

# Adatelemzés és adatbányászat MSc

1. téma

## IR és adatalemzési alapok

# Információ szintjei

Kell-e targonca?

Mennyi targonca kell?

Mennyi kezelők kellenek?

Mit szállítsanak ma?

Mikorra vigyék át?

Hova kell menni most?

Mit kell vinni?

Merre kell menni?

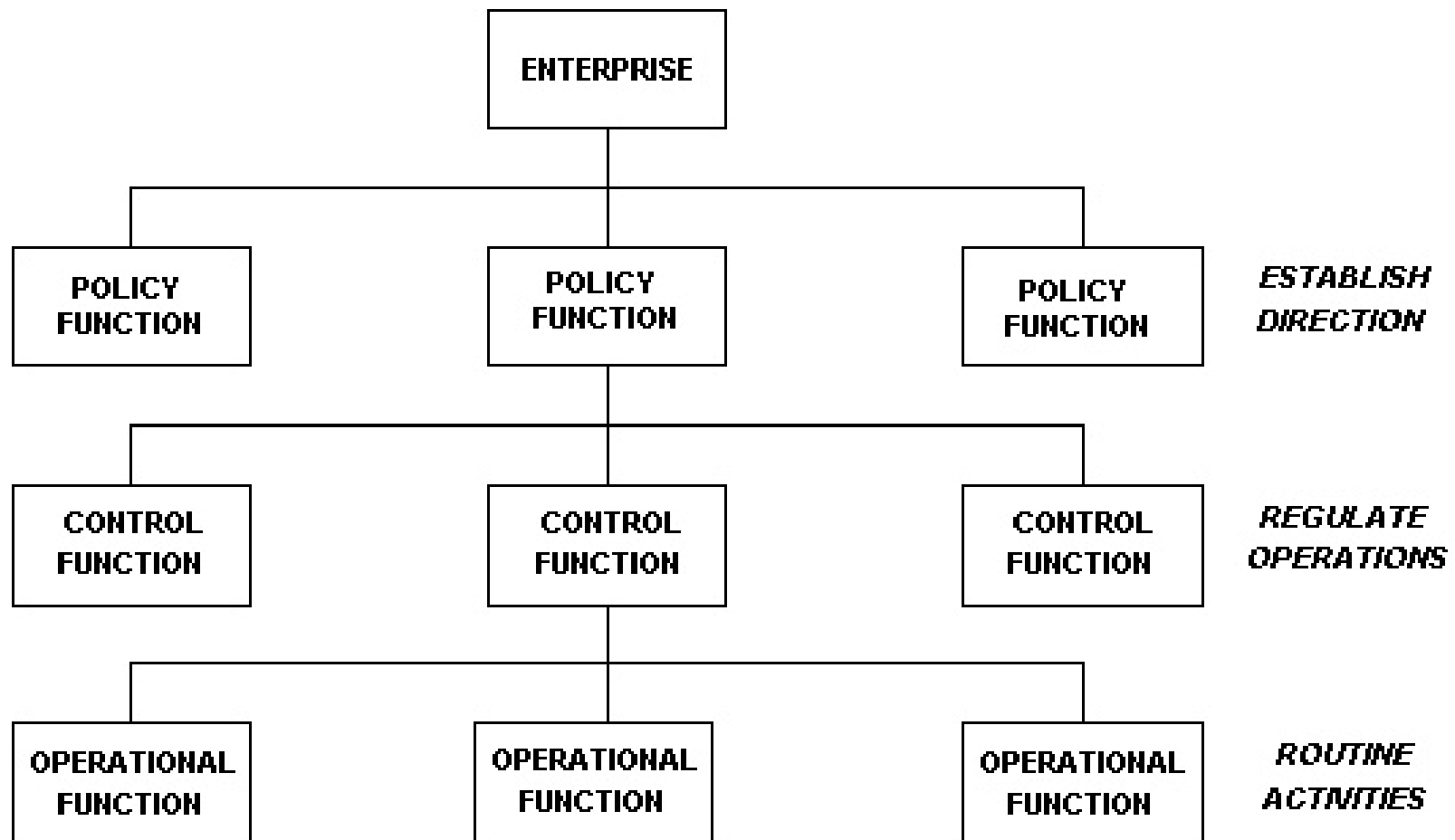
Hogyan kell beindítani?

Hogyan kell fékezni?

?

