



SQL search for keywords and statistics

Markus Schatten, PhD

Faculty of Organization and Informatics, University of Zagreb Pavlinska 2, 42000 Varaždin markus.schatten@foi.hr

08.04.2011.













Markus Schatten, PhD (FOI)

SQL search for keywords and statistics

Outline

Introduction

- Pattern Matching in PostgreSQL
- 3 POSIX Regular Expressions
- 4 Statistics and aggregation
- 5 Assignment

Introduction	Pattern Matching in PostgreSQL	POSIX Regular Expressions	Statistics and aggregation	Assignment
Tools				

PostgreSQL 8.4

Introduction	Pattern Matching in PostgreSQL	POSIX Regular Expressions	Statistics and aggregation	Assignment
Tools				

- PostgreSQL 8.4
- OpenOffice.org Base
 - Connect with host=localhost dbname=enron port=5432 username ubuntu

"Good" databases often do not need pattern matching

- "Good" databases often do not need pattern matching
- A curious case:
 - field name "data"
 - type varchar(200)
 - structure: first 10 characters internal user ID (numeric), then user's surname ended by a '\$' sign, then user's name again ended by a '\$' sign, one character indicator of user type ...

- "Good" databases often do not need pattern matching
- A curious case:
 - field name "data"
 - type varchar(200)
 - structure: first 10 characters internal user ID (numeric), then user's surname ended by a '\$' sign, then user's name again ended by a '\$' sign, one character indicator of user type ...
- Databases with lots of textual information (web applications, document management systems, ...)

- "Good" databases often do not need pattern matching
- A curious case:
 - field name "data"
 - type varchar(200)
 - structure: first 10 characters internal user ID (numeric), then user's surname ended by a '\$' sign, then user's name again ended by a '\$' sign, one character indicator of user type ...
- Databases with lots of textual information (web applications, document management systems, ...)

o ...

Outline



Pattern Matching in PostgreSQL

- 3 POSIX Regular Expressions
- 4 Statistics and aggregation
- 5 Assignment

The LIKE expression

Standard SQL pattern matching expression

The LIKE expression

- Standard SQL pattern matching expression
- Uses simple wildcard characters to match strings
 - % matches any 0 or more characters
 - _ matches exactly one character

The LIKE expression

- Standard SQL pattern matching expression
- Uses simple wildcard characters to match strings
 - % matches any 0 or more characters
 - matches exactly one character
- Usefull for keyword search and (very) simple patterns

Introduction	Pattern Matching in Postgr	eSQL	POSIX Regular Expressi		Statistics and aggregation	Assignment
SELECT	'e-Discovery'	LIKE	'%e%very%'			
SELECT	'e-Discovery'	LIKE	'%Disco%'			
SELECT	'e-Discovery'	LIKE	'is	_'		
SELECT	'e-Discovery'	NOT	LIKE '_e%'			
SELECT	'e-Discovery'	NOT	LIKE ' %asy%'			
SELECT	'e-Discovery'	LIKE	'%over_'			

SELECT * FROM emails WHERE subject LIKE '%risk%'	
SELECT * FROM emails WHERE "From" LIKE '%@enron.com%'	

The ILIKE expression

PostgreSQL specific pattern matching expression (not standard!)

The ILIKE expression

- PostgreSQL specific pattern matching expression (not standard!)
- Same functionality as LIKE but case insensitive

SELECT *)
FROM emails	
WHERE subject ILIKE '%risk%'	
SELECT *	
FROM emails	
WHERE recipients NOT ILIKE '%@enron.com%'	

LIKE and ILIKE

Shortcuts:

- ~~ = LIKE
- !~~ = NOT LIKE
- ~~* = ILIKE
- !~~★ = NOT ILIKE

LIKE and ILIKE

Shortcuts:

- ~~ = LIKE
- !~~ = NOT LIKE
- ~~* = ILIKE
- !~~* = NOT ILIKE
- Advantage very fast

LIKE and ILIKE

Shortcuts:

- ~~ = LIKE
- !~~ = NOT LIKE
- ~~* = ILIKE
- !~~* = NOT ILIKE
- Advantage very fast
- Dissadvantage very limited

The SIMILAR TO expression

• "a curious cross between LIKE notation and common regular expression notation"

The SIMILAR TO expression

- "a curious cross between LIKE notation and common regular expression notation"
- Uses SQL standard definition of regular expressions

The SIMILAR TO expression

- "a curious cross between LIKE notation and common regular expression notation"
- Uses SQL standard definition of regular expressions
- Like the LIKE expression always match the whole string

