

Eddig: többváltozós térben bonyolult kérdések voltak:

- adatok mintázata \rightarrow klaszterek
- csoportok - osztályok - klaszifikáció
- folytonos változó modellezése

Néhány egyszerű kérdés egváltozós térben:

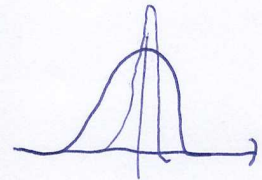
Mit mérünk?

nyers adatok - scatter plot

túlsó adat: átlag, sűrűs, min-max, kvantilisok, histogram, boxplot, $\bar{y}(s)$ forma

Mennyi az értéke - átlagra való kérdések:

levezethető: $\sigma_{\bar{y}}^2 = \frac{\sigma_y^2}{N} \rightarrow \frac{1}{\sqrt{N}}$ a sűrűsége



átlagnál $\frac{1}{\sqrt{N}}$ -re csökken a sűrűsége

t-eloszlás: \bar{y} átlagokból μ valószínűségi átlagra

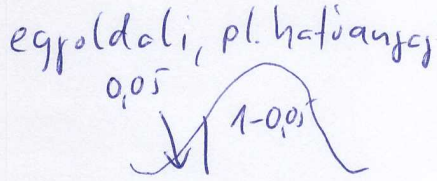
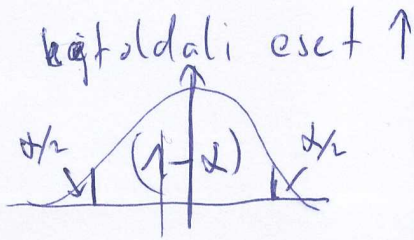
mintából $\frac{\bar{y} - \mu}{\sigma/\sqrt{N}}$ $\nu = N - 1$ szabadsági fokú
t-eloszlást követ

Gosset 1908 (student)

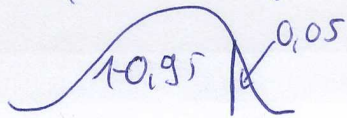
intervallumot adjunk meg ahová pl. 95%-esik

$1 - 0,95 = 0,05 = \alpha$ szignifikancia szint mellett

$$\bar{y} \pm \frac{s \cdot t^{-1}\left(\frac{\alpha}{2}, N = N - 1\right)}{\sqrt{N}}$$



pl. szennyező



lehet 0,01; 0,005; 0,001...

Statistikai tesztek (próba, hipotézis vizsgálata...)

döntési célra használjuk: állításról mondjuk meg, hogy igaz-e? NEM! Többször: csak annyit mondunk, hogy

nem nem igaz valószínűleg...

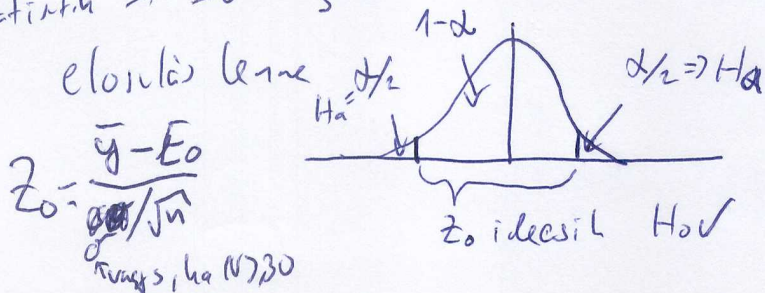
- állítás, amire kíváncsi vagyunk (populációra! nem a mintára)
- elméleti modell felállítás
- elméleti modellhez paraméterek (mérési, mérésóra)
- adott eset \Leftrightarrow elméleti modell valószínűségi összehasonlítása

