

Full-Text Search in Django with PostgreSQL



#pycon8 - PyDatabase - Firenze, 2017-04-08

Q Paolo Melchiorre |

- **Computer Science Engineer**
- **Backend Python Developer *since 2006***
- **Django Developer *since 2011***
- **Senior Software Engineer @ TwentyTab**
- **Remote Worker**
- **PostgreSQL user, not a DBA**

Q Goal

To show how we have used Django Full Text Search and PostgreSQL in a real project.

Q Motivation

To implement Full-Text Search using only Django and PostgreSQL functionalities, without resorting to external tools.

Q Agenda

- **Full-Text Search**
- **Existing Solutions**
- **PostgreSQL FTS**
- **Django Support**
- **Concerti@Roma**
- **What's next**
- **Conclusions**
- **Questions**

Q Full-Text Search |

“... Full-Text Search refers to techniques for searching a single computer-stored document or a collection in a full text database ...”*

-- **Wikipedia**

* **FTS = Full-Text Search**

Q Features of a FTS |

- **Stemming**
- **Ranking**
- **Stopwords**
- **Multiple languages support**
- **Accent support**
- **Indexing**
- **Phrase search**

Q Tested Solutions |

Lucene



elastic

Solr 

Q Elasticsearch |

Project: Snap Market (~500,000 mobile users)

Issues:

- Management problems
- Patching a Java plugin

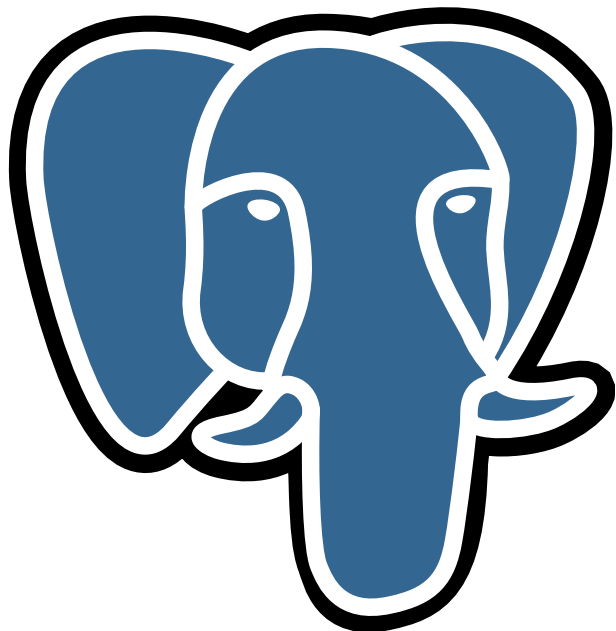
```
@@ -52,7 +52,8 @@ public class DecompondTokenFilter ... {  
-         posIncAtt.setPositionIncrement(0);  
+         if (!subwordsonly)  
+             posIncAtt.setPositionIncrement(0);  
         return true;  
}
```

Q Apache Solr

Project: GoalScout (~25,000 videos)

Issues:

- Synchronization problems
- All writes to PostgreSQL and reads from Solr



Q Existing Solutions |

PROS

- **Full featured solutions**
- **Resources** (*documentations, articles, ...*)

CONS

- **Synchronization**
- **Mandatory use of driver** (*haystack, bungiesearch...*)
- **Ops oriented: focus on system integrations**

Q FTS in PostgreSQL |

- **FTS Support *since* version 8.3**
- **TSVECTOR to represent text data**
- **TSQUERY to represent search predicates**
- **Special Indexes (GIN, GIST)**
- **Phrase Search *since* version 9.6**

Q What are documents |

“... a document is the unit of searching in a full-text search system; for example, a magazine article or email message ...”

-- PostgreSQL documentation

Q Django Support |

- Module `django.contrib.postgres`
- FTS Support *since version 1.10*
- BRIN and GIN indexes *since version 1.11*
- Dev oriented: focus on programming

Q Making queries |

```
class Blog(models.Model):
    name = models.CharField(max_length=100)
    tagline = models.TextField()

class Author(models.Model):
    name = models.CharField(max_length=200)
    email = models.EmailField()

class Entry(models.Model):
    blog = models.ForeignKey(Blog)
    headline = models.CharField(max_length=255)
    body_text = models.TextField()
    pub_date = models.DateField()
    authors = models.ManyToManyField(Author)
```


Q Standard queries

```
>>> Author.objects.filter(name__contains='Terry')  
[<Author: Terry Gilliam>, <Author: Terry Jones>]
```

```
>>> Author.objects.filter(name__icontains='Erry')  
[<Author: Terry Gilliam>, <Author: Terry Jones>,  
<Author: Jerry Lewis>]
```