% matches any 0 or more characters

- % matches any 0 or more characters
- _ matches exactly one character

- % matches any 0 or more characters
- _ matches exactly one character
- denotes alternation (either of two alternatives)

- % matches any 0 or more characters
- _ matches exactly one character
- denotes alternation (either of two alternatives)
- * denotes repetition of the previous item zero or more times

- % matches any 0 or more characters
- _ matches exactly one character
- denotes alternation (either of two alternatives)
- * denotes repetition of the previous item zero or more times
- + denotes repetition of the previous item one or more times

- % matches any 0 or more characters
- _ matches exactly one character
- denotes alternation (either of two alternatives)
- denotes repetition of the previous item zero or more times
- + denotes repetition of the previous item one or more times
- () parentheses can be used to group items into a single logical item

- % matches any 0 or more characters
- _ matches exactly one character
- denotes alternation (either of two alternatives)
- denotes repetition of the previous item zero or more times
- + denotes repetition of the previous item one or more times
- () parentheses can be used to group items into a single logical item
- [...] a bracket expression specifies a character class

	Introduction	Pattern Matching in PostgreSQL	POSIX Regular Expressions	Statistics and aggregation	Assignment
ĺ	SELECT	'never' SIMILAR TO '	never'		
	SELECT	'gonna' NOT SIMILAR	TO 'gon'		
	SELECT	'give' SIMILAR TO '%	5(i w)%'		
	SELECT	'you up' NOT SIMILAR	TO ' (you_) (up)	,	

```
SELECT 'email: mschatte@foi.hr'
SIMILAR TO 'email: *%@foi.hr'
SELECT 'email: mschatte@foi.hr'
SIMILAR TO 'email: *%@foi.hr'
```

SELECT	'foooooo'	SIMILAR TO	′f(o)*′
--------	-----------	------------	---------

SELECT 'fooooo' SIMILAR TO 'f(o)+'

```
SELECT 'f' SIMILAR TO 'f(o) *'
```

SELECT 'f' NOT SIMILAR TO 'f(\circ)+'

[0-9] matches any digit (0-9)

- [0-9] matches any digit (0-9)
- [a-z] matches any lowercase character

- [0-9] matches any digit (0-9)
- [a-z] matches any lowercase character
- [A-Z] matches any uppercase character

- [0-9] matches any digit (0-9)
- [a-z] matches any lowercase character
- [A-Z] matches any uppercase character
- [.?#] matches either the '.', the '?' or the '#' character
Bracket expressions

- [0-9] matches any digit (0-9)
- [a-z] matches any lowercase character
- [A-Z] matches any uppercase character
- [.?#] matches either the '.', the '?' or the '#' character
- [a-z3-6.,] matches a lowercase character, a digit from 3 to 6, the '' or the '' character

Bracket expressions

- [0-9] matches any digit (0-9)
- [a-z] matches any lowercase character
- [A-Z] matches any uppercase character
- [.?#] matches either the '.', the '?' or the '#' character
- [a-z3-6.,] matches a lowercase character, a digit from 3 to 6, the '' or the '' character

...

SELECT '27 October 2001'
SIMILAR TO
'[0-9]+ [A-Z][a-z]+ [1-2][0-9][0-9][0-9]'
SELECT *
FROM emails
WHERE subject
SIMILAR TO
'% [0-9]+ [A-Z][a-Z]+ [1-2][0-9][0-9][0-9]%'

Extracting strings with SUBSTRING/3

• The substring function with three parameters, substring(string from pattern for escape-character), provides extraction of a substring that matches an SQL regular expression pattern

Extracting strings with SUBSTRING/3

- The substring function with three parameters, substring(string
 from pattern for escape-character), provides extraction of a
 substring that matches an SQL regular expression pattern
- pattern pattern as with SIMILAR TO

SELECT substring(

```
'An important date was 27th October 2001 since ...' from '% #"[0-9]+th [A-Z][a-z]+ [1-2][0-9][0-9][0-9]#"%' for '#' )
```

SELECT substring(

```
subject
from '% #"[0-9]+th [A-Z][a-z]+ [1-2][0-9][0-9][0-9]#"%'
for '#' ), emails.*
FROM emails
WHERE subject
SIMILAR TO
    '% [0-9]+th [A-Z][a-z]+ [1-2][0-9][0-9][0-9]%'
```