Példa: z-próba: $n > 30$ minta y_i , E_0

$H_0: E = E_0$ $H_a: E \neq E_0$ (H_0 állításának ellentét eseménye)

mintából: $\bar{y} \rightarrow \sigma$ s hisztogram, majd

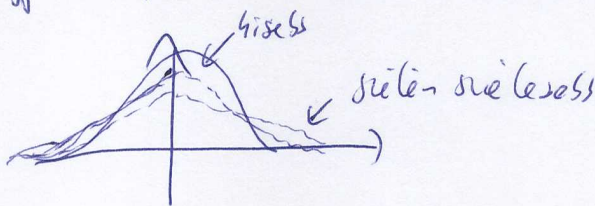
próba statisztika $\rightarrow z_0 = \frac{\bar{y} - E_0}{s/\sqrt{n}}$ hisztogram, E_0 közepes, ha std. normál



terület alulról integrálva
 $0,025 \leq z_0 \leq 0,975$ $H_0 \checkmark$
 +össi H_a

egymintás t-próba

ugyanaz, csak a fv. a t-elosulás $\approx N-1$ gye



$$t_0 = \frac{\bar{y} - E_0}{s/\sqrt{n}}$$

kétmintás t-próba

két mérési sor $H_0: E_1 = E_2$ $H_a: E_1 \neq E_2$

$$d = \bar{y}_1 - \bar{y}_2 \approx 0$$

$$t_0 = \frac{d - E(d)}{s_d} \text{ vizsgálat}$$

két populáció
 átlaga egyenlő-e?
 (pontosabban...)

ún. p-érték itt többször $0,05 \leq p \Rightarrow H_0$
 $> \Rightarrow H_a$

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

ANOVA - egyfaktor (analysis of variance \Leftarrow átlag miatt) ANOVA

alaphérdés: ha többcsoportban kaptunk eredményeket, fűszere az átlagok szempontjából, melyikben

additív modell $y_{ij} = \mu + \alpha_j + \epsilon_{ij}$ \leftarrow egyéni bizonytalanság

\uparrow mért \uparrow teljes \uparrow csoportvárható értékek ellen
 várható érték

$\mu + \alpha_j = \mu_j$

$H_0: \mu_1 = \mu_2 = \mu_3 = \dots$ vagy $\alpha_1 = \alpha_2 = \alpha_3 = \dots = 0$

H_1 : létezik \neq

\bar{y} \leftarrow teljes átlag \bar{y}_j \leftarrow csoport átlag n elem egy-egy csoportban g csoport

$$SS_T = \sum_i^n \sum_j^g (y_{ij} - \bar{y})^2 = SS_{\text{csoport belül}} + SS_{\text{csoport között}} =$$

$$= \sum_i^n \sum_j^g (y_{ij} - \bar{y}_j)^2 + \sum_j^g n(\bar{y}_j - \bar{y})^2$$

\leftarrow mért ~~szórás~~ \downarrow variánciá

szabadsági fokok: SS_T $n \cdot g - 1$ SS_{csoport} $g \cdot (n - 1)$ SS_{csk} $g - 1$

\downarrow

$$s_T^2 = SS_T / (n \cdot g - 1)$$
 ~~s_{csoport}~~

$$s_{\text{csoport}}^2 = \frac{SS_{\text{csoport}}}{g(n-1)}$$

$$s_{\text{csk}}^2 = \frac{SS_{\text{csk}}}{g-1}$$

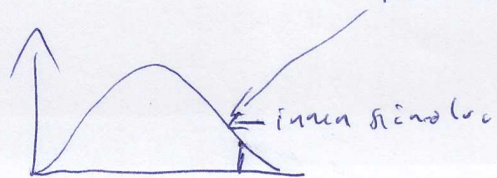
\rightarrow Mindegyik Feloszlással írható le páronként osztva (normál eloszlás...)

$$F = \frac{s_{\text{csk}}^2}{s_{\text{csoport}}^2} \Rightarrow F \sim e^2 \quad \begin{matrix} \nu_1 = (g-1) \\ \nu_2 = g(n-1) \end{matrix} \leftarrow \text{null. h. közti } F_{\nu_1, \nu_2}$$

táblázatban:

forrás	négyzetösszeg	ν	var	F	$P(H_0)$
csoport között	SS_{csk}	$g-1$	$\frac{s_{\text{csk}}^2}{s_{\text{csoport}}^2}$	$\frac{s_{\text{csk}}^2}{s_{\text{csoport}}^2}$	rejtett
csoport belül	SS_{csoport}	$g(n-1)$	s_{csoport}^2		\uparrow
teljes	SS_T	$g \cdot n - 1$			pontosabban vizsgáljuk! $P(X > F)$

~~további régi táblák:~~
 egyenlő a $P > 0,9$ is! miért?