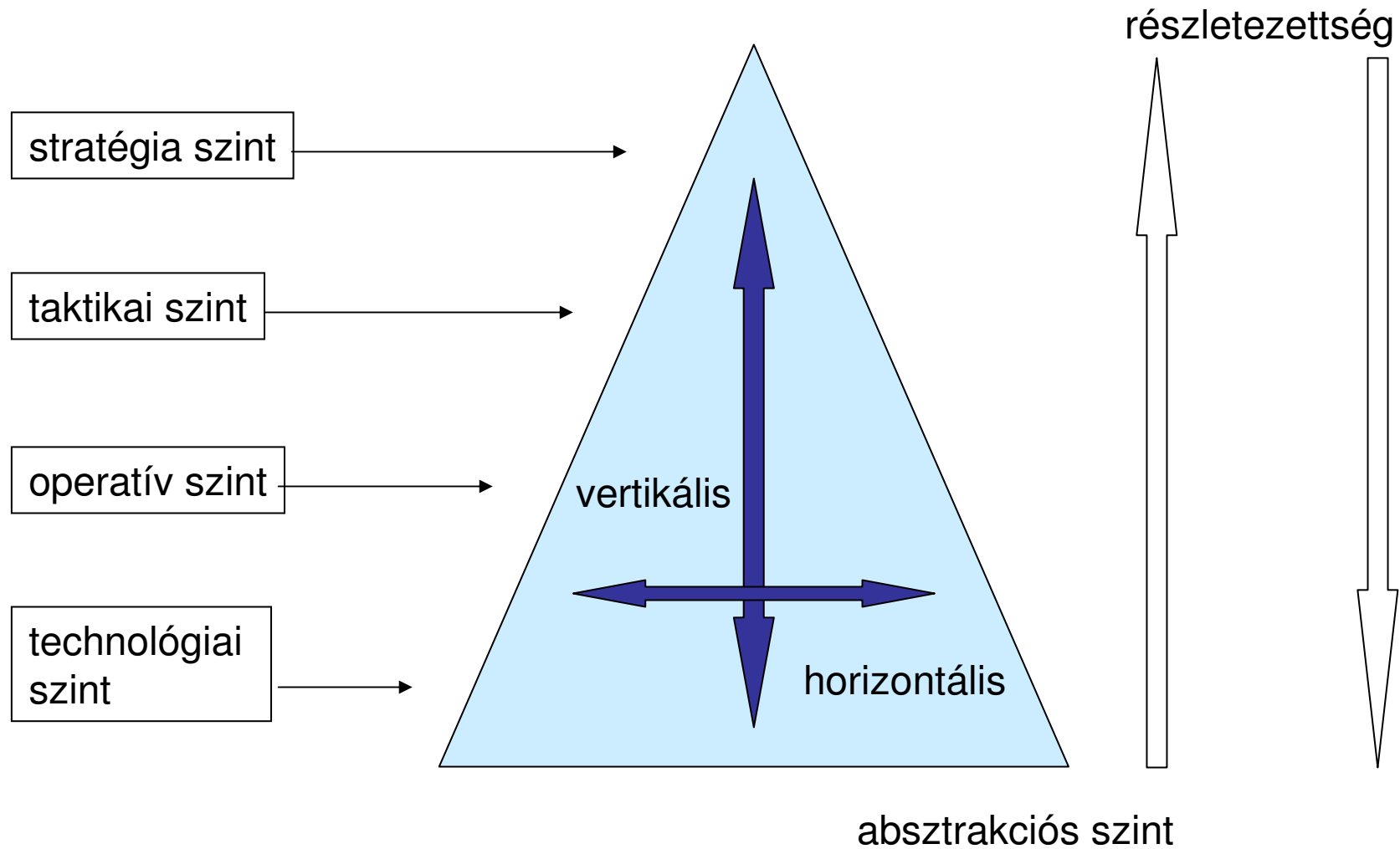


# Vállalati irányítási szint modulok



# VIR információ szintjei

Információ áramlások funkció nézetben



# Vállalati információs rendszer (VIR)

VIR: Olyan információs (informatikai) rendszer, melynek célja, a vállalati információkezelés hatékonyságjavítása, az információk integrált kezelése.

A VIR előnyei:

információ gyors elérése

információ pontossága

piaci versenyhelyzet javítása

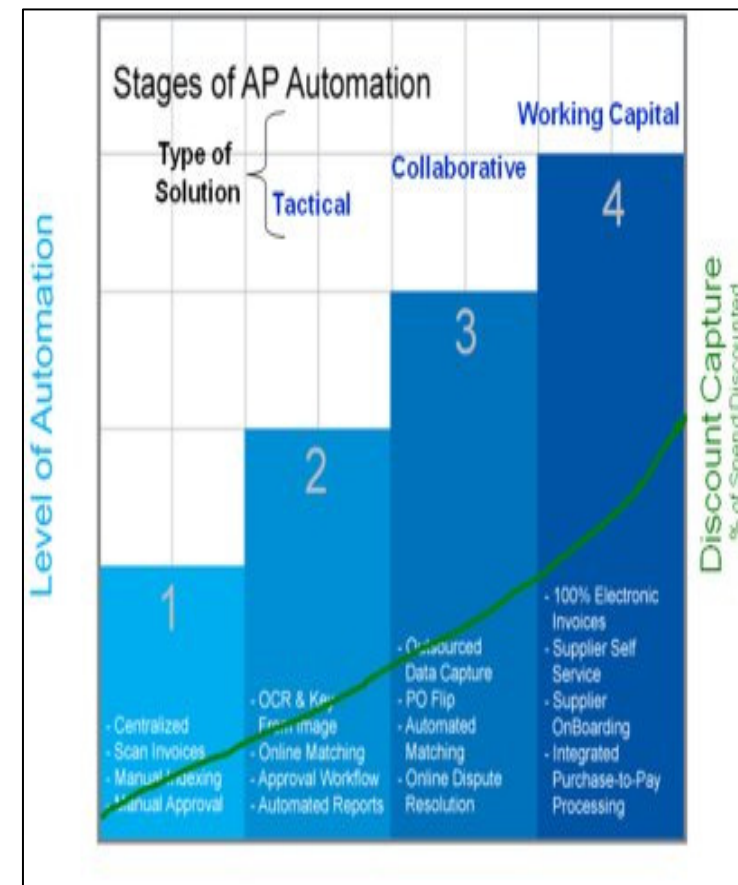
termelékenység növelése

ügyvitel javítása

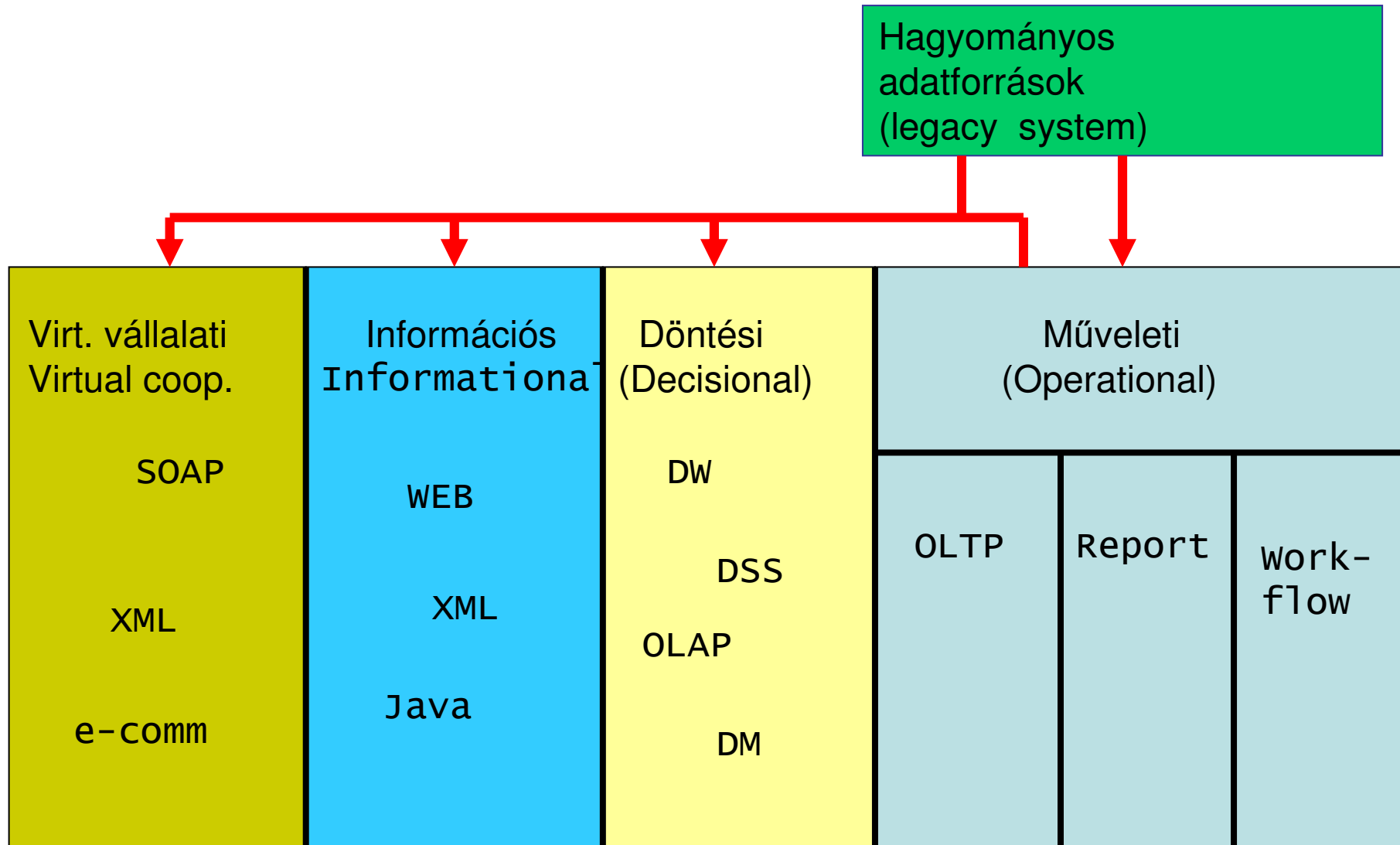
technológia fejlesztése

hatékonyabb kapcsolat a partnerekkel

hatékonyabb döntéshozatal , tervezés



# Információs rendszerek típusai



# OLTP jellemzői

adatmódosítás

aktuális állapot

nagy konkurencia

konzisztencia

rövid tranzakciók

homogenitás

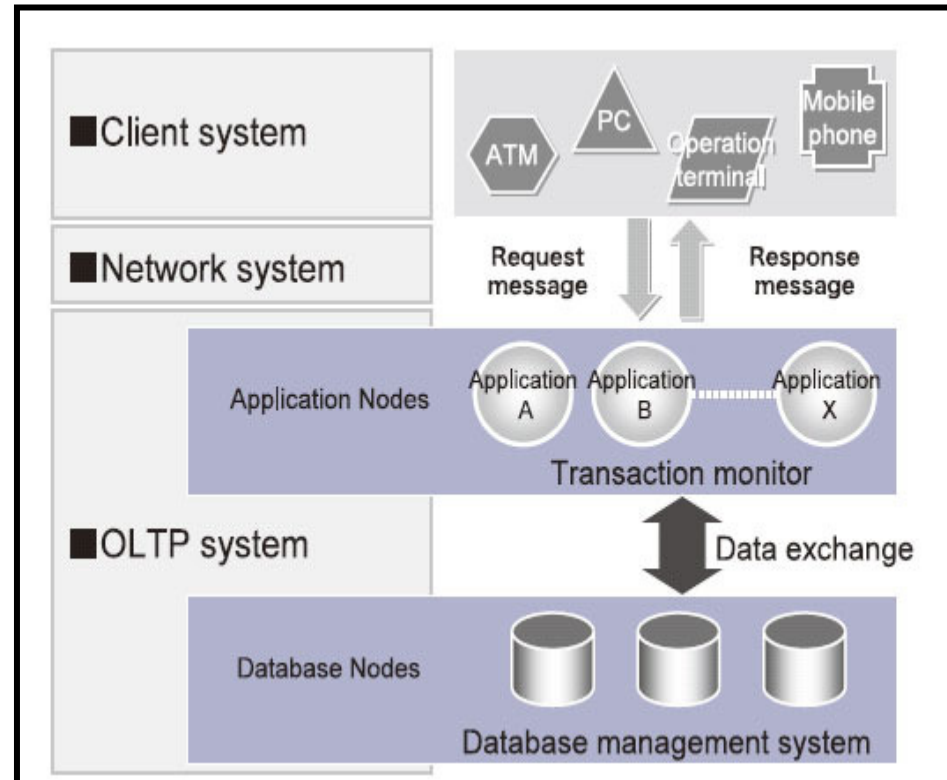
normalizált

relációs és OO

SQL felület

ACID elvek

adatvesztés elleni védelem



# OLAP jellemzői

adatelekérdezés

korábbi állapotok

kis konkurencia

betöltés

hosszú tranzakciók

heterogenitás

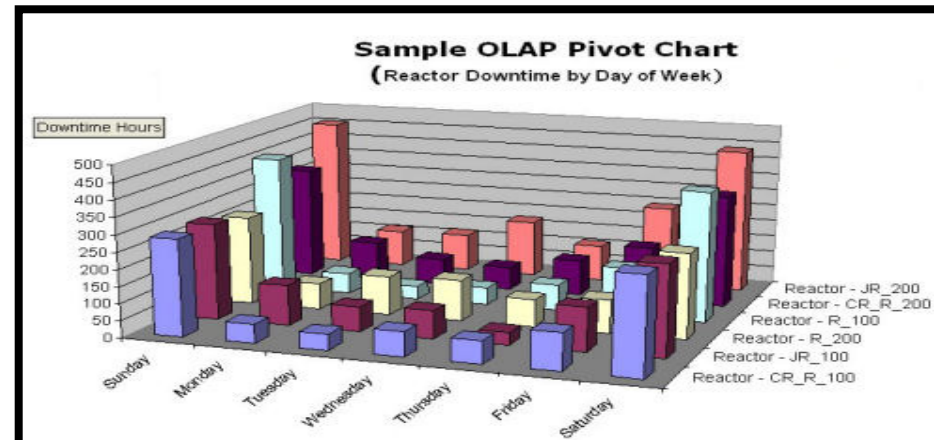
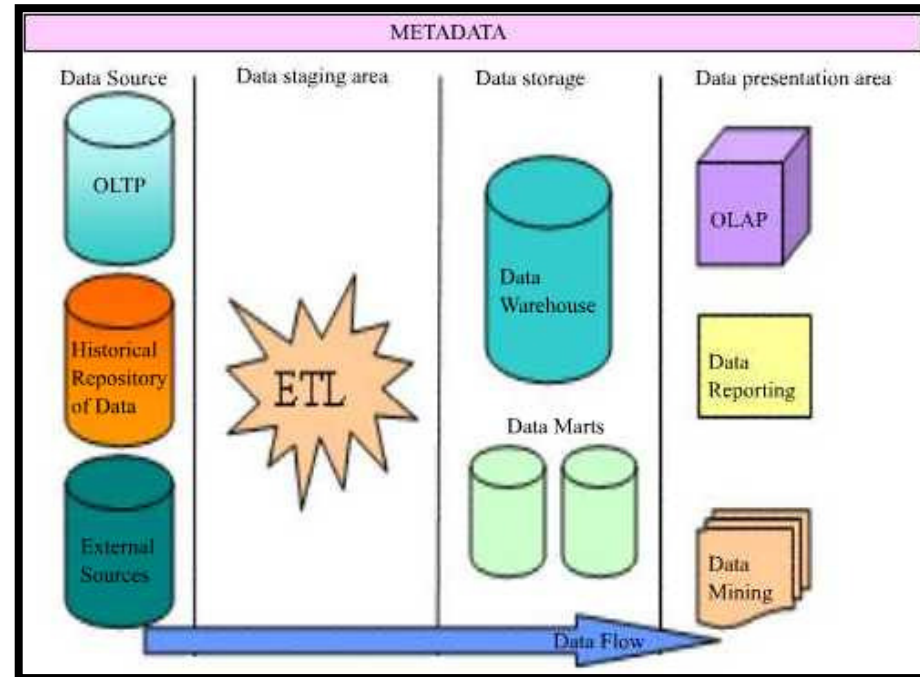
adatkocka

