## Discovering python search engine

José Manuel Ortega - Pycones 2017

## Agenda

- Introduction to search engines
- ElasticSearch, whoosh, django-hystack
- ElasticSearch example
- Other solutions & tools
- Conclusions

## Search engines

## Search engines

- Document based
- A document is the unit of searching in a full text search system.
- A document can be a json or python dictionary

## Core concepts

## Core concepts

- Index: Named collection of documents that have similar characteristics(like a database)
- **Type:**Logical partition of an index that contains documents with common fields(like a **table**)
- **Document:**basic unit of information(like a **row**)
- **Mapping:**field properties(datatype,token extraction).

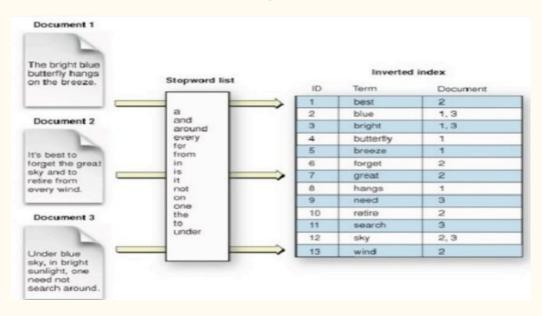
  Includes information about how fields are stored in the index

## Core concepts

- Relevance are the algorithms used to rank the results based on the query
- Corpus is the collection of all documents in the index
- Segments: Sharded data storing the inverted index. Allow searching in the index in a efficient way

#### Inverted index

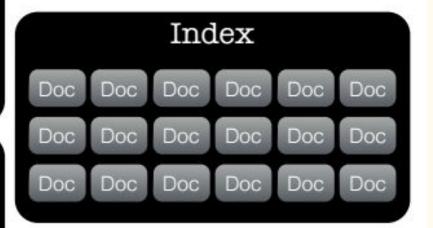
- Is the heart of the search engine
- Each inverted index stores position and document IDs



## Cluster

## Index Doc Doc

# Index Doc Doc















## ElasticSearch



Search & Analyze Data in Real Time







stackoverflow

- Open source search server based on Apache Lucene
- Written in Java
- Cross-platform
- Communications with the search server is done through HTTP REST API
- curl -X<GET|POST|PUT|DELETE>
   http://localthost:9200/<index>/<type\_document>/id

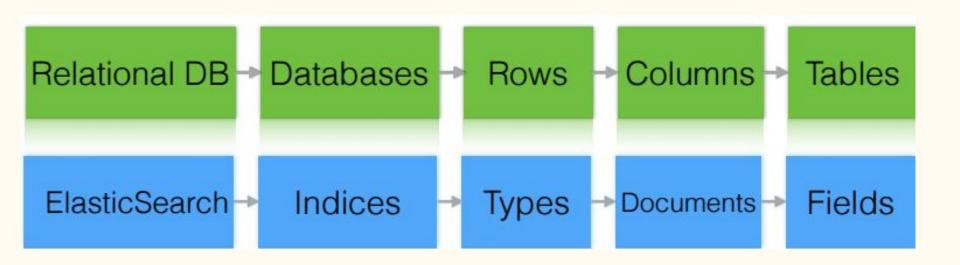
- You can add a document without creating an index
- ElasticSearch will create the index,mapping type and fields automatically
- ElasticSearch will infer the data types based on the document's data

```
"took": 4.
                                                  Information about the
        "timed_out": false,
                                                 execution of the request.
        "_shards": {
           "total": 5,
           "successful": 5,
                                                 Object with information about the search
           "failed": 0
                                                    results, including the actual results.
        "hits": {
10
                                                     Total number of documents that
           "total": 2, <
                                                            match the query.
           "max_score": 0.095891505.
11
12
           "hits": [
13
                                                       Array with search hits.
14
                 "_index": "movies",
15
                 "_type": "movie",
16
                 "_{d": "5",
                                                               Meta data about the hit.
17
                 "_score": 0.095891505.
18
                 "_source": { -
                    "title": "Kill Bill: Vol. 1",
19
                                                             The document that produced the hit.
                    "director": "Quentin Tarantino",
20
```

#### Metadata Fields

- Each document has metadata associated with it
- \_index:Allows matching documents based on their indexes.
- \_type:Type of the document
- \_id:Document id(not indexed)
- \_uid:\_type + \_id(indexed)
- \_source:contains the json passed in creation time of the index or document(not indexed)
- version

#### ElasticSearch vs Relational DB



Port to connect to. Protocol used, HTTP is Elasticsearch listens supported out-of-the-box. to 9200 by default. Type name http://localhost:9200/get-together/group/1 Hostname of the Index name Document ID Elasticsearch node to connect to. Use localhost if Elasticsearch is on the local machine.

