

# Test Kolmogorova

- Neparametarski test (nezavisan od raspodele obeležja).
- Primjenjuje se za obeležja koja imaju neprekidne raspodele.
- Nulta hipoteza  $H_0$  je da je raspodela  $F(x)$  jednaka raspodeli  $F_0(x)$ , a alternativna hipoteza je da je  $F(x)$  različita od  $F_0(x)$ .
- Test-statistika, tj. statistika Kolmogorova je

uzoračka funkcija  
raspodele

$$D_n = \sup_{-\infty < x < \infty} |F_n^*(x) - F_0(x)|$$

- Kolmogorov je pokazao da za neprekidne funkcije raspodela važi

$$\lim_{n \rightarrow \infty} P[\sqrt{n} D_n < \lambda] = K(\lambda) = \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2\lambda^2}, \quad \lambda > 0$$

$$K(\lambda) = 0, \quad \lambda \leq 0$$

# Test Kolmogorova, nastavak

- Neka je realizovana vrednost statistike Kolmogorova

$$d_n = \sup_{-\infty < x < \infty} |F_n^*(x) - F_0(x)|$$

- Kritična oblast je

$$C = [d_{n,\alpha}, \infty)$$

određuje se iz  
tablica

- Hipotezu  $H_0$  odbacujemo (za dati prag značajnosti i za dati uzorak), ako je

$$d_n > d_{n,\alpha}$$

# Poređenje neparametarskih testova

- $\chi^2$  – test se odnosi na sve raspodele. Test Kolmogorova samo za neprekidne raspodele.
- U  $\chi^2$  – testu mogu figurisati i raspodele sa nepoznatim parametrima.
- Kod  $\chi^2$  – testa se upoređuju empirijske i teorijske frekvencije, a kod testa Kolmogorova empirijska i teorijska funkcija raspodele.
- U  $\chi^2$  – testu se vrši grupisanje podataka i samo je važno koliko ih ima po pojedinim intervalima, a ne i koji su. Time se gubi deo informacije o uzorku.

# Testiranje hipoteze o homogenosti

- Da li nezavisni uzorci obima  $n$  za obeležja  $X$  i  $Y$  opisuju isti proces, tj., da li potiču iz iste raspodele?
- Hipoteza homogenosti je oblika  $H_0(F_X(x) = F_Y(y))$ , gde su  $F_X(x)$  i  $F_Y(y)$  funkcije raspodele obeležja  $X$  i  $Y$ .
- Alternativna hipoteza je  $H_1(F_X(x) \neq F_Y(y))$ .
- **Test Kolmogorova-Smirnova:**
- Prost slučajan uzorak obima  $n_1$  obeležja  $X$  i prost slučajan uzorak obima  $n_2$  obeležja  $Y$ . Izračunaju se uzoračke funkcije raspodele  $F_{n_1}^*(x) \quad F_{n_2}^*(x)$

# Kolmogorov-Smirnov test

- Test-statistika je

$$D_m = \sup_{-\infty < x < \infty} |F_{n_1}^*(x) - F_{n_2}^*(x)|$$

$$m = \frac{n_1 n_2}{n_1 + n_2}$$

- Iz tablica se, za dati prag značajnosti  $\alpha$ , odredi broj  $d_{m,\alpha}$  koji je granica kritične oblasti

$$C = [d_{m,\alpha}, \infty)$$

- Hipotezu o jednakosti raspodела obeležja  $X$  i  $Y$  odbacujemo, ako  $d_{m,\alpha} \in C$ .

# Regresija i korelacija

- Podaci iz uzorka se mogu koristiti za utvrđivanje međusobne zavisnosti obeležja. Neka obeležja su više, a neka manje povezana.
- Na osnovu pojma uslovnog matematičkog očekivanja se rešava problem određivanja funkcionalne zavisnosti dve slučajne promenljive (obeležja).
- U teoriji verovatnoće moguće je odrediti raspodelu  $Y$  ( $Y=f(X)$ ) ako je poznata raspodela  $X$

$$F_Y(y) = P(Y < y) = P(f(X) < y) = P(X < f^{-1}(y)) \quad g_Y(y) = \frac{dF_Y}{dy}$$

ako funkcija  $f$  ima inverznu,  
i ako je ona rastuća

# Uslovno matematičko očekivanje

- Ako je  $(X, Y)$  diskretna 2D slučajna promenljiva sa vrednostima  $(x_i, y_j)$ ,  $i=1,\dots,r$ ,  $j=1,\dots,s$  i zakonom raspodele  $p_{ij}=P[X=x_i, Y=y_j]$ ,  $i=1,\dots,r$ ,  $j=1,\dots,s$ , tada je

$$E(Y / X = x_i) = \sum_{j=1}^s y_j P[Y = y_j / X = x_i] = \sum_{j=1}^s y_j \frac{P[Y = y_j, X = x_i]}{P[X = x_i]}$$