Q Unaccented query |

```
>>> Author.objects.filter(name__unaccent__icontains='Helen')  
[<Author: Helen Mirren>, <Author: Helena Bonham Carter>,  
<Author: H  l  ne Joy>]
```

Q Trigram similar |

```
>>> Author.objects.filter(
    name__unaccent__lower__trigram_similar='Hélène')
[<Author: Helen Mirren>, <Author: Helena Bonham Carter>, <Author:
Hélène Joy>]
```

Q The search lookup |

```
>>> Entry.objects.filter(body_text__search='Cheese')  
[<Entry: Cheese on Toast recipes>, <Entry: Pizza Recipes>]
```

Q SearchVector

```
>>> from django.contrib.postgres.search import SearchVector
>>> Entry.objects.annotate(
...     search=SearchVector('body_text', 'blog__tagline'),
... ).filter(search='Cheese')
[<Entry: Cheese on Toast recipes>, <Entry: Pizza Recipes>]
```

Q SearchQuery |

```
>>> from django.contrib.postgres.search import SearchQuery
>>> SearchQuery('potato') & SearchQuery('ireland')
# potato AND ireland
>>> SearchQuery('potato') | SearchQuery('penguin')
# potato OR penguin
>>> ~SearchQuery('sausage')
# NOT sausage
```

Q SearchRank

```
>>> from django.contrib.postgres.search import (  
    SearchQuery, SearchRank, SearchVector)  
>>> vector = SearchVector('body_text')  
>>> query = SearchQuery('cheese')  
>>> Entry.objects.annotate(  
    rank=SearchRank(vector, query)).order_by('-rank')  
[<Entry: Cheese on Toast recipes>, <Entry: Pizza recipes>]
```

Q Search configuration |

```
>>> from django.contrib.postgres.search import (
    SearchQuery, SearchVector)
>>> Entry.objects.annotate(
...     search=SearchVector('body_text', config='french'),
... ).filter(search=SearchQuery('œuf', config='french'))
[<Entry: Pain perdu>]

>>> from django.db.models import F
>>> Entry.objects.annotate(
...     search=SearchVector('body_text', config=F('blog__lang')),
... ).filter(search=SearchQuery('œuf', config=F('blog__lang')))
[<Entry: Pain perdu>]
```


Q Weighting queries |

```
>>> from django.contrib.postgres.search import (  
        SearchQuery, SearchRank, SearchVector)  
>>> vector = SearchVector('body_text', weight='A') +  
        SearchVector('blog__tagline', weight='B')  
>>> query = SearchQuery('cheese')  
>>> Entry.objects.annotate(rank=SearchRank(  
        vector, query)).filter(rank__gte=0.3).order_by('rank')
```

Q SearchVectorField |

```
>>> Entry.objects.update(search_vector=SearchVector('body_text'))  
>>> Entry.objects.filter(search_vector='cheese')  
[<Entry: Cheese on Toast recipes>, <Entry: Pizza recipes>]
```

Q Concerti@Roma |

The numbers of the project:

- > 11,000 bands
- ~ 1,000 venues
- ~ 1,500 shows
- > 100 festivals
- ~ 15,000 user/month

Q Concerti@Roma 2 |

CONCERTI A ROMA

CONCERTI VENUES BANDS

7-4-2017

Electro Hip Hop Jazz Metal Pop Punk Reggae Rock

Clear

Search

LOG IN

VENERDÌ 07 APRILE

Rock Pop	Le Luci della Centrale Elettrica + Flavio Giurato + Colombre @ Atlantico	👍
Rock	Peter Hook & the Light @ Teatro Quirinetta	👍
Rock	Area 765 @ Mentelocale 2.0	👍
Rock Pop	So Does your Mother + The Shiny Moss @ Contestaccio	👍
Rock	Boddah (Nirvana Tribute) + Clamôr + Closed Speech + Nessuno (acoustic) + Sara Pik (acoustic) + Pocaroba (acoustic) @ Defrag	👍
Hip Hop	Jack the Smoker @ Zoobar	👍
Pop	Carlo Palermo (release party) + guest @ Magazzino 33	Cancelled

Tweet di @ConcertiaRoma

Concerti a Roma @ConcertiaRoma

CONCERTI DAL 07/04 al 09/04 su www.concertiaroma.com

Per [maggiore](#), date scrivete per messaggio privato su <https://twitter.com/concertiaroma>

Tutti i concerti di Roma - conc...
Tutti i concerti a Roma e provincia.
concertiaroma.com