## Creating an Index

```
curl -XPUT 'localhost:9200/myindex'-d {
    "settings":{..}
    "mappings":{..}
```

```
curl -XPUT http://localhost:9200/blog -d '{
    "settings" : {
        "index" : {
            "number_of_shards" : 1,
            "number_of_replicas" : 0
    "mapping": {
        "posts": {
            "properties": {
                 "title": "string",
                 "published": "date"
```

```
curl -XPOST http://localhost:9200/blog/posts -d '{
    "title": "Discover python search engine",
    "published": "2017-09-22"
curl -XPOST http://localhost:9200/blog/posts/_search?
pretty= true -d '{
    "query": {
        "match": {
            "title": "python"
```

```
import json
import requests
url = 'http://localhost:9200/pycones/talk/1'
document = {
  'title': "Discovering python search engine",
  'description' : "Full text search"
requests.post(url, data=json.dumps(document))
```

## Searching a document

```
curl -XGET 'localhost:9200/myindex/mydocument/_search?q=elasticSearch'
curl -XGET 'localhost:9200/myindex/mydocument/_search?pretty' -d{
   "query":{
                                     Query DSL
        "match":{
           " all":"elasticSearch"
```

```
import requests
import json
def search(uri, term):
    """Simple Elasticsearch Query"""
    query = json.dumps({
        "query": {
            "match": {
                "title": term
    response = requests.get(uri, data=query)
    results = json.loads(response.text)
    return results
```

## Searching a document

- Search can get much more complex
  - Multiple terms
  - Multi-match(math query on specific fields)
  - Bool(true,false)
  - Range
  - RegExp
  - GeoPoint,GeoShapes

### ElasticSearch python client

- The official low-level client is elasticsearch-py
  - o pip install elasticsearch

```
Collecting elasticsearch

Downloading elasticsearch-5.4.0-py2.py3-none-any.whl (58kB)

100% | 61kB 1.5MB/s

Collecting urllib3<2.0,>=1.8 (from elasticsearch)

Downloading urllib3-1.22-py2.py3-none-any.whl (132kB)

100% | 133kB 1.7MB/s

Installing collected packages: urllib3, elasticsearch

Successfully installed elasticsearch-5.4.0 urllib3-1.22
```

## ElasticSearch-py API

```
from elasticsearch import Elasticsearch
es= Elasticsearch()
document = {"title":"elasticSearch","descripion":"python"}
##create inded if doesnt already exist
response = es.index(index="myindex",doc type="mydocument",id=1,body=document)
print(response)
##search for "python" in all fields in all documents
query = {'query':{'match':{' all':'python'}}}
response search = es.search(index="myindex",doc type="mydocument",body=query)
print(response search)
```