- Ako je  $(X, Y)$  neprekidna 2D slučajna promenljiva sa gustinom raspodele  $g(x, y)$  i ako je  $g_X(x)$  gustina raspodele slučajne promenljive  $X$ , tada je

$$E(Y / X = x) = \int_{-\infty}^{\infty} y \frac{g(x, y)}{g_X(x)} dy$$

# Regresija

- Slučajna promenljiva  $E(Y/X) = R(X)$  se naziva **regresija**.
- To je funkcija koja najbolje opisuje zavisnost  $Y$  od  $X$  u smislu da je srednje kvadratno odstupanje minimalno

$$E(Y-h(X))^2$$

ako je  $h(X)= E(Y/X) = R(X)$ .

- Da bismo odredili  $R(X)$ , potrebno je poznavanje zajedničke raspodele slučajnih promenljivih  $X$  i  $Y$ .
- Ako slučajna promenljiva  $(X, Y)$  ima dvodimenzionalnu normalnu raspodelu, tada je

$$E(Y / X) = aX + b$$

# Metoda najmanjih kvadrata

- Ako zajednička raspodela nije poznata, moguće je na osnovu podataka iz uzorka  $(x_i, y_i), i=1,..,n$  odrediti zavisnost  $Y$  od  $X$  u obliku  $Y=f(X, a_1, a_2, \dots, a_m)$ .
- Funkciju  $f$  zovemo regresija, a parametre  $a_1, \dots, a_m$  određujemo iz uslova da suma kvadrata odstupanja bude minimalna

$$S(a_1, \dots, a_m) = \sum_{j=1}^n (y_j - f(x_j, a_1, \dots, a_m))^2$$

- Rešava se sistem jednačina (normalni sistem)

$$\frac{\partial S}{\partial a_1} = 0, \dots, \frac{\partial S}{\partial a_m} = 0$$

# Određivanje parametara

- Neka su  $(\hat{a}_1, \hat{a}_2, \dots, \hat{a}_m)$  jedno rešenje normalnog sistema.  
Proveravamo da li je za te vrednosti postignut lokalni minimum f-je

$$d^2S = \frac{\partial^2 S}{\partial a_1^2} da_1^2 + \dots + \frac{\partial^2 S}{\partial a_m^2} da_m^2 + 2 \frac{\partial^2 S}{\partial a_1 \partial a_2} da_1 da_2 + \dots + 2 \frac{\partial^2 S}{\partial a_{m-1} \partial a_m} da_{m-1} da_m$$

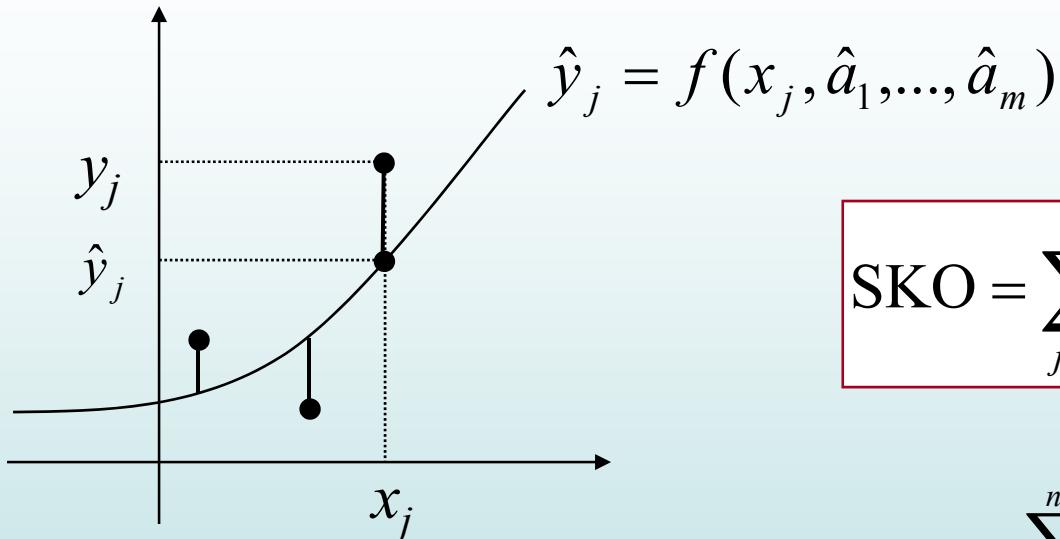
$$d^2S(\hat{a}_1, \hat{a}_2, \dots, \hat{a}_m) > 0$$

tu je minimum funkcije

- Ako ima više tačaka u kojima se postiže lokalni minimum, među svima treba odrediti onu koja ima najmanju vrednost.

