

90 perc STATISZTIKA

Bioanal., ka
2018/2019 II
febr
Tóth György
Kémia SA-ba

http://tothgyorgy.web.elte.hu ← letölthető innen
Vegyes hallgatói háttér ⇒ alapfogalmak is megjelölve

(elmélet, anyag - gyakorlati rész támogatására)

Valószínűségsszámítás: események sokaságának tárgyalása
valószínűségi alapon (determinisztikus til
összetett) (populáció)

Statistika: minta tulajdonságai alapján mit mondhatunk a
sokaságról

Várható érték és becslés

VSI-ben. $E(x) = \int_{-\infty}^{\infty} x \cdot f(x) dx$
↑ val. sűrűség fv.
folytonos

diszkrét val. v. t. eloszlás
 $E(x) = \sum_i x_i \cdot p(x_i)$
↑ p_i kockázat
 $p(x_i) = 1/6$
 $E(x) = 3,5$

minta: N darab y_1, y_2, \dots, y_n

átlag $\bar{y} = \frac{\sum_{i=1}^n y_i}{N}$

kitöltés adatnál:
median

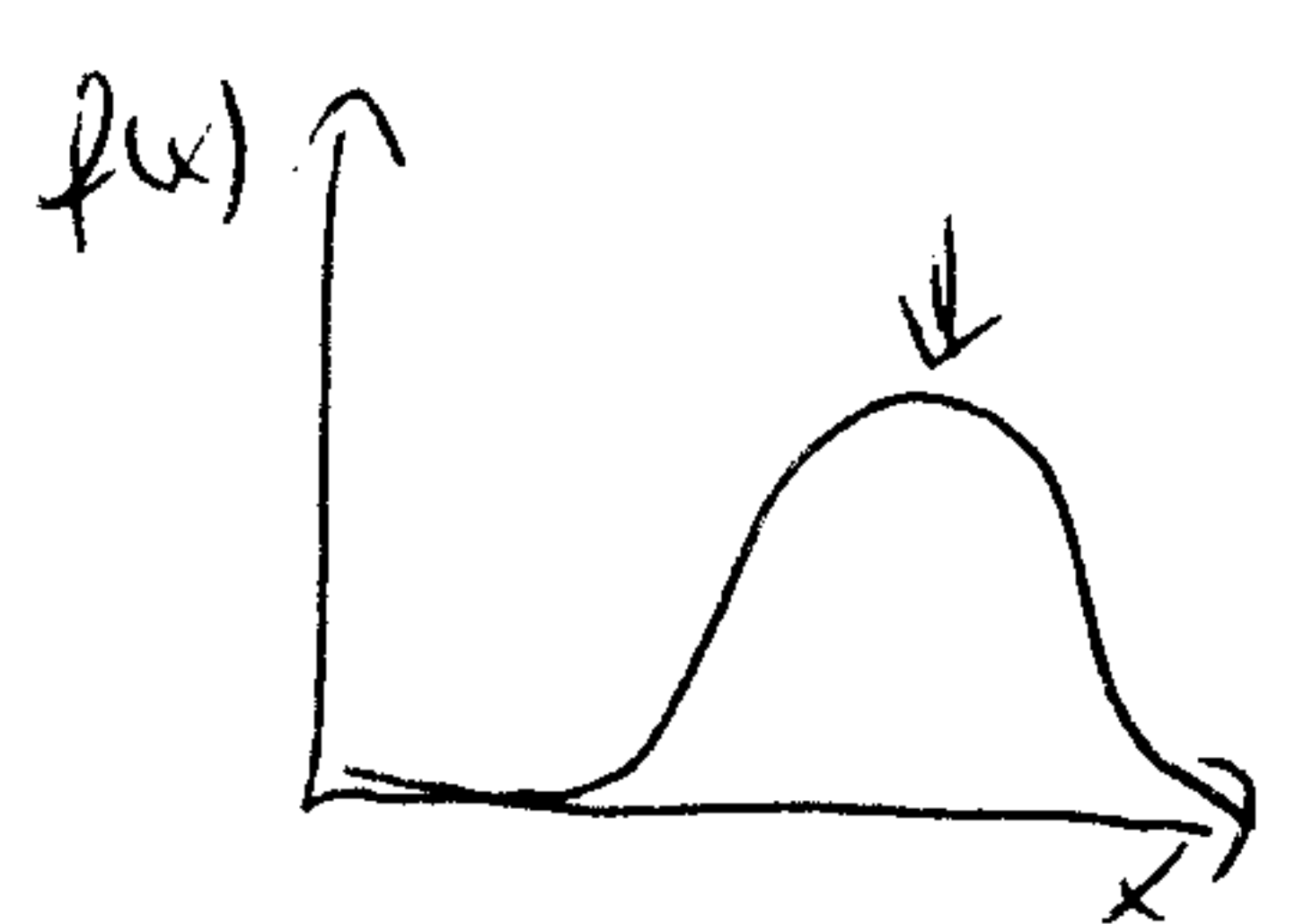
- 1) sorbarendezés
 $y_1 \leq y_2 \leq y_3 \leq \dots \leq y_n$
- 2) középső, vagy középső átlag

robosztus becslés: sokféle adatra jó, hibásokra is (nem értékes)