## ElasticSearch-py API

```
{'_index': 'myindex', '_type': 'mydocument', '_id': 'l', '_version': 9, 'result': 'updated', '_sh
ards': {'total': 2, 'successful': 1, 'failed': 0}, 'created': False}
{'took': 4, 'timed_out': False, '_shards': {'total': 5, 'successful': 5, 'failed': 0}, 'hits': {'
total': 2, 'max_score': 0.25811607, 'hits': [{'_index': 'myindex', '_type': 'mydocument', '_id':
'2', '_score': 0.25811607, '_source': {'title': 'elasticSearch', 'descripion': 'python'}}, {'_index': 'myindex', '_type': 'mydocument', '_id': 'l', '_score': 0.25811607, '_source': {'title': 'elasticSearch', 'descripion': 'python'}}}}
```

```
import csv
from elasticsearch import ElasticSearch
#by default connect with localhost:9200
es = ElasticSearch()
#create an index in es,ignore status code 400
es.indices.create(index='myindex',ignore=400)
with open('data.csv') as file:
  reader = csv.DictReader(file)
  for line in reader:
    es.index(index='myindex',doc_type,body=line)
```

```
def check_index_existence(index_name):
   Function for checking the existence of an index
    :param index_name: Name of the index to be created
    :param doc_type: Name of document type to be created
    :return: False if index does not exist else True
   if not es.indices.exists(index=index_name):
       return False
   else:
       return True
```

```
def index_document(index_name, doc_type, doc_id):
   Function for index a document
    :param index_name: Name of the index
    :param doc_type: Name of document type
    :param doc_id: Unique id to be set for the document
    :return:
   doc1 = {
            'name': 'Elasticsearch Essentials',
            'category': ['Big Data', 'search engines',
              'Analytics']
   es.index(index=index_name, doc_type=doc_type, body=doc1,
     id=doc_id)
```

```
def get_document(index_name, doc_type, doc_id):
   Function for fetching a single document
    :param index_name: Name of the index
    :param doc_type: Name of document type
    :param doc_id: Document id to be retrieved
    :return:
    response = es.get(index=index_name, doc_type=doc_type,
     id=doc_id, ignore=404)
   print response
   print response.get('_source').get('category')
```

```
def check_document_existence(index_name, doc_type, doc_id):
    Function for checking if a document exists ot not in the
      given index
    :param index_name: Name of the index
    :param doc_type: Name of document type
    :param doc_id: Document id whose existence need to be
      checked
    :return:
    if not es.exists(index=index_name, doc_type=doc_type, id
      =doc_id):
        return False
    else:
        return True
```

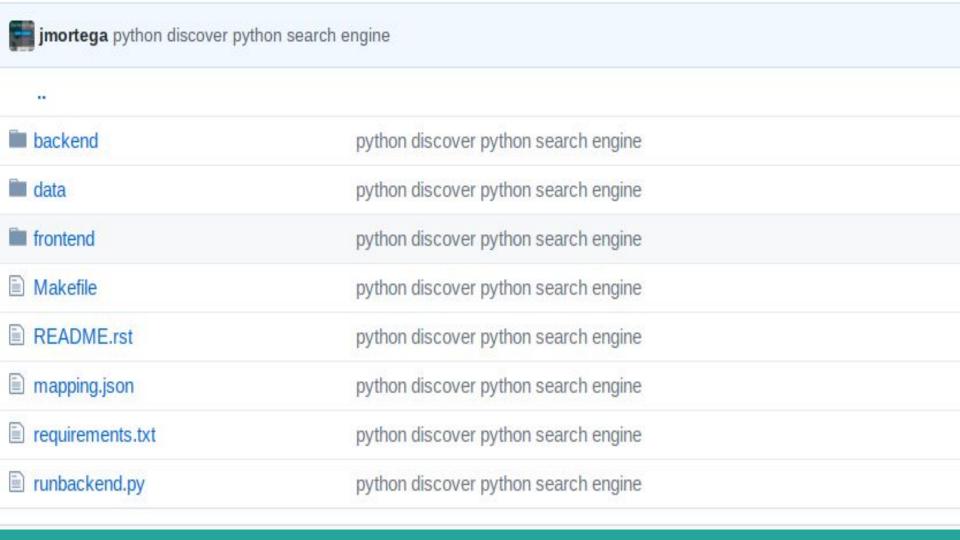
```
def search_documents(index_name, doc_type, query=None):
   Function for performing search using Query-DSL queries
   :param index_name: Name of the index
   :param doc_type: Name of document type
   :param query: Query in json format
   if query is None:
       query = {
        "query":{
          "match":{"text":"Elasticsearch"}
   response = es.search(index=index_name, doc_type=doc_type
     , body=query, size=2, request_timeout=20)
   for hit in response['hits']['hits']:
       print json.dumps(hit.get('_source'))
```