Rengeleg teszt köthető a χ^2 eloszláshoz

$\chi^2 = \sum x_i^2$ $x_i \in N(0, 1)$ $\nu = \text{degrees of freedom} = N$



$E(\chi^2) = \nu$ $\max \chi^2 \rightarrow \nu - 2, \text{ ha } \nu > 2$
 $\sigma^2(\chi^2) = 2\nu$

pl: Illeszkedés vizsgálása egy adott eloszláshoz: H_0 : minta A-ból származik
 H_1 : nem

A felosztás A_i részekre, histogramm képzése

$\chi^2_0 = \sum_{i=1}^k \frac{(y_i - e_i)^2}{e_i}$

y_i : mintából A_i -be esők száma
 e_i : N.P.C(A_i) $N = \sum_{i=1}^k y_i$

elvetés, ha $\chi^2_0 \geq \chi^2_{\alpha, \nu}$

$\nu = k - 1 - r$
 r = paraméterek száma az eloszlásban

Kontingencia táblázat

mért:

	dadog	püste	
férfi	32	28	→ 60
nő	18	22	→ 40
	↓	↓	↓
y_{ij}	50	50	100

függetlenség esetén várt

	d	p	
f	30	30	→ 60
n	20	20	→ 40
	↓	↓	
e_{ij}	50	50	

H_0 : A populációk ugyanazok a kategóriákban
 H_1 : nem

vagy: H_0 : Sor és Oszlop kategóriák függetlenek egymástól

H_1 : függnek

$\chi^2_0 = \sum_i \sum_j \frac{(y_{ij} - e_{ij})^2}{e_{ij}}$ $e_{ij} = \frac{\sum_k y_{ik} \cdot \sum_l y_{lj}}{\sum_k \sum_l y_{kl}}$

elvetés, ha $\chi^2_0 \geq \chi^2_{\alpha, \nu}$

$\nu = (\text{sorok} - 1)(\text{oszlopok} - 1)$

Másik gyakorlat: t-eloszlás

TS4

párosított adatokra, pl.: meredekségre

$$y = ax + b \text{ modellre } (y = ax + b + \varepsilon)$$

\uparrow normál eloszl.

$$S_{yy} = \sum y_i y_i - \bar{y} \bar{y}$$

$$H_0: a = a_0$$

$$H_a: \neq, \text{ vagy } >, \text{ vagy } <$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{\sum x_i y_i - \bar{x} \bar{y}}{\sum x_i x_i - \bar{x} \bar{x}}$$

$$t_0 = \frac{a - a_0}{s_{xy} / \sqrt{S_{xx}}}$$

$$s_{xy}^2 = (S_{yy} - a S_{xy}) / (N - 2)$$

elvetés:

$$|t_0| \geq t_{\alpha/2, n-2}, \text{ vagy } t_0 > t_{\alpha/2, n-2}, \text{ vagy } t_0 < -t_{\alpha/2, n-2}$$

Korrelációs e.h.-ra

$$H_0: r = 0$$

$$H_a: r \geq 0 \text{ vagy } r > 0 \text{ vagy } r < 0$$

$$t_0 = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

$$\text{elvetés: } |t_0| \geq t_{\alpha/2, n-2}, t_0 \geq t_{\alpha/2, n-2}, t_0 \leq -t_{\alpha/2, n-2}$$

Élvezik $r = r_0$ és $r_1 = r_2$ változatis, másokgy

tengelymetre

$$t_0 = \frac{b - b_0}{\text{RMSE} \cdot \sqrt{\frac{1}{N} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2}}}$$

 $\gg n-2$ t-vel összehasonlítva

Ha nem normál eloszlás:

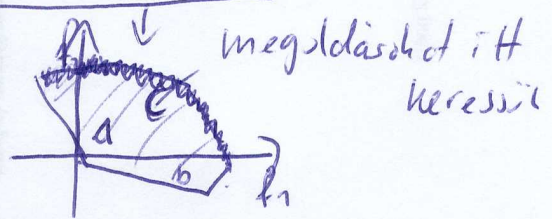
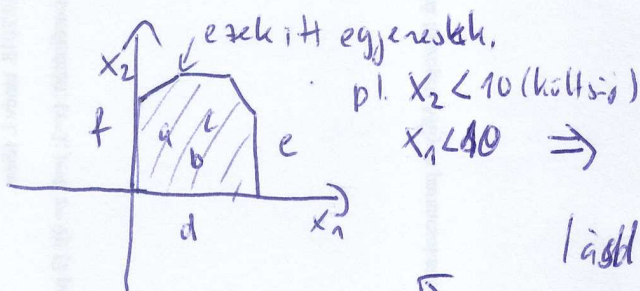
páros t-teszt \leftrightarrow Mann-Whitney u-tesztanova \leftrightarrow Kruskal-Wallis testegymintás \leftrightarrow Wilcoxon^o előjelrang teszt

Több kritériumú döntések (Multi-criteria Decision Making) MCDM/A

Optimális döntés keresése, ha több szempont van.

$x_1 \dots x_m$ kritérium $f_1(x) \dots f_n(x)$ döntési függvény
(min. v. max)

elfogadható teret definiálható a kritériumok terében,
pl. 2db kritérium \leftarrow döntési fu.-ek terében



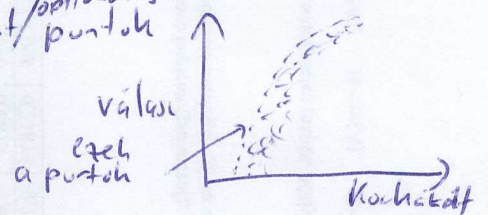
lásd: [Multicriteria decision analysis wikipedia](https://en.wikipedia.org/wiki/Multicriteria_decision_analysis)

Hol vannak a pontok \Rightarrow a, b, c van belül
non-dominált pont (nem leuralható):
ha nincs olyan, ami mindkétben jobb:

\rightarrow b, c-értéke a legrosszabb
az origótól, ha pl. f_1, f_2 max
a kezdés, b vagy c

pl. ~~van~~ $a + a$ b, c \rightarrow

\approx Pareto hatékonyság = pontok, ahol
nem javítható valahelyen tovább a jölele,
amikor, hogy a többieknek romlana.
Pareto efficient/optimalis pontok



Derringer módszer MCDM-re:

$x_i \in T$ -dik kritériumra egy adott esetben (nagy legyen a jó)

x_{max} : legjobb x_{min} : legrosszabb $d_i = \frac{x_i - x_{min}}{x_{max} - x_{min}}$ vagy

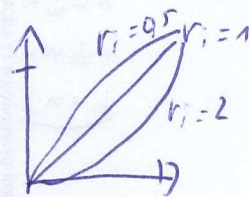
lehet úgy is, hogy x_{min}, x_{max} az értelmes határ $x_j < x_{min} \rightarrow d_j = 0$
és akkor ha $x_j < x_{min} \rightarrow d_j = 0$ $x_j > x_{max} \rightarrow d_j = 1$

adott esetre

$$D = \sqrt{\prod_{i=1}^m d_i^{r_i}} \quad \text{számolás}$$

$r_i = 1$ lineáris
 $r_i = 0,5, r_i = 2$

MCDM/2



max D a legjobb döntés

Topsis algoritmus

- $X^{n \times m}$ mátrix, n eset, m kritérium (X_{ij} nagy lehet $\frac{z_j}{z_j^+}$)