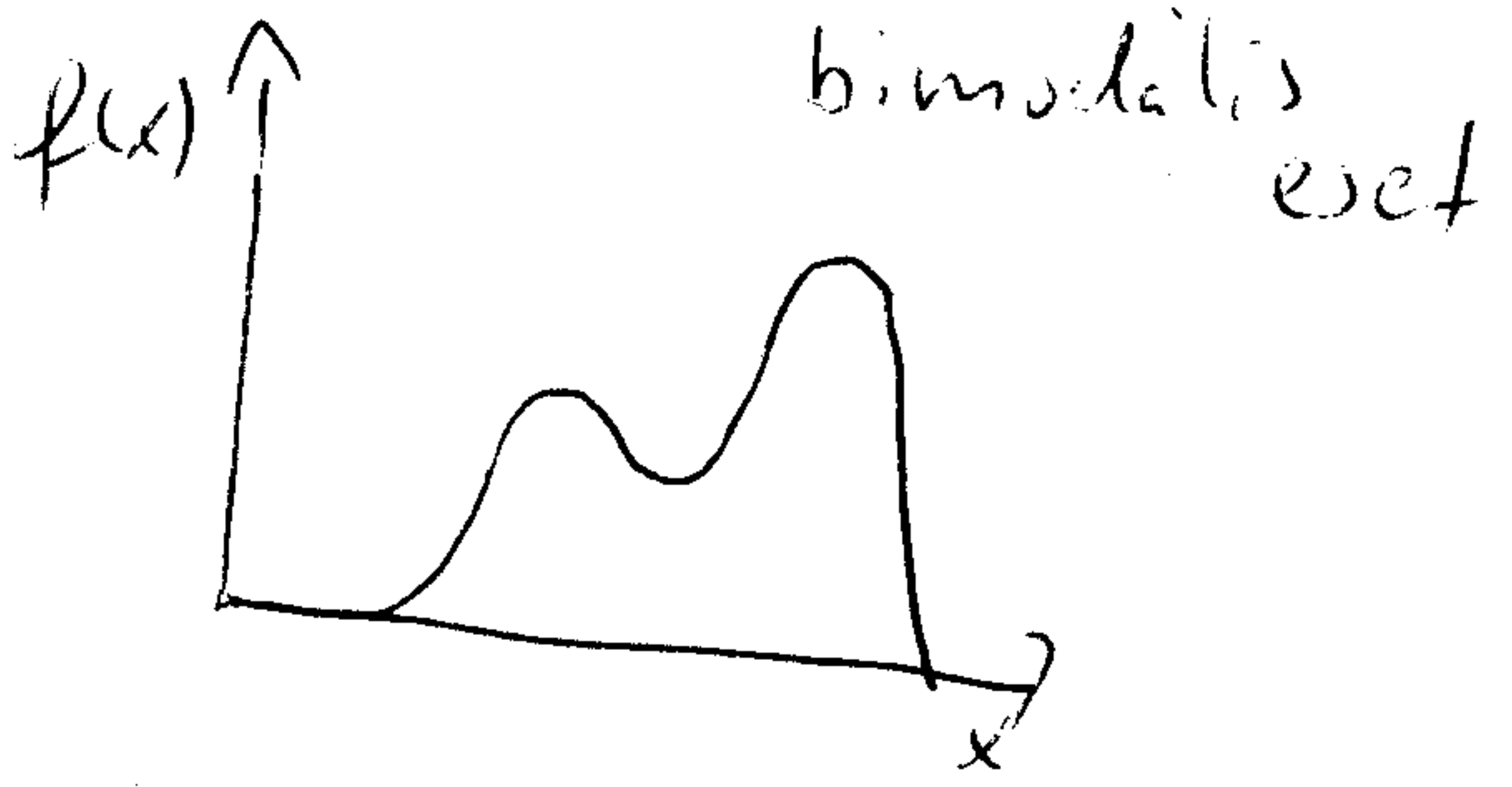
Modus: leggyakoribb



diszkrét val. v. t. eloszlás

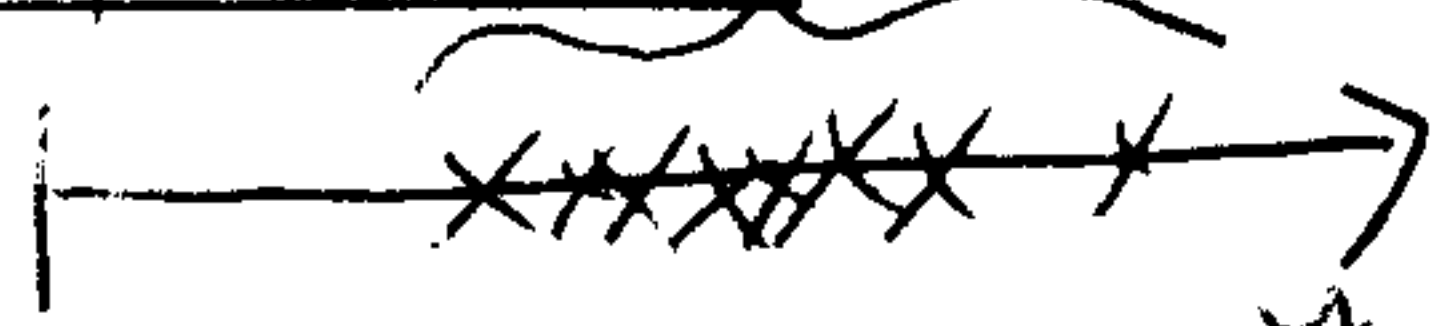


$f(x)$: val. sűrűség függvény, integrálja 1



bimodális eset

Terjedelem?



↑ helye: várható érték

Variancia, szórásnégyzet:
$$\sigma^2 = \int_{-\infty}^{\infty} (x - E(x))^2 f(x) dx$$

$$\sigma^2 = \sum (x_i - E(x))^2 p(x_i)$$

mintából:

estimated variance

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2$$

$s = \sqrt{s^2}$ becslt szórási (standard deviation)

variációs koefficiens (coefficient of variation)

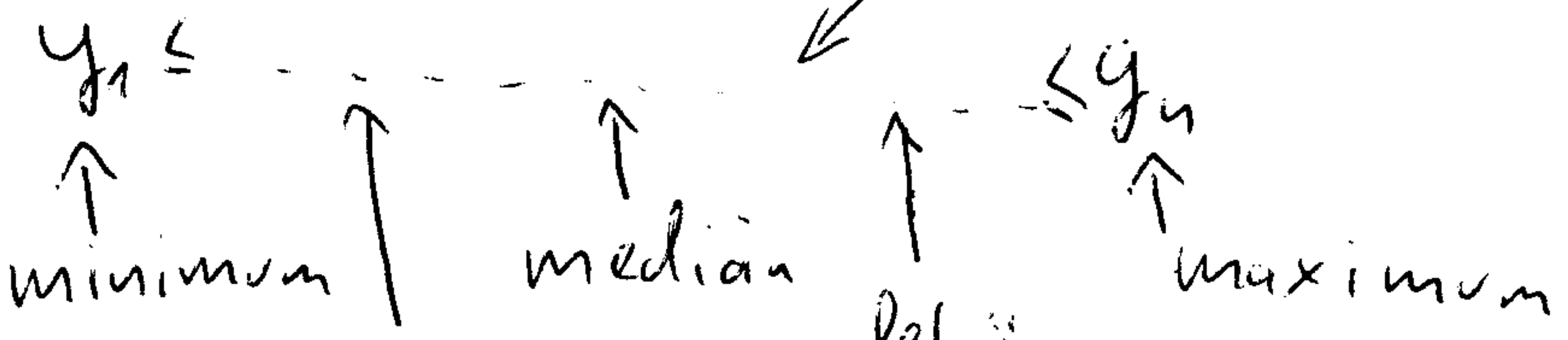
$$c.o.v = \frac{s}{\bar{y}}$$

más ha k vagy °, 0 körüli a good!