## Geo queries

- Elastic search supports two types of geo fields
  - geo\_point(lat,lon)
  - geo\_shapes(points,lines,polygons)
- Perform geographical searches
  - Finding points of interest and GPS coordinates

```
def create_index_with_geopoint_mapping(index_name,doc_type):
    Function to create index with nested mapping
    :param index_name: Name of the index to be created
    :param doc_type: Name of document type to be created
    doc_mapping = {
      "properties": {
       "location": {
            "type": "geo_point"
    body = dict()
    mapping = dict()
    mapping[doc_type] = doc_mapping
    body['mappings'] = mapping
    es.indices.create(index=index_name, body = body)
```

```
def find_by_distance(index_name, doc_type):
   query = {
      "query": {
          "filter": {
            "geo_distance": {
              "distance": "12km",
              "location": {
                "lat": 28.67,
                "lon": 77.42
   response = es.search(index=index_name, doc_type=doc_type
       body=query)
```





### List of Talks

Title	Author	
Modelo de Redes Neuronales Recurrentes para el Análisis y Predicción de Series Temporales	J.C. González-Avella y J. M. Tuduri	Delete
Discovering python search engine	Jose Manuel Ortega	Delete
Usando Python y SocketlO para aplicaciones en tiempo real	Oscar Ramírez	Delete
High-impact refactors while keeping the lights on	Diego Muñoz	Delete
Mi coche ya es mayor de edad, creo que puede conducir él solito.	Ricardo Guerrero Gómez-Olmedo	Delete
City attractiveness seen through Twitter	Antònia Tugores	Delete
Factories, what the hell?	Miguel Jiménez	Delete
Machine Learning en producción	Luis Mesas	Delete





## Talks Home

machine learning

Title

Machine Learning en producción

Insert New Talk Full Talk List

Author

Luis Mesas



# Machine Learning en producción

Author:	Luis Mesas
Туре:	Talk
Description:	Una vez tenemos el modelo diseñado por los data scientist hay que ponerlo en pro, pero ninguna de las herramientas que han usado son viables en producción y tienen problemas de escalabilidad ¿Que hacemos ahora?. En esta charla veremos arquitecturas aplicadas a 2 casos de uso típicos: arquitectura de ejecución única reutilizable y arquitectura de uso contínuo con bucle de feedback.

Go to Search

See All Talks

Insert New Talk

#### Insert New Talk

	- 1			
0.7		w	tit	
		IIN.	LILL	

Machine learning y datascience con python con sklearn y pyspark

#### Author:

author

#### Type:

**✓** Talk **✓** Workshop

#### Talk Description:

description



Cancel

Create basic database (use data from data/pycones2017.data):

make index

Run backend service:

make backend

Run frontend service:

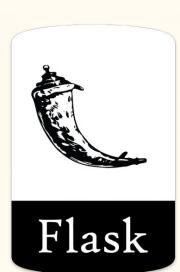
make frontend

Point your browser to:

http://localhost:8000

and search for a talk.





```
def get(self):
    print("Call for: GET /talks")
    url = config.es_base_url['talks']+'/_search'
    query = {
        "query": {
            "match_all": {}
        "size": 100
    resp = requests.post(url, data=json.dumps(query))
    data = resp.json()
    talks = []
    for hit in data['hits']['hits']:
        talk = hit[' source']
        talk['id'] = hit[' id']
        talks.append(talk)
    return talks
```

# Whoosh

- Pure-python full-text indexing and searching library
- Library of classes and functions for indexing text and then searching the index.
- It allows you to develop custom search engines for your content.
- Mainly focused on index and search definition using schemas
- Python 2.5 and Python 3

## Schema

```
class whooshSCHEMA(fields.SchemaClass):
  title = fields.TEXT(stored=True, sortable=True)
  content = fields.TEXT(spelling=True)
  date = fields.DATETIME(stored=True)
  summary = fields.STORED
  url=fields.ID(stored=True, unique=True))
WHOOSH\_SCHEMA = fields.Schema(
  title=fields.TEXT(stored=True,sortable=True),
  content=fields.TEXT(spelling=True),
  date = fields.DATETIME(stored=True),
  summary = fields.STORED,
  url=fields.ID(stored=True, unique=True))
```