modulokból áll

nincs szabvány

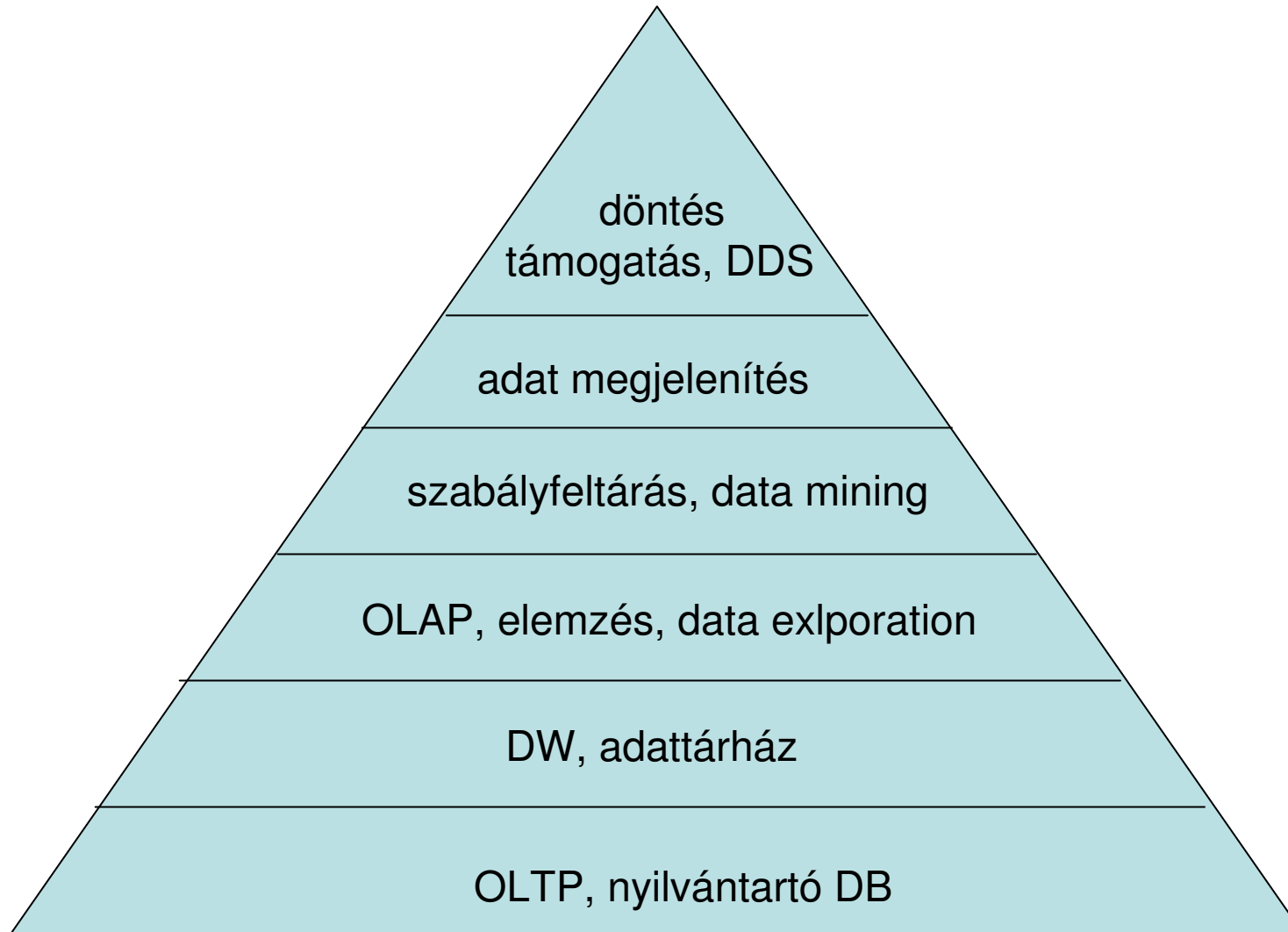
nem normalizált

adatvesztés elleni védelem



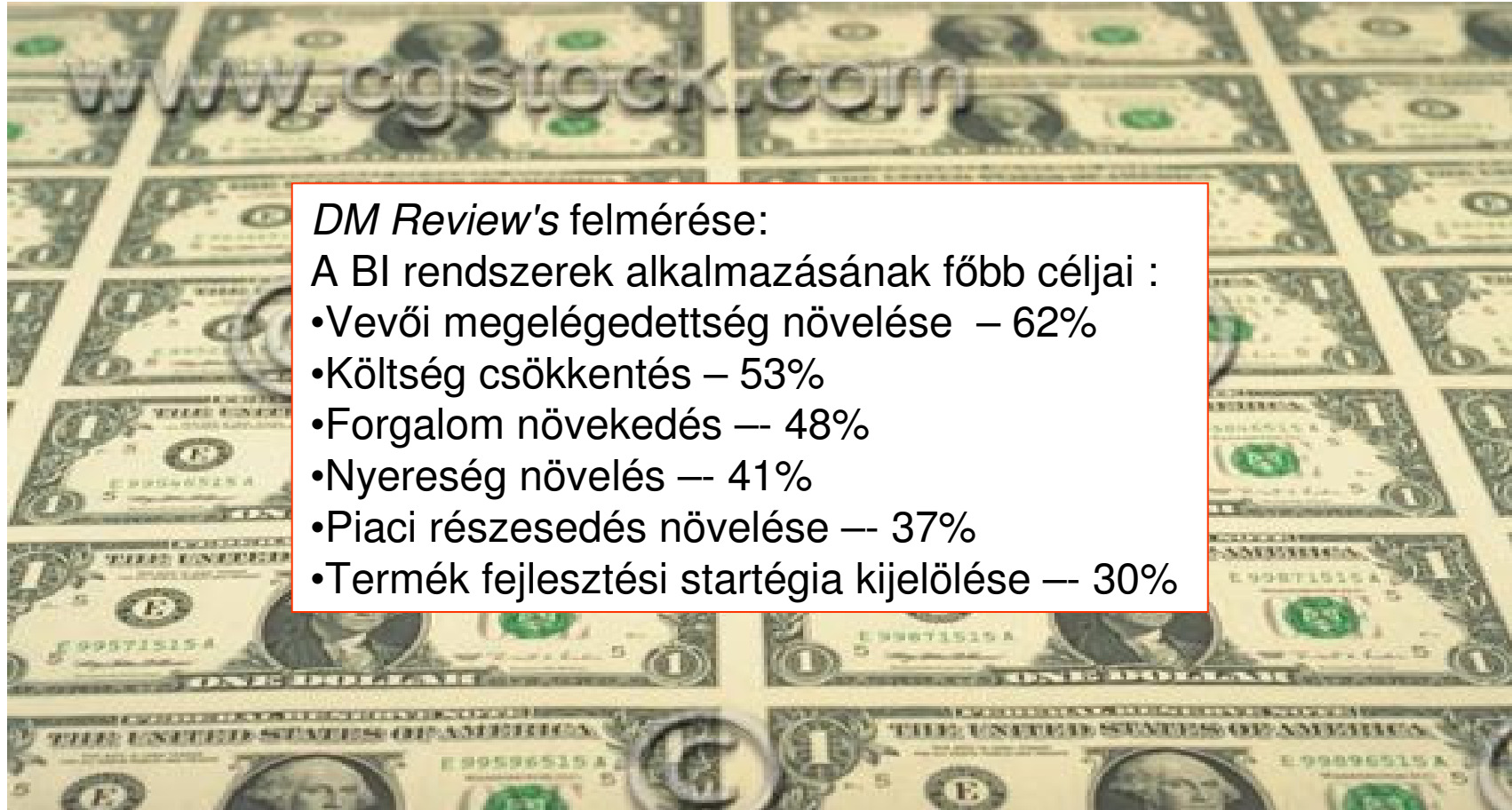


# DSS, OLAP rendszerek szerepe



Információ feldolgozási szintek

# OLAP rendszerek célja



*DM Review's* felmérése:

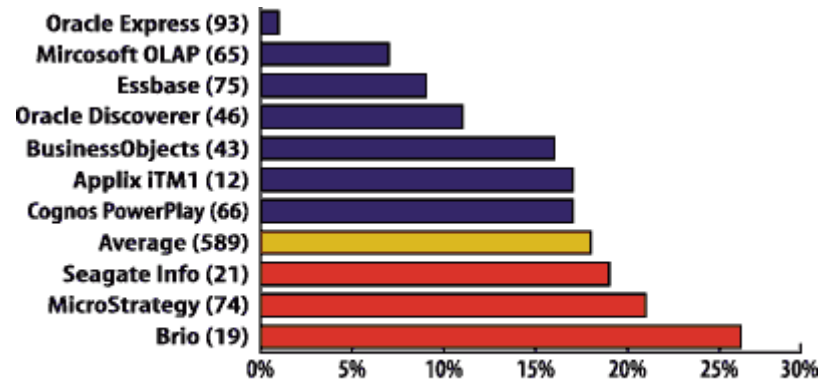
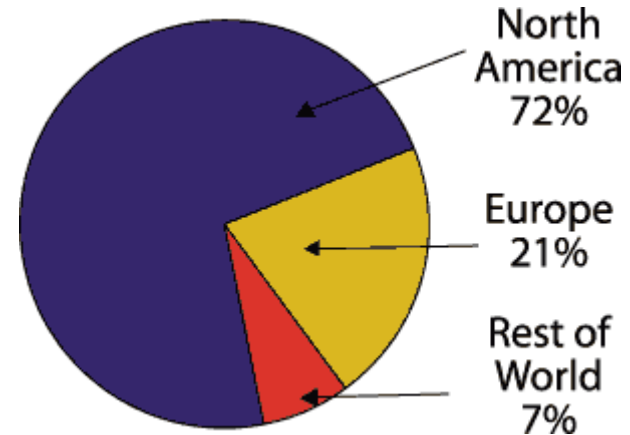
A BI rendszerek alkalmazásának főbb céljai :

- Vevői megelégedettség növelése – 62%
- Költség csökkentés – 53%
- Forgalom növekedés – 48%
- Nyereség növelés – 41%
- Piaci részesedés növelése – 37%
- Termék fejlesztési startégia kijelölése – 30%

# OLAP/DW rendszerek piaca

- évi 3 milliárd dolláros piac
- több szereplős
- főbb területek:
  - ERP
  - CRM
  - Sales analysis
  - planning
- 85% sikeres projekt
  - MicroStrategy, MS, Oracle
- megbízhatóság
- kb 2/3 részben kihasználtak
- nem dominál a Web-OLAP
- piaci részesedés:
  - MS
  - Oracle
  - Essbase
  - MicroStrategy

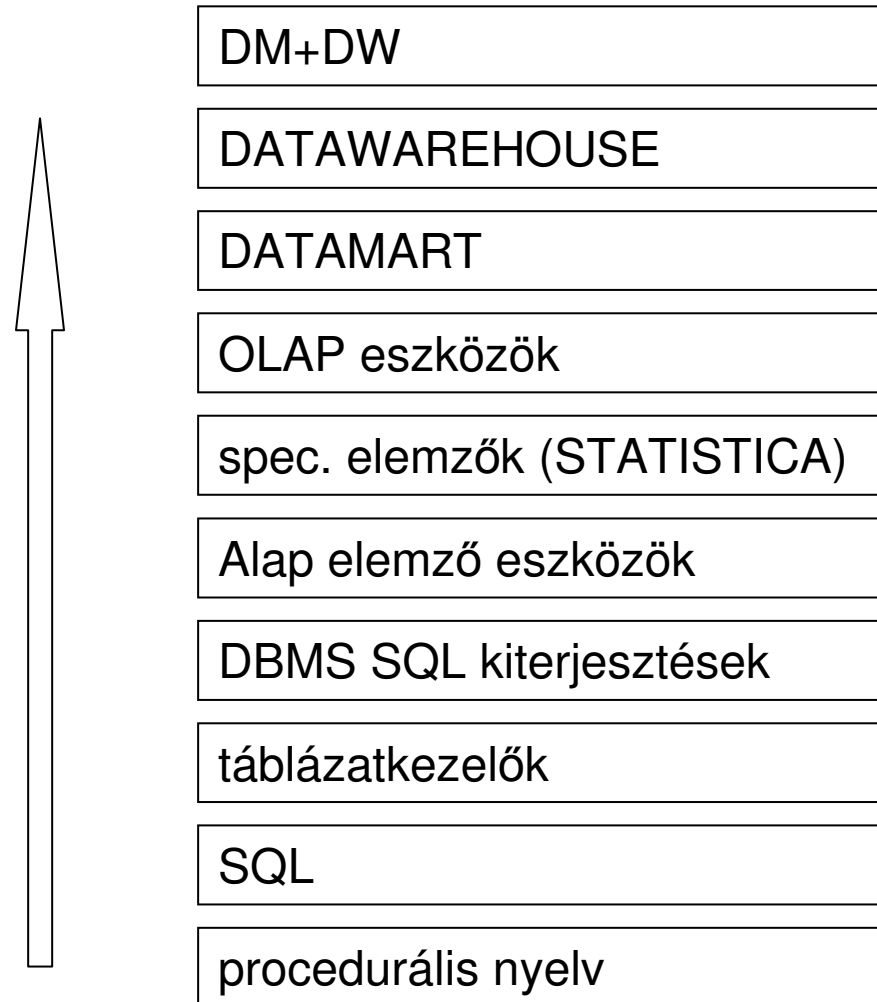
OLAP  
elterjedtsége



Shorter bars mean fewer query performance problems reported

teljesítmény viszonyok

# Adatelemzési eszközök



összetettség

# Adatelemzés SQL-ben

szórás

```
SELECT  
SQRT(AVG(POWER(x.ar – (SELECT AVG(y.ar) FROM auto_eladas),2)))  
FROM auto_eladas;
```

Gyakorisági tábla  
(évenként a Fiat autók relatív aránya)

```
SELECT  
ev, SUM(CASE WHEN tipus='FIAT' THEN 1 ELSE 0 END) / SUM(1)  
FROM auto_eladas  
GROUP BY ev;
```

