

به نام خدا



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دستور کار آزمایشگاه پایگاه داده

جلسه چهارم

Window functions ، Crosstab ، CUBE ، ROLLUP

استاد درس

دکتر شهریاری

ROLLUP

ROLLUP از جمله دستوراتی است که به کمک آن میتوان پرس وجوهایی را طراحی و اجرا کرد که ماهیت گزارشی دارند و میتوانند به عنوان گزارش های مدیریتی در نظر گرفته شوند که قرار است مبدأ تصمیم سازی های مدیریتی باشند. برای استفاده از آن، عبارت ROLLUP را بعد از عبارت BY GROUP می آوریم و فیلدهایی که می خواهیم بر اساس آن ها گزارش تهیه شود را به عنوان مؤلفه های ROLLUP در نظر می گیریم. در این صورت ROLLUP موجب اضافه شدن یک سطر بصورت مجموع مقادیر محاسبه شده برای تابع AGGREGATE براساس مؤلفه تعریف شده در ROLLUP ، خواهد شد.

فرض کنید جدول فروشی داریم که حاوی نام برند، بخش و تعداد است و می خواهیم بدانیم که علاوه بر این که هر برند در هر بخش چه تعدادی فرخته است، می خواهیم این مورد را نیز اضافه کنیم که یک برند در همه بخش ها و همه برند ها در همه بخش ها چه فروشی داشته اند. پرس و جوی زیر با استفاده از ROLLUP این گزارش را می دهد.

```
1. SELECT
2. case grouping(brand)
3. when 0 then brand
4. when 1 then 'All brands'
5. end as brand,
6. case grouping(segment)
7. when 0 then segment
8. when 1 then 'All segment'
9. end as segment,
10. SUM (quantity)
11. FROM
12. sales
13. GROUP BY
14. ROLLUP (brand, segment)
15. ORDER BY
16. brand,
17. segment;
```

نتیجه:

	brand character varying	segment character varying	sum bigint
1	ABC	All segment	300
2	ABC	Basic	200
3	ABC	Premium	100
4	All brands	All segment	700
5	XYZ	All segment	400
6	XYZ	Basic	300
7	XYZ	Premium	100

عبارت GROUPING(brand) به منظور بررسی گروه بندی شدن یا نشدن هر رکورد استفاده شده است. در صورتی که نتیجه یک رکورد حاصل گروه بندی بر اساس فیلد brand باشد مقدار یک و در غیر این صورت مقدار صفر بر می گرداند.

CUBE

این دستور در واقع مشابه ROLLUP است اما بصورت پیشرفته‌تر، گزارش‌های دقیق‌تری را تولید می‌کند. این دستور در واقع به ازای هر گروه‌بندی ممکن از داده‌ها، یک مقدار متناسب با تابع AGGREGATE استفاده شده در پرس و جو را به مجموعه نتایج اضافه می‌کند.

به عنوان مثال، ROLLUP پرس و جو مثال قبل را با CUBE جایگزین می‌کنیم.

```
1. SELECT
2.   case grouping(brand)
3.   when 0 then brand
4.   when 1 then 'All brands'
5.   end as brand,
6.   case grouping(segment)
7.   when 0 then segment
8.   when 1 then 'All segment'
9.   end as segment,
10.  SUM (quantity)
11. FROM
12.   sales
13. GROUP BY
14.   cube (brand, segment)
15. ORDER BY
16.   brand,
17.   segment;
```

نتیجه :

	brand character varying	segment character varying	sum bigint
1	ABC	All segment	300
2	ABC	Basic	200
3	ABC	Premium	100
4	All brands	All segment	700
5	All brands	Basic	500
6	All brands	Premium	200
7	XYZ	All segment	400
8	XYZ	Basic	300
9	XYZ	Premium	100

همانطور که مشخص است دو ردیف به جدول اضافه شده است که شامل ردیف ۵ و ۶ می‌شود. این دو ردیف گروه‌بندی brand با تک تک segment ها می‌باشد.

تفاوت CUBE و ROLLUP

تفاوت این دو دستور در نوع انتخاب گروه بندی می باشد. در CUBE همه زیرمجموعه های ممکن را در نظر می گیرد ولی در ROLLUP ترتیب و تقدم مهم است. به عنوان مثال گروه بندی که بواسطه CUBE برای مجموعه (c1,c2,c3) انجام می شود بصورت زیر است.

```
(c1, c2, c3)
(c1, c2)
(c2, c3)
(c1,c3)
(c1)
(c2)
(c3)
()
```

و گروه بندی که با ROLLUP صورت می گیرد بصورت زیر است.

```
(c1, c2, c3)
(c1, c2)
(c1)
()
```

سوال: اگر بخواهیم گروه بندی انجام دهیم که خروجی آن گزارشی از میزان فروش سالیانه، ماه و روزانه در سال ها و ماه ها و روز های خاص باشد باید از کدام دستور استفاده کرد.

Window Functions

معمولا از این نوع توابع روی مجموعه ای از سطرهای یک جدول، در جهت اعمال عملیات های محاسباتی، ارزیابی داده ها، رتبه بندی و غیره ... استفاده میگردد. به بیان ساده تر به وسیله Window Function ها میتوان، سطرهای یک جدول را گروه بندی نمود. و روی گروه ها از توابعی مثل توابع تجمعی (Functions Aggregate) استفاده کرد. این نوع توابع از قابلیت و انعطاف پذیری زیادی برخوردار هستند و به وسیله آنها میتوان نتایج (خروجی) بسیار مفیدی را از پرس و جو ها بدست آورد، معمولا از این نوع توابع در داده کاوی و گزارش گیری ها استفاده میگردد. کلمه "Window" در Window Function، به مجموعه سطرهایی اشاره میکند، که محاسبات و ارزیابی و ... روی آنها اعمال میگردد.

فرض کنید جدول کالایی به شکل زیر وجود دارد، که حاوی نام کالا، قیمت و گروه کالا می باشد.

	product_id [PK] integer	product_name character varying (255)	price numeric (11,2)	group_id integer
1	1	Microsoft Lumia	200.00	1
2	2	HTC One	400.00	1
3	3	Nexus	500.00	1
4	4	iPhone	900.00	1
5	5	HP Elite	1200.00	2
6	6	Lenovo Thinkpad	700.00	2
7	7	Sony VAIO	700.00	2
8	8	Dell Vostro	800.00	2
9	9	iPad	700.00	3
10	10	Kindle Fire	150.00	3
11	11	Samsung Galaxy Tab	200.00	3

مطابق با توضیحات داده شده می‌خواهیم پرس و جویی را اجرا کنیم که در نهایت به ازای هر کالا مشخص کند که میانگین قیمت گروهی که کالا در آن است چه مقدار است. برای این کار نیاز است محدوده window را به اندازه‌ای داشته باشیم که در آن شناسه گروه (group_id) در آن یکسان باشد در این محدوده تابع میانگین را اجرا کنیم. پرس و جوی زیر با استفاده از windowing این کار را انجام می‌دهد. (جدول product_groups حاوی نام گروه‌ها به همراه شناسه آن‌ها می‌باشد)

```

1. SELECT
2.   product_name,
3.   price,
4.   group_name,
5.   AVG (price) OVER (
6.     PARTITION BY group_name
7.   )
8. FROM
9.   products
10.  INNER JOIN
11.    product_groups USING (group_id);

```

نتیجه:

	product_name character varying (255)	price numeric (11,2)	group_name character varying (255)	avg numeric
1	HP Elite	1200.00	Laptop	850.0000000000000000
2	Lenovo Thinkpad	700.00	Laptop	850.0000000000000000
3	Sony VAIO	700.00	Laptop	850.0000000000000000
4	Dell Vostro	800.00	Laptop	850.0000000000000000
5	Microsoft Lumia	200.00	Smartphone	500.0000000000000000
6	HTC One	400.00	Smartphone	500.0000000000000000
7	Nexus	500.00	Smartphone	500.0000000000000000
8	iPhone	900.00	Smartphone	500.0000000000000000
9	iPad	700.00	Tablet	350.0000000000000000
10	Kindle Fire	150.00	Tablet	350.0000000000000000
11	Samsung Galaxy Tab	200.00	Tablet	350.0000000000000000

همانطور که مشخص است مقدار میانگین صرفاً در محدوده تعیین شده محاسبه شده است. تابع avg بر روی محدوده ای اجرا شده است که با کلیدواژه over بیان شده و به هر مجموعه ای که توسط over انتخاب می شود یک محدوده یا پنجره گفته می شود. کلیدواژه partition بیان گر گروه بندی یا پارتیشن بندی بوسیله یک ستون مشخص است.

window function ها همیشه بعد از Join ، where ، Group by و Having انجام می شوند و پس از آن order by اجرا می شود.

حال این نکته حائز اهمیت است که می توان این محدوده را هوشمندانه تر انتخاب کرد. برای این کار می توان ساختار window function ها را بررسی کرد، که به صورت زیر است.

```
1. window_function(arg1, arg2,..) OVER (  
2.   [PARTITION BY partition_expression]  
3.   [ORDER BY sort_expression [ASC | DESC] [NULLS {FIRST | LAST }]])
```

تابع window_function یک تابع تجمعی است که می تواند ورودی های مختلفی داشته باشد. و همانطور که در قبل گفته شد کلیدواژه over محدوده مورد نظر را انتخاب می کند. و بعد از پارتیشن بندی می توان مشخص کرد که در هر پارتیشنی ردیف ها با چه ترتیبی در خروجی ظاهر شوند. در خصوص مقادیری که مقدار نداشته باشند یعنی NULLS، می توان مشخص کرد که در ابتدا یا در انتها ظاهر شوند، این کار با order By و FIRST یا LAST مشخص می شود.

توابعی که می توان استفاده کرد انواع مختلف دیگری نیز دارد که در لینک زیر برخی از آن ها نشان داده شده است.

<https://www.postgresql.org/docs/9.1/functions-window.html>

بطور مثال توابع ROW_NUMBER() ، RANK() ، DENSE_RANK() عددی را برای هر ردیف از یک محدوده یا پنجره متناسب با ترتیب آن ها انتخاب می کند. با این تفاوت که ROW_NUMBER() رتبه بندی را بصورت پی در پی انجام می دهد. پرس و جوی زیر مثالی از این تابع است.

```
1. SELECT  
2.   product_name,  
3.   group_name,  
4.   price,  
5.   ROW_NUMBER () OVER (  
6.     PARTITION BY group_name  
7.     ORDER BY  
8.       price  
9.   )  
10. FROM  
11.   products  
12. INNER JOIN product_groups USING (group_id);
```

نتیجه :

	product_name character varying (255)	group_name character varying (255)	price numeric (11,2)	row_number bigint
1	Sony VAIO	Laptop	700.00	1
2	Lenovo Thinkpad	Laptop	700.00	2
3	Dell Vostro	Laptop	800.00	3
4	HP Elite	Laptop	1200.00	4
5	Microsoft Lumia	Smartphone	200.00	1
6	HTC One	Smartphone	400.00	2
7	Nexus	Smartphone	500.00	3
8	iPhone	Smartphone	900.00	4
9	Kindle Fire	Tablet	150.00	1
10	Samsung Galaxy Tab	Tablet	200.00	2
11	iPad	Tablet	700.00	3

تابع RANK() مقدار یکسانی برای مقادیر یکسان در نظر می‌گیرد. پرس وجوی زیر مثالی از آن است.

```

1. SELECT
2.     product_name,
3.     group_name,
4.     price,
5.     RANK () OVER (
6.         PARTITION BY group_name
7.         ORDER BY
8.             price
9.     )
10. FROM
11.     products
12. INNER JOIN product_groups USING (group_id);

