

9. GRUPAREA DATELOR

- Care sunt funcțiile de grup și cum sunt ele folosite
- Cum se pot grupa datele într-o tabelă
- Care sunt regulile de folosire a clauzei GROUP BY
- Cum se pot filtra grupurile folosind clauza HAVING
- Care este diferența dintre clauzele WHERE și HAVING
- În ce ordine se execută clauzele WHERE, GROUP BY și HAVING

9.1. FUNCȚII DE GRUP

- Într-un capitol anterior am discutat despre funcțiile singulare, adică despre funcțiile care operează la un moment dat asupra unei singure înregistrări.
- Este acum momentul să discutăm despre funcțiile de grup, care returnează o singură valoare pentru un **grup sau set** de linii dintr-un tabel. Puteți calcula cea mai mare valoare dintr-un set de valori, puteți determina numărul de înregistrări ce respectă o anumită condiție etc.
- Pentru exemplificarea acestor funcții vom folosi tabela **VOTURI** și tabela **JUDEȚE** care conțin următoarele date

Tabela **JUDEȚE**

Cod_județ	Judet	Număr_alegători
B	București	1750192
IS	Iași	650029
SB	Sibiu	363380

Tabela **VOTURI**

Județ	Candidat	Număr voturi
B	1	347016
B	2	1552
B	3	1374
IS	1	196508
IS	2	1038
IS	3	1267
SB	1	65084
SB	2	561
SB	3	533
B	4	96744
B	5	25656
B	6	13361
IS	4	35784
IS	5	5558
IS	6	4094
SB	4	19937
SB	5	4323
SB	6	2366
B	7	25937
B	8	4619
B	9	4323
IS	7	3682
IS	8	1291
IS	9	327
SB	7	4225
SB	8	765
SB	9	3797
B	10	2037
B	11	22687
B	12	514366
IS	10	1312
IS	11	3781
IS	12	12184
SB	10	660
SB	11	3768
SB	12	105993
SB	13	100
B	13	(null)
IS	13	(null)

- Principalele funcții de grup.
- **COUNT(x)** – determină numărul de valori ale lui **x**. Funcția, ca de altfel toate funcțiile de grup ignoră câmpurile completate cu **NULL**, adică va număra doar valorile nenule ale lui **x**.
- De exemplu, comanda
**SELECT COUNT(JUDET), COUNT(numar_voturi)
FROM voturi**
- va afișa numărul total de înregistrări din tabelă, **39** (câmpul **JUDET** nu are nici o valoare **NULL**) precum și numărul de linii pentru care câmpul **numar_voturi** este nenul, adică **37**, ultimele două linii din tabel având valoare **null** în câmpul **numar_voturi**.

COUNT(JUDET)	COUNT(NUMAR_VOTURI)
39	37



- Funcția **COUNT** poate fi folosită în combinație cu clauza **DISTINCT**, pentru a număra doar valorile distincte dintr-un domeniu. De exemplu dacă dorim să știm pentru câte județe avem rezultatele votării în tabela noastră, vom folosi comanda:

```
SELECT count(distinct judet)
```

```
FROM voturi
```

- Se va obține valoarea **3**, întrucât avem doar **3** județe înregistrate (București, Iași, Sibiu).

COUNT(DISTINCTJUDET)
3



- Exemplu:

```
SELECT count(distinct candidat), count(candidat)  
FROM voturi
```

- Evident primul apel de funcție afișează valoarea **13** , deoarece există **13** candidați pentru care au fost exprimate voturi, iar a doua comandă afișează valoarea **39**, adică exact numărul de linii din tabel deoarece toate liniile au completat câmpul **candidat**.

COUNT(DISTINCTCANDIDAT)	COUNT(CANDIDAT)
13	39



- **MAX(x)** – determină valoarea maximă a valorilor expresiei **x**.
- Să vedem de exemplu cum putem afla care este cel mai număr de voturi exprimate pentru un candidat într-un județ.

```
SELECT MAX(numar_voturi)  
FROM voturi
```

MAX(NUMAR_VOTURI)
514366



- Se poate observa pe tabelul cu datele din tabela voturi că acest maxim a fost obținut în București de către candidatul având codul **12**.
- Totuși această informație nu este foarte relevantă pentru că și populația din București este mult mai mare decât în celelalte județe. Ar trebui să putem determina numărul de voturi primite de către un candidat raportat la numărul de alegători (persoane cu drept de vot). SQL ne permite să aplicăm funcțiile de grup nu doar pe câmpuri din baza de date ci și pe expresii, ca în exemplul următor:

```
SELECT max(100*numar_voturi/numar_alegatori)  
FROM voturi v, judete j  
WHERE v.judet=j.cod_judet
```