# Create index and insert document

```
#To create an index basically you need a writer object
ix = index.create in("index",schema=WHOOSH SCHEMA)
writer = ix.writer()
writer.add document(title="pycones 2017",content="python conference",
  date = datetime(2017, 9, 22),
  summary = "discovering python search engine",
  url="http://pycones.es")
writer.add document(title="python 2017",content="pycones2017",
  date = datetime(2017, 9, 22),
  summary = "discovering python search engine",
  url="http://pycones.es")
  iter.commit()
```

# Searching single field

```
#searching in the index by a single field
from whoosh import qparser

queryParser = qparser.QueryParser("content",schema = ix.schema)
query = queryParser.parse("python")
with ix.searcher() as searcher:
    results = searcher.search(query)
    print(results)
    for result in results:
        print(result)
```

```
<Top 1 Results for Term('content', 'python') runtime=0.000558944000658812>
<Hit {'date': datetime.datetime(2017, 9, 22, 0, 0), 'summary': 'discovering py
thon search engine', 'title': 'pycones 2017', 'url': 'http://pycones.es'}>
```

# Searching multiple field

```
<Top 2 Results for Or([Term('title', 'python'), Term('content', 'python')]) ru
ntime=0.000842152001496288>
<Hit {'date': datetime.datetime(2017, 9, 22, 0, 0), 'summary': 'discovering py
thon search engine', 'title': 'python 2017', 'url': 'http://pycones.es'}>
<Hit {'date': datetime.datetime(2017, 9, 22, 0, 0), 'summary': 'discovering py
thon search engine', 'title': 'pycones 2017', 'url': 'http://pycones.es'}>
```

# Django-haystack



More Like This

Faceting

**Spatial Search** 

Stored (non-indexed) fields

Highlighting

**Spelling Suggestions** 

Boost

#### Find the needle you're looking for.



Download



Documentation

Search doesn't have to be hard. Haystack lets you write your search code once and choose the search engine you want it to run on. With a familiar API that should make any Djangonaut feel right at home and an architecture that allows you to swap things in and out as you need to, it's how search ought to be.

Haystack is BSD licensed, plays nicely with third-party apps without needing to modify the source and supports Solr, Elasticsearch, Whoosh and Xapian.

#### Haystack 2.6.0!

Posted on 2017/01/04

Changelog is available here:

http://django-

haystack.readthedocs.io/en/v2.6.0/changelog.html.

#### Haystack 2.0.0!

Posted on 2013/05/12 by Daniel

After two (far too long) years,

- Multiple backends (you have a Solr & a Whoosh index, or a master Solr & a slave Solr, etc.)
- An Elasticsearch backend
- Big query improvements
- Geospatial search (Solr & Elasticsearch only)
- The addition of Signal Processors for better control
- Input types for improved control over queries
- Rich Content Extraction in Solr

Branch: master ▼ django-hayst	ack / haystack / backends /	Create new file	Find file	History
acdha Whoosh: prevent more_like	e_this from hitting an uninitialized variable	Latest commit	7d42273 7	days ago
•				
initpy	PEP-8		4 mo	nths ago
elasticsearch2_backend.py	Refactoring the code in pull request #1336 . This pull request is to		9 mo	nths ago
elasticsearch_backend.py	fix: contains filter, add endswith filter		a	year ago
simple_backend.py	Avoid unsafe default value on backend clear() methods	unsafe default value on backend clear() methods 2 years ago		
solr_backend.py	Merge branch '1504-solr-6-by-default'	4 months ago		
whoosh_backend.py	Whoosh: prevent more_like_this from hitting an uninitialized variable	zed variable 7 days ago		

```
# Add Haystack to INSTALLED_APPS
INSTALLED_APPS = settings.INSTALLED_APPS + (
    'haystack',
HAYSTACK_CONNECTIONS = {
    'default': {
        'ENGINE': 'haystack.backends.solr_backend
          .SolrEngine',
        'URL': 'http://localhost:9001/solr/example'
    'elasticsearch': {
        'ENGINE': 'haystack.backends.elasticsearch_backend
          .ElasticsearchSearchEngine',
        'URL': 'http://localhost:9200',
        'INDEX_NAME': 'example_project'
    'whoosh': {
        'ENGINE': 'haystack.backends.whoosh_backend
```

```
from haystack import indexes
from haystack import site
from blog.models import BlogEntry
class BlogEntryIndex(indexes.SearchIndex,
indexes.Indexable):
    text = indexes.CharField(document=True, use_template
      =True)
    author = indexes.CharField(model_attr='author')
    pub_date = indexes.DateField(model_attr='pub_date')
   def get_queryset(self):
    return BlogEntry.objects.filter(pub_date= datetime
      .datetime.now(), status=BlogEntry.LIVE_STATUS)
site.register(BlogEntry, BlogEntryIndex)
```