# Suma kvadrata odstupanja

- Provera adekvatnosti izabranog modela – pomoću sume kvadrata odsupanja (SKO)



$$\text{SKO} = \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

- Indeks krivolinijske korelacije  $R^2 = 1 - \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y}_n)^2}$   $\bar{y}_n = \frac{1}{n} \sum_{j=1}^n y_j$   
 $R^2 \in [0,1]$
- Bolji izbor  $f$  što je  $R^2$  bliže 1.

# Linearna regresija

- Linearna zavisnost veličina  $X$  i  $Y$  je najjednostavniji oblik zavisnosti ( $Y=aX+b$ ).
- Treba odrediti koeficijente  $a$  i  $b$ .
- Ako raspodela za  $(X, Y)$  nije poznata, tada se  $a$  i  $b$  mogu odrediti iz uslova da

$$E(Y - (aX + b))^2$$

bude minimalno. Iz sistema jednačina

$$\frac{\partial S(a,b)}{\partial a} = 0 \quad \frac{\partial S(a,b)}{\partial b} = 0$$

dobijaju se koeficijenti

$$a = \frac{E(XY) - E(X)E(Y)}{D(X)}$$

$$b = E(Y) - aE(X)$$

# Linearna regresija

- Dat je uzorak  $(x_i, y_i)$ ,  $i=1,\dots,n$ . Neka je zavisnost oblika  $f(x,a,b) = a x + b$ . Koeficijenti  $a$  i  $b$  se nalaze iz uslova

$$\min_{a,b \in R} S(a,b) = \min_{a,b \in R} \sum_{j=1}^n (y_j - (ax_j + b))^2$$

$$\frac{\partial S(a,b)}{\partial a} = 0 \quad \frac{\partial S(a,b)}{\partial b} = 0 \quad \text{Sistem normalnih jednačina}$$

$$\begin{aligned} a \sum_{j=1}^n x_j + nb &= \sum_{j=1}^n y_j \\ a \sum_{j=1}^n x_j^2 + b \sum_{j=1}^n x_j &= \sum_{j=1}^n x_j y_j \end{aligned}$$

$\hat{a}$  i  $\hat{b}$   
rešenje sistema

# Određivanje parametara

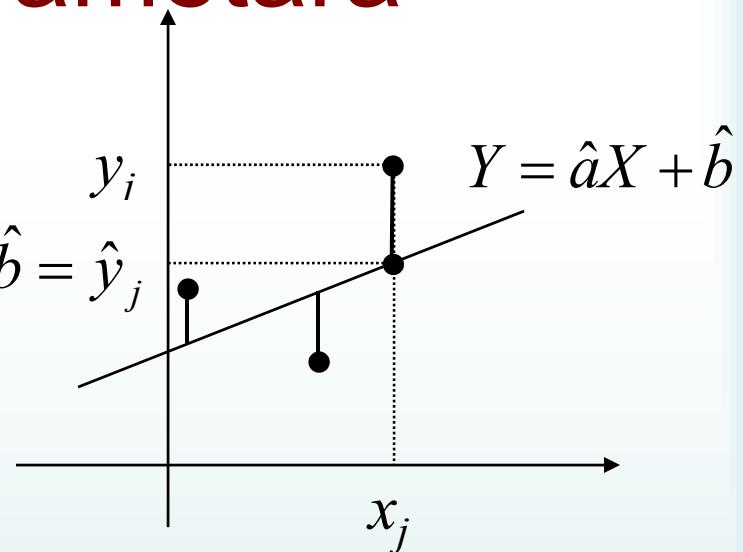
$$\Delta = \left( \sum_{j=1}^n x_j \right)^2 - n \sum_{j=1}^n x_j^2$$

$$\hat{a} = \frac{\sum_{j=1}^n x_j \sum_{j=1}^n y_j - n \sum_{j=1}^n x_j y_j}{\Delta}$$

$$\hat{b} = \frac{\sum_{j=1}^n x_j \sum_{j=1}^n x_j y_j - \sum_{j=1}^n y_j \sum_{j=1}^n x_j^2}{\Delta}$$

$$\hat{b} = \frac{1}{n} \left( \sum_{j=1}^n y_j - \hat{a} \sum_{j=1}^n x_j \right)$$

$$\bar{y}_n = \frac{1}{n} \sum_{j=1}^n y_j \quad \bar{x}_n = \frac{1}{n} \sum_{j=1}^n x_j$$