Szórásrendezett mintánál

tetszőleges = percentilis

$y_{max} - y_{min} = \text{range (terjedelem)}$

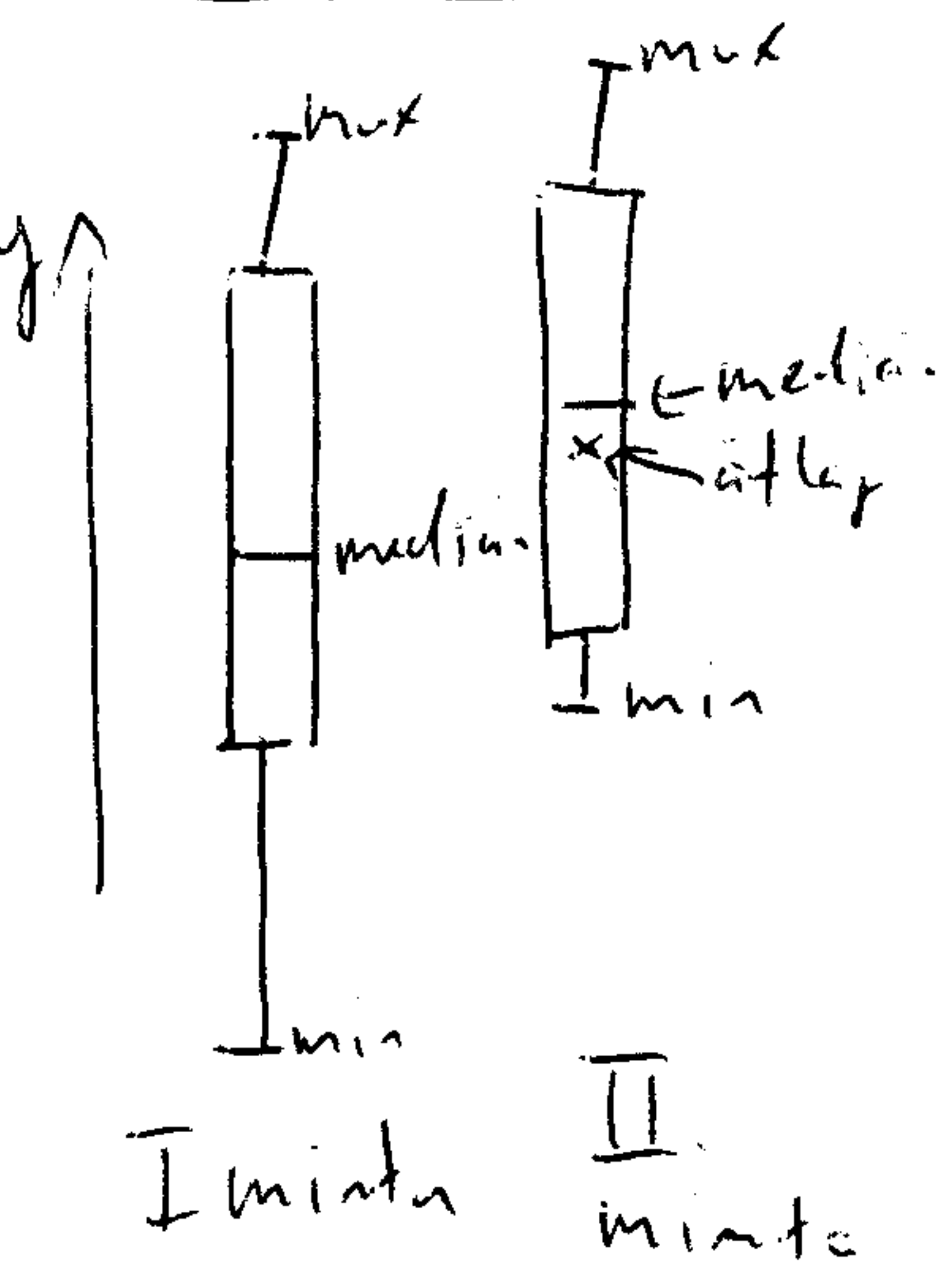


$y_{3/4} - y_{1/4} = \text{interkvartilis távolság} = d_{int}$

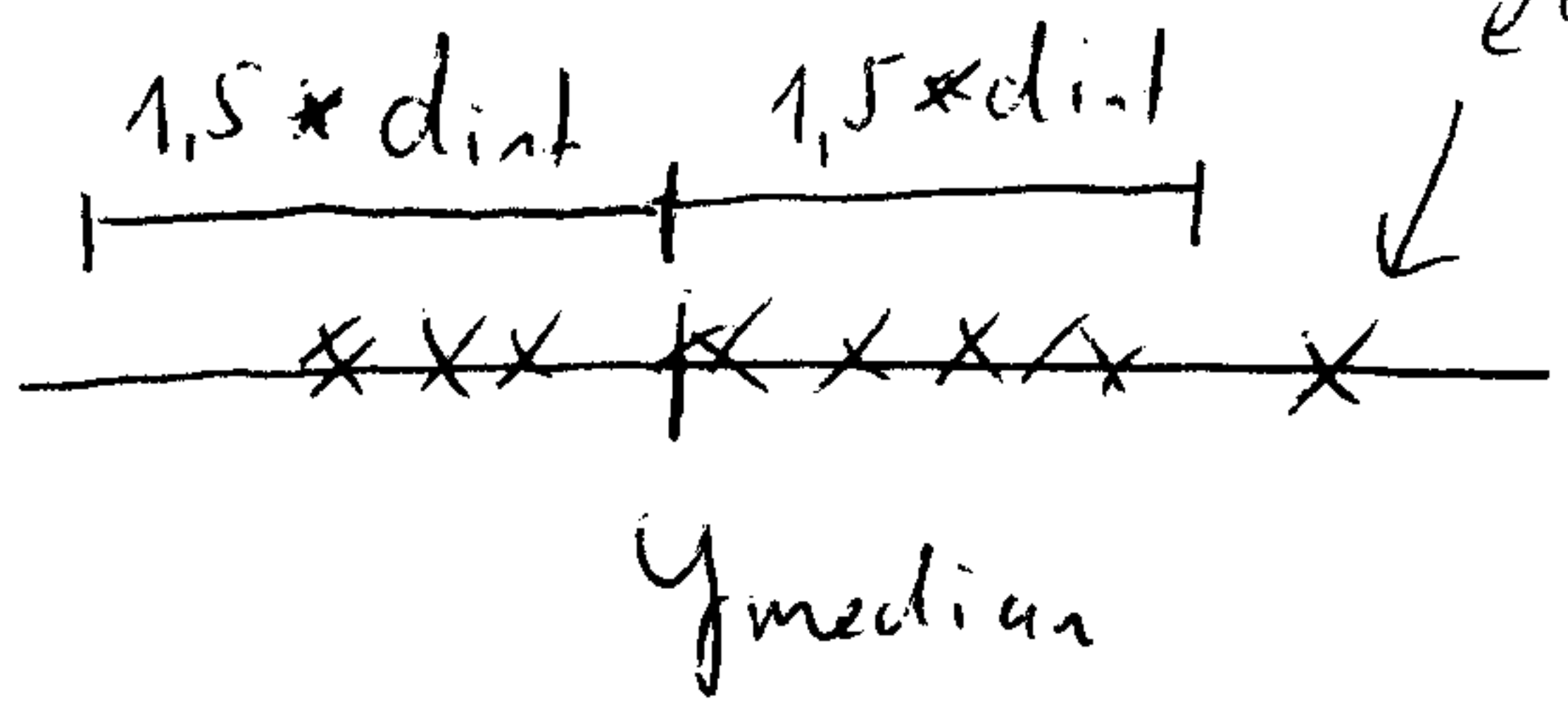
alsó kvartilis első kvartilis

miért jó:

boxplot:



hílygi adat eldobása:



eldobható hílygi adat!

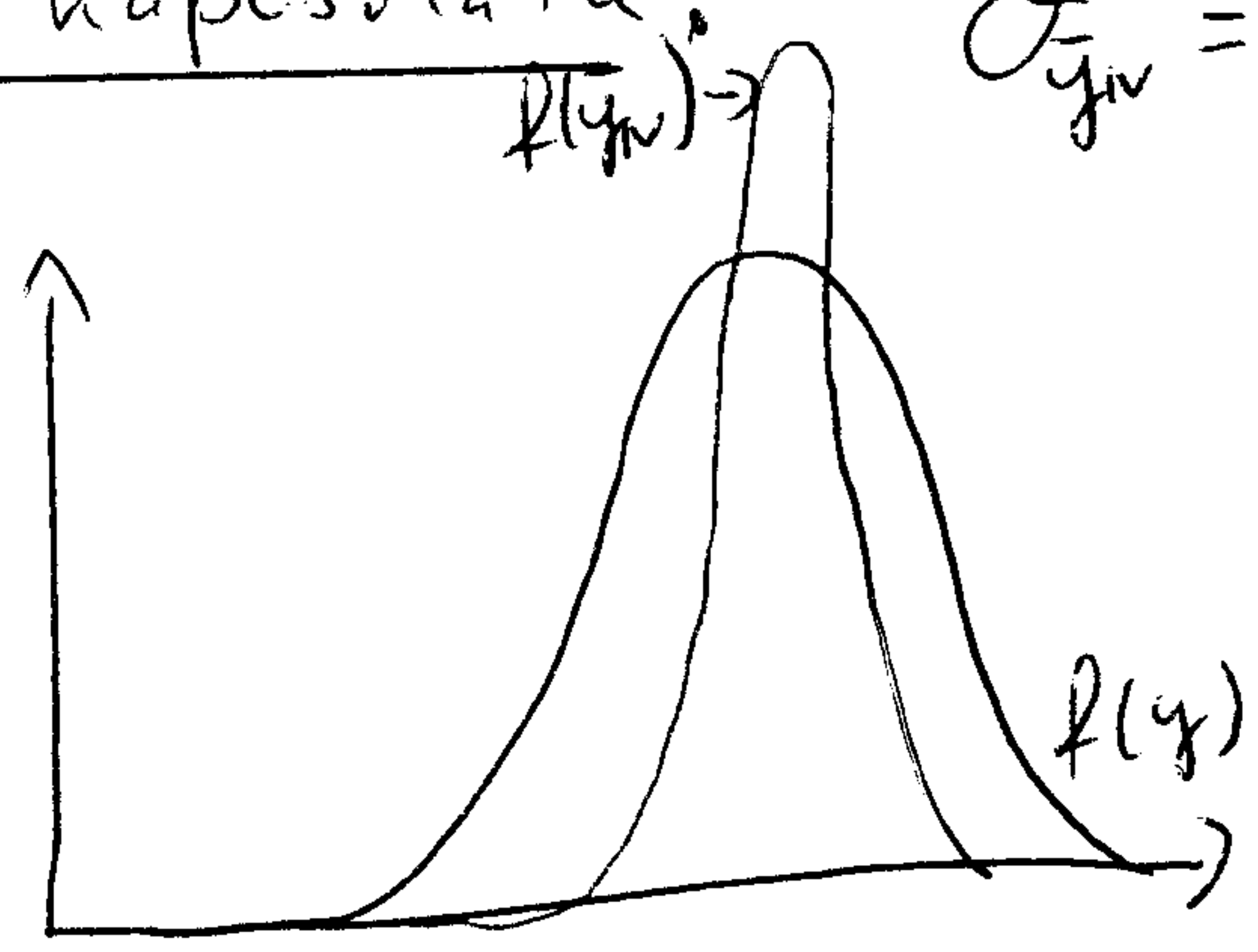
Analitikus mér, hogy adatot adjon meg (becsüljön)

- pontbecslés \rightarrow 1 szám
- intervallumbecslés \rightarrow adott valószínűséggel ebben az intervallumban van (pl. 95%, 99%)

