

201709 evaluate C...

%pyspark

FINISHED

```
# Zeppelin notebook to demonstrate evaluation of CommonCrawl-derived domain vectors by
# classify domains according to high-level topic in the DMOZ dataset. Currently configu
# Bill's domain hex feature vectors from the 'Bill 6' notebook, and to use only Pyspark
# All cells should complete in less than a few minutes on an m4.2xlarge cluster.
# End-to-end run-time: approx 30 mins, with nfiles=128.
# NOTE: Should we really be trying to predict domain links instead? Or predicting bad d
# PJ - 20 Sept 2017
```

```
import boto
from pyspark.sql.types import *
```

```
# Import the DMOZ domain category dataset as an RDD
# (downloaded from https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910
```

```
dmoz_labels=sc.textFile('s3://billsdata.net/CommonCrawl/DMOZ/dmoz_domain_category.csv')
header = dmoz_labels.first() # extract header
dmoz_labels = dmoz_labels.filter(lambda row: row != header).map(lambda row: row.replace
dmoz_labels.take(3)
```

```
[[u'sdcastroverde.com', u'Top/World/Galego/regional/Galicia/Lugo/municipalities/Castrove
rde'], [u'www.232analyzer.com', u'Top/Computers/Hardware/Test_Equipment/Analyzers'], [u'
zschachwitz-tischtennis.de', u'Top/World/Deutsch/Sport/ball_Sports/table_tennis/Teams/Ge
rmany/Saxony']]
```

%pyspark

FINISHED

```
# Convert our labels RDD into a Spark DataFrame with a schema - neither column can be N
schema=StructType([StructField("domain", StringType(), False), StructField("categories"
dmoz_labels_df=spark.createDataFrame(dmoz_labels,schema)
dmoz_labels_df.printSchema()
print(dmoz_labels_df.count())
dmoz_labels_df.show(1)
dmoz_labels_df.cache()
```

root

```
 |-- domain: string (nullable = false)
 |-- categories: string (nullable = false)
```

2488259

```
+-----+-----+
|          domain|          categories|
+-----+-----+
|sdcastroverde.com|Top/World/Galego/...|
+-----+-----+
```

only showing top 1 row

DataFrame[domain: string, categories: string]

%pyspark

FINISHED

```
# Make a dictionary of short domains (removing www. prefix) to top-level category label
prefix="www."
dmoz_labels_clean=dmoz_labels_df.rdd.map(lambda row: ((row['domain'][len(prefix):] if row
                                                         row['categories'].split("/")[1].:

dmoz_labels_df.unpersist()
schema=StructType([StructField("domain", StringType(), False), StructField("category", :
dmoz_labels_clean_df=spark.createDataFrame(dmoz_labels_clean,schema)
dmoz_labels_clean_df.show(2)
dmoz_labels_clean_df.cache()
```

domain	category
lsdcastroverde.com	World
232analyzer.com	Computers

only showing top 2 rows
DataFrame[domain: string, category: string]

%pyspark

FINISHED

```
# Summarize categories in the DMOZ data
dmoz_labels_clean_df.groupBy('category').count().show()
```

category	count
Recreation	46095
World	1273970
Science	28138
Home	6952
Computers	45194
Sports	34890
Health	24218
Society	82079
Shopping	54062
Reference	21663
Games	10246
Arts	66721
Business	148144
Regional	642176
News	3711

%pyspark

FINISHED

```
# Load domain feature vectors from s3, in the following format:
# (u'www.angelinajolin.com', [4.30406509320417, 0.02702702702702703, 0.0, 0.13513513513513513])
```

```
# Load feature vectors from WAT files (from 'Bill 6' notebook) as an RDD:
#nfiles=128 # (takes about 5 mins for 128 files)
#inputURI = "s3://billsdata.net/CommonCrawl/domain_hex_feature_vectors_from_%d_WAT_file:

# Load Tom V's LDA topic-based vectors
inputURI = "s3://billsdata.net/CommonCrawl/lda_topic_vectors/tom_lda_vecs_sample3.txt"

features_rdd = sc.textFile(inputURI).map(eval)
features_rdd.take(3)
import pyspark.sql.types as typ
schema=StructType([StructField("domain", StringType(), False), StructField("vector", ArrayType(DoubleType(), False)),])
features_df=spark.createDataFrame(features_rdd,schema)
features_df.cache()
print("Nr domains:", features_df.count())
features_df.show(1)
features_df.printSchema()
```

('Nr domains:', 9763)

domain	vector
l10kidsin2010.blog...	[6.80862914E-5, 0...

only showing top 1 row

None

root

```
-- domain: string (nullable = false)
-- vector: array (nullable = true)
-- element: double (containsNull = false)
```

```
%pyspark
```

FINISHED

```
# Spark.ML classifiers require VectorUDT type, rather than Array, so we need to convert
from pyspark.ml.linalg import Vectors, VectorUDT
from pyspark.sql.functions import udf

vectorize=udf(lambda vs: Vectors.dense(vs), VectorUDT())
features_df = features_df.withColumn("vec", vectorize(features_df['vector'])).drop('vec')
features_df.show(1)
features_df.printSchema()
```

domain	vec
l10kidsin2010.blog...	[6.80862914E-5,0....]

