

# Tom 2 Wiki Topic ...

%pyspark

FINISHED

```
# PySpark CommonCrawl Topic Modelling
# Tom V / Paul J - 13/2/2018

# SET THE spark.driver.maxResultSize PROPERTY TO 16g

import boto
from boto.s3.key import Key
from gzipstream import GzipStreamFile
from pyspark.sql.types import *
import warc
import ujson as json
from urlparse import urlparse
from langdetect import detect_langs
import pyclld2 as cld2

#wetlist = sc.textFile("s3://commoncrawl/crawl-data/CC-MAIN-2017-04/wet.paths.gz") #
# April(?) 2017 Crawl - 57000 files
# Latest blog/documentation: http://commoncrawl.org/2017/10/october-2017-crawl-archive
# -now-available/
wetlist = sc.textFile("s3://commoncrawl/crawl-data/CC-MAIN-2017-43/wet.paths.gz") #
# October 2017 Crawl - 89100 files (10.58TB)

wetlist.cache()

def unpack(uri):
    conn = boto.connect_s3(anon=True, host='s3.amazonaws.com')
    bucket = conn.get_bucket('commoncrawl')
    key_ = Key(bucket, uri)
    file_ = warc.WARCFile(fileobj=GzipStreamFile(key_))
    return file_

def detect(x):
    try:
        return detect_langs(x)[0].lang # Maybe we can get away with looking at less
        characters, or do something less expensive?
    except Exception as e:
        return None

def detect2(x):
    try:
        isReliable, textBytesFound, details = cld2.detect(x)
        return details[0][1]
    except Exception as e:
        print(e)
        return None

def process_wet(id_, iterator):
    for uri in iterator:
```

```

file = unpack(uri)
for record in file: # Approx 53k web pages per WET file
    try:
        #url = record.rec_headers.get_header('WARC-Target-URI')
        #yield record, record.content_stream().read().decode('utf-8')
        url = record.url

        # TODO: Limit number of bytes read per record e.g. read(200000)

        domain = None if not url else urlparse(url).netloc
        text = record.payload.read().decode('utf-8') #.limit(100) # TODO:
            Limit this read to ensure max length (for improving
            parallizability)
        lang = detect2(text[:300]) # Use PyCLD2, not langdetect, which was
            killing performance!
        yield domain, url, text, lang
    except Exception as e:
        yield e

def process_wet_simple(id_, iterator):
    count=0
    for uri in iterator:
        file = unpack(uri)
        for record in file:

```

```

%pyspark
detect2("this is a test")

'en'

```

READY

```

%pyspark

# PARAMETER - number of input files
nfiles = 1024 # Total 89100

# PARAMETER - slices / partitions of input
files = sc.parallelize(wetlist.take(nfiles)) #, numSlices=nfiles/32) # TODO: Try
    numSlices=nfiles/32 etc, or just default!

# Should parallelize
print(files.getNumPartitions())
rdd=files.mapPartitionsWithIndex(process_wet)

print(str(rdd))
docs = rdd.toDF(["host", "url", "text","lang"]) # "lang"
#docs.cache()

320
PythonRDD[179] at RDD at PythonRDD.scala:48

```

FINISHED

```

%pyspark

```

FINISHED

```
# Filter for English only
docs_en = docs.filter(docs.lang == 'en')
```

## Load saved vectors from Wikipedia model (created by python Wikipedia Text Processing.ipynb)

FINISHED

```
%pyspark
from pyspark.ml import Pipeline, PipelineModel
from pyspark.ml.feature import RegexTokenizer, CountVectorizer, StopWordsRemover
from pyspark.ml.clustering import LocalLDAModel

textModel = PipelineModel.load('s3://billsdata.net/CommonCrawl/wikipedia/text_model')
ldaModel = LocalLDAModel.load('s3://billsdata.net/CommonCrawl/wikipedia/lda_model')
```

```
%pyspark
```