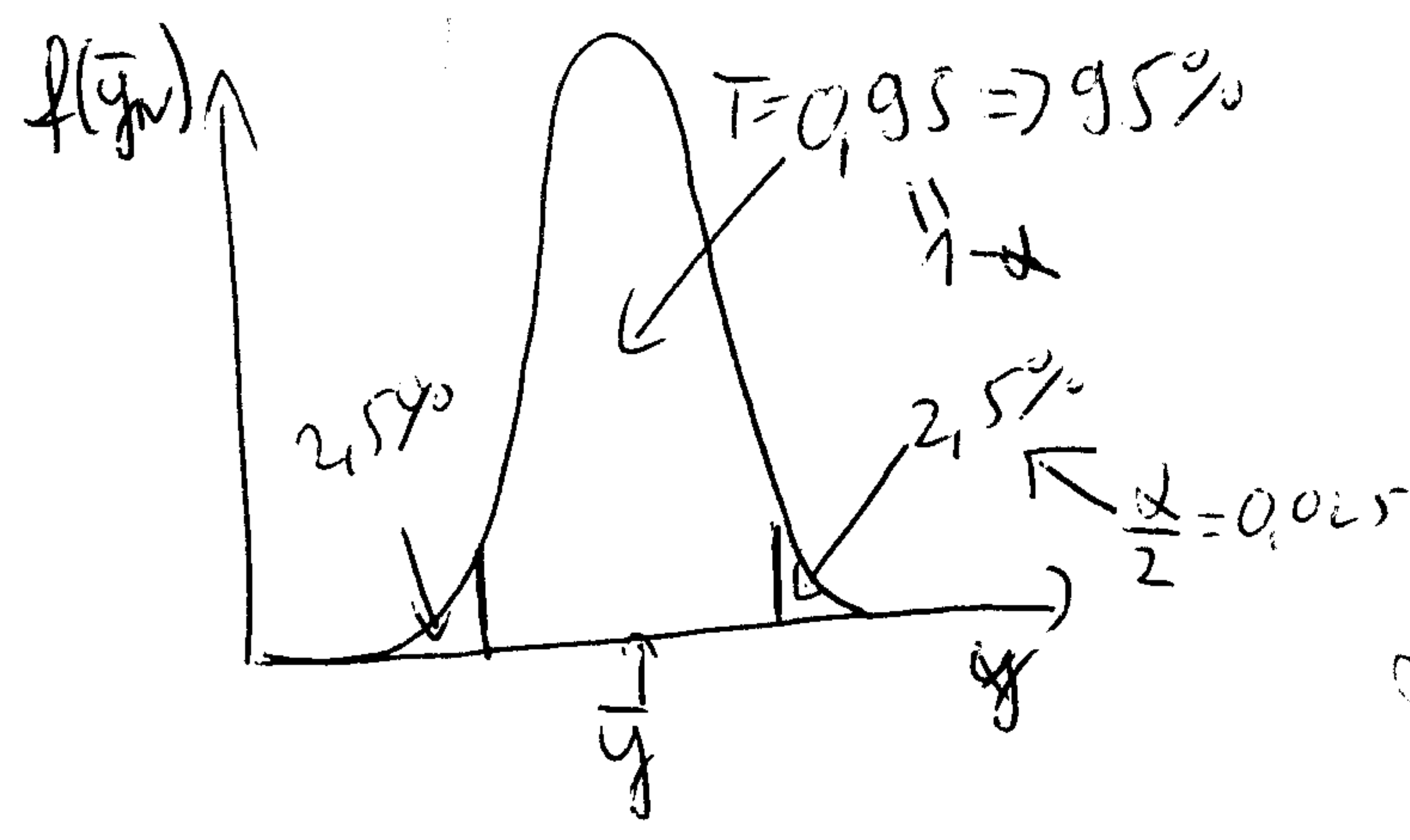
σ_y és $\sigma_{\bar{y}_n}$ kapcsolata:

$$\sigma_{\bar{y}_n} = \frac{\sigma_y}{\sqrt{N}}$$

\Rightarrow több mérés, szűk az intervallum



Konfidencia intervallum várható értékre: $\bar{y} \pm ?$

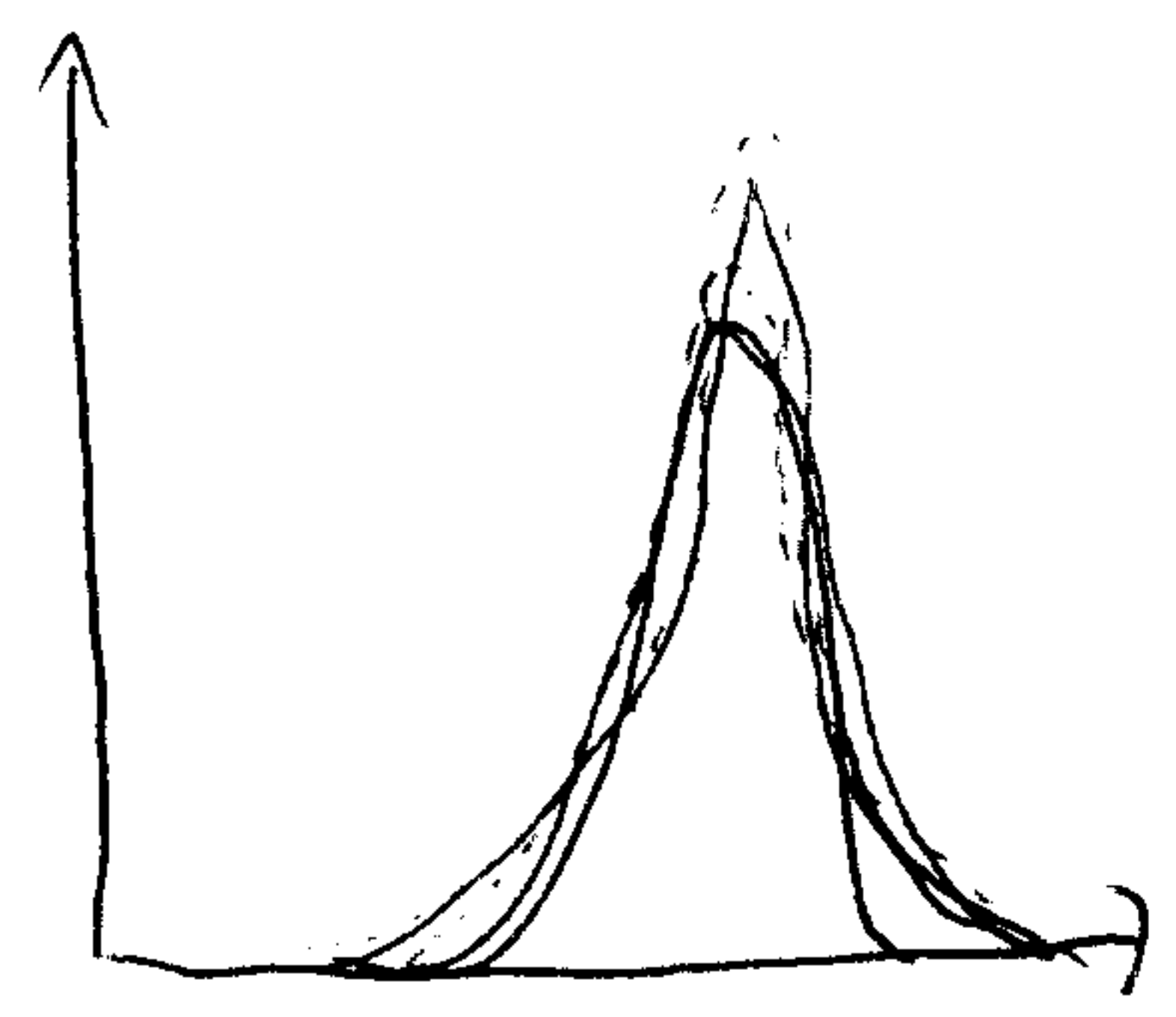


szimmetrikus, ha normál eloszlás és $N \geq 30$

$$\bar{y} \pm \frac{1,96 \cdot s}{\sqrt{N}}$$

α szignifikancia szint

hisz mintánál gond, hogy s és \bar{y} is ugyan abból van becsülve, normál eloszlás \rightarrow student eloszlás (t-eloszlás)



$$\bar{y} \pm \frac{|t_{\alpha/2}^{-1}| \cdot s}{\sqrt{N}}$$

γ szabadsági fok

$$\nu = N - 1$$

$\bar{y} \pm$ \neq $y(s)$
 \uparrow \uparrow
 ez a Konfidencia intervallum ide sorost soka irni

