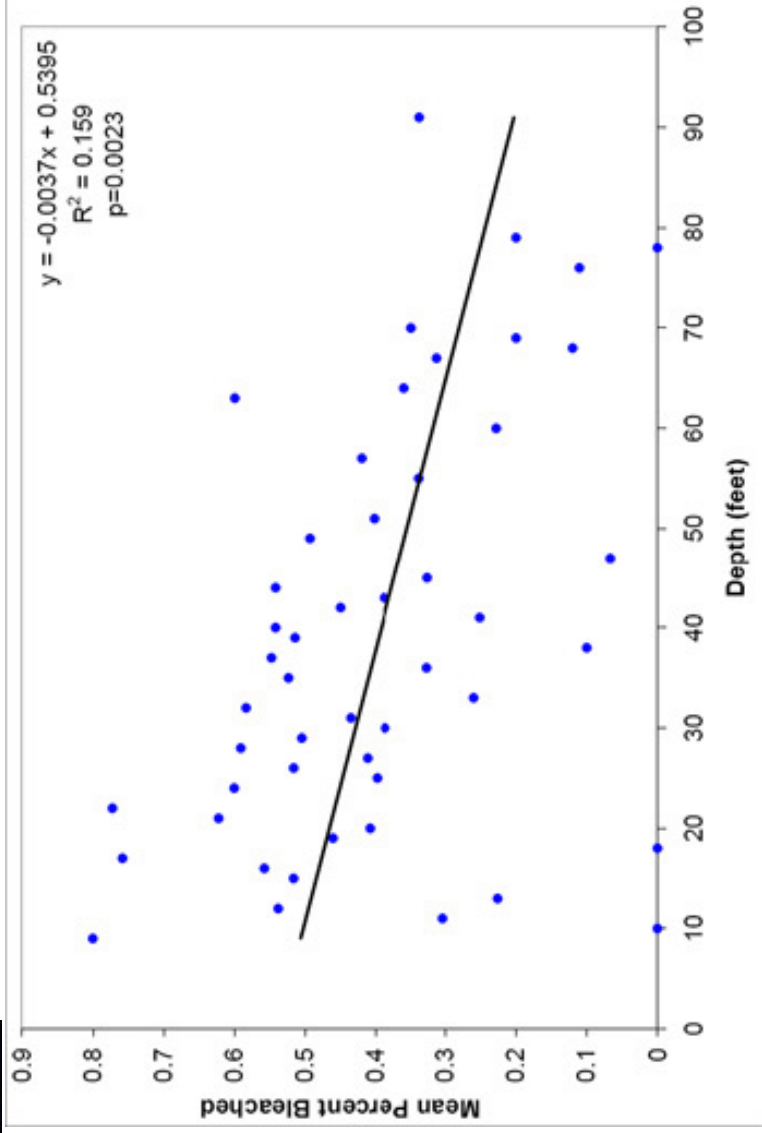
lineáris függvények esetében:

$E(Y - aX - b)^2$  minimumhelye:

$$a = \frac{\text{cov}(X, Y)}{D^2 X} = \frac{\text{corr}(X, Y)DY}{DX}$$

$$b = EY - aEX$$







# Lineáris modell

---

- $Y_i = aX_i + b + \varepsilon_i$  ( $X_i$  a magyarázó változó értéke,  $\varepsilon_i$  független, azonos eloszlású hiba.  $E(\varepsilon_i) = 0$ , általában feltesszük, hogy normális eloszlásúak)
- $a, b$  a becsülendő együtthatók



# Megoldás

---

$$\hat{a} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \hat{b} = \bar{y} - \hat{a}\bar{x}$$



# Szórások

---



$$D(\hat{a}) = \frac{\sigma}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, D(\hat{b}) = \sigma \sqrt{\frac{1}{\bar{x}} + \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2}}$$

Az  $x^*$  pontban előrejelzett érték  $\hat{a}x^* + \hat{b}$

és ennek szórása

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

A szórásbecslésnél  $\sigma$  helyett

$$\text{annak becslt értékét használjuk: } \hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{a}x_i - \hat{b})^2}{n-2}$$