```

نتیجه:

	product_name character varying (255)	group_name character varying (255)	price numeric (11,2)	rank bigint
1	Sony VAIO	Laptop	700.00	1
2	Lenovo Thinkpad	Laptop	700.00	1
3	Dell Vostro	Laptop	800.00	3
4	HP Elite	Laptop	1200.00	4
5	Microsoft Lumia	Smartphone	200.00	1
6	HTC One	Smartphone	400.00	2
7	Nexus	Smartphone	500.00	3
8	iPhone	Smartphone	900.00	4
9	Kindle Fire	Tablet	150.00	1
10	Samsung Galaxy Tab	Tablet	200.00	2
11	iPad	Tablet	700.00	3

سوال: تابع فرق تابع DENSE_RANK() با تابع RANK() چیست؟

تابع FIRST_VALUE و LAST_VALUE اولین و آخرین مقدار در هر محدوده را انتخاب می‌کند. در پرس و جوی زیر کمترین قیمت در هر گروه برای هر ردیف مشخص شده است.

```
1. SELECT
2.   product_name,
3.   group_name,
4.   price,
5.   FIRST_VALUE (price) OVER (
6.     PARTITION BY group_name
7.     ORDER BY
8.       price
9.   ) AS lowest_price_per_group
10. FROM
11.   products
12. INNER JOIN product_groups USING (group_id);
```

نتیجه:

	product_name character varying (255)	group_name character varying (255)	price numeric (11,2)	lowest_price_per_group numeric
1	Sony VAIO	Laptop	700.00	700.00
2	Lenovo Thinkpad	Laptop	700.00	700.00
3	Dell Vostro	Laptop	800.00	700.00
4	HP Elite	Laptop	1200.00	700.00
5	Microsoft Lumia	Smartphone	200.00	200.00
6	HTC One	Smartphone	400.00	200.00
7	Nexus	Smartphone	500.00	200.00
8	iPhone	Smartphone	900.00	200.00
9	Kindle Fire	Tablet	150.00	150.00
10	Samsung Galaxy Tab	Tablet	200.00	150.00
11	iPad	Tablet	700.00	150.00

ORDER BY

با استفاده از ORDER BY می‌توان محدوده ای را درون پنجره یا محدوده پارتیشن بندی شده ایجاد کرد. به طور مثال می‌خواهیم در هر گروه قیمت کالا را با قیمت کالای قبلی خود جمع کنیم. در این حالت علاوه بر این که پارتیشن بندی باید متناسب با گروه باشد، باید مشخص کرد که در زمان محاسبه هر ردیف فقط کالای قبلی و کالای فعلی را در نظر بگیرد. پرس و جوی زیر نشان دهنده همین حالت است.

```
1. SELECT
2.   product_name,
3.   group_name,
4.   price,
5.   sum (price) OVER (
6.     PARTITION BY group_name
7.     ORDER BY
8.       price ROWS BETWEEN 1 PRECEDING
```



```

9.          AND CURRENT ROW
10.      ) AS highest_price_per_group
11. FROM
12.      products
13. INNER JOIN product_groups USING (group_id);

```

نتیجه:

	product_name character varying (255)	group_name character varying (255)	price numeric (11,2)	highest_price_per_group numeric
1	Sony VAIO	Laptop	700.00	700.00
2	Lenovo Thinkpad	Laptop	700.00	1400.00
3	Dell Vostro	Laptop	800.00	1500.00
4	HP Elite	Laptop	1200.00	2000.00
5	Microsoft Lumia	Smartphone	200.00	200.00
6	HTC One	Smartphone	400.00	600.00
7	Nexus	Smartphone	500.00	900.00
8	iPhone	Smartphone	900.00	1400.00
9	Kindle Fire	Tablet	150.00	150.00
10	Samsung Galaxy Tab	Tablet	200.00	350.00
11	iPad	Tablet	700.00	900.00