Outline



- Pattern Matching in PostgreSQL
- POSIX Regular Expressions
 - 4 Statistics and aggregation
 - 5 Assignment

Introduction	Pattern Matching in PostgreSQL	POSIX Regular Expressions	Statistics and aggregation	Assignment
Atoms				

(re) (where re is any regular expression) matches a match for re, with the match noted for possible reporting

Atoms

- (re) (where re is any regular expression) matches a match for re, with the match noted for possible reporting
- (?:re) as above, but the match is not noted for reporting (a "non-capturing" set of parentheses)

Atoms

- (re) (where re is any regular expression) matches a match for re, with the match noted for possible reporting
- (?:re) as above, but the match is not noted for reporting (a "non-capturing" set of parentheses)
 - . matches any single character (similar to $_$ in <code>LIKE</code>)

Atoms

- (re) (where re is any regular expression) matches a match for re, with the match noted for possible reporting
- (?:re) as above, but the match is not noted for reporting (a "non-capturing" set of parentheses)
 - . matches any single character (similar to $_$ in <code>LIKE</code>)
- [chars] a bracket expression, matching any one of the chars

 $\label{eq:k} $$ (where k is a non-alphanumeric character) matches that character taken as an ordinary character, e.g., $$ matches a backslash character $$$

Atoms - cont.

- κ (where k is a non-alphanumeric character) matches that character taken as an ordinary character, e.g., κ matches a backslash character
- \c where c is alphanumeric (possibly followed by other characters) is an escape

Atoms - cont.

- $\label{eq:k} $$ (where k is a non-alphanumeric character) matches that character taken as an ordinary character, e.g., $$ ackslash character $$ backslash ch$
- \c where c is alphanumeric (possibly followed by other characters) is an escape
 - { when followed by a character other than a digit, matches the left-brace character {; when followed by a digit, it is the beginning of a bound (see below)

Atoms - cont.

- $\label{eq:k} $$ (where k is a non-alphanumeric character) matches that character taken as an ordinary character, e.g., $$ ackslash character $$ backslash ch$
- \c where c is alphanumeric (possibly followed by other characters) is an escape
 - { when followed by a character other than a digit, matches the left-brace character {; when followed by a digit, it is the beginning of a bound (see below)
 - x where x is a single character with no other significance, matches that character

* a sequence of 0 or more matches of the atom

- \star a sequence of 0 or more matches of the atom
- + a sequence of 1 or more matches of the atom

- \star a sequence of 0 or more matches of the atom
- + a sequence of 1 or more matches of the atom
- ? a sequence of 0 or 1 matches of the atom

- * a sequence of 0 or more matches of the atom
 + a sequence of 1 or more matches of the atom
 ? a sequence of 0 or 1 matches of the atom
- a sequence of 0 or 1 matches of the atom
- $\{m\}\$ a sequence of exactly m matches of the atom

* a sequence of 0 or more matches of the atom
+ a sequence of 1 or more matches of the atom
? a sequence of 0 or 1 matches of the atom
{m} a sequence of exactly m matches of the atom
{m, } a sequence of m or more matches of the atom

- $\star\,$ a sequence of 0 or more matches of the atom
- + a sequence of 1 or more matches of the atom
- ? a sequence of 0 or 1 matches of the atom
- $\{m\}$ a sequence of exactly m matches of the atom
- $\{m_{\text{r}},\}~$ a sequence of m or more matches of the atom
- $\{m, n\}$ a sequence of m through n (inclusive) matches of the atom; m cannot exceed n

Constraints

- matches at the beginning of the string
- \$ matches at the end of the string

~ matches regular expression, case sensitive

- ~ matches regular expression, case sensitive
- ~* matches regular expression, case insensitive

- ~ matches regular expression, case sensitive
- ~* matches regular expression, case insensitive
- $! \sim$ does not match regular expression, case sensitive

- ~ matches regular expression, case sensitive
- ~* matches regular expression, case insensitive
- !~ does not match regular expression, case sensitive
- $! \sim \star \$ does not match regular expression, case insensitive

SELECT

'Never gonna let you down' ~* '^never'

SELECT

'Never gonna run around and desert you' !~ '^gonna'

SELECT

'Never gonna make you cry' ~ 'cry\$'