FINISHED

```
# Test the models - for debugging only
import numpy as np
import pandas as pd

X=ldaModel.topicsMatrix().toArray()
vocab = np.array(textModel.stages[2].vocabulary)

topicLabels = [' '.join(vocab[np.argsort(X[:,i])[:, :-1]][:5])] for i in range(100)]

def score_topics(text):
    df = sqlContext.createDataFrame(pd.DataFrame({'text':[text]}))
    vec = textModel.transform(df)
    scores = ldaModel.transform(vec).select('topicDistribution').collect()[0]
    .topicDistribution.toArray()
    return pd.Series(dict(zip(topicLabels, scores)))

# Try it on an arbitrary sentence
school students education university college      0.001276
season team first teams cup                       0.001276
series book published books novel                 0.001261
series show television also episode               0.001292
ship ships two navy war                           0.001220
social one may also people                        0.001240
space earth light solar star                      0.001223
species found also large may                      0.001261
station line railway service train                 0.001261
team season coach football first                  0.001253
tom oliver ghost haiti kay                        0.001183
ukrainian ukraine dog dogs stamps                 0.001198
university research professor published science   0.001323
war union soviet communist political              0.001213
water company construction new coal               0.001240
world olympics championships summer women         0.430250
zealand new grand auckland prix                   0.001216
Length: 100, dtype: float64
```

`%pyspark`

FINISHED

```
# Now score pages from our WET files
docs_en.show(5)
vec = textModel.transform(docs_en)
vec.show(5)
```

```
+-----+-----+-----+-----+
|          host|          url|          text|lang|
+-----+-----+-----+-----+
|          null|          null|Software-Info: ia...| en|
|1000daysofwriting...|http://1000daysof...|1000 Days of Writ...| en|
|100unhappydays.bl...|http://100unhappy...|100 Unhappy Days:...| en|
|          10in30.com|http://10in30.com...|LearnOutLoud_300x...| en|
|123-free-download...|http://123-free-d...|MusicBoxTool - [3...| en|
+-----+-----+-----+-----+
```

only showing top 5 rows

```
+-----+-----+-----+-----+-----+
|          host|          url|          text|lang|          word|
+-----+-----+-----+-----+-----+
|          null|          null|Software-Info: ia...| en|[software, info, ..
|1000daysofwriting...|http://1000daysof...|1000 Days of Writ...| en|[1000, daysof, writ...
|100unhappydays.bl...|http://100unhappy...|100 Unhappy Days:...| en|[100, unhappy, days...
|          10in30.com|http://10in30.com...|LearnOutLoud_300x...| en|[learn, outloud, 300...
|123-free-download...|http://123-free-d...|MusicBoxTool - [3...| en|[musicboxtool, -, [3...
+-----+-----+-----+-----+-----+
```

`%pyspark`

FINISHED

```
# Create topic distribution vectors and tidy up
scores = ldaModel.transform(vec)
scores2 = scores.drop("url").drop("text").drop("lang").drop("words").drop("filtered")
              .drop("vec")
```

```
+-----+-----+-----+
|          host| topicDistribution|
+-----+-----+
|          null|[0.13311231517264...|
|1000daysofwriting...|[3.26238961502545...|
|100unhappydays.bl...|[5.81230792144733...|
|          10in30.com|[8.42785330446217...|
|123-free-download...|[6.44166735802645...|
+-----+-----+
```

only showing top 5 rows

`%pyspark`

RUNNING 0%

```
# Save these vectors to disc, so we can just load them later
scores2.write.parquet('s3://billsdata.net/CommonCrawl/topic_model_%d_files')
```

Started 20 minutes ago.

**Load saved scores from nfiles of WET files**

FINISHED

%pyspark

# Restart here

nfiles=1024

```
scores2 = spark.read.parquet('s3://billsdata.net/CommonCrawl/topic_model_%d_files
/cc_page_wiki_topic_scores' % nfiles)
```

```
+-----+-----+
|          domain|  topicDistribution|
+-----+-----+
|          null|[0.12615256587902...|
|1000daysofwriting...|[3.26238961502545...|
|100unhappydays.bl...|[5.81230792144733...|
|          10in30.com|[8.42785330446217...|
|123-free-download...|[6.44166735802645...|
+-----+-----+
```