import datetime

```
from haystack import indexes
from haystack import site
from blog.models import BlogEntry
class BlogEntryIndex(RealTimeSearchIndex):
    text = indexes.CharField(document=True, use_template
      =True)
    author = indexes.CharField(model_attr='author')
    pub_date = indexes.DateField(model_attr='pub_date')
  def get_queryset(self):
    return BlogEntry.objects.filter(pub_date= datetime
      .datetime.now(), status=BlogEntry.LIVE_STATUS)
site.register(BlogEntry, BlogEntryIndex)
```

import datetime

```
#Create search_sites.py inside my Django site, at the same
import haystack
haystack.autodiscover()
from haystack.query import SearchQuerySet
def search entries(request):
    entries = SearchQuerySet.autocomplete
    (content auto=request.POST.get('search text',''))
    return render to response('ajax search.html',{'entries'.entries})
```

- Create the index
  - Run ./manage.py rebuild\_index to create the new search index.
- Update the index
  - ./manage.py update\_index will add new entries to the index.
  - ./manage.py rebuild\_index will recreate the index from scratch.

- Pros:
  - Easy to setup
  - Looks like Django ORM but for searches
  - Search engine independent
  - Support 4 engines (Elastic, Solr, Xapian, Whoosh)
- Cons:
  - Poor SearchQuerySet API
  - Difficult to manage stop words
  - Loose performance, because extra layer
  - Django Model based

# Other solutions

## Other solutions

Xapian

- https://xapian.org
- https://docs.djangoproject.com/en/1.11/ref/contrib/pos tgres/search/
- https://www.postgresql.org/docs/9.6/static/textsearch.
   html