only showing top 1 row

None

root

```
-- domain: string (nullable = false)
-- vec: vector (nullable = true)
```

%pyspark

FINISHED

```
# Filter embeddings for only those vectors that have entries in the DMOZ dictionary (i.e.
#common_domains_df= features_df.join(dmoz_labels_clean_df, features_df.domain == dmoz_labels_clean_df)
common_domains_df=features_df.join(dmoz_labels_clean_df, ["domain"]) # doesn't create embeddings
common_domains_df.cache()
features_df.unpersist()
dmoz_labels_clean_df.unpersist()
print("Number of labelled domains = " + str(common_domains_df.count()))
common_domains_df.show(3)
common_domains_df.printSchema()
```

Number of labelled domains = 647

```
+-----+-----+-----+
|      domain|      vec| category|
+-----+-----+-----+
|roanokeisland.com|[0.0599172335,4.5...| Regiona|
|      vietbao.vn|[1.04846403E-4,8....| World|
|      tv.adobe.com|[0.472358937,4.30...|Computers|
+-----+-----+-----+
```

only showing top 3 rows

root

```
-- domain: string (nullable = false)
-- vec: vector (nullable = true)
-- category: string (nullable = false)
```

%pyspark

FINISHED

```
# Create numeric indexes for our classes
from pyspark.ml.feature import IndexToString, StringIndexer, VectorIndexer
labelIndexer = StringIndexer(inputCol="category", outputCol="indexedCategory").fit(common_domains_df)

# Split into training and test sets using spark.ML API
domains_train, domains_test = common_domains_df.randomSplit([0.7,0.3],seed=42)
domains_test.groupBy('category').count().show()
```

```
+-----+-----+
| category|count|
+-----+-----+
|Recreation|    7|
|    World|   12|
|   Science|    9|
|    Home|    5|
|Computers|   24|
|   Sports|   11|
|   Health|    5|
| Society|   16|
| Shopping|    6|
|Reference|   13|
|    Games|    5|
|    Arts|   20|
|Business|   10|
|Regional|   35|
|    News|    1|
```

%pyspark

FINISHED

```
# Create a pipeline and fit a RandomForest Classifier using spark.ml
from pyspark.ml import Pipeline
from pyspark.ml.classification import RandomForestClassifier

# Build our RF classifier
rf = RandomForestClassifier(labelCol="indexedCategory", featuresCol="vec", numTrees=10)

# Convert indexed labels back to original labels
labelConverter = IndexToString(inputCol="prediction", outputCol="predictedCategory", labelCol="indexedCategory")

# Define and run the full pipeline to train the model and make predictions
pipeline = Pipeline(stages=[labelIndexer, rf, labelConverter])
model=pipeline.fit(domains_train)
predictions=model.transform(domains_test)
print(predictions.take(1))
predictions.select("predictedCategory", "category", "vec").show(5)

# FYI, Equivalent code in sklearn
#from sklearn.ensemble import RandomForestClassifier
#rf = RandomForestClassifier(max_depth=2, random_state=0)
#rf.fit(X_train, y_train)
#print(classification_report(y_test, rf.predict(X_test)))
```

```
[Row(domain=u'411.info', vec=DenseVector([0.7023, 0.001, 0.0011, 0.001, 0.001, 0.0011, 0.2896, 0.001, 0.0009, 0.001]), category=u'Reference', indexedCategory=5.0, rawPrediction=DenseVector([3.442, 0.2423, 1.3674, 1.9328, 0.23, 1.3183, 0.1013, 0.3289, 0.1688, 0.3291, 0.1651, 0.144, 0.0455, 0.0911, 0.0935]), probability=DenseVector([0.3442, 0.0242, 0.1367, 0.1933, 0.023, 0.1318, 0.0101, 0.0329, 0.0169, 0.0329, 0.0165, 0.0144, 0.0045, 0.0091, 0.0094]), prediction=0.0, predictedCategory=u'Regional')]
```

```
+-----+-----+-----+
|predictedCategory| category|          vec|
+-----+-----+-----+
|          Regional|Reference|[0.70233998,0.001...|
|          Regional|      Arts|[0.167306443,1.85...|
|          Regional|Regional|[0.868170457,1.52...|
|          Regional|Reference|[0.241388187,2.36...|
|          Regional|Business|[0.417959319,0.04...|
+-----+-----+-----+
```

only showing top 5 rows

%pyspark

FINISHED

```
# Select (prediction, true label) and compute test error
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
evaluator1 = MulticlassClassificationEvaluator(labelCol="indexedCategory", predictionCol="prediction", metricName="accuracy")
evaluator2 = MulticlassClassificationEvaluator(labelCol="indexedCategory", predictionCol="predictedCategory", metricName="accuracy")

accuracy = evaluator1.evaluate(predictions)
f1=evaluator2.evaluate(predictions)
```

```
print("Accuracy=%g, F1=%g" % (accuracy, f1))
```

```
Accuracy=0.213904, F1=0.131569
```

```
%pyspark
```

READY