- normalítás:

$$r_{ij} = \frac{x_{ij}}{\sqrt{\sum_{k=1}^n x_{kj}^2}} \quad i=1, 2, \dots, n$$

- w_j eredeti súlyokkal $w_j = \frac{w_j}{\sum w_k}$ normált súlyokkal

- legrossabb és legjobb elvi pont összehasonlítása

$A_w \leftarrow$ hipotetikus pont, összes kritériumnál a legrossabb az érv

$A_b \leftarrow$ hipotetikus pont, összes kritériumnál a legjobb az érv

- n eset ~~szere~~ távolság meghatározása A_w -től és A_b -től

pl: $d_{iw} = \sqrt{\sum_{j=1}^m (t_{ij} - t_{w_j})^2}$ $d_{ib} = \dots$

$$S_{iw} = \frac{d_{iw}}{d_{iw} + d_{ib}}$$

\in használható a legrossabbhoz

S_{iw} rangsorba vehető, nagy a legjobb döntés

$S_{iw} = 1$ elvi legjobb

$S_{iw} = 0$ elvi legrossabb

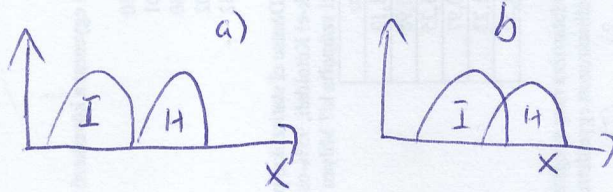
Egy döntés a létrehozásra (lehetőleg előző!)

$f(\underline{w}, \underline{x}_+)$ alapján pl. max-ra

ROC görbe (Receiver-Operating Characteristic)

ROC

Klassifikáció 1 változó szerint:



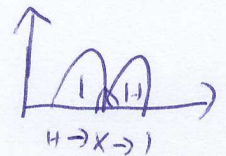
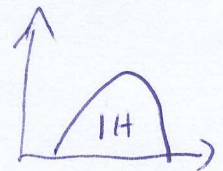
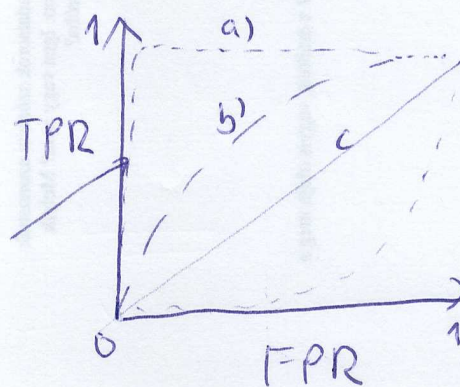
ha x -legyen a változóval egy görbe, hanem két függő x -től:

valós/döntés igaz hamis
 igaz TP FN
 hamis FP TN

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{TN + FP}$$

a küszöbön a döntési határ változik



AUC = görbe alatti terület \Rightarrow 1 tökéletes a klassz.

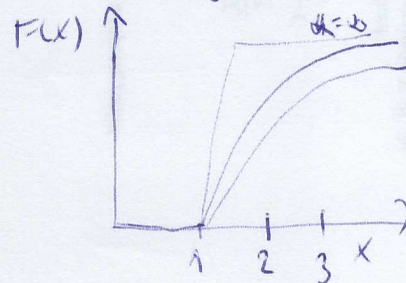
Pareto-eloszlás

statisztika: 80%-a a javaknak 20%-nál a

Pareto index társadalomnál... Illyen eloszlást sok minden

↓ követi
 $x \geq x_m$
 $x < x_m$

$$F(x) = \begin{cases} 1 - \left(\frac{x_m}{x}\right)^\alpha & x \geq x_m \\ 0 & x < x_m \end{cases}$$



kumulatív!

Lehet pl. fah torony is így kategorizálható \rightarrow
 (bemenő változó)



mit juttatunk...

Box-Cox transformáció:

Szeretjük, ha a hiba/bizonytalanság kromoskedantikus,
vagyis nem függ pl. y értéktől. átalakítsuk át eredetire:

$$\sigma_y = y^\lambda \Rightarrow y^* = y^\lambda \Rightarrow \text{var}(y^*) = \text{konstans}$$

λ	λ	$y^* =$
2	-1	$1/y$
1,5	-0,5	$1/\sqrt{y}$
1	0	$\ln y$
0,5	0,5	\sqrt{y}
0	1	y (nincs transformáció)