MAX(100*NUMAR_VOTURI/NUMAR_ALEGATORI)
30.2306512478674028389502622190702260976



- Prin această comandă am obținut cel mai mare procent de voturi obținut de către un candidat într-un județ. Acest procent a fost obținut raportat la totalul persoanelor cu drept de vot și a fost obținut de către candidatul cu codul 1 în județul Iași:

```
SELECT 100*numar_voturi/numar_alegatori,  
j.judet, v.candidat  
FROM voturi v, judete j  
WHERE v.judet=j.cod_judet
```

100*NUMAR_VOTURI/NUMAR_ALEGATORI	JUDET	CANDIDAT
19.8273103750902758097397314123250477662	Bucuresti	1
.088675985263331108815489957673215281523	Bucuresti	2
.078505672520500607933301032115333631967	Bucuresti	3
30.2306512478674028389502622190702260976	Iasi	1
...

- În acest moment nu știm încă să scriem o comandă pentru a afișa județul și candidatul pentru care s-a obținut valoarea maximă.
- **MIN(x)** – determină valoarea minimă a valorilor expresiei **x**.
- **SUM(x)** – determină suma valorilor expresiei **x**.



- Cum aflăm oare numărul total de voturi valabil exprimate în județul Sibiu? Foarte simplu:

```
SELECT sum(numar_voturi)
```

```
FROM voturi
```

```
WHERE judet='SB'
```

SUM(NUMAR_VOTURI)
212112



- **AVG(x)** – determină media valorilor expresiei **x**. De exemplu, putem afla procentul mediu obținut un candidat în toate județele:

```
SELECT avg(100*numar_voturi/numar_alegatori)  
FROM voturi v, judete j  
WHERE (candidat=12) and (v.judet=j.cod_judet)
```

- Comanda afișează media procentelor obținute în fiecare județ de către candidatul cu codul **12**:

AVG(100*NUMAR_VOTURI/NUMAR_ALEGATORI)
20.1440450845973468926087992135771906663



- Am dori să afișăm un tabel cu procentele obținute de toți candidații, însă vom vedea cum realizăm acest lucru într-un paragraf următor.
- După cum am precizat la funcția **COUNT**, funcțiile de grup, deci și **AVG** ignoră valorile **NULL**. Așadar dacă vom rula comanda:

```
SELECT avg(numar_voturi)
```

```
FROM voturi
```

```
WHERE candidat=13
```

- vom obține valoarea **100**, deși în baza de date există 3 linii pentru candidatul **13**, și doar o linie are completat câmpul **numar_voturi** cu valoarea **100**.



- Dacă dorim să obținem valoarea **33.333**, adică **100/3**, vom scrie:

```
SELECT AVG(NVL(numar_voturi,0))
```

```
FROM voturi
```

```
WHERE candidat=13
```

- adică înlocuim valorile **null** cu valoarea **0**, pentru ca acestea să intre în calculul mediei.



- **STDEV(x)** – funcție statistică definită ca fiind abaterea pătratică a expresiei date. Cu cât valoarea funcției este mai mică cu atât valorile expresiei **x** sunt mai apropiate de medie.
- **VARIANCE(x)** – este o funcție statistică care calculează dispersia expresiei **x**. Se definește ca pătratul abaterii medii pătratice.



- Funcțiile **COUNT**, **MIN**, **MAX** pot fi aplicate și datelor de tip șir de caractere sau dată calendaristică, celelalte funcții fiind aplicabile doar valorilor numerice.
- De exemplu comanda următoare va afișa data celei mai vechi angajări, data celei mai recente angajări, numărul de date de angajare, și numărul de date distincte de angajare din tabela **employees**:

```
select min(hire_date), max(hire_date),  
count(distinct hire_date), count(hire_date)  
from employees
```

MIN(HIRE_DATE)	MAX(HIRE_DATE)	COUNT(DISTINCT HIRE_DATE)	COUNT(HIRE_DATE)
17-JUN-87	29-JAN-00	19	20

9.2. GRUPAREA DATELOR. CLAUZA GROUP BY

- Uneori am putea dori să grupăm liniile dintr-o tabelă și să obținem anumite informații despre grupurile respective.
- De exemplu am dori să calculăm numărul total de voturi obținut de fiecare candidat în toată țara. Cu ceea ce am învățat până acum, am putea rula o comandă de forma celei de mai jos pentru fiecare candidat în parte:

```
SELECT sum(numar_voturi)
```

```
FROM voturi
```

```
WHERE candidat=1
```

- însă această metodă nu este convenabilă, întrucât am dori să obținem un tabel cu toate aceste date.