# Adatelemzés SQL-ben

korreláció

(a Fiat és Opel autók havi darabszámai között)

```
SELECT CORR (dbF,dbO) FROM
( SELECT
ev, SUM(CASE WHEN tipus='FIAT' THEN 1 ELSE 0 END) / SUM(1) dbF,
SUM(CASE WHEN tipus='OPEL' THEN 1 ELSE 0 END) / SUM(1) dbO
FROM auto_eladas
GROUP BY ev )
WHERE dbF * dbO IS NOT NULL;
```

differentia képzés

Havi összforgalom változás

```
SELECT T1.ho, T1.db – T2.db dbdiff FROM
( SELECT ho, SUM(ertek) db FROM auto_eladas GROUP BY ho ) T1
INNER JOIN
( SELECT ho, SUM(ertek) db FROM auto_eladas GROUP BY ho ) T2
ON (T1.ho = T2.ho + 1 and T1.ev = T2.ev) OR
(T1.ho =1 and T2.ho =1 and T1.ev = T2.ev + 1)
```

# Adatelemzés SQL-ben

keresztreferencia tábla  
(autótípusok éves darabszámai)

```
SELECT
  SUM (CASE WHEN ev = 2006 THEN db ELSE 0 END) e06,
  SUM (CASE WHEN ev = 2007 THEN db ELSE 0 END) e07,
  SUM (CASE WHEN ev = 2008 THEN db ELSE 0 END) e08,
  SUM (CASE WHEN ev = 2009 THEN db ELSE 0 END) e09,
  ---
FROM auto_eladas
GROUP BY típus
```

differentia maximuma  
Legnagyobb havi összforgalom változás

```
SELECT MAX(dbdiff) FROM
(SELECT T1.ho, T1.db - T2.db dbdiff FROM
( SELECT ho, SUM(ertek) db FROM auto_eladas GROUP BY ho ) T1
INNER JOIN
( SELECT ho, SUM(ertek) db FROM auto_eladas GROUP BY ho ) T2
ON (T1.ho = T2.ho + 1 and T1.ev = T2.ev) OR
(T1.ho =1 and T2.ho =1 and T1.ev = T2.ev + 1))
```

# Adatelemzési módszerek

## Regresszió

Feladat: adott mérési pontra legjobban illeszkedő görbe megkeresése

Adottak: mérési pontok, függvényosztály (paraméteresen)

Feladat: a mérési pontokra legjobban illeszkedő paraméterek meghatározása

Optimalizálási feladat:

Célfüggvény: illeszkedési hiba: eltérések négyzetösszege

Optimalizálási módszerek:

Derivált zérushelye

Gradiens módszer

Sztohasztikus keresés

$$\{(x_i, y_i)\}$$

$$\{f(\bar{p}_i, x)\}$$

$$E(\bar{p}_i) = \sum_i (f(\bar{p}_i, x_i) - y_i)^2$$



# Adatelemzési módszerek

## Lineáris regresszió

többváltozós lineáris regresszió: a mérési pontokat legjobban közelítő függvény meghatározása

egy függő változó feltételes várható érték becslésére szolgál

$$E(y|x_1, x_2, \dots) = F(x_1, x_2, \dots, \alpha_1, \alpha_2, \dots)$$
$$y = F(x_1, x_2, \dots, \alpha_1, \alpha_2, \dots) + \varepsilon$$

lineáris regresszió : a paraméterekben lineáris az F függvény

$$y = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 + \varepsilon$$

$$y = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 + \alpha_4 x_1^2 + \alpha_5 x_2^2 + \varepsilon$$

a feltétel szerint  $\varepsilon$  egy 0 várható értékű, azonos paraméterű normál eloszlású

# Adatelemzési módszerek

## Lineáris regresszió

a paraméterek várható értékének meghatározása a legkisebb négyzetek elvével történik  
elemi esetre:

$$y_i = \alpha_1 x_{i1} + \alpha_2 x_{i2} + \alpha_3 + \varepsilon_i$$

$$\varepsilon_i = y_i - (a_1 x_{i1} + a_2 x_{i2} + a_3)$$

$$E(\varepsilon_i) = 0$$

$$\sum \varepsilon_i^2 \rightarrow \text{minimális}$$

a szélsőérték szükséges feltétele a deriváltak zérus értéke

$$\partial_{a_1} \sum_i (y_i - (a_1 x_{i1} + a_2 x_{i2} + a_3))^2 = 0$$

$$\partial_{a_2} \sum_i (y_i - (a_1 x_{i1} + a_2 x_{i2} + a_3))^2 = 0$$

$$\partial_{a_3} \sum_i (y_i - (a_1 x_{i1} + a_2 x_{i2} + a_3))^2 = 0$$

# Adatelemzési módszerek

## Lineáris regresszió

egy változós esetre:

$$\partial_{a_1} \sum_i (y_i - (a_1 x_i + a_2))^2 = 0$$

$$\partial_{a_2} \sum_i (y_i - (a_1 x_i + a_2))^2 = 0$$

$$\partial_{a_1} \sum_i (y_i^2 + a_1^2 x_i^2 + a_2^2 + 2 a_1 a_2 x_i - 2 y_i a_1 x_i - 2 y_i a_2) = 0$$

$$\partial_{a_1} \sum_i (a_1^2 x_i^2 + 2 a_1 a_2 x_i - 2 y_i a_1 x_i + a_2^2 - 2 y_i a_2 + y_i^2) = 0$$

$$\partial_{a_2} \sum_i (a_2^2 + 2 a_1 a_2 x_i - 2 y_i a_2 + a_1^2 x_i^2 - 2 y_i a_1 x_i + y_i^2) = 0$$

$$a_1 \sum_i x_i^2 + a_2 \sum_i x_i - \sum_i y_i x_i = 0$$

$$a_2 n + a_1 \sum_i x_i - \sum_i y_i = 0$$

$$a_1 = (n \sum_i x_i y_i - \sum_i x_i \sum_i y_i) / (n \sum_i x_i^2 - \sum_i x_i \sum_i x_i)$$

$$a_2 = (\sum_i y_i - a_1 \sum_i x_i) / n$$

# Adatelemzési módszerek

## Lineáris regresszió

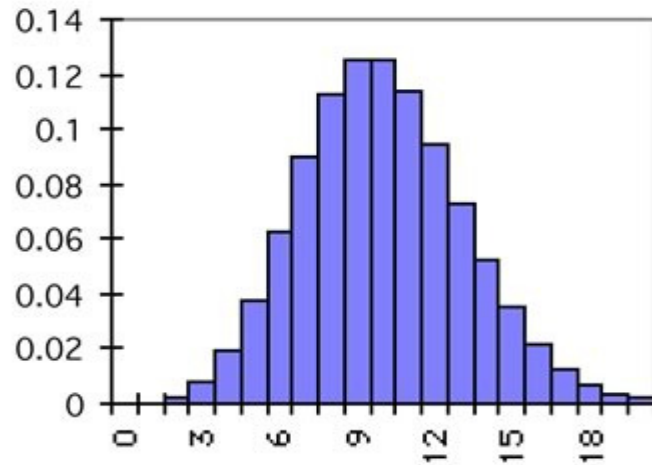
Minta  $\{(2.,4.), (4.,6.2), (6.,4.5)\}$  és  $f(a,b,x) = ax+b$

Excel megvalósítás:

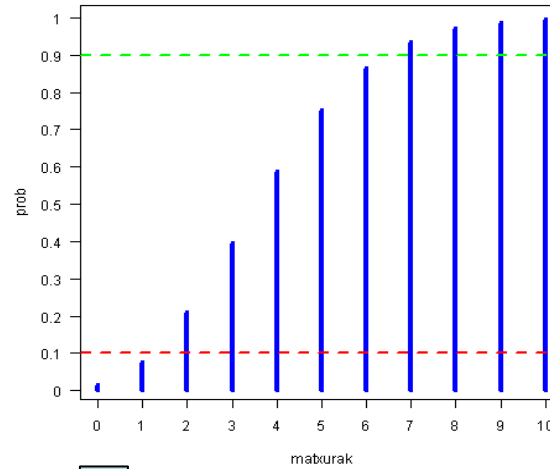
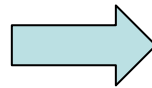
- elemi lépésekben
- regressziós görbével

# Adatelemzési módszerek

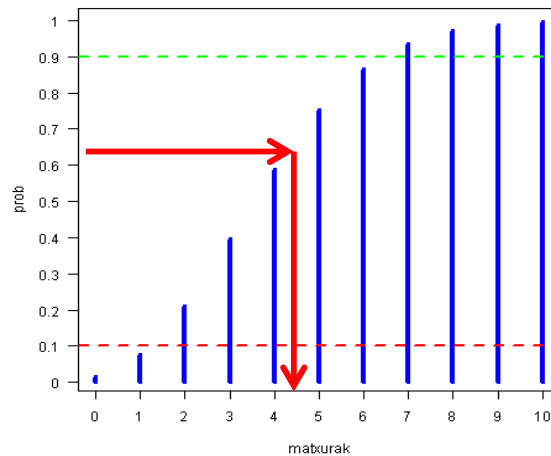
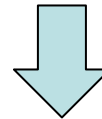
Minta generálása adott eloszlás függvényhez



sűrűség fv.



eloszlás fv.



érték

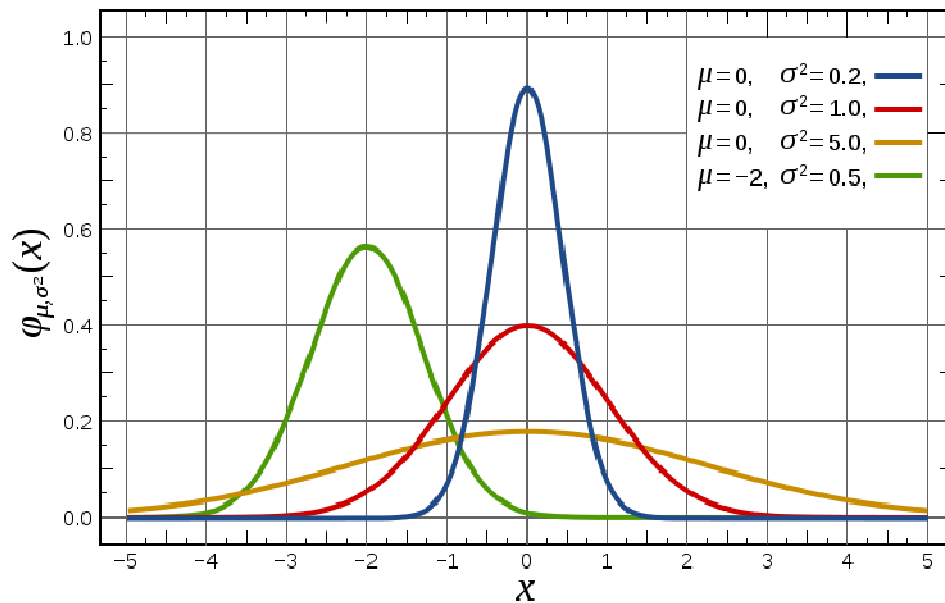
# Adatelemzési módszerek

## Normál eloszlás

Centrális határeloszlás tétele: Nagy n-re a mért empirikus átlagok normális eloszlást mutatnak

A  $x = (a - a') / \sigma$  változó  
N(0,1) normál eloszlású lesz

A N(0,1) eloszlás esetén az  
 $|x| > \approx 2.8$  pontok „lehetetlen”  
eseményeknek tekintetők



$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

várható érték

szórásnégyzet

# Adatelemzési módszerek

## Zajok szűrése

Zajok kiszűrésének módszere: az elosztást  $N(0,1)$ -re hozva a 2.5-nél nagyobb abszolút értékek zajoknak tekinthetők

Példa: 46, 48, 38, 45, 47, 58, 44, 45, 43, 44

46	0,2	458	0,04	25,73333	0,039426
48	2,2	45,8	4,84	5,072803	0,433685
38	-7,8		60,84		-1,53761
45	-0,8		0,64		-0,1577
47	1,2		1,44		0,236556
58	12,2		148,84		2,404982
44	-1,8		3,24		-0,35483
45	-0,8		0,64		-0,1577
43	-2,8		7,84		-0,55196
44	-1,8		3,24		-0,35483
<b>a</b>	<b>a-a'</b>	<b>a'</b>	<b>(a-a')<sup>2</sup></b>	<b><math>\sigma</math></b>	<b>x</b>



# Adatelemzési módszerek

## Véletlenszerűség ellenőrzése

Wald-Wolfowitz teszt: figyeli a sorozatok (runs) eloszlását (ne legyen se túl kevés, se túl sok sorozat)

Induló adatsor: mérési értékek

Lépések:

- $\bar{a}$  átlag kiszámítása
- $\text{sig}(a - \bar{a})$  -val helyettesítjük  $a$ -kat
- $n^+$ ,  $n^-$  (elemek db),  $R$  (sorozatok száma) meghatározása
- $a'' = 1 + 2n^+n^- / (n^+ + n^-)$
- $\sigma^2 = (a'' - 1)(a'' - 2) / (n^+ + n^- - 1)$ .
- $z = (R - a'') / \sigma$
- ha  $|z| > z_0$  akkor nem véletlen a sorozat ( $\sim 2.5$ )



# Adatelemzési módszerek

## Véletlenszerűség ellenőrzése

Példa: 3, 5, 12, 7, 9, 8, 21, 17, 87, 22, 18, 24

Excel megvalósítás:

- elemi lépésekben
- regressziós görbével

# Adatelemzések statisztikai háttere

A statisztikusok is ideális világból indulnak ki.