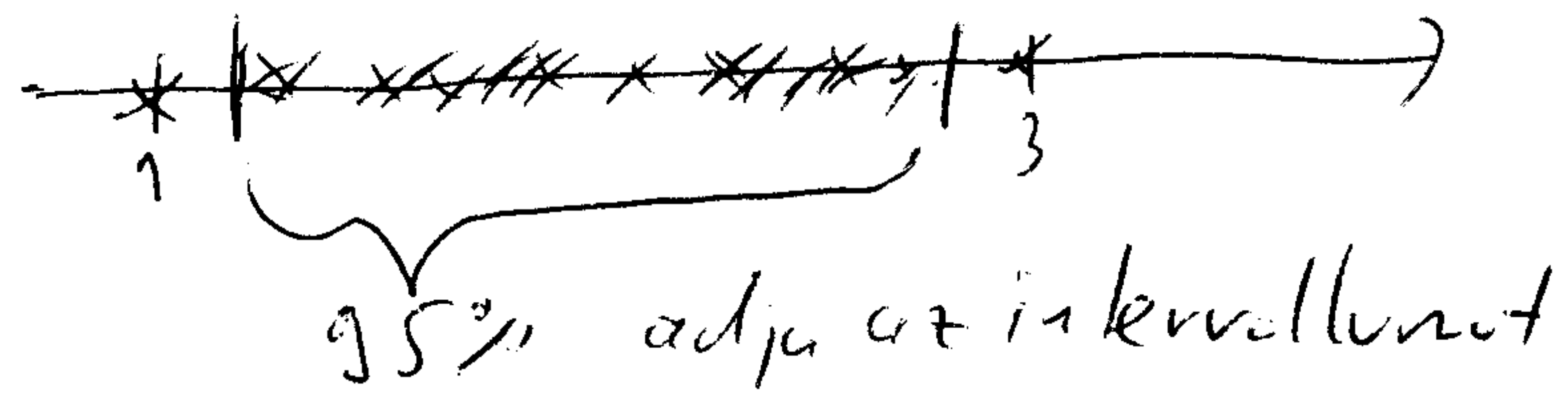
Konfidencia intervallum ^{magyar} nem normalis elosulasra.

pl. bootstrap módszer: minta \Rightarrow populacio, visszateveses ujra mintaveveleses:

$y_1=1, y_2=2, y_3=3 \Rightarrow 27$ új minta

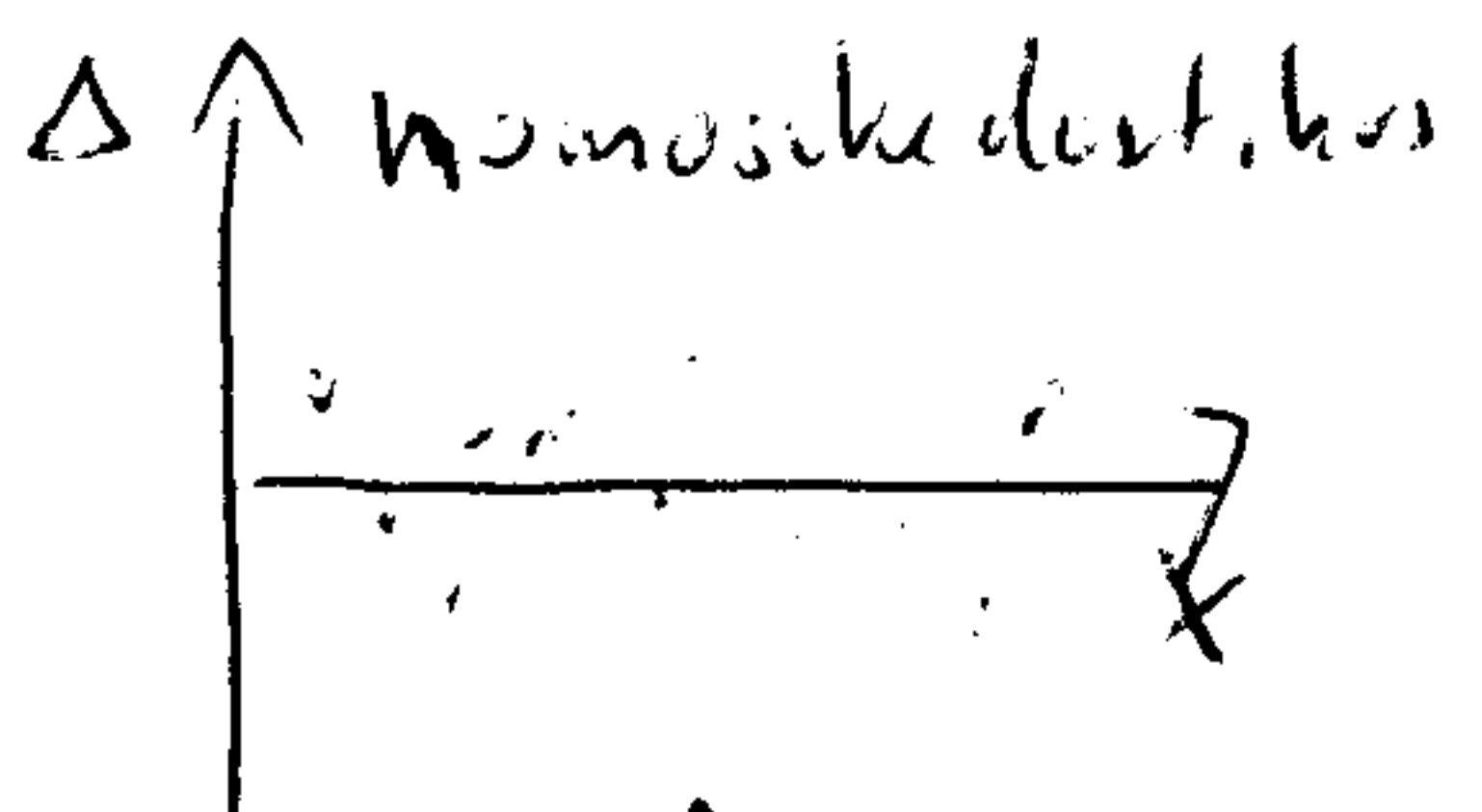
- 111
 - 222
 - 333
 - 123
 - 122
 - 322
 - 211
 - 311
 - 133
 - 233
- } 27 db

ezek $y_{a,b}$



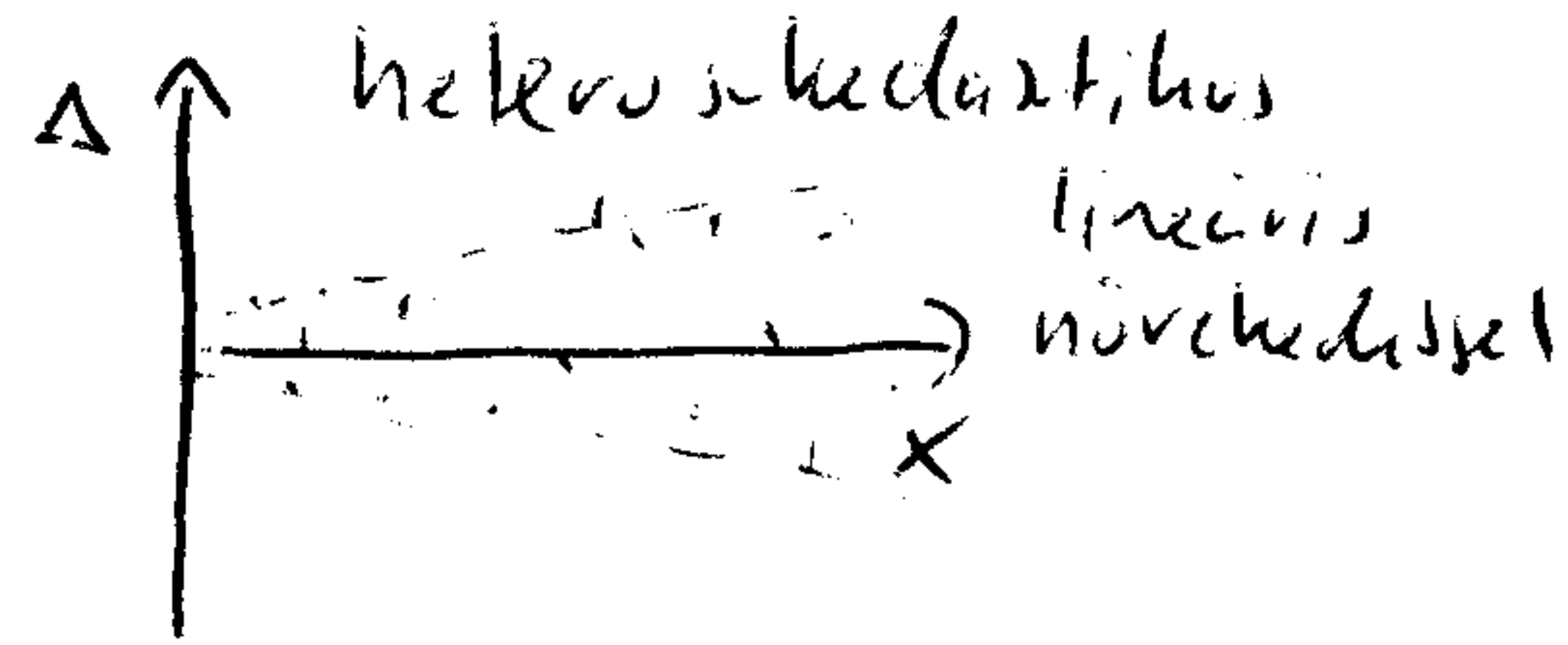
Statistika: döntéseket hozhatunk
statisztikai tesztek (próba, hipotézis vizsgálata)

pl: Validálás során tesztatlanság vizsgálható.
 x pontszám bemeiert - y meghatározott koncentráció $\Delta = y - x$