SUM(NUMAR_VOTURI)
608608



- O astfel de grupare a datelor se poate face folosind clauza **GROUP BY**.

```
SELECT candidat,  
sum(numar_voturi)  
AS "TOTAL VOTURI"  
FROM voturi  
GROUP BY candidat
```

CANDIDAT	TOTAL VOTURI
1	608608
2	3151
3	3174
4	152465
5	35537
6	19821
7	33844
8	6675
9	8447
10	4009
11	30236
12	632543
13	100



- Se observă că pentru fiecare grup de înregistrări s-a obținut câte o singură valoare, adică pentru fiecare candidat am obținut o sumă a tuturor voturilor primite.
- De exemplu candidatul cu codul 1 a obținut în București **347016** voturi, la Iași **196508** voturi iar la Sibiu **65084** voturi, în total **608608** voturi.
- Clauza **GROUP BY** poate fi folosită și fără funcții de grup, doar pentru a afișa liniile grupate după anumit criteriu:

```
SELECT candidat,numar_voturi FROM voturi  
GROUP BY candidat, numar_voturi
```



- Procentul mediu obținut de către fiecare candidat.

SELECT

candidat,AVG(100*numar_voturi/numar_alegatori)

FROM voturi v, judete j WHERE v.judet=j.cod_judet

GROUP BY candidat

CANDIDAT	AVG(100*NUMAR_VOTURI/NUMAR_ALEGATORI)
1	22.6562295618455989756920853154424336476
2	.13424833638246597421520567348011821838
3	.14003282051378316208662701165544554467
4	5.50638342911377223943294320779193667412
5	1.17019960040031154685700409684022462069
6	.68144301477468349708451524754117789898
7	1.07036088696521333741348008106908880327
8	.224347948245650587054654794284450480961
9	.447406194157174323863379832964865398693
10	.166617475745310934099468805148008754076
11	.97161839160269742583080767121400220691
12	20.1440450845973468926087992135771906663
13	.027519401177830370411139853596785733942



REGULI DE FOLOSIRE A CLAUZEI **GROUP BY**

- În clauza **GROUP BY** nu se acceptă aliasele coloanelor, comanda următoare va genera o eroare

```
SELECT department_id As Departament,  
       job_id, MAX(salary)
```

```
FROM employees GROUP BY Departament,  
       job_id
```



- Toate câmpurile care apar în select în afara funcțiilor de grup trebuie să apară în clauza **GROUP BY** ca în exemplele de mai jos:

```
SELECT department_id, job_id, MAX(salary)  
FROM employees GROUP BY department_id, job_id
```

- sau

```
SELECT department_id, department_name, max(salary)  
FROM employees NATURAL JOIN departments  
GROUP BY department_id, department_name
```

- sau

```
SELECT upper(last_name), sum(salary)  
FROM employees GROUP BY last_name
```

- Observați în ultimul exemplu că deși în clauza **SELECT** câmpului **last_name** îi este aplicată o funcție (simplă nu de grup), în clauza **GROUP BY**, **last_name** poate să apară fără funcția respectivă. (*se pot confunda funcțiile singulare cu cele de grup*)

- Nu se pot folosi funcții de grup în clauza **WHERE**. De aceea următoarea comandă nu va putea fi rulată ea generând o eroare:

```
SELECT * FROM voturi
```

```
WHERE numar_voturi=max(numar_voturi)
```

- Pentru a putea afla ce candidat/candidați au obținut cele mai multe voturi vom folosi o subinterogare (asupra acestui subiect vom reveni în capitolul următor) astfel:

```
SELECT * FROM voturi
```

```
WHERE numar_voturi =
```

```
(SELECT max(numar_voturi) from voturi)
```



- În clauza **GROUP BY** pot să apară și alte coloane care nu apar în **SELECT**

```
SELECT MAX(salary) FROM employees  
GROUP BY departments
```



- funcțiile de grup pot fi imbricate ca în exemplul următor, în care am determinat cel mai mare număr total de voturi obținut de către un candidat.