توابع LEAD و LAG

کاربرد این توابع هنگامی است که نیاز به ردیف های بعدی یا قبلی در یک پارتیشن وجود دارد. بطور مثال پرس وجوی زیر در هر گروه نشان می دهد که قیمت کالای قبلی چه مقدار بوده است. ورودی های این دو تابع به ترتیب ، عملیات محاسباتی که بر روی یک ستون انجام می شود، تعداد ردیفی که باید قبل تر یا بعدتر باید به آن دسترسی داشته باشد و مقدار پیشفرض در صورت نبود ردیف قبلی یا بعدی. پرس و جوی زیر مثالی از این مورد است.

```

1. SELECT
2.     product_name,
3.     group_name,
4.     price,
5.     LAG (price, 1) OVER (
6.         PARTITION BY group_name
7.         ORDER BY
8.             price
9.     ) AS prev_price,
10.    price - LAG (price, 1) OVER (
11.        PARTITION BY group_name
12.        ORDER BY
13.            price
14.    ) AS cur_prev_diff
15. FROM
16.     products
17. INNER JOIN product_groups USING (group_id);

```

نتیجه:

	product_name character varying (255)	group_name character varying (255)	price numeric (11,2)	prev_price numeric	cur_prev_diff numeric
1	Sony VAIO	Laptop	700.00	[null]	[null]
2	Lenovo Thinkpad	Laptop	700.00	700.00	0.00
3	Dell Vostro	Laptop	800.00	700.00	100.00
4	HP Elite	Laptop	1200.00	800.00	400.00
5	Microsoft Lumia	Smartphone	200.00	[null]	[null]
6	HTC One	Smartphone	400.00	200.00	200.00
7	Nexus	Smartphone	500.00	400.00	100.00
8	iPhone	Smartphone	900.00	500.00	400.00
9	Kindle Fire	Tablet	150.00	[null]	[null]
10	Samsung Galaxy Tab	Tablet	200.00	150.00	50.00
11	iPad	Tablet	700.00	200.00	500.00

Crosstab

یکی افزونه های پایگاه داده PostgreSQL توابعی هستند که به آن ها tablefunc گفته می شود. ویژگی این توابع این است که در خروجی آن ها یک جدول است. یکی از این توابع CROSSTAB نام دارد. این تابع در واقع جایگزین PIVOT TABLE ها در این پایگاه داده می باشد. برای استفاده از این تابع باید افزونه tablefunc را بصورت زیر فعال کرد.

```
CREATE EXTENSION IF NOT EXISTS tablefunc;
```

به کمک این دستور میتوان جای سطر و ستون را در گزارش عوض کرد به عبارت بهتر به جای آنکه چندین سطر داشته باشیم، جهت راحت شدن تجزیه و تحلیل گزارش، تعدادی ستون به گزارش اضافه می کنیم.

ساختار کلی این تابع بصورت زیر است.

```
crosstab(text source_sql, text category_sql)
```

ورودی اول یعنی source_sql، پرس و جویی است که اطلاعات اصلی از آن نشئت می گیرد، این پرس و جو باید یک ستون row_name، یک ستون category و یک ستون value داشته باشد. (نام ستون ها در اینجا با کاربریشان تناسب دارد). ستون row_name حتما باید اولین ستون باشد. و category, value باید حتما آخرین ستون ها باشند. هر تعداد ستونی می تواند بعد از row_name باشد.

ورودی دوم یعنی category_sql باید پرس و جویی باشد که تنها یک ستون، حداقل یک ردیف داشته باشد و بدون ردیف تکراری.

در نهایت خروجی این تابع به این صورت است که به ازای هر مقدار یکتا در row_name یک ردیف ایجاد می کند و مقدار row_name به همراه ستون های اضافی (به غیر از دو ستون آخر) را بدون تغییر در آن ردیف می گذارد. مقدار value نیز با توجه

به این که در ستون category چه مقداری وجود دارد، در ستونی که متعلق به آن category است قرار می‌گیرد. اگر مقدار category با هیچکدام از مقادیر خروجی پرس و جوی category_sql یکسان نباشد آن ردیف نادیده گرفته می‌شود.