↑
átlagos tesztatlanság vizsgálata

$H_0: E(y-x) = 0 = E_0$



↑
átlagos visszatérés vizsgálata

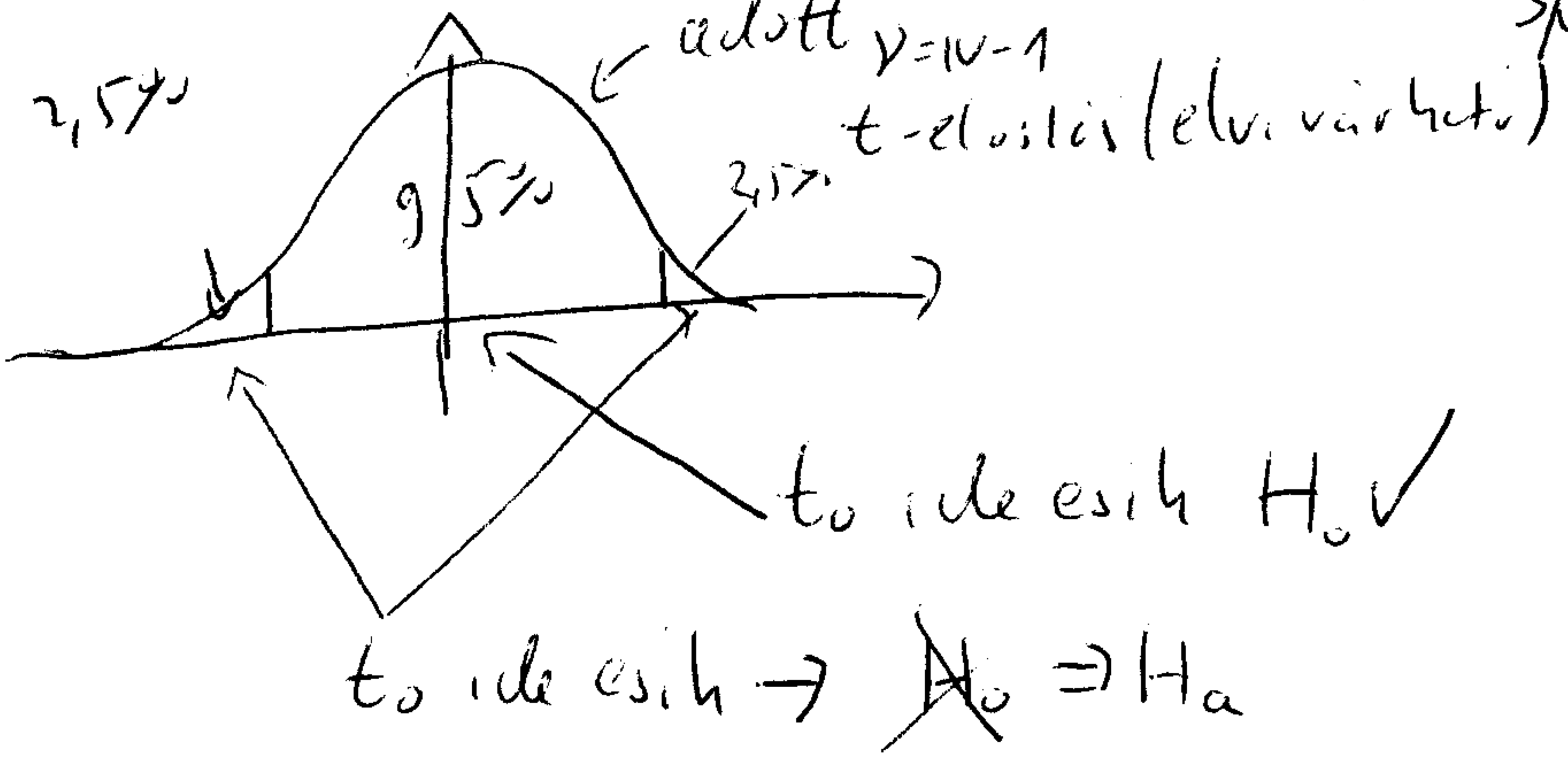
$H_0: E(\frac{y}{x}) = 1 = E_0$

módszer: egy mintás t-próba

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

EXCEL:
 T.ELOSUL($t_0, N-1, 1$)

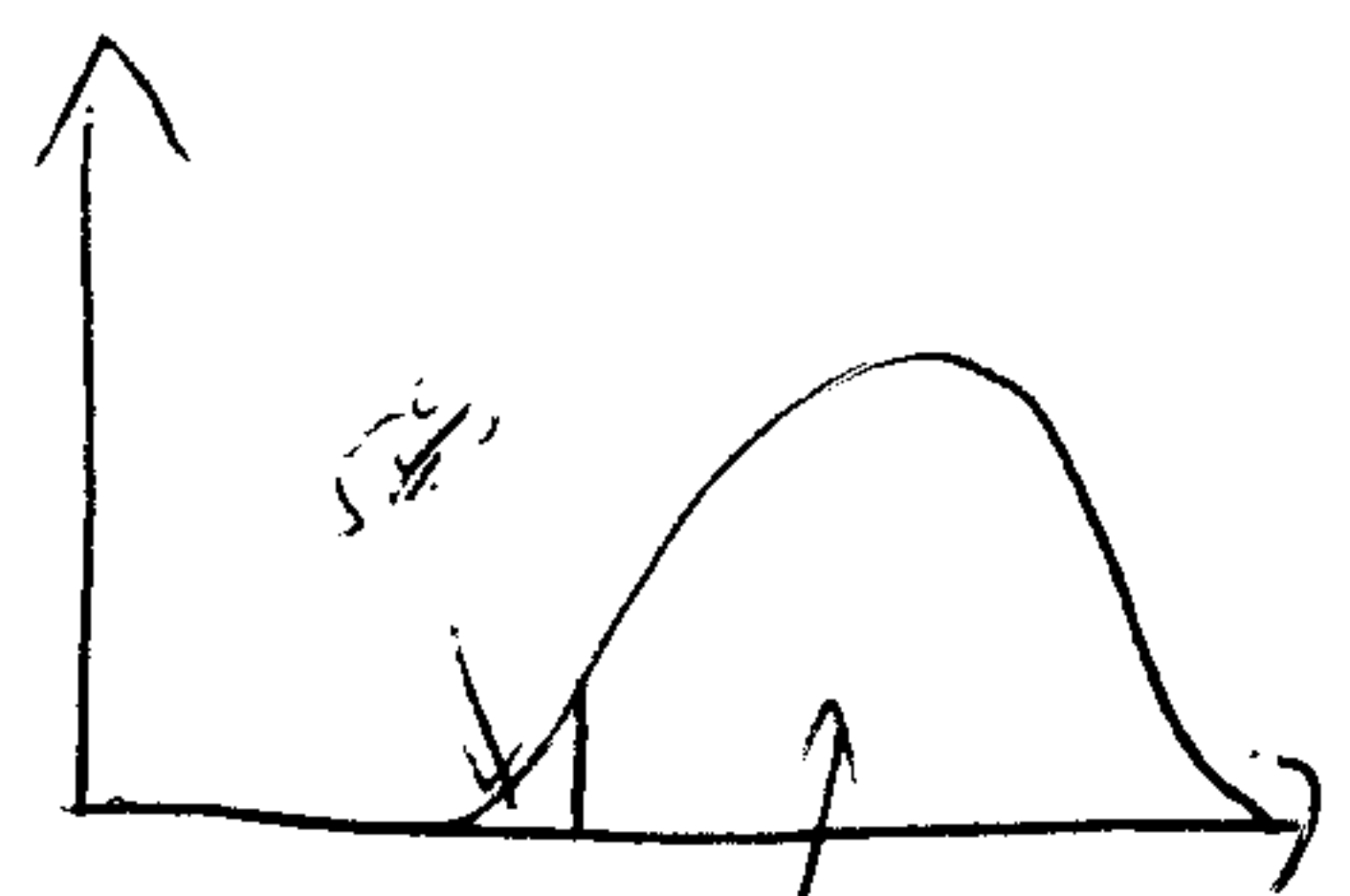
↓
 kiválasztás



ha s helyett σ

↓
 a vagy z próba
 a neve, normal elositas

lehet egyoldali eset



ide essen a határangaj



ide essen a szemérvé

A nagyominusos statisztikai proba H_0 -nak kedveznek (a gyártónak)

