

Слабо-структурированные данные в СУБД PostgreSQL (мастер-класс)

Александр Коротков, Иван Панченко

Postgres Professional

2017

1. Материалы

Docker

```
docker run --rm -p 127.0.0.1:9000:5432 waaeer/semitutorial
```

VirtualBox

```
https://pgconf.ru/media/2017/03/14/json-tutorial96.ova
```

Руки

```
https://pgconf.ru/2017/93658
```

Там же ссылки на документацию.

2. Jsonb

```
# \d company
```

```
Table "public.company"
```

```
Column | Type | Modifiers
```

```
-----+-----+-----
js      | jsonb |
```

```
# SELECT js FROM company limit 5;
```

```
-----+-----+-----
{"_id": {"$oid": "52cdef7c4bab8bd675297d8a"}, "name": "Wetpaint"}, {"_id": {"$oid": "52cdef7c4bab8bd675297d8b"}, "name": "AdventNe"}, {"_id": {"$oid": "52cdef7c4bab8bd675297d8c"}, "name": "Zoho"}, {"_id": {"$oid": "52cdef7c4bab8bd675297d8d"}, "ipo": null, "nan"}, {"_id": {"$oid": "52cdef7c4bab8bd675297d8e"}, "ipo": {"pub_day"
```

```
(5 rows)
```

```
# SELECT jsonb_pretty(js) FROM company LIMIT 1;
-[ RECORD 1 ]+-----
jsonb_pretty | {
              |     "_id": {
              |         "$oid": "52cdef7c4bab8bd675297d8a"
              |     },
              |     "name": "Wetpaint",
              |     "image": {
              |         "available_sizes": [
              |             [
              |                 150,
              |                 75
              |             ]
              |         ]
              |     }
              | }
.....
```

```
# SELECT js -> 'name',
       js -> 'acquisition' -> 'price_amount',
       js #> '{acquisition,price_amount}'
FROM company
LIMIT 5;
```

?column?	?column?	?column?
"Wetpaint"	30000000	30000000
"AdventNet"	NULL	NULL
"Zoho"	NULL	NULL
"Digg"	500000	500000
"Facebook"	NULL	NULL

(5 rows)

-> вытаскивает тип jsonb

```
# SELECT js -> 'name' FROM company
   WHERE js -> 'name' ILIKE '%paint%';
ERROR: operator does not exist: jsonb ~>* unknown
```

Приведение значения поля к строке

```
# SELECT (js -> 'name')::text, js ->> 'name'
   FROM company LIMIT 5;
   text      | ?column?
```

```
-----+-----
"Wetpaint"   | Wetpaint
"AdventNet"  | AdventNet
"Zoho"       | Zoho
"Digg"       | Digg
"Facebook"   | Facebook
(5 rows)
```



```
# SELECT
  js ->> 'name',
  js ->> 'number_of_employees'
FROM company
WHERE
  js ->> 'name' ILIKE '%paint%' AND
  (js ->> 'number_of_employees')::int >= 10;
?column? | ?column?
-----+-----
Wetpaint | 47
(1 row)
```

```
# SELECT
    js -> 'offices' -> 0 ->> 'city',
    js #>> '{offices,0,city}'
FROM company LIMIT 5;
?column? | ?column?
-----+-----
Seattle  | Seattle
Pleasanton | Pleasanton
Pleasanton | Pleasanton
San Francisco | San Francisco
Menlo Park | Menlo Park
(5 rows)
```

```

# CREATE UNIQUE INDEX company_id_idx
ON company ((js -> '_id' -> '$oid'));
CREATE INDEX

# EXPLAIN ANALYZE SELECT js -> 'name' FROM company
  WHERE js -> '_id' -> '$oid' = '52cdef7c4bab8bd675297d8a';
                                                    QUERY PL

-----

Index Scan using company_id_idx on company (cost=0.29..8.31 rows=4)
  Index Cond: (((js -> '_id'::text) -> '$oid'::text) = '52cdef7c4bab8bd675297d8a')
Planning time: 0.084 ms
Execution time: 0.156 ms
(4 rows)
  
```

```
# SELECT jsonb_array_elements(js -> 'competitions')
FROM company
WHERE js -> '_id' ->> '$oid' = '52cdef7c4bab8bd675297d8a';
      jsonb_array_elements
```

```
{ "competitor": { "name": "Wikia", "permalink": "wikia" } }
{ "competitor": { "name": "JotSpot", "permalink": "jotspot" } }
{ "competitor": { "name": "Socialtext", "permalink": "socialtext" } }
{ "competitor": { "name": "Ning by Glam Media", "permalink": "ning" } }
{ "competitor": { "name": "Soceeo", "permalink": "soceeo" } }
{ "competitor": { "name": "Yola", "permalink": "yola" } }
{ "competitor": { "name": "SocialGO", "permalink": "socialgo" } }
{ "competitor": { "name": "IslamNor", "permalink": "islamnor" } }
(8 rows)
```

```
# SELECT (jsonb_each(js #> '{offices, 0}')).*
FROM company
WHERE js -> '_id' ->> '$oid' = '52cdef7c4bab8bd675297d8a';
```

key	value
city	"Seattle"
address1	"710 - 2nd Avenue"
address2	"Suite 1100"
latitude	47.603122
zip_code	"98104"
longitude	-122.333253
state_code	"WA"
description	""
country_code	"USA"

(9 rows)

- ▶ Для скалярных a и b , $a @> b$ тогда и только тогда, когда $a = b$.
- ▶ Если a и b объекты, то для каждого ключа в b должен существовать ключ в a , содержащий соответствующее значение.
- ▶ Если a и b массивы, то для каждого элемента в b должен существовать содержащий его элемент в a .

