

BRANDON WONG – SG DAT 1

I. INTRODUCTION

II. APPROACH

III. CHALLENGES

IV. NEXT STEPS

V. Q&A

INTRODUCTION

\$40,000 prize

Timeline: 18/01/16 – 25/04/16 1,269 teams

Predict Relevance of Search Results



Home Depot: Home Depot is an American retailer of home improvement and construction products and services



Competition Basics: Improve customer's shopping experience by developing a model that can accurately predict the relevance of search results



INTRODUCTION (CONT)

Data Files

File Name	Available Formats
sample_submission.csv	.zip (226.76 kb)
train.csv	.zip (2.51 mb)
test.csv	.zip (4.74 mb)
product_descriptions.csv	.zip (34.77 mb)
attributes.csv	.zip (27.21 mb)
relevance_instructions	.docx (105.01 kb)

Train and
Test have
similar
columns

BUT

relevance
score not
provided
in Test.

scikit-learn algorithm cheat-sheet

START

classification

- get more data
 - >50 samples
 - predicting a category
 - <100K samples
 - do you have labeled data
 - YES
 - SVC
 - Ensemble Classifiers
 - Naive Bayes
 - NO
 - SGD Classifier
 - Linear SVC
 - Text Data
 - KNeighbors Classifier
 - YES
 - SGD Classifier
 - Linear SVC
 - NO
 - SGD Classifier
 - Linear SVC
 - NO
 - SGD Classifier
 - Linear SVC
 - NO
 - SGD Classifier
 - Linear SVC

regression

- get more data
 - >50 samples
 - predicting a category
 - <100K samples
 - do you have labeled data
 - YES
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
 - NO
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
 - YES
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
 - NO
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
 - NO
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
 - NO
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors
- NO
 - SGD Regressor
 - Lasso
 - ElasticNet
 - SVR(kernel='rbf')
 - EnsembleRegressors

clustering

- get more data
 - >50 samples
 - predicting a category
 - <100K samples
 - do you have labeled data
 - YES
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
 - NO
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
 - YES
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
 - NO
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
 - NO
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
 - NO
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM
- NO
 - Spectral Clustering
 - GMM
 - KMeans
 - MiniBatch KMeans
 - MeanShift
 - VBGMM

dimensionality reduction

- get more data
 - >50 samples
 - predicting a category
 - <100K samples
 - do you have labeled data
 - YES
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
 - NO
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
 - YES
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
 - NO
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
 - NO
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
 - NO
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation
- NO
 - Randomized PCA
 - Isomap
 - Spectral Embedding
 - LLE
 - kernel approximation

Back

scikit learn

I have:

- **Labeled data (Classification?)**
- **Mid-scale volume, classifiers (Random Forest?)**

I need:

- Root words / Stems (Snowball Stemmer / PyStemmer?)

APPROACH (CONT)

id	product_uid	product_title	search_term	relevance
2	100001	Simpson Strong	angle bracket	3
3	100001	Simpson Strong	l bracket	2.5
9	100002	BEHR Premium	deck over	3
16	100005	Delta Vero 1-Ha	rain shower head	2.33
17	100005	Delta Vero 1-Ha	shower only faucet	2.67
18	100006	Whirlpool 1.9 cu	convection otr	3
20	100006	Whirlpool 1.9 cu	microwave over s	2.67
21	100006	Whirlpool 1.9 cu	microwaves	3
23	100007	Lithonia Lightin	emergency light	2.67
27	100009	House of Fara 3	mdf 3/4	3

Relevance is a number between 1 (not relevant) and 3 (most relevant)

E.g. Search for Steel Saw

Steel Saw (R = 3)



Steel Nails (R = 2)

Shovel (R = 1)

Each pair was (*search_term*, *product*) evaluated by at least 3 human raters.

The provided relevance scores are the average value of the ratings

LEADERBOARD

#	Δ1w	Team Name <small>* in the money</small>	Score <small>?</small>	Entries	Last Submission UTC (Best - Last Submission)
1	—	Turing test <small>👤 *</small>	0.44014	113	Tue, 29 Mar 2016 16:06:47 (-4d)
2	—	. <small>*</small>	0.44222	69	Fri, 01 Apr 2016 07:00:18 (-29.9h)
3	↑11	Alex&Andreas&Nurlan <small>👤 *</small>	0.44268	79	Thu, 31 Mar 2016 20:32:12
954	↓51	Ganapriya Kalavagunta	0.48720	1	Sun, 13 Mar 2016 16:38:06
955	new	Brandon Wong	0.48721	2	Fri, 01 Apr 2016 19:15:11
<div>Your Best Entry ↑ You improved on your best score by 0.00000. You just moved up 87 positions on the leaderboard.</div> <div> Tweet this!</div>					
956	↓52	 Iqbal Hossain	0.48721	7	Fri, 12 Feb 2016 17:01:41 (-6.1d)

CHALLENGES

1. Not trying to predict the true relevancy of the product as a response to a search query
2. Instead, build program to mimic human raters, assuming they are the most efficient method of assessing relevancy
3. Have to teach the models/machines to act like humans? Need to “create a search system auditor that can help measure the efficacy of changes in algorithms preferably in real time”

NEXT STEPS

1. Check out winning strategies on Kaggle, improve on them
2. Keep competing in Kaggle competitions
3. Try out Natural Language Processing

Q&A