only showing top 5 rows

%pyspark

FINISHED

```
# Aggregate page-scores per Host for now (will be same process for aggregating host
-scores to PLD-scores) - first count and sum vectors
```

```
scores3=scores2.rdd.map(lambda x: (x['domain'], (1,x['topicDistribution']
))).reduceByKey(lambda acc, val: (acc[0]+val[0], acc[1]+val[1]))
```

```
# Next, divide by the total to create averaged vectors, and convert back to a
dataframe
```

```
scores4=scores3.map(lambda x: (x[0], (x[1][1]/x[1][0]))).toDF(["host",
"averageTopicDistribution"])
```

```
+-----+-----+
|          host|averageTopicDistribution|
+-----+-----+
|lgattinawritercram...| [2.15245888670342...|
|www.stratifiedaut...| [6.52478581426803...|
|          www.edcast.com| [3.20044544708958...|
|www.indianyellowp...| [0.00651007834318...|
|          www.huskers.com| [2.66916339495641...|
+-----+-----+
```

only showing top 5 rows

```
{'averageTopicDistribution': <class 'pyspark.ml.linalg.VectorUDT'>, 'host': <class 'pysp
ark.sql.types.StringType'>}
```

%pyspark

FINISHED

```
# Just playing - code to help understand the different libraries and vector types!
```

```
import pyspark.mllib.linalg as mllib
```

```
import pyspark.ml.linalg as ml
df = sc.parallelize([
    (mllib.DenseVector([1, ]), ml.DenseVector([1, ])),
    (mllib.SparseVector(1, [0, ], [1, ]), ml.SparseVector(1, [0, ], [1, ]))
]).toDF(["mllib_v", "ml_v"])
df.show()
{s.name: type(s.dataType) for s in df.schema}

+-----+-----+
|      mllib_v|      ml_v|
+-----+-----+
|      [1.0]|      [1.0]|
|(1,[0],[1.0])|(1,[0],[1.0])|
+-----+-----+
{'ml_v': <class 'pyspark.ml.linalg.VectorUDT'>, 'mllib_v': <class 'pyspark.mllib.linalg.VectorUDT'>}
```

%pyspark

FINISHED

```
# Enrich each row with the corresponding PLD (using code from Paul J)
# TODO: Should just pickle the Bloom filter and load it in!
saveURI="s3://billsdata.net/CommonCrawl/hyperlinkgraph/cc-main-2017-aug-sep-oct
        /domaingraph/vertices/"
pld_df=spark.read.load(saveURI)
pld_df.show(3)
```

```
+---+-----+
| ID|  PLD|
+---+-----+
| 0|aaa.1|
| 1|aaa.2|
| 2|aaa.3|
+---+-----+
only showing top 3 rows
DataFrame[ID: string, PLD: string]
```

%pyspark

FINISHED

```
# Next, we'll construct a local dictionary from of all the PLDS (key is the PLD, value
  is the ID)
# This is our truth-table of known PLDs that we'll use when counting hosts
# Create a bloom filter using a pure python package (might be a little slow)
from pybloom import BloomFilter
pld_bf = BloomFilter(capacity=94000000, error_rate=0.005) # was 91M

for row in pld_df.rdd.collect(): # limit(100000000) # TODO: Still bad (and exceeds
    spark.driver.maxResultSize with all rows)!
    pld_bf.add(row['PLD'])

print(pld_df.rdd.take(3))
print(pld_df.rdd.take(3)[2]['PLD'])
#pld_bf.add(pld_df.rdd.take(3)[2]['PLD'])
```

```

print("aaa.aaa" in pld_bf) # Should be True

import sys
print(sys.getsizeof(pld_bf))
print(len(pld_bf)) # Should match number of items entered

# Broadcast the bloom filter so it's available on all the slave nodes - we don't need
  to change
# it any more so it's fine being immutable.
pld_bf_distrib=sc.broadcast(pld_bf)

print("aaa.aaa" in pld_bf) # Should be true
print("aaa.aaa.1" in pld_bf) # Should be false

[Row(ID=u'0', PLD=u'aaa.1'), Row(ID=u'1', PLD=u'aaa.2'), Row(ID=u'2', PLD=u'aaa.3')]
aaa.3
True
64
93110180
True
False
True
False