Примеры.

arg1	arg2	result
1	1	true
1	"1"	false
{"a": 1, "b": 2}	{"a": 1}	true
{"a": 1, "b": 2}	{"a": 2}	false
[1, 2]	[1]	true
[{"a": 1, "b": 2}]	[{"a": 1}, {"b": 2}]	true
[{"a": 1}, {"b": 2}]	[{"a": 1, "b": 2}]	false

```
$ git clone https://github.com/plv8/plv8.git
$ cd plv8
$ make USE_PGXS=1
$ sudo make USE_PGXS=1 install
$ make USE_PGXS=1 installcheck
$ psql DB -c "CREATE EXTENSION plv8;"
```

Написать на plv8 функцию, которая для компании выводит список стран, где у неё есть офисы.

```
offices: [
  { country_code : "RUS" }
]
```



```
CREATE FUNCTION extract_countries(company jsonb)
RETURNS text[]
AS $$
    var tmp = {}, result = [];
    for(var i = 0; i < company.offices.length; i++)
        tmp[company.offices[i].country_code] = true;
    for(var i in tmp)
        result.push(i)
    return result;
$$ LANGUAGE plv8 IMMUTABLE STRICT;
```

```
CREATE INDEX company_countries_idx ON company
USING gin(extract_countries(js));
# EXPLAIN ANALYZE SELECT * FROM company
WHERE extract_countries(js) @> '{RUS}';
```

QUERY

```
Bitmap Heap Scan on company (cost=8.73..348.48 rows=94 width=8)
  (actual time=0.038..0.089 rows=50 loops=1)
  Recheck Cond: (extract_countries(js) @> '{RUS}'::text[])
  Heap Blocks: exact=48
   -> Bitmap Index Scan on company_countries_idx (cost=0.00..8.73)
        Index Cond: (extract_countries(js) @> '{RUS}'::text[])
Planning time: 0.057 ms
Execution time: 0.113 ms
(7 rows)
```

Написать на plv8 агрегат, который который посчитает частоты различных размеров картинок.

```
image: { available_sizes : [
  [800,600], [1024,768]
]}
```

```

CREATE FUNCTION image_size_agg_sfunc(state jsonb,
                                     company jsonb) RETURNS jsonb
AS $$
    if (!company.image) return state;
    var sizes = company.image.available_sizes
    for (var i = 0; i < sizes.length; i++)
    {
        var sizeKey = sizes[i][0][0] + 'x' +
                     sizes[i][0][1];
        if (sizeKey in state) state[sizeKey]++;
        else state[sizeKey] = 1;
    }
    return state;
$$ LANGUAGE plv8 IMMUTABLE STRICT;
  
```

```

CREATE AGGREGATE image_size_agg (jsonb) (
  SFUNC = image_size_agg_sfunc,
  STYPE = jsonb,
  INITCOND = '{}');

# SELECT image_size_agg(js) FROM company;
-[ RECORD 1 ]-----+-----
image_size_agg | {"1x1": 15, "32x8": 3, "150x7": 1, "150x8": 5,

```

```
$ git clone https://github.com/postgrespro/jsquery.git
$ cd jsquery
$ make USE_PGXS=1
$ sudo make USE_PGXS=1 install
$ make USE_PGXS=1 installcheck
$ psql DB -c "CREATE EXTENSION jsquery;"
```

- ▶ GIN jsonb_ops: индексируем ключи отдельно, значения отдельно (операторы @>, ?, ?&, ?|).
- ▶ GIN jsonb_path_ops: индексируем hash(path + value) (операторы @>).
- ▶ jsonb_path_value_ops: индексируем hash(path) + hash(value) (операторы @>, jsonb @@ jsquery).
- ▶ jsonb_value_path_ops: индексируем hash(value) + bloom(path) (операторы @>, jsonb @@ jsquery).
- ▶ Индексы по отдельным ключам (js -> 'key').

```
-- Оператор @>
# SELECT js->'name' FROM company
   WHERE js @> '{"offices":[{"country_code": "RUS"}]}' ;

-- Оператор jsquery
# SELECT js->'name' FROM company
   WHERE js @@ 'offices.#.country_code = "RUS"' ;

-- Разворачивание jsonb в rowset
# SELECT js->'name' FROM company
   WHERE EXISTS (
     SELECT 1
     FROM jsonb_array_elements(js->'offices') e1
     WHERE e1 ->> 'country_code' = 'RUS') ;
```


Найти компании, конкуренты Wikia, у которых есть офис в США.