# Veza parametara, $\rho_{X,Y}$ i uzor. disp.

- Koeficijent korelaciјe

$$\rho_{X,Y} = \frac{E(XY) - E(X)E(Y)}{\sqrt{D(X)}\sqrt{D(Y)}}$$

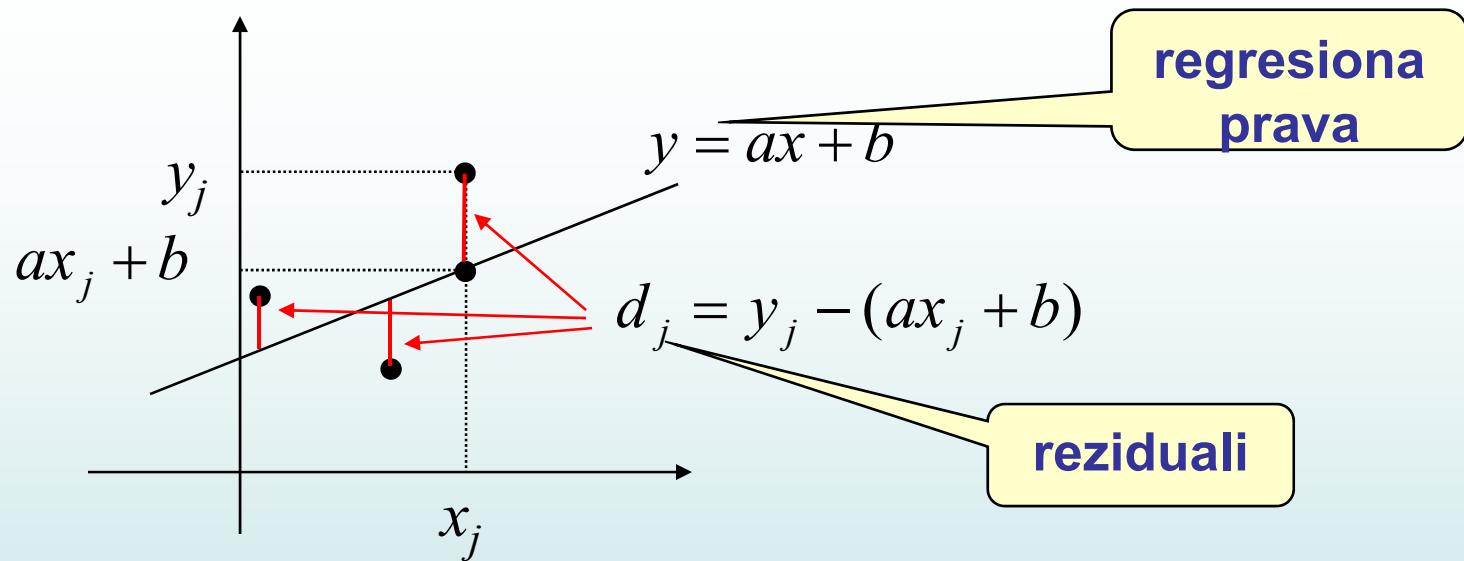
- Uzorački koeficijent korelaciјe

$$E(X) \rightarrow \bar{x}_n \quad E(Y) \rightarrow \bar{y}_n \quad E(X, Y) \rightarrow \bar{x_n y_n} = \frac{1}{n} \sum_{j=1}^n x_j y_j$$

$$r_{X,Y} = \frac{\frac{1}{n} \sum_{j=1}^n x_j y_j - \frac{1}{n} \sum_{j=1}^n x_j \frac{1}{n} \sum_{j=1}^n y_j}{\sqrt{\left[ \frac{1}{n} \sum_{j=1}^n x_j^2 - \left( \frac{1}{n} \sum_{j=1}^n x_j \right)^2 \right] \left[ \frac{1}{n} \sum_{j=1}^n y_j^2 - \left( \frac{1}{n} \sum_{j=1}^n y_j \right)^2 \right]}}$$

$$\hat{a} = \frac{\sqrt{\bar{S}_{nY}^2}}{\sqrt{\bar{S}_{nX}^2}} r_{XY}$$

# Reziduali, regresija



- Od svih pravih  $y=ax+b$  biramo onu za koju je zbir kvadrata odsečaka minimalan.