**SELECT max(sum(numar_voturi)) FROM voturi
GROUP BY candidat**

MAX(SUM(NUMAR_VOTURI))
632543



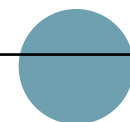
9.3. SELECTAREA GRUPURILOR. CLAUZA HAVING

- De multe ori nu ne interesează să afișăm toate grupurile de obținute prin folosirea clauzei **GROUP BY**. Pentru a filtra grupurile folosim clauza **HAVING**. Așa cum am văzut în exemplele anterioare putem folosi clauza **GROUP BY** fără clauza **HAVING** însă clauza **HAVING** poate fi folosită doar atunci când este prezentă clauza **GROUP BY**.
- Presupunem că dorim să afișăm toți candidații care au obținut un procent în alegeri mai mare de 5% din numărul total de persoane cu drept de vot. Pentru aceasta procedăm astfel:
 - folosim clauza **GROUP BY** pentru a grupa liniile după candidați și calculăm pentru fiecare candidat procentul obținut:

```
SELECT candidat,  
100*sum(numar_voturi)/sum(numar_alegatori)  
FROM voturi v JOIN judete j  
ON v.judet=j.cod_judet  
GROUP BY candidat
```



CANDIDAT	100*SUM(NUMAR_VOTURI)/SUM(NUMAR_ALEGATORI)
1	22.0222817982769582150245277809640393096
2	.114017906347551618341432066351112190219
3	.114850153839139586358522811360974322994
4	5.51689625238954537938001904037522059082
5	1.28589474385050519231973067023785271463
6	.717216414381091915945898123499014510416
7	1.22463409153492128567039887451191398469
8	.24153269592824723974264012786216244675
9	.305651937454068080015892308621975458831
10	.145064356251137555674643336719012621576
11	1.09407978937625221585894635296484550411
12	22.8883619596316544971578748162270892216
13	.003618467354730295726481500042878838154



- Folosim clauza **HAVING** pentru a filtra grupurile care se vor afișa

SELECT candidat,

100*sum(numar_voturi)/sum(numar_alegatori)

FROM voturi v **JOIN** judete j

ON (v.judet=j.cod_judet)

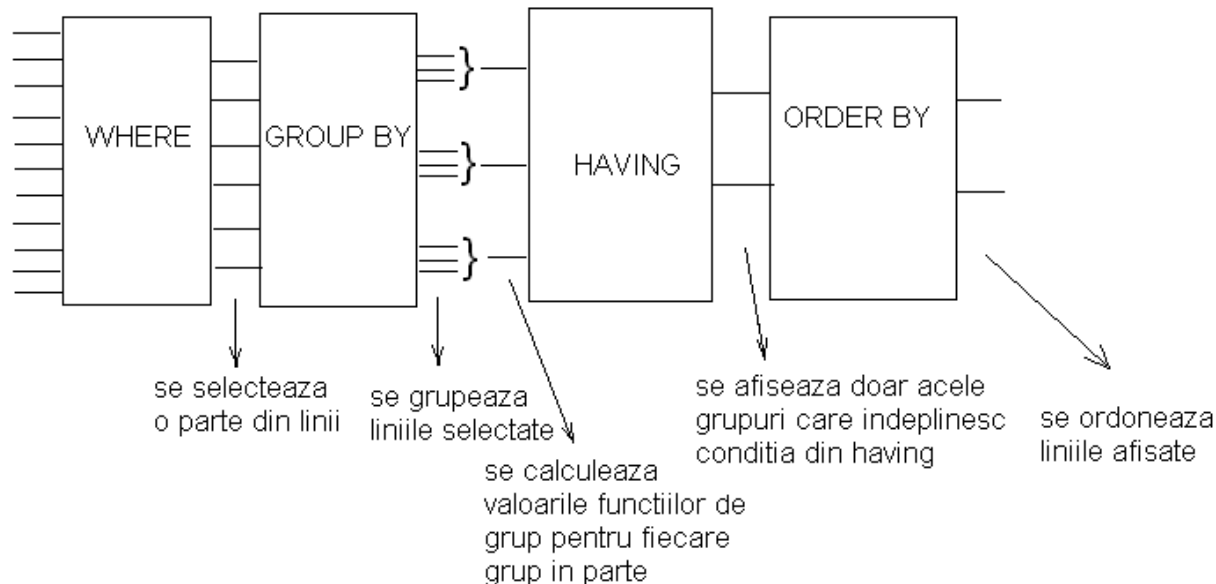
GROUP BY candidat

HAVING **100*sum(numar_voturi)/sum(numar_alegatori)>5**

CANDIDAT	100*SUM(NUMAR_VOTURI)/SUM(NUMAR_ALEGATORI)
1	22.0222817982769582150245277809640393096
4	5.51689625238954537938001904037522059082
12	22.8883619596316544971578748162270892216