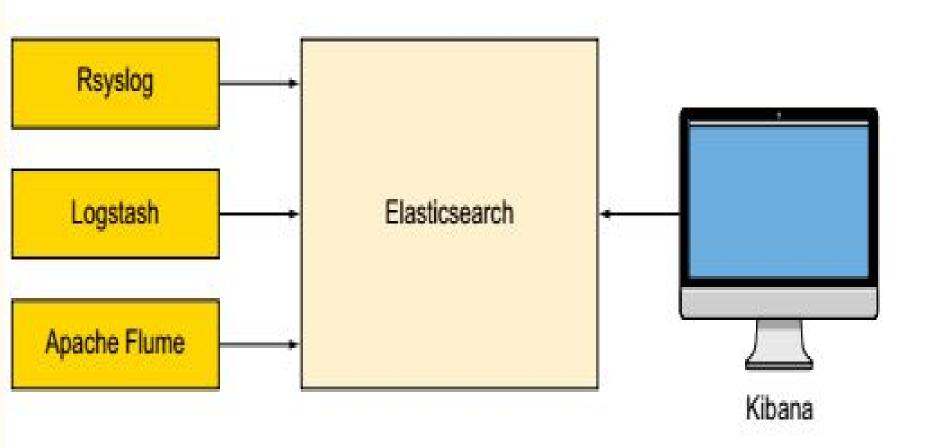


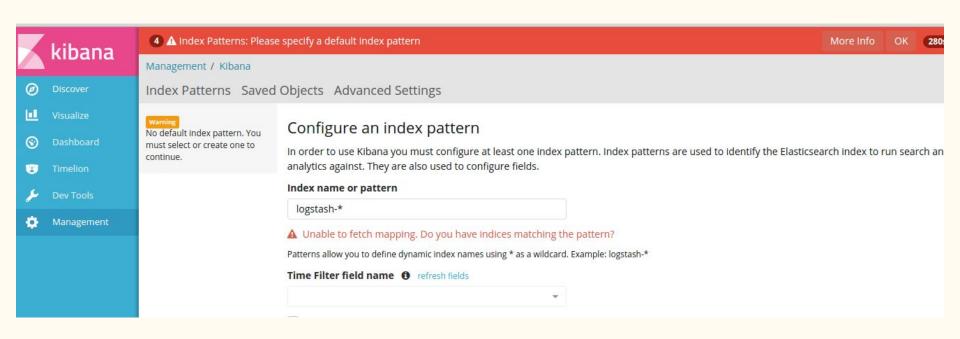


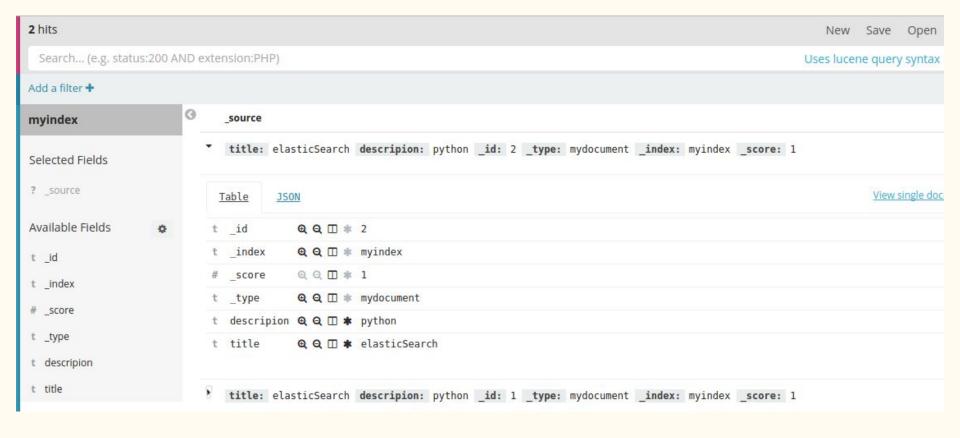
# pysolr

```
import pysolr
# Setup a Solr instance. The timeout is optional.
solr = pysolr.Solr('http://localhost:8983/solr/', timeout=10)
# How you'd index data.
solr.add([
        "id": "doc 1",
        "title": "A test document",
        " doc": [
            { "id": "doc 1 1", "title": "document 1 1" }
results = solr.search('document')
```

# Other tools

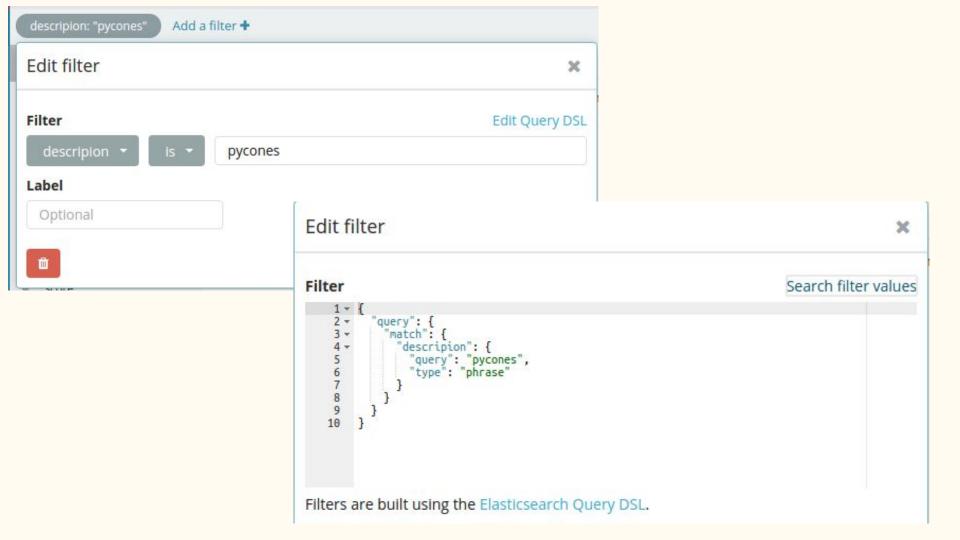






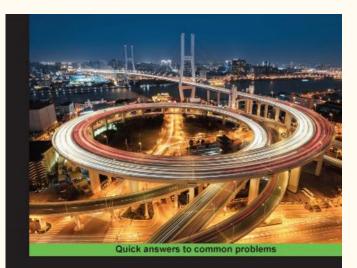
#### Console

```
GET _search
                                                                        1 - {
                                                                              "took": 10.
 2 -
       "query": {
                                                                              "timed_out": false,
                                                                        4 +
                                                                              "_shards": {
         "multi_match":
 4
                                                                                "total": 6.
 6
               "query": "pycones",
                                                                                "successful": 6.
               "fields":["title", "descripion"]
                                                                                "failed": 0
                                                                        8 -
 8 4
 9
                                                                        9 +
                                                                              "hits": {
                                                                                "total": 2,
10 - }
                                                                       10
11 4 }
                                                                       11
                                                                                "max_score": 0.25811607,
                                                                       12 -
                                                                                "hits":
                                                                       13 -
                                                                                      _index": "myindex",
                                                                       14
                                                                       15
                                                                                      _type": "mydocument",
                                                                                    "_id": "2",
                                                                       16
                                                                       17
                                                                                      score": 0.25811607.
                                                                       18 -
                                                                                      source": {
                                                                       19
                                                                                      "title": "pycones 2017",
                                                                       20
                                                                                      "descripion": "discoverting python search engine"
                                                                       21 4
                                                                       22 4
                                                                                  },
                                                                       23 -
                                                                       24
                                                                                      _index": "myindex",
                                                                       25
                                                                                      type": "mydocument",
                                                                       26
                                                                                      id": "1",
                                                                       27
                                                                                      score": 0.25811607.
                                                                       28 -
                                                                                      _source": {
                                                                       29
                                                                                      "title": "python conference",
                                                                       30
                                                                                      "descripion": "pycones 2017"
                                                                       31 4
                                                                       32 -
                                                                       33 *
                                                                       34 -
                                                                       35 4 }
```



## References

- http://elasticsearch-py.readthedocs.io/en/master/
- https://whoosh.readthedocs.io/en/latest
- http://django-haystack.readthedocs.io/en/master/
- http://solr-vs-elasticsearch.com/
- https://wiki.apache.org/solr/SolPython
- https://github.com/django-haystack/pysolr



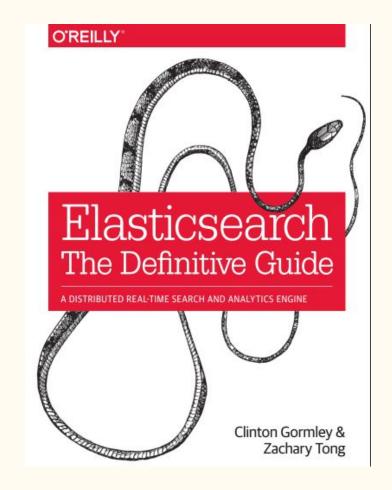
#### ElasticSearch Cookbook

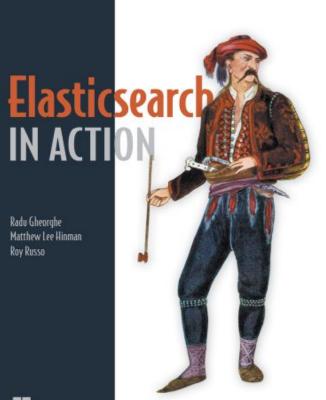
Second Edition

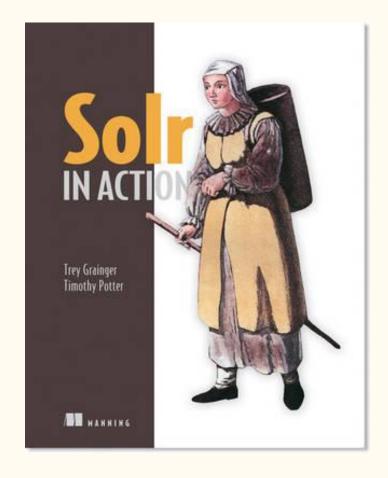
Over 130 advanced recipes to search, analyze, deploy, manage, and monitor data effectively with ElasticSearch

Alberto Paro









ME HANNING

Roy Russo

# Thanks!

# jmortega.github.io

@jmortegac