فرض شود جدولی فروشی مثل زیر وجود دارد.

	year integer	month integer	qty integer
1	2007	1	1000
2	2007	2	1500
3	2007	7	500
4	2007	11	1500
5	2007	12	2000
6	2008	1	1000

و پرس و جوی زیر را بر روی آن اجرا می‌کنیم.

```

1. select * from crosstab(
2.   'select year, month, qty from sales order by 1',
3.   'select m from generate_series(1,12) m'
4. ) as (
5.   year int,
6.   "Jan" int,
7.   "Feb" int,
8.   "Mar" int,
9.   "Apr" int,
10.  "May" int,
11.  "Jun" int,
12.  "Jul" int,
13.  "Aug" int,
14.  "Sep" int,
15.  "Oct" int,
16.  "Nov" int,
17.  "Dec" int
18. );

```

نتیجه:

	year integer	Jan integer	Feb integer	Mar integer	Apr integer	May integer	Jun integer	Jul integer	Aug integer	Sep integer	Oct integer	Nov integer	Dec integer
1	2007	1000	1500	[null]	[null]	[null]	[null]	500	[null]	[null]	[null]	1500	2000
2	2008	1000	[null]	[null]	[null]	[null]	[null]	[null]	[null]	[null]	[null]	[null]	[null]

همانطور که مشخص است، خروجی به شکلی است که انگار جای سطر و ستون category عوض شده است. این نوع پرس و جوی در گزارش گیری بسیار مورد استفاده است.

توجه: در تمارینی که مطرح می‌شود، از جداولی استفاده شده است که در پایگاه داده rental DVD وجود دارد. این پایگاه داده یک نوع پایگاه داده آموزشی است. برای دسترسی به آن می‌توانید کل پایگاه داده را دانلود کرده و با استفاده از فایل tar، آن را در پایگاه داده خود بازیابی کنید. مشخصات کامل این پایگاه داده نمونه به همراه model ER آن در لینک زیر وجود دارد.

<https://www.postgresqltutorial.com/postgresql-sample-database>

تمرین

۱. با استفاده از crosstab پرس و جویی بنویسید که مشخص کند تعداد فیلم اجاره ای مربوط به هر store در هر

category چه تعداد بوده است. بخشی از خروجی به شکل زیر است:

۲. پرس و جویی بنویسید که مجموع amount اجاره فیلم هارا بر اساس روز، ماه و سال را با استفاده از ROLLUP گزارش دهد. بخشی از خروجی به شکل زیر است.

	year double precision	mmonth double precision	dday double precision	sum numeric
1	2005	6	14	41.89
2	2005	6	15	1179.97
3	2005	6	16	1191.11
4	2005	6	17	1158.19
5	2005	6	18	1284.99
6	2005	6	19	1283.92
7	2005	6	20	1223.09
8	2005	6	21	986.69
9	2005	6	[null]	8349.85
10	2005	7	5	128.73

۳. در سوال قبل بجای ROLLUP از CUBE استفاده کنید و توضیح دهید چه چیزی را نشان می‌دهد.

۴. با استفاده از window functions پرس و جویی بنویسید که در جدول film، فیلم ها را بر اساس length رتبه بندی کند. این رتبه بندی باید در هر category باشد.

نکات مهم هنگام تمرین:

- کلیه مستند تکلیف‌ها (کد و گزارش و خروجی‌ها) باید در قالب یک فایل فشرده با نام `DBlab_#studentnumber_#assignmentnumber` (عدد آخر شماره سری تمرین است) در سامانه **courses** آپلود شوند.
- در صورت دستی نوشتن تکالیف تشریحی و ناخوانا بودن آن‌ها باعث کسر نمره می‌شود.
- در صورت بروز اشکال در مورد هر تمرین تا ۲۴ ساعت قبل از اتمام زمان ارسال، به آدرس ایمیل بزنید.
- تاخیر در ارسال تکالیف به ازای هر روز مشمول کسر نمره خواهد شد.
- در صورت کشف تقلب براساس مقررات آموزشی برخورد خواهد شد.
- حداکثر تعداد صفحات باید ۱۰ صفحه باشد.

موفق باشید