106.6
RADIO ROCK

▪ Python2.7 / Django 1.7 / PostgreSQL 9.1 / Like

Q Concerti@Roma 3 |

The screenshot shows the website's header with navigation links: COLLABORA CON NOI, REGISTRATI, ACCEDI, and CREA UN EVENTO. Below the header is a search bar labeled 'CONCERTI A ROMA' with the subtitle 'Cerca i concerti di oggi nella Capitale' and a search input field 'Cerca nel sito'. A date filter is set to '09/29/2016'. Genre filters include JAZZ, ELECTRO, HIP HOP, METAL, POP, PUNK, REGGAE, ROCK, and a CLEAR button. The main content area is titled 'Concerti Di Oggi A Roma' and lists four events:

- Tedio + Aharo**: Gio 6 Apr, Na Wishlist, Inizio concerti: 23:00:00. Genres: ELECTRO, ROCK. Button: Scheda Concerto →. Notification: Ricevi notifiche su questo concerto.
- Lamorivostr**: Gio 6 Apr, Na Cosetta, Inizio concerti: 22:00:00. Genre: POP. Button: Scheda Concerto →. Notification: Ricevi notifiche su questo concerto.
- Marina Re**: Gio 6 Apr, Auditorium (Parco della Musica), Inizio concerti: 21:00:00. Genre: POP. Button: Scheda Concerto →. Notification: Ricevi notifiche su questo concerto.
- The Real Kids + Transe**: Gio 6 Apr. Genres: PUNK, ROCK. Button: Scheda Concerto →. Notification: Ricevi notifiche su questo concerto.

On the right side, there is a poster for a performance by 'DESTRAGE' at 'SABATO TRAFFICO LIVE CLUB'. The poster features a child's face and text: 'NO SON MUSIC, TIME TO KILL RECORDS, EROCKS PRODUCTION PRESENTA DESTRAGE', 'SUBLIMINAL FEAR', 'ONNET', and 'SABATO TRAFFICO LIVE CLUB'.

▪ Python3.6 / Django1.11 / PostgreSQL 9.6 / FTS

Q C@R Manager

```
class BandManager(models.Manager):  
    def search(self, text):  
        vector = (  
            SearchVector('nickname', weight='A', config=LANG) +  
            SearchVector('genres__name', weight='B', config=LANG)+  
            SearchVector('description', weight='D', config=LANG)  
        )  
        query = SearchQuery(text, config=LANG)  
        rate = SearchRank(vector, query)  
        return self.get_queryset().annotate(rate=rate).filter(  
            search=query).annotate(search=vector).distinct(  
                'id', 'rate').order_by('-rate', 'id')
```

Q C@R Test Setup |

```
class BandTest(TestCase):  
    def setUp(self):  
        metal, _ = Genre.objects.get_or_create(name='Metal')  
        doom, _ = Genre.objects.get_or_create(name='Doom')  
        doomraiser, _ = Contact.objects.get_or_create(  
            nickname='Doom raiser', description='Lorem...')  
        doomraiser.genres.add(doom)  
        forgotten_tomb, _ = Contact.objects.get_or_create(  
            nickname='Forgotten Tomb', description='Lorem...')  
        forgotten_tomb.genres.add(doom)  
        ....
```

Q C@R Test Method

```
class BandTest(TestCase):  
    def setUp(self):  
        ...  
  
    def test_band_search(self):  
        band_queryset = Band.objects.search(  
            'doom').values_list('nickname', 'rate')  
        band_list = [  
            ('Doom raiser', 0.675475),  
            ('The Foreshadowin', 0.258369),  
            ('Forgotten Tomb', 0.243171)]  
        self.assertEqual(list(  
            OrderedDict(band_queryset.items()), band_list))
```


Q What's next

- **Multiple language ranking**
- **Search suggestions**
- **SearchVectorField with triggers**
- **JSON/JSONB Full-Text Search**
- **RUM indexing**

Q Conclusions

Conditions to implement this solution:

- No extra dependencies
- Not too complex searches
- Easy management
- No need to synchronize data
- PostgreSQL already in your stack
- Python-only environment

Q Resources

- <https://docs.djangoproject.com/en/1.11/ref/contrib/postgres/search/>
- <https://www.postgresql.org/docs/9.6/static/textsearch.html>
- <https://github.com/damoti/django-tsvector-field>
- https://en.wikipedia.org/wiki/Full-text_search

Q Acknowledgements |

- 20tab



twentytab

- Marc Tamlyn `django.contrib.postgres`

Q Thank you

 BY - SA (Attribution-ShareAlike)

- <https://creativecommons.org/licenses/by-sa/4.0/>

 Slides

- <https://speakerdeck.com/pauloxnet>



twentytab

Q Social Time



PYCON OTTO





twentytab

Q Questions ?



Q Contacts



<https://twitter.com/pauloxnet>



<https://linkedin.com/in/paolomelchiorre>



<https://github.com/pauloxnet>