Реализовать решение с помощью:

- ▶ оператора @>,
- ▶ jquery,
- ▶ разворачивания jsonb в rowset.

```
{ competitions: [ { competitor : { name : "EE" }}, ...] }
```

```
-- Оператор @>
# SELECT js->'name' FROM company
WHERE
  js @>
    '{"competitions":[{"competitor":
      {"name": "Wikia"}
    ]}]' AND
  js @>
    '{"offices":[{"country_code": "USA"}]}' ;
```

```
-- jsquery
# SELECT js->'name' FROM company
WHERE js @@
      'competitions.#.competitor.name = "Wikia" AND
      offices.#.country_code = "USA"');
```

```

-- Разворачивание jsonb в rowset
SELECT js->'name' FROM company
WHERE
  EXISTS (
    SELECT 1
    FROM jsonb_array_elements(js->'competitions') e1
    WHERE e1 -> 'competitor' ->> 'name' = 'Wikia'
  ) AND EXISTS (
    SELECT 1
    FROM jsonb_array_elements(js->'offices') e1
    WHERE e1 ->> 'country_code' = 'USA'
  );
  
```

Найти компании, в которые Trinity Ventures инвестировала в 2007 или 2008 году.

Реализовать решение с помощью:

- ▶ оператора @>,
- ▶ jquery,
- ▶ разворачивания jsonb в rowset.

```
funding_rounds: [
  { funded_year : 2007, investments: [
    {financial_org: { name: "000"}}
  ]},...
```

```

-- Оператор @>
# SELECT js->'name' FROM company
WHERE
  js @>
    '{"funding_rounds":[{"funded_year": 2007,
      "investments":[{"financial_org":
        {"name":"Trinity Ventures"}}
      ]}
    ]}' OR
  js @>
    '{"funding_rounds":[{"funded_year": 2008,
      "investments":[{"financial_org":
        {"name":"Trinity Ventures"}}
      ]}
    ]}' ;
  
```

```
-- jsquery
# SELECT js->'name' FROM company
WHERE
  js @@
    'funding_rounds.#(
      funded_year IN (2007,2008) AND
      investments.#.
      financial_org.name = "Trinity Ventures")';
```

```

-- Разворачивание jsonb в rowset
# SELECT js->'name' FROM company
WHERE EXISTS (
  SELECT 1
  FROM jsonb_array_elements(js->'funding_rounds') e1
  WHERE e1 -> 'funded_year' IN ('2007','2008') AND EXISTS
  (SELECT 1
   FROM jsonb_array_elements(e1 -> 'investments') e2
   WHERE e2 -> 'financial_org' ->>
         'name' = 'Trinity Ventures'));
  
```


Найти компании, которые в 2008 году привлекли более 100 000 000 USD инвестиций.

Реализовать решение с помощью:

- ▶ jquery,
- ▶ разворачивания jsonb в rowset.

```
funding_rounds: [
  { funded_year : 2007, raised_amount : 1000000,
    raised_currency_code: "USD"
  }, ....
]
```

```
-- jquery
# SELECT js -> 'name'
FROM company
WHERE js @@
      'funding_rounds.# (funded_year = 2008 AND
                        raised_currency_code = "USD" AND
                        raised_amount > 100000000)';
```

```
-- Разворачивание jsonb в rowset
# SELECT js->'name' FROM company
  WHERE EXISTS (
    SELECT 1
    FROM jsonb_array_elements(js->'funding_rounds') e1
    WHERE
      e1 -> 'funded_year' = '2008' AND
      e1 ->> 'raised_currency_code' = 'USD' AND
      e1 -> 'raised_amount' > '100000000');
```

Найти компании, у которых суммарный объём инвестиций был более 1 000 000 000 USD.

Реализовать решение с помощью разворачивания jsonb в rowset.

```
funding_rounds: [
  { funded_year : 2007, raised_amount : 1000000,
    raised_currency_code:"USD"
  }, ....
]
```

```
-- Разворачивание jsonb в rowset
SELECT js->'name' FROM company
WHERE (
  SELECT SUM((e1 ->> 'raised_amount')::numeric)
  FROM jsonb_array_elements(js->'funding_rounds') e1
  WHERE
    e1 ->> 'raised_currency_code' = 'USD') > 1000000000;
```

Выдать названия компаний в порядке количества упомянутых лиц.

Реализовать решение с помощью разворачивания jsonb в rowset.

```
relationships: [
  { person : { last_name : "Ivanov" } },
  ....
]
```

```

SELECT js->>'name',
       (SELECT COUNT      (DISTINCT r->'person'->>'last_name')
        FROM jsonb_array_elements(js->'relationships') r
       ) n
FROM company ORDER BY n DESC;
  
```

Конкурс Выдать названия компаний в порядке количества упомянутых лиц, и сам список этих лиц для каждой компании в виде массива.

Реализовать решение с помощью разворачивания `jsonb` в `rowset`.

Спасибо за внимание!