SELECT

```
'Never gonna say goodbye' ! * 'SAY$'
```

linte	-	Ч		~	1		
II IU	υ		u			וט	

SELECT '27-9-2001' ~ '[0-9]{1,2}\-[0-9]{1,2}\-[1-2][0-9]{3}'
SELECT * FROM emails
WHERE
subject ~ '[0-9]{1,2}\-[0-9]{1,2}\-[1-2][0-9]{3}'

substring(*string* from *pattern*) where pattern is a regular expression extracts the matching substring

substring(*string* from *pattern*) where pattern is a regular expression extracts the matching substring

regexp_replace(source, pattern, replacement) replaces the matching substring with replacement

substring(*string* from *pattern*) where pattern is a regular expression extracts the matching substring

regexp_replace(source, pattern, replacement) replaces the matching substring with replacement

regexp_matches(string, pattern) returns all matching substrings

substring(*string* from *pattern*) where pattern is a regular expression extracts the matching substring

regexp_replace(source, pattern, replacement) replaces the matching substring with replacement

regexp_matches(string, pattern) returns all matching substrings

substring(*string* from *pattern*) where pattern is a regular expression extracts the matching substring

regexp_replace(source, pattern, replacement) replaces the matching substring with replacement

regexp_matches(string, pattern) returns all matching substrings

regexp_split_to_array(string, pattern) same as above, but
 returns array

```
SELECT substring( 'Never gonna tell a lie' from 'gonna(.*)lie'
)
SELECT regexp_replace( '... and hurt you', 'hurt', 'rickroll' )
SELECT regexp_matches( '27 Sep 2001', '([0-9]{1,2})
    ([A-Z][a-z]{2}) ([1-2][0-9]{3})' )
SELECT regexp_split_to_table( 'joza@foi.hr; ivek@foi.hr;
    bara@foi.hr', '\; ' )
```

Outline

Introduction

- Pattern Matching in PostgreSQL
- 3 POSIX Regular Expressions
- 4 Statistics and aggregation
 - 5 Assignment

Aggregate functions

avg finds the average of a numerical attribute
avg finds the average of a numerical attribute sum finds the sum of a numerical attribute

avg finds the average of a numerical attribute sum finds the sum of a numerical attribute min find the smalest value of an attribute

avg finds the average of a numerical attribute sum finds the sum of a numerical attribute min find the smalest value of an attribute max find the highest value of an attribute

avg finds the average of a numerical attribute
 sum finds the sum of a numerical attribute
 min find the smalest value of an attribute
 max find the highest value of an attribute
 count find the number of values of some attribute

avg finds the average of a numerical attribute
sum finds the sum of a numerical attribute
min find the smalest value of an attribute
max find the highest value of an attribute
count find the number of values of some attribute
variance find the variance of some numerical attribute

avg finds the average of a numerical attribute
sum finds the sum of a numerical attribute
min find the smalest value of an attribute
max find the highest value of an attribute
count find the number of values of some attribute
variance find the variance of some numerical attribute
stddev find the standard deviation of some numerical attribute

<pre>SELECT count(*) FROM emails</pre>			
SELECT min(date), FROM emails	<pre>max(date)</pre>		

GROUP BY and HAVING

Aggregate functions are usually used with a GROUP BY clause to group aggregate values by some given attribute. HAVING is used for constraints on aggragate expressions.

```
Introduction
```

```
SELECT date, count(*)
FROM emails
GROUP BY date
ORDER BY 2 DESC
```

```
SELECT custodianname, count(*)
FROM emails
GROUP BY 1
HAVING count(*) > 100
ORDER BY 2
```

Outline

Introduction

- 2) Pattern Matching in PostgreSQL
- 3 POSIX Regular Expressions
- 4 Statistics and aggregation
- 5 Assignment

Project definition

Keyword/pattern definition

Define a set of keywords and data possibly found in email headers (subject, from, to, cc, etc.) which for you as an investigator would be messages of importance. What kind of data is important to an investigative process?

- e.g. keywords: risk, plan, transaction, signature etc.
- e.g. data: email addresses, domains, dates, account numbers, etc.

Keyword extraction

Create queries that will find all the emails sent **to** Enron email addresses (@enron.com) with the keywords and data defined in the previous task.

Project definition - cont.

Mentioned dates

Create queries that will find all emails that mention a date. Extract these dates.

• Think of all possible ways to write a date (e.g. September 27th 2001; 2001-9-27; 27 Sep 2001 etc.)

Communication structure

Try to create a query that will connect emails to their replies (e.g. subject 'Re: Hello' is a reply to subject 'Hello')

- Familiarize your self with the TEXTCAT/2 function which concatenates strings (to construct a pattern)
- Familiarize your self with the LIMIT and OFFSET modifiers in SELECT (to keep computation reasonable)