- Bineînțeles că putem folosi clauzele **WHERE**, **GROUP BY** și **HAVING** împreună. În acest caz, clauza **WHERE** va filtra mai întâi liniile din tabelă, liniile rămase vor fi grupate apoi conform criteriului dat de clauza **GROUP BY** și în final sunt afișate doar acele grupuri care respectă condiția dată de clauza **HAVING**.
- Atenție! Trebuie făcută distincția clară dintre clauzele **WHERE** și **HAVING**. Clauza **WHERE** acționează asupra liniilor în timp ce **HAVING** acționează la nivel de grup.



Ordinea de executare a clauzelor comenzii **SELECT**

- Să vedem de exemplu cum se evaluează comanda următoare

SELECT candidat,

100*sum(numar_voturi)/sum(numar_alegatori)

FROM voturi v JOIN judete j

ON (v.judet=j.cod_judet)

WHERE numar_voturi>15000

GROUP BY candidat

HAVING

100*sum(numar_voturi)/sum(numar_alegatori)>5

CANDIDAT	100*SUM(NUMAR_VOTURI)/SUM(NUMAR_ALEGATORI)
1	22.0222817982769582150245277809640393096
4	5.51689625238954537938001904037522059082
12	29.3512120713181287412967242185267405132

- Observați însă mai întâi că prin adăugarea clauzei **WHERE**, rezultatele obținute diferă puțin de cele obținute anterior, aceasta pentru că la calculul procentului obținut de către candidatul 12 de exemplu nu mai este inclusă următoarea linie din tabelă

JUDET	CANDIDAT	NUMAR_VOTURI
IS	12	12184



- Comanda se evaluează astfel:
 - Mai întâi sunt filtrate liniile din tabelă

```
SELECT candidat, numar_voturi, numar_alegatori  
FROM voturi v JOIN judete j  
ON (v.judet=j.cod_judet)  
WHERE numar_voturi>15000
```

- au fost afișate doar 11 linii din totalul de 39 câte are tabela.

CANDIDAT	NUMAR_VOTURI	NUMAR_ALEGATORI
1	347016	1750192
1	196508	650029
1	65084	363380
4	96744	1750192
5	25656	1750192
4	35784	650029
4	19937	363380
7	25937	1750192
11	22687	1750192
12	514366	1750192
12	105993	363380



- Liniile obținute la pasul anterior sunt grupate pe candidați și se aplică funcțiile de grup

SELECT candidat,

100*sum(numar_voturi)/sum(numar_alegatori)

FROM voturi v JOIN judete j

ON (v.judet=j.cod_judet)

WHERE numar_voturi>15000

GROUP BY candidat

CANDIDAT	100*SUM(NUMAR_VOTURI)/SUM(NUMAR_ALEGATORI)
1	22.0222817982769582150245277809640393096
4	5.51689625238954537938001904037522059082
5	1.46589631309022095861482625906186292704
7	1.48195169444266686169288855165604687943
11	1.29625778200334591861921434905427518809
12	29.3512120713181287412967242185267405132



- În final sunt afișate doar acele linii obținute la pasul anterior care îndeplinesc condiția din clauza **HAVING**.

SELECT candidat,

100*sum(numar_voturi)/sum(numar_alegatori)

FROM voturi v JOIN judete j

ON (v.judet=j.cod_judet)

WHERE numar_voturi>15000

GROUP BY candidat

HAVING

100*sum(numar_voturi)/sum(numar_alegatori)>5

CANDIDAT	100*SUM(NUMAR_VOTURI)/SUM(NUMAR_ALEGATORI)
1	22.0222817982769582150245277809640393096
4	5.51689625238954537938001904037522059082
12	29.3512120713181287412967242185267405132



În această lecție am învățat despre:

- Care sunt funcțiile de grup și cum sunt ele folosite
- Cum se pot grupa datele într-o tabelă
- Care sunt regulile de folosire a clauzei GROUP BY
- Cum se pot filtra grupurile folosind clauza HAVING
- Care este diferența dintre clauzele WHERE și HAVING
- În ce ordine se execută clauzele WHERE, GROUP BY și HAVING

SFÂRȘIT