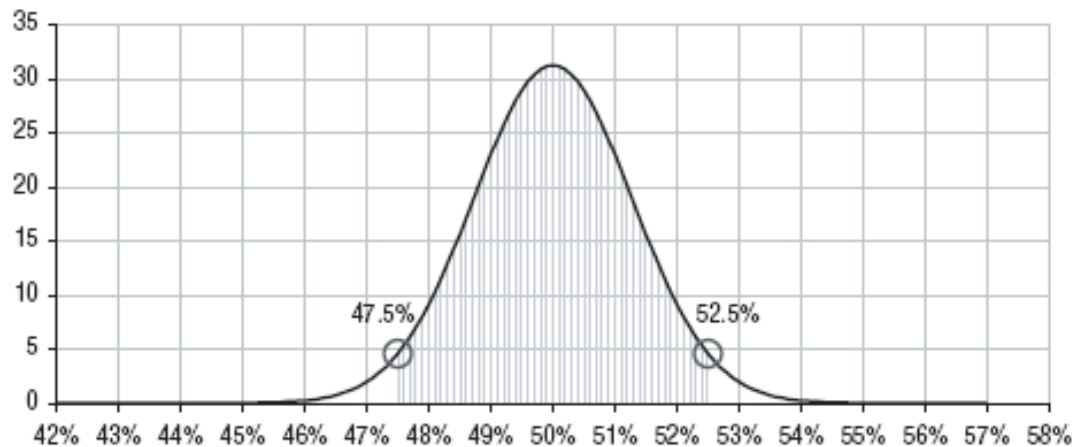
Null-hipotézis elve: a mérési eltérés csak a véletlen műve

A vizsgálat tárgya:

- milyen paraméterű az ideális eloszlás?
- mennyire teljesül a null-hipotézis?

A mérési adatokon próbákat lehet végrehajtani a hipotézis ellenőrzésére, a hipotézis konfidencia szintjének megállapítására

Az elemzés megadja, hogy milyen konfidencia értékkel fog a paraméter egy megadott konfidencia intervallumba esni.



# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-eloszlás

Feltétel: legyenek  $X_i$  független normál eloszlású változók,  $(\mu, \sigma)$  paraméterekkel.

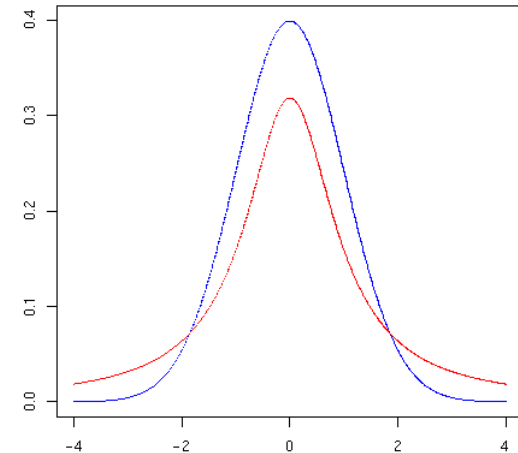
Ekkor a minta átlaga ( $n$ : mintaszám): 
$$\bar{X} = \frac{\sum X_i}{n}$$

minta szórásnégyzete: 
$$S^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$$

Ekkor normál(0,1) eloszlású: 
$$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

Illetve Chi-négyzet eloszlású: 
$$\frac{(n-1)S^2}{\sigma^2}$$

Emiatt T Student eloszlású lesz: 
$$\frac{\bar{X} - \mu}{S / \sqrt{n}}$$



# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-teszt

Egyváltozós eset.

- null hipotézis: az eloszlás várható értéke:  $\mu$

- feladat: a tapasztalati eloszlás illeszkedik-e?

- vizsgált eloszlás:  $t = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}}$

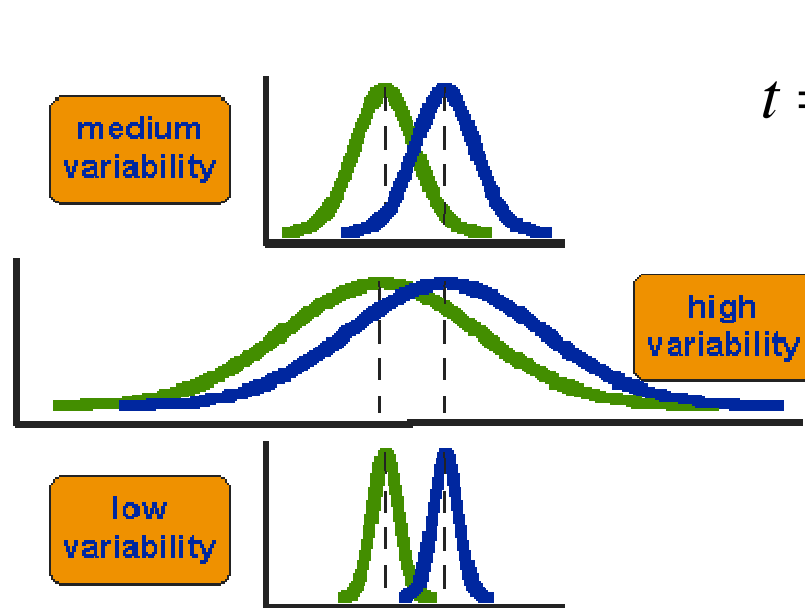
- függetlenségi tényező:  $n - 1$

# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-teszt

Kétváltozós eset.

A feladat adott kontroll és mérési eloszlás mellett eldönteni, hogy a mérési eloszlás mennyire illeszkedik a kontrollra



$$t = \frac{\overline{X}_1 - \overline{X}_2}{\sqrt{\frac{(n_1 - 1)S^2_1 + (n_2 - 1)S^2_2}{n_1 + n_2 - 2} \cdot \frac{n_1 + n_2}{n_1 n_2}}}$$

Függetlenségi tényező:  $n_1 + n_2 - 2$

# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-teszt

T-tábla használata:

- az oszlop jelöli a konfidenciát
- a sor jelöli a függetlenségi értéket
- a cella adja meg az előírt maximum t értékek  
(ha a tábla érték nagyobb mint a számított, akkor megtartjuk a hipotézist)

FD	0.1	0.05	0.01
5	2.02	2.57	4.03
6	1.94	2.45	3.71
7	1.89	2.37	3.50
9	1.83	2.26	2.68
20	1.72	2.09	2.85
30	1.70	2.04	2.75



# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-teszt

Adott az alábbi adatsor: 483, 502, 498, 496, 502, 483, 494, 491, 505, 486.

Kérdés: tekintető-e 5%-os kockázat mellett a eloszlás  $m=500$  várható értékűnek?

$$t = \frac{\overline{X} - \mu}{\frac{S}{\sqrt{n}}}$$

$X=494$ ,  $S = 8.05$ ,  $\mu=500$ ,  $n=10$ ,  $df=9$

$t=2.36$ , tablazat:2.26

megoldás: nem fogadható el a hipotézis

# Adatelemzések statisztikai háttere

## Hipotézis vizsgálat, T-teszt

Adott az alábbi adatsor, két eltérő helyen élő egyedhalmaz súlyértékei:

X1: 52; 57; 62; 55; 64; 57; 56; 55

X2: 41; 34; 33; 36; 40; 25; 31; 37; 34; 30; 38.

Kérdés: tekintető-e azonosnak a két eloszlás 5%-os kockázat mellett?

$$t = \frac{\overline{X}_1 - \overline{X}_2}{\sqrt{\frac{(n_1 - 1)S^2_1 + (n_2 - 1)S^2_2}{n_1 + n_2 - 2} \cdot \frac{n_1 + n_2}{n_1 n_2}}}$$

X1=57.25, X2=34.45

n1=8, n2=11

S1<sup>2</sup>=15.36, S2<sup>2</sup>=21.87

t=11.12, táblázat=2.11

megoldás: nem fogadható el a hipotézis