Validáció: két labor összehasonlítása

kétmintás t-próba

$H_0: E_1 = E_2$ $H_a: E_1 \neq E_2$

mivel s_1 sokkal nagyobb, mint s_2 , ezért feltételezzük az H_0 -t, hogy nem egyenlő varianciájuk. ~~hát~~ (variancia becslés mindkettőnél)

labor-1	labor-2
-	-
-	-
-	-
-	-

\Rightarrow EXCEL P-értékkel kapunk pl. adatelemzés (van párosít $H_{1,1}$)

$0,05 \leq p\text{-value}$ $H_0 \checkmark$

$0,05 > p\text{-value}$ H_a

Regressió:

egyváltozós \rightarrow egyenes illesztés $x \rightarrow y$

többváltozós \rightarrow sík, hipersík illesztése $\underline{x} \rightarrow y$

N x_i, y_i pár. $y_i = ax_i + b + \epsilon_i$

ϵ_i homoskedasztikus \rightarrow nem kell súlyozni, a pontok w_i

$w_i = 1/\sigma_i^2$

leghibesebb négyzetek módszere:

OK, ha $cov(\epsilon_i, \epsilon_j) = 0$

$Q = \sum_{i=1}^N (y_i - ax_i - b)^2$

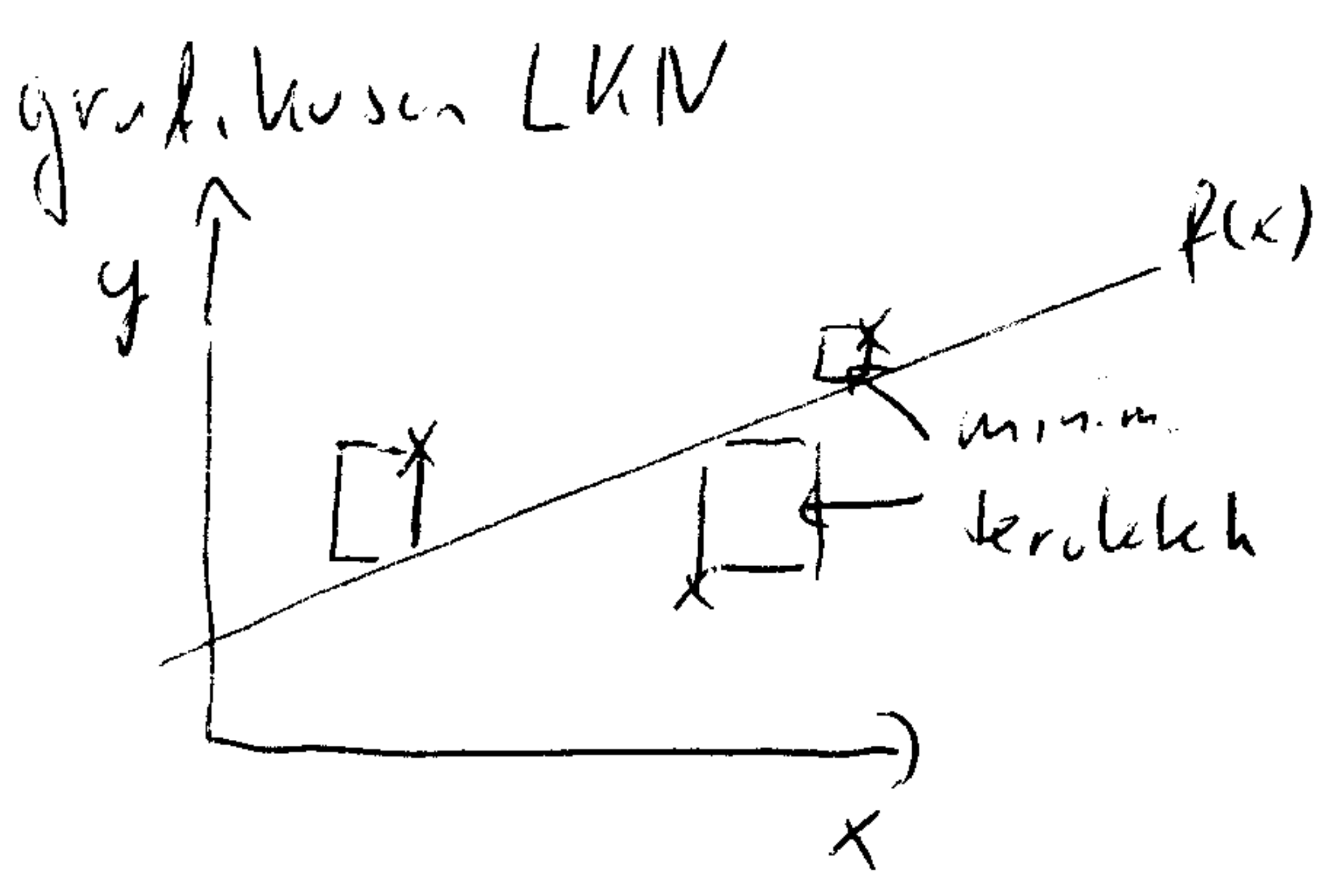
$E(\epsilon_i) = 0$
normal eloszlás

$\min Q \Rightarrow \frac{\partial Q}{\partial a} = 0, \frac{\partial Q}{\partial b} = 0 \Rightarrow$

maximum likelihood becslés:

$a = \frac{\sum(x_i y_i) - \bar{x} \bar{y}}{\sum(x_i x_i) - \bar{x} \bar{x}}$

$b = \bar{y} - a \bar{x}$



eredmény
 \hat{y}_i számolt, y_i mért

$\bar{\hat{y}}_i = \bar{y}_i$

$a \pm b \cdot t$ pl. 95% intervallummal
 ↑ hivatalosan = 0, kell tengelymetret!

residuális: $r_i = y_i - \hat{y}_i$ $S_r^2 = \sum r_i^2 / (N-2) \Rightarrow S_a^2, S_b^2$ ebből

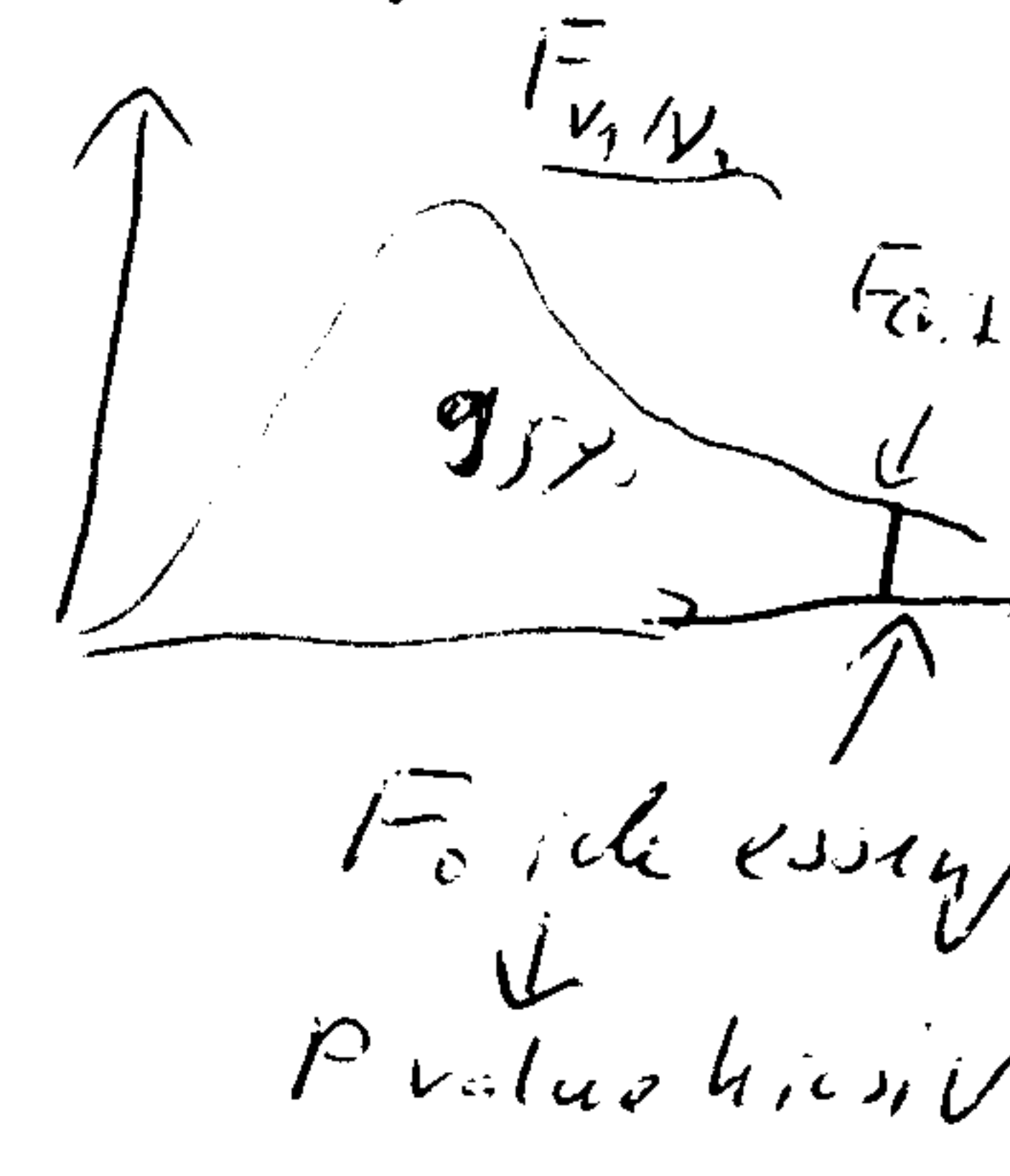
↑
 ábrátolása, főleg standardizálva segít → közösi értékek