```

%pyspark

FINISHED

```

from pyspark.sql.functions import udf

# Returns a Boolean to say whether PLD is a hostname in itself
def is_a_pld(hostname):
    #if hostname in pld_lookup_table:
    #if pld_lookup_table.filter(lambda a: a == hostname).count()>0:
    if hostname in pld_bf_distrib.value:
        return True
    else:
        return False

# Define a function to do the hostname->pld conversion, if the pld exists in our
  dictionary
def convert_hostname(hostname):
    # Return hostname as-is, if this is already a PLD
    #if hostname in pld_lookup_table:
    #if pld_lookup_table.filter(lambda a: a == hostname).count()>0:
    if hostname in pld_bf_distrib.value:
        return hostname
    # Otherwise we're going to have to split it up and test the parts
    try:
        parts=hostname.split('.')
        if (len(parts)>4 and is_a_pld('.'.join(parts[0:4]))):
            return '.'.join(parts[0:4])
        if (len(parts)>3 and is_a_pld('.'.join(parts[0:3]))):
            return '.'.join(parts[0:3])
        if (len(parts)>2 and is_a_pld('.'.join(parts[0:2]))):
            return '.'.join(parts[0:2])
        if (len(parts)>1 and is_a_pld('.'.join(parts[0:1]))):

```

```

        return '.'.join(parts[0:1])
    return "ERROR" # Couldn't find a corresponding PLD - this should never happen!
except:
    return "ERROR"

```

```
udf_convert_hostname = udf(convert_hostname, StringType())
```

```
# Test
```

```
print(convert_hostname("aaa.aaa"))
```

```
aaa.aaa
```

```
True
```

```
%pyspark
```

FINISHED

```
# Function to reverse hostnames
```

```
from pyspark.sql.functions import udf
```

```
def reverse_domain(domain):
```

```
    try:
```

```
        ret = '.'.join(reversed(domain.split('.')))
```

```
        return ret
```

```
    except:
```

```
        return "NODOMAIN"
```

```
print(reverse_domain("com.facebook"))
```

```
udf_reverse_domain = udf(reverse_domain, StringType())
```

```
# Convert hosts in Topic DF to PLDs using convert_hostname function from Paul 5.
```

```
scores5=scores4.withColumn("pld_rev",udf_reverse_domain(udf_convert_hostname
```

```
(udf_reverse_domain("host")))) # Reverse the hostnames prior to lookup, then back
```

```
facebook.com
```

```

+-----+-----+-----+
|          host|averageTopicDistribution|          pld_rev|
+-----+-----+-----+
|lgattinawritercram...| [2.15245888670342...|          blogspot.com|
|www.stratifiedaut...| [6.52478581426803...| stratifiedauto.com|
|      www.edcast.com| [3.20044544708958...|          edcast.com|
|www.indianyellowp...| [0.00651007834318...| indianyellowpages...|
|      www.huskers.com| [2.66916339495641...|          huskers.com|
|jeofurry.blogspot...| [1.68556714565171...|          blogspot.com|
|www.brickmodeldes...| [4.93337221161091...| brickmodeldesign.com|
|      us.petvalu.com| [2.40757190496072...|          petvalu.com|
|www.bobgarontrain...| [3.31587146684675...| bobgarontraining.com|
|      www.homefocus.iel| [0.01917307556468...|          homefocus.iel|
+-----+-----+-----+

```

```
only showing top 10 rows
```

```
%pyspark
```

FINISHED

```
# Now we can aggregate page-scores per PLD, using a map-reduce similar to the host
  aggregation above - first count and sum vectors
```

```
scores6=scores5.rdd.map(lambda x: (x['pld_rev'], (1,x['averageTopicDistribution']
```



```

    )).reduceByKey(lambda acc, val: (acc[0]+val[0], acc[1]+val[1]))
# Next, divide by the total to create averaged vectors, and convert back to a
# dataframe
scores7=scores6.map(lambda x: (x[0], (x[1][1]/x[1][0]))).toDF(["pld_rev",
    "averageTopicDistribution"])

```

```

+-----+-----+
|          pld_rev|averageTopicDistribution|
+-----+-----+
|digitalhorizonson...| [0.00144483565419...|
|kitchencollection...| [0.01798544424072...|
|          ieee.org| [3.84272202396751...|
|          museumca.org| [1.67996725189244...|
|thepsoriasisprogr...| [3.37113604892104...|
+-----+-----+

```

only showing top 5 rows

12243

9141

```
%pyspark
```

FINISHED

```

# Save pld topic distributions in parquet format for Tom to play with (and to figure
# out how to create a PLD topic summary from this).
# TODO: Maybe a numpy argmax to get the index of the 'top' topic for each PLD with a

```

```
%pyspark
```

READY