regresszió variancia analízise:

	V	összetör
modell (regr.)	p	$MSS = \sum (\hat{y}_i - \bar{y})^2$
maradék (resid.)	$N-p-1$	$RSS = \sum (\hat{y}_i - y_i)^2$
teljes (total)	$N-1$	$TSS = \sum (y_i - \bar{y})^2$
		$(TSS = MSS + RSS)$

p = meredekség száma, tengelymetr.
 nincs benne
 variancia
 $MSS/p = S_M^2$
 $RSS/(N-p-1) = S_e^2$
 $F_0 = S_M^2 / S_e^2$ P-value

F eloszlás két
 variancia alapján
 a számosságok



R^2 determinációs
~~variancia~~ együttható

$R^2 = 1 - \frac{RSS}{TSS}$

$R^2 \in [0, 1]$ (itt)

\hat{y}, \hat{y} vektorok egy N -dim térben
 betart mindig \cos^2 -ja

ha különböző N -ű modellek → R^2 adj, lehet negatív
 ↓
 de et csak átskálázás $[-(N, n), 1]$

ha nem kényszer nélküli regresszió, pl $b \neq 0$ egy más modell.

$R^2 \in (-\infty, 1], \bar{y}_0 \neq \bar{y}, TSS \neq MSS + RSS$

RMSE - root mean square error

$$RMSE = \sqrt{\frac{\sum (\hat{y}_i - y_i)^2}{N}}$$

szemléletes, átlagos hiba (átlag "négyzetes hiba")
nem statisztikai, de jó

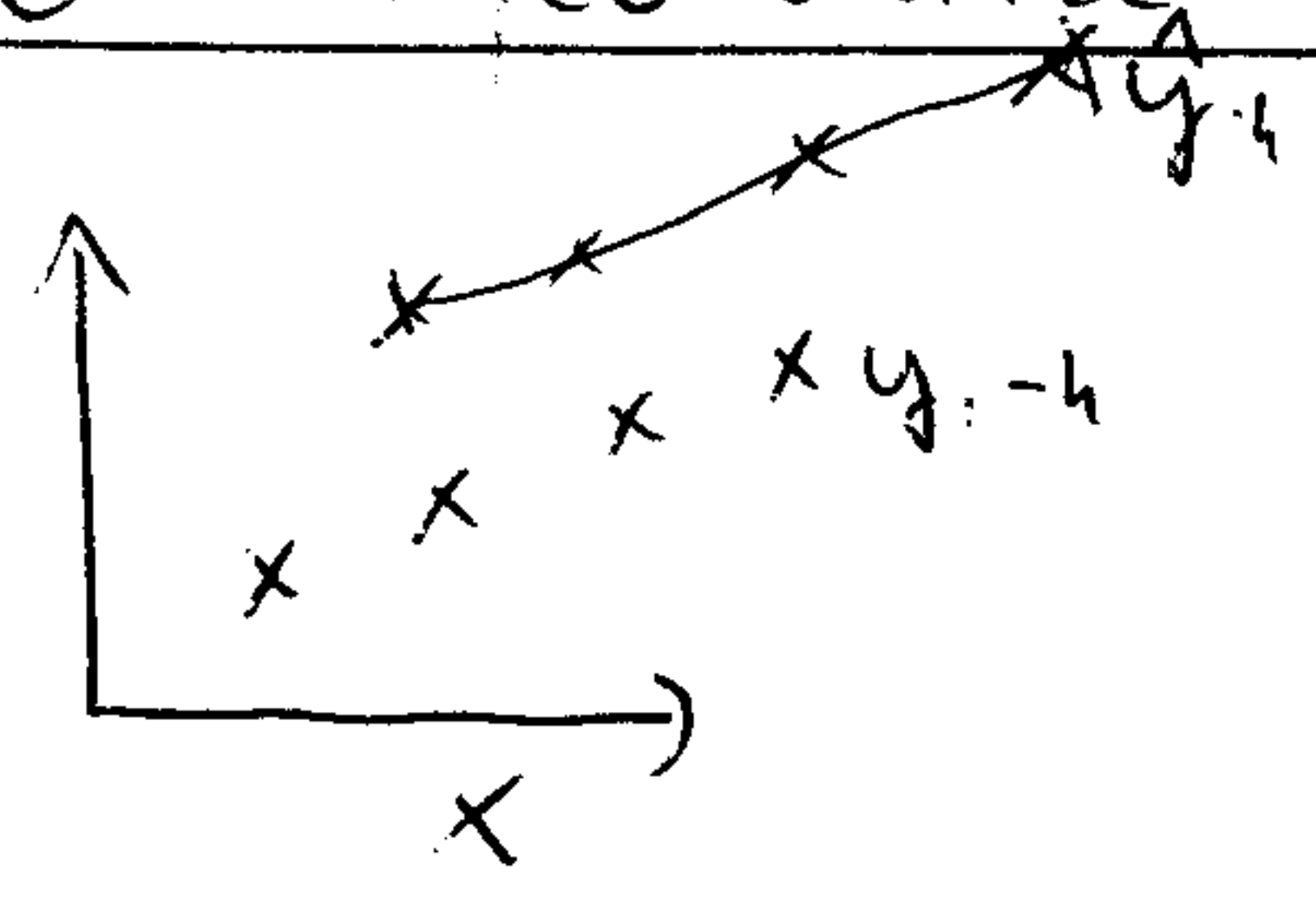
$Q^2, (R^2_{cv})$
↑
cross-validation

$\hat{y}_{j_{i=1}}$ ← olyan modellben számolt érték, ahol ez pont
nem volt benne

test-train (test/training) felosztás

pl. 80% train → ~~ered~~ modell elkezdésére 20%-on
≈ R^2_{test} helyett, $RMSE_{test}$

CCC Concordance Correlation Coefficient (újabb)



csak az irány azonos ⇒ $R^2 = 1$
vektor \neq irányított síkba

↑
ezek egyszerre kellene.

CCC $\in [-1, 1]$ mint korrelációs
eh csak irányított
síkhasonló, nem
vektorok

Feladat: illeszkedés jósa: $F, R^2, RMSE, CCC$

robustasság: Q^2 (belső validáció)

predikció test-train felosztás
(külső validáció)

Pelda:

y = (2,1; 2,2; 2,2; 2,3; 2,1; 2,1) EXCEL

átlag	2,11428	(átlag)
medián	2,1	(medián)
módusz	2,1	(módusz)
becs. szórás	0,10694	(szórás)

$2,11428 \pm 0,09887 \Rightarrow \underline{2,11 \pm 0,10} \checkmark$
 $\Rightarrow 2,114 \pm 0,099$

egymintás t-teszt próba ~~2,5-re~~ $E_0 = 2,25$ -re

$t_0 = \frac{2,5 - \bar{y}}{s/\sqrt{n}} \rightarrow T.ELOSIL(t_0, \nu = n-1, 1)$

↓
 3,358 0,9923

↓
 $0,025 \leq \leq 0,975$ ide esik $H_0: E = E_0$

↑ elosztás pr: