

The Python Way: Predicting Entity Popularity on Twitter



Pedro Saleiro

DSPT#2
Porto, 2016

Reputation

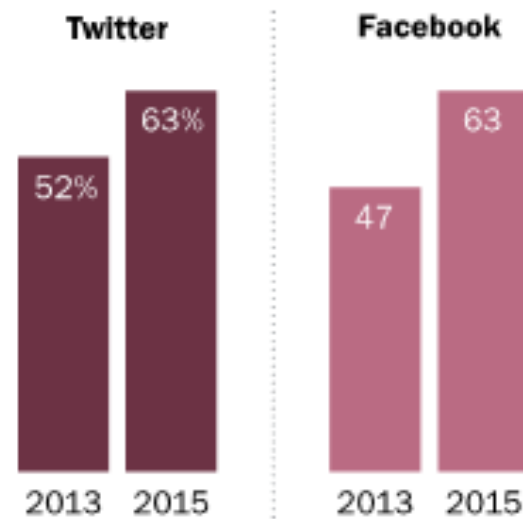
[[Van Riel et al., 2007](#)] define reputation as “overall assessments of organizations by their stakeholders”



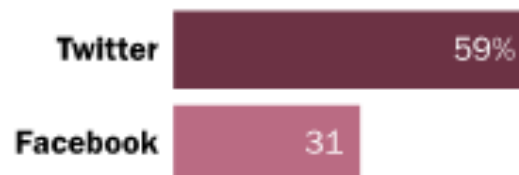
Social Media & Online News

Facebook and Twitter News Use is on the Rise

% of ___ users who get news there



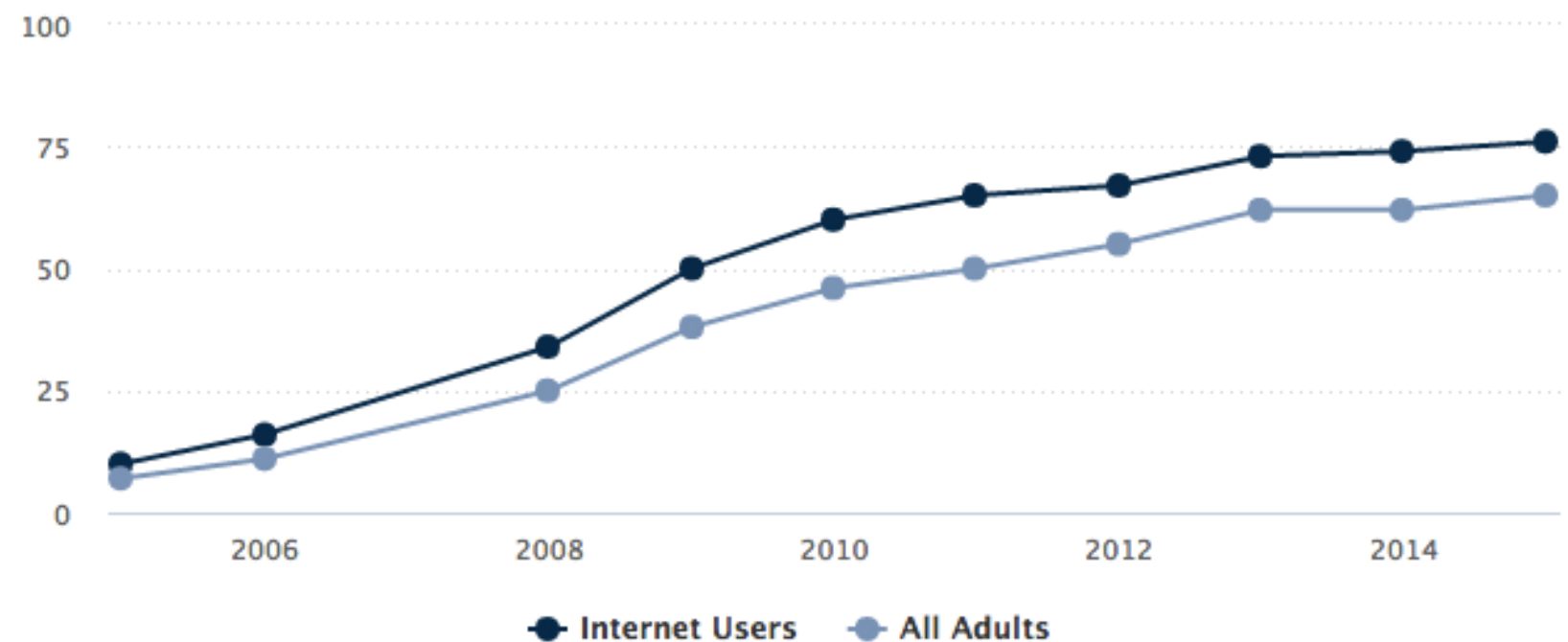
Of those who get news from ___ in 2015, percent who have kept up with a news event as it was happening



Social Media and News Survey, March 13-15 & 20-22, 2015. Q2, Q4, Q7, Q11.

PEW RESEARCH CENTER

% of all American adults and internet-using adults who use at least one social networking site



Online Reputation Monitoring

Tracking what is said about a given entity on Social Media

Early ORM systems focused on counting entity mentions on Social Media

Implies collecting, cleaning, filtering, mining, exploring and analysing large streams of unstructured text data

Current Systems focus in NER, NED, Polarity Classification and Visualisation (Social Media Analytics)

POPmine

Framework for ORM

Data

Tweets

100K “portuguese” users panel

24/7 since 2012

228M tweets by today

News

Sapo news crawler

55 online news outlets since 2010

6M news articles

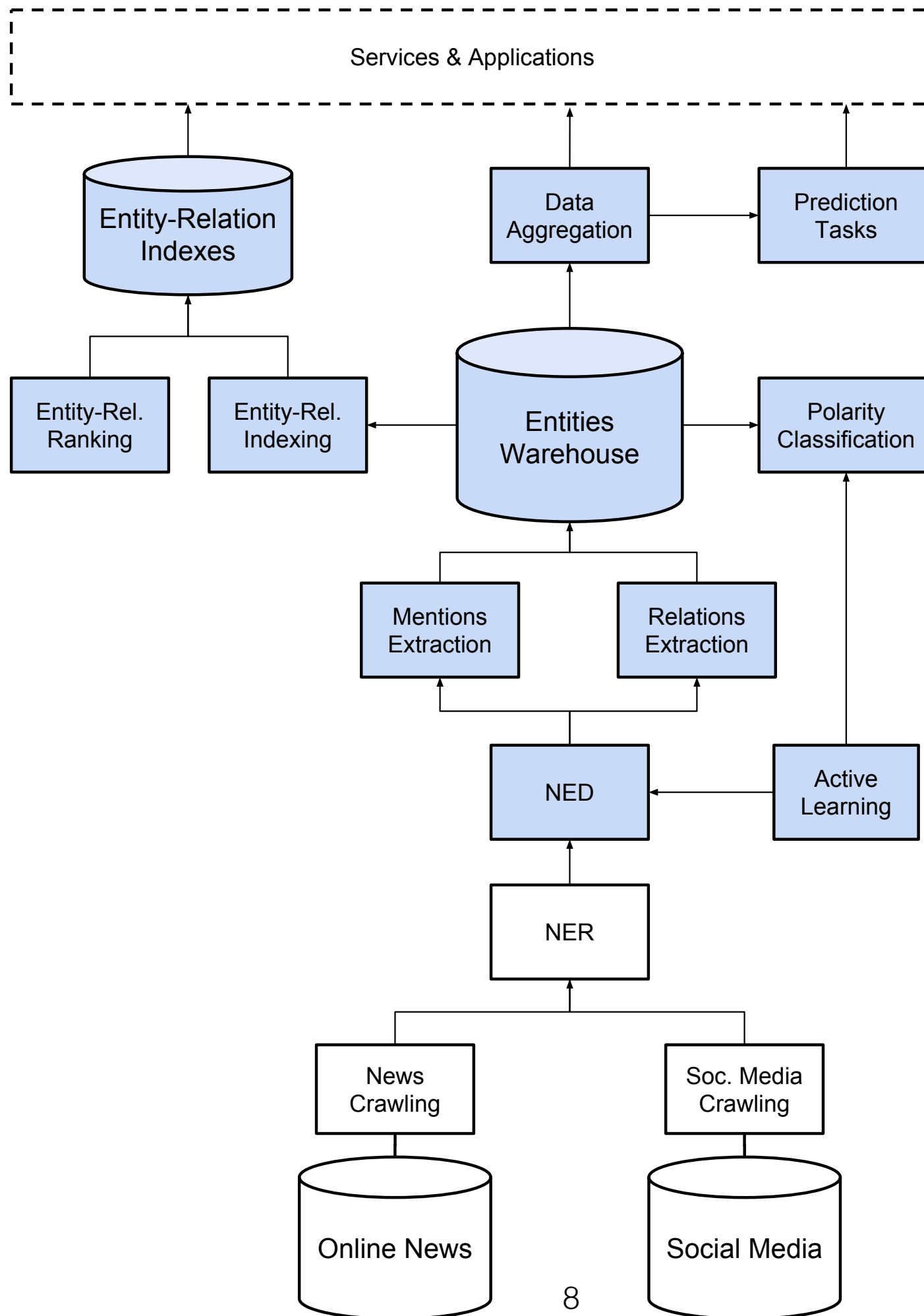
Named Entities

Filters using Natural Language Processing, Information Retrieval and Machine Learning.

“Passos” is Passos Coelho based on remaining words, hashtags, links.

Use of SapoLabs Verbetes (names of personalities and their lexical variations).

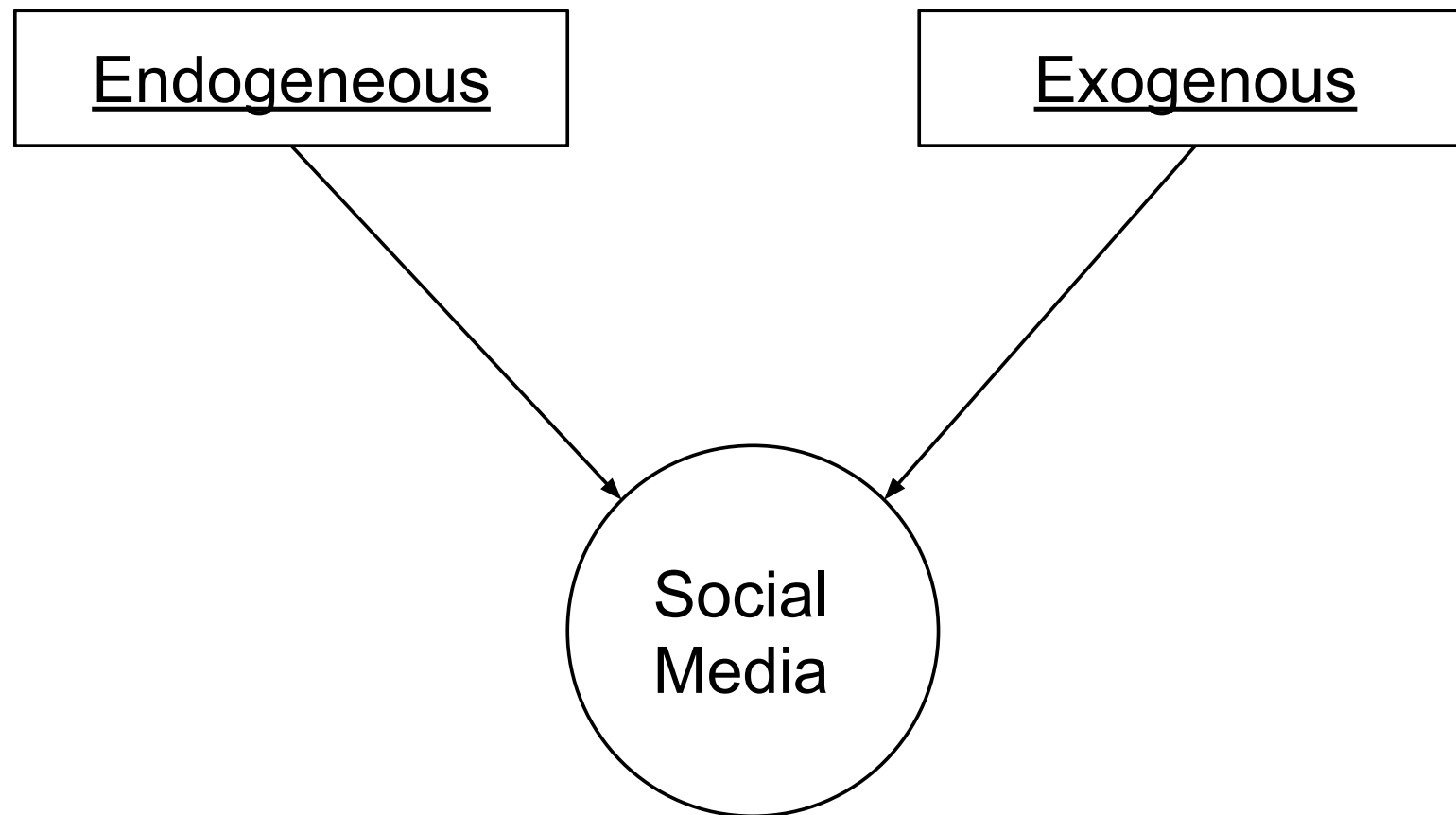
RepLab Filtering competition: 91% accuracy.



Prediction Task

Predict entity popularity on Twitter?

Social Media Attention



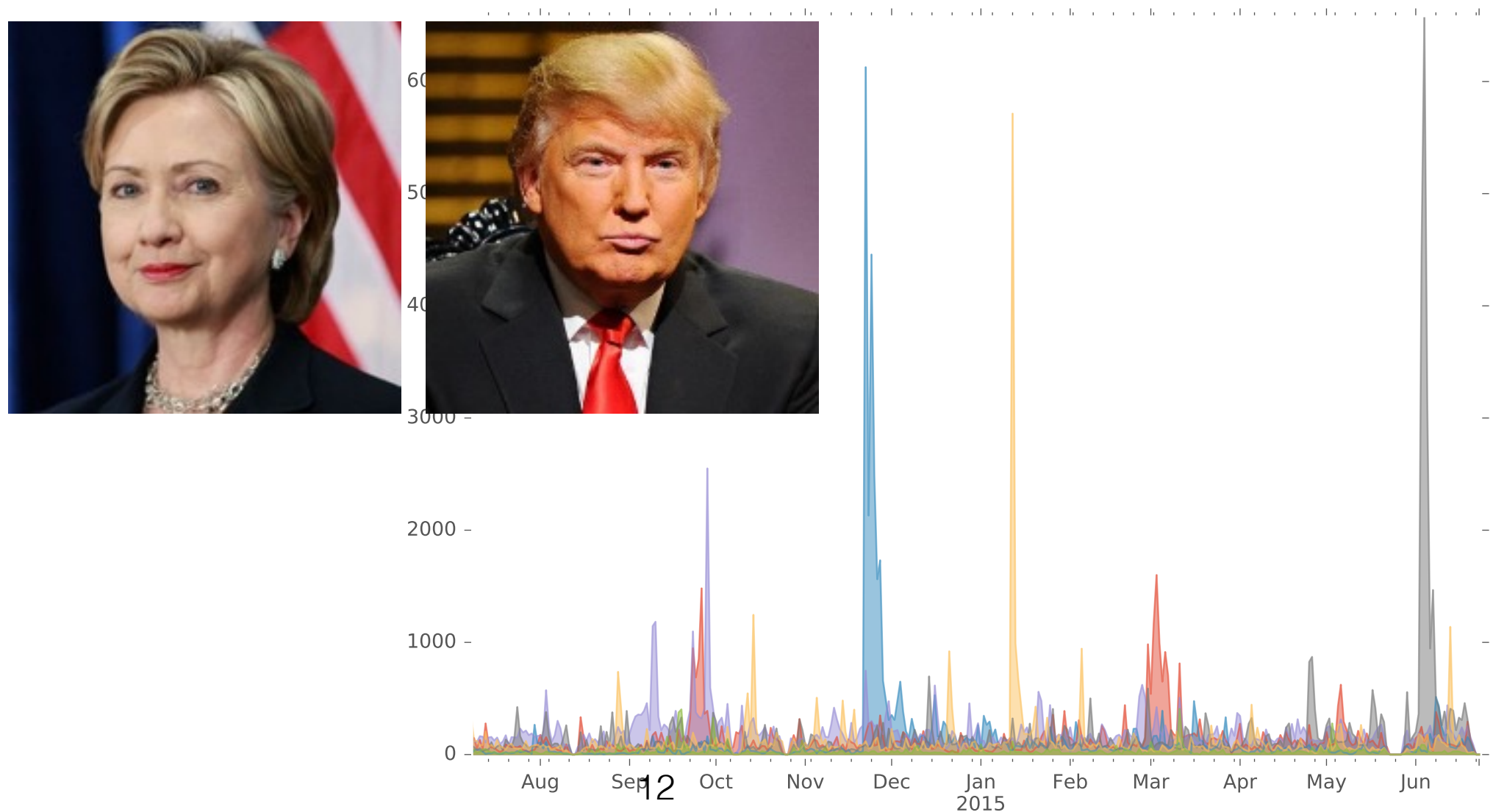
Social Media Attention

Endogenous: Followers network, Influence, Hashtags,

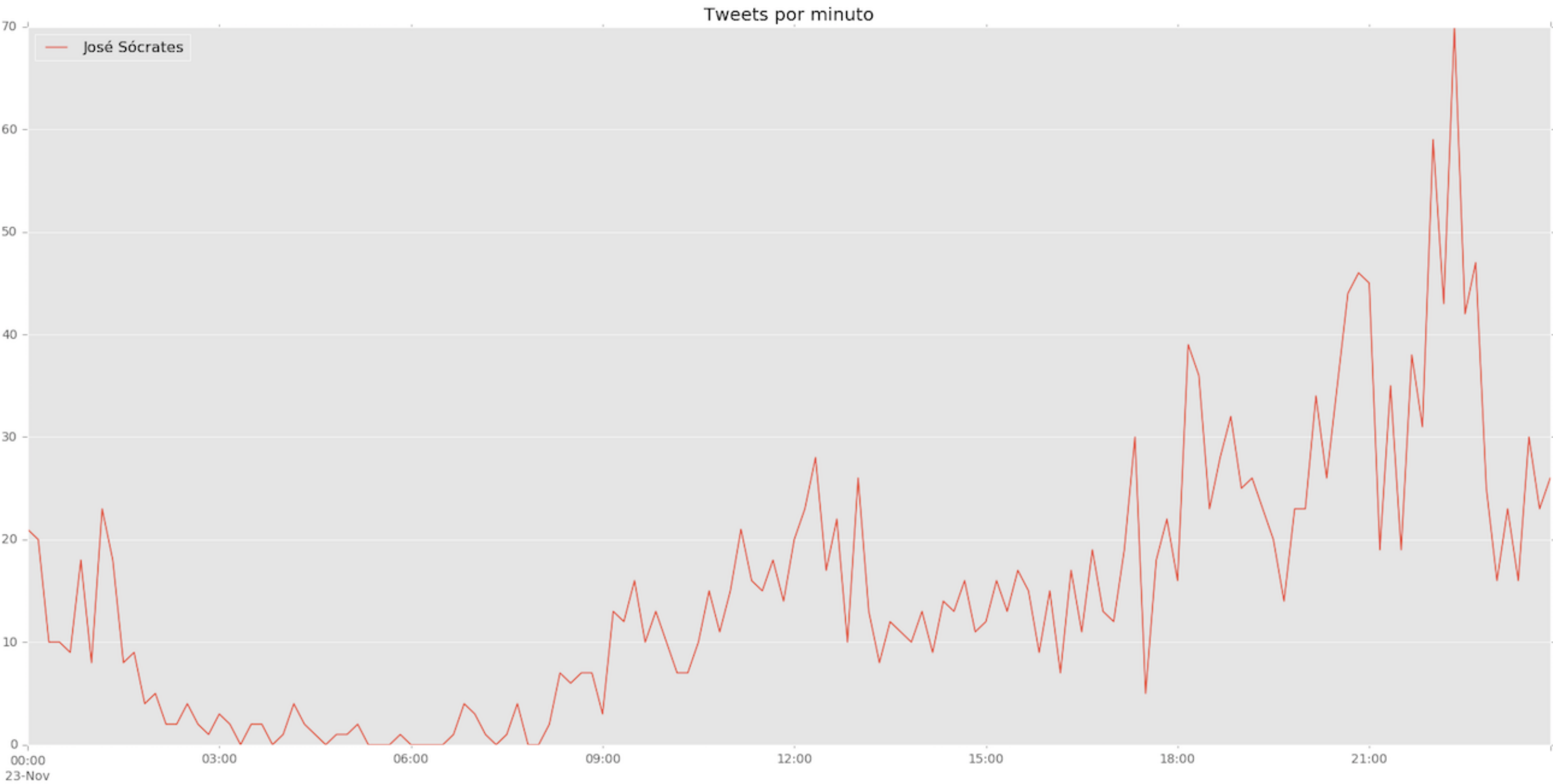
Exogenous: Online Media, TV, Traditional Media,
Friends, Family, Personality, Macroeconomy,...

Goal

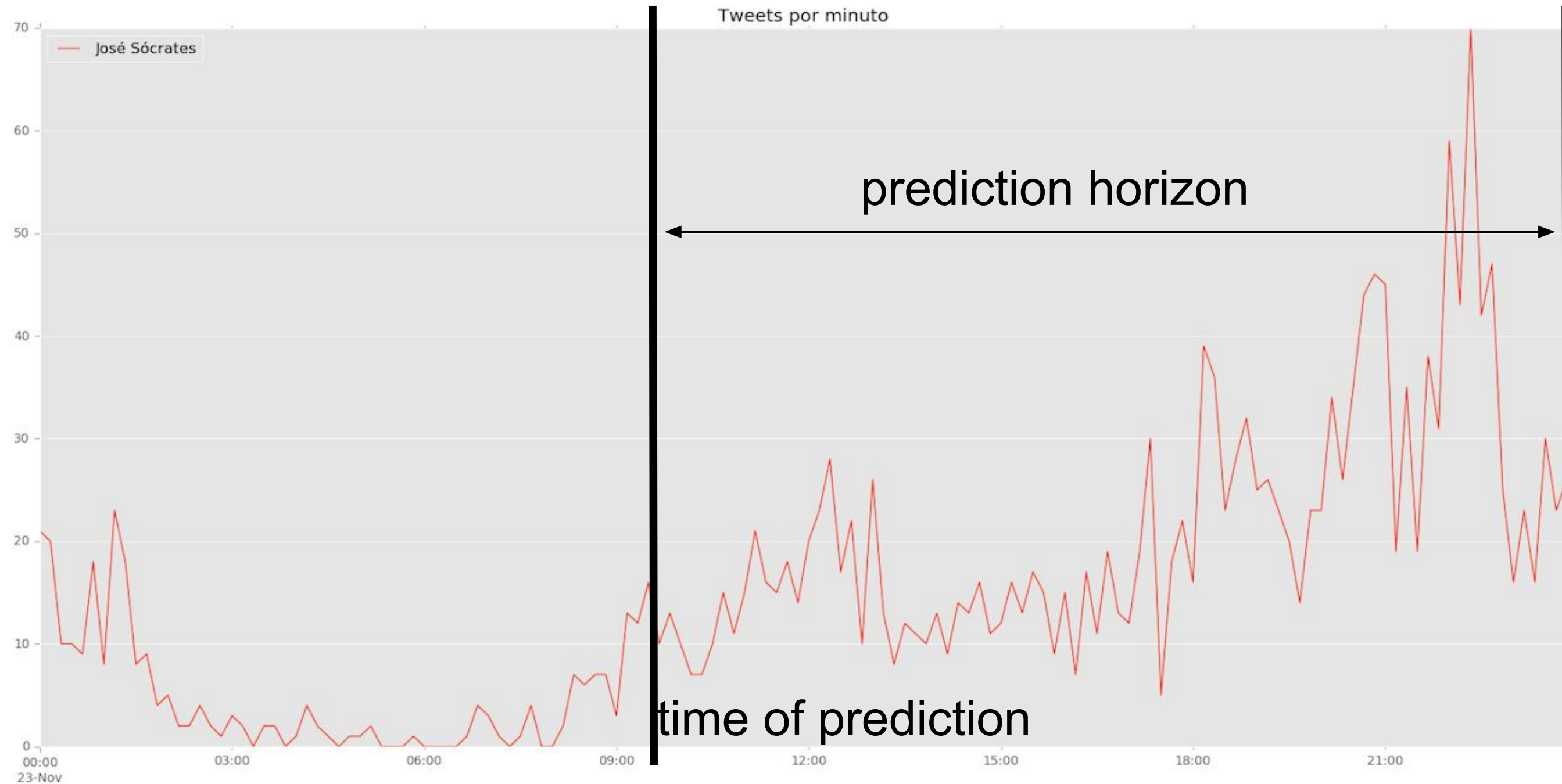
Predict if an entity (e.g. politician) will be *frequently mentioned on Twitter* in the hours *following appearing in the news*



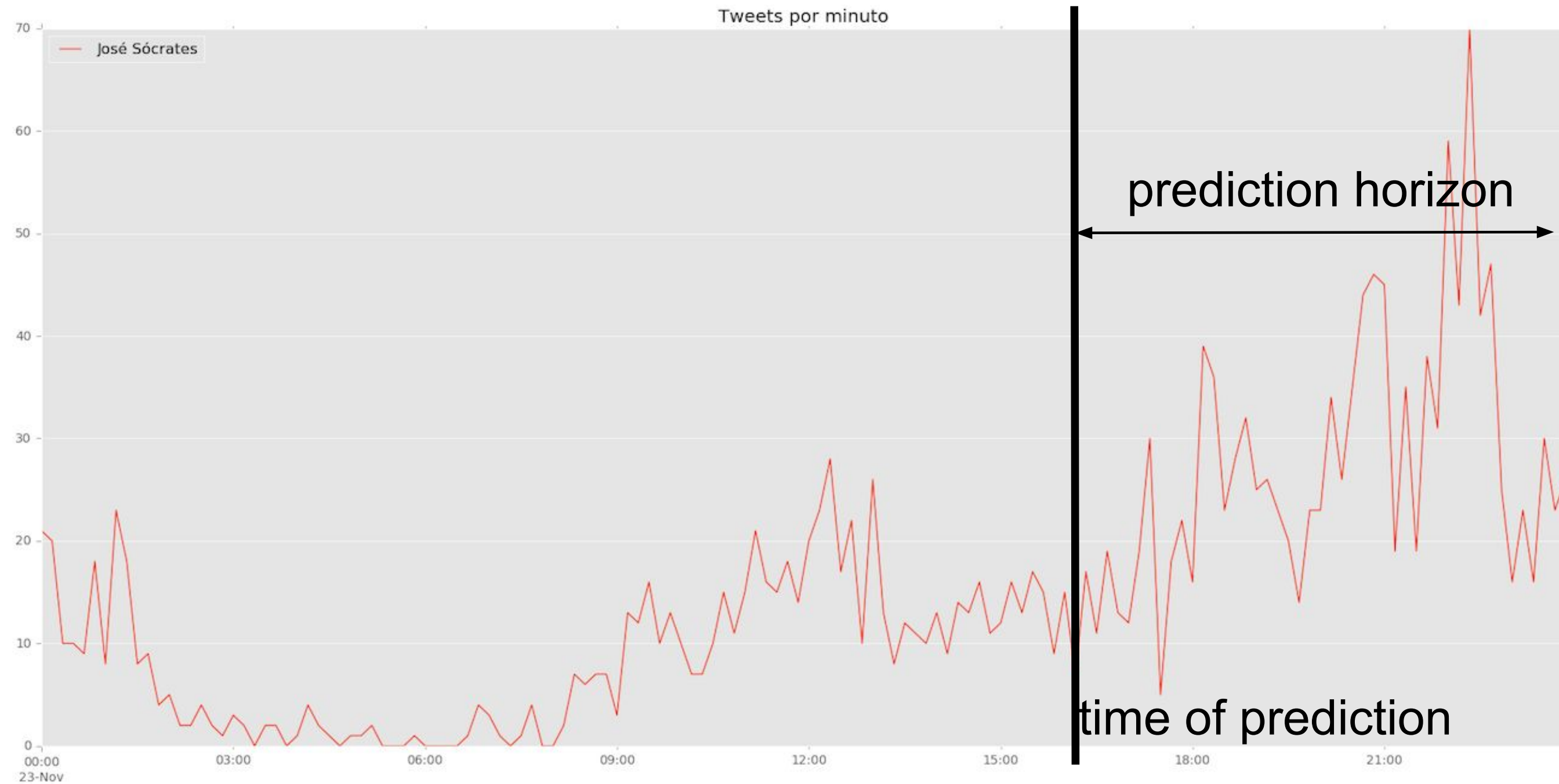
Example



Example



Example



Formalisation

Given a set of entities $E = \{e_1, e_2, \dots, e_i, \dots\}$,

a daily stream of social media messages $S = \{s_1, s_2, \dots, s_i, \dots\}$,

a daily stream of online news articles $N = \{n_1, n_2, \dots, n_i, \dots\}$,

a discrete function $f_m(e_i, S)$ representing mentions of an entity e_i on the social media stream S ,

a daily time frame $T = [t_p, t_{p+h}]$, where the time t_p is the time of prediction and t_{p+h} is the prediction horizon time.

We want to learn a target popularity function

$$f_p(e_i, N, T) = \sum_{t=t_p}^{t=t_{p+h}} f_m(e_i, S)$$

Approach

Supervised learning approach

News features:

- *signal*
- *textual*
- *semantic*
- *sentiment*

Multiple prediction horizons

- *e.g. impact until 24:00 of news published at 8am*

Features - Signal

total mentions of entity in the news

between midnight and time of prediction

same for previous day (lagged)

total mentions of entity in news titles

average news length

number of different news outlets mentioning entity

weekday/weekend

Features - Textual

Create a meta-document of news titles mentioning entity between midnight and time of prediction

Everyday

Calculate TF-IDF for each day (compare with others)

Create a topic model using Latent Dirichlet Allocation

assign topic probability for each day

Features - Sentiment

Use a sentiment lexicon

Count occurrences in news titles mentioning entity

Count positives, negatives and neutrals

Calculate aggregate functions (e.g. pos/neg)

Calculate a TF-IDF score of adjectives in the news

Features - Semantic

Take advantage of journalists tags

Semantic categories (e.g. Internal Politics)

Calculate TF-IDF of semantic tags

Create a BOW representation of named entities

Calculate a TF-IDF of named entities

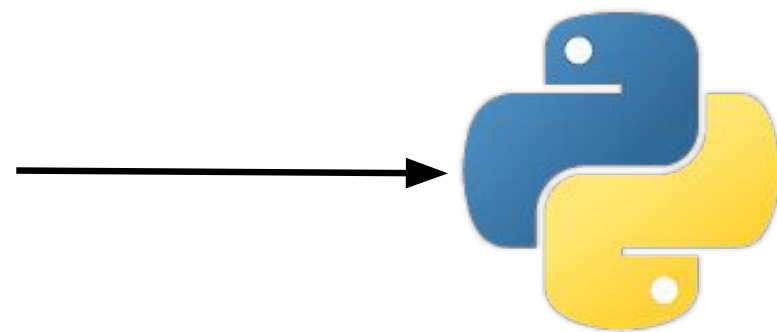
Prediction

Regression or classification?

Implementation

The Python Way

db
entity
horizon
train period
test period
features list



model
predictions
evaluation

Tweets & News DB

```
db.tweets.find_one()
```

```
{u'_id': u'144424904090062849',  
  u'datasource_s': u'mysql',  
  u'description': u'Jornalista redactor da RTP',  
  u'followers_count': 89,  
  u'friends_count': 269,  
  u'language_s': u'pt',  
  u'location': u'',  
  u'name': u'Carlos Santos Neves',  
  u'profile_image_url': u'http://a0.twimg.com/profile_images/1206504715/IMG_0639_normal.JPG',  
  u'retweet_count': 0,  
  u'screen_name': u'csantosneves',  
  u'source': u'web',  
  u'status_in_reply_to_status_id': u'0',  
  u'status_in_reply_to_user_id': u'0',  
  u'text': u'A\xed est\xed Paulo Portas: "A Europa, em todo o caso, n\xeo se faz a dois, faz-se a 27" - http://t.co/aa4XZYGk',  
  u'time_zone': u'Lisbon',  
  u'tokenized_text': u'A\xed est\xed Paulo Portas : " A Europa , em todo o caso , n\xeo se faz a dois , faz-se a 27" - http://t.co/aa4XZYGk',  
  u'tweet_date': datetime.datetime(2011, 12, 7, 14, 35, 55),  
  u'urls': [u'http://t.co/aa4XZYGk'],  
  u'user_created_at': datetime.datetime(2010, 10, 26, 9, 54, 22),  
  u'user_id': u'207933284',  
  u'version': 10}
```

Tweets & News DB

```
db.news.find_one()
```

```
{u'_id': u'6066313',  
  u'content': u"O avanço Karim Benzema mostrou pontaria afinada no triunfo sobre os LA Galaxy (3-1), terminando o encontro com dois golos na folha de marcadores. A partida está inserida na Guinness Cup, competição amigável realizada na digressão dos merengues pelos Estados Unidos.\nO fã da primeira parte foi apontado por Angel Di Maria. O extremo argentino não se fez rogado e inaugurou a contenda, logo aos 15 minutos.\nA segunda metade também começou com feio para os merengues, que aumentaram vantagens por intermédio de Karim Benzema, aos 51'. Jose Villarreal ainda reduziu para a antiga forma de David Beckham (63'), mas o avançado francês matou definitivamente o encontro, com um bis (75').\nNa mesma competição está inserido o Chelsea de José Mourinho que, esta madrugada, com golos de Oscar e Eden Hazard derrotou o Inter de Milão. Os italianos jogaram com menos um elemento desde a expulsão de Campagnaro, aos 58'.",  
  u'link': u'http://www.futebol365.pt/noticias/artigo.asp?id=90615&utm_source=rss&utm_medium=feed&utm_campaign=noticias_xml',  
  u'numComments': u'0',  
  u'occurrences': [u'José Mourinho'],  
  u'pubdate': datetime.datetime(2013, 8, 2, 9, 30, 37),  
  u'source': u'Futebol 365',  
  u'tags': u'',  
  u'title': u'Jogos Amigáveis: Benzema em destaque na vitória sobre os LA Galaxy'}
```

Tweets & News DB

```
db.tweets_mentions.find_one({'target_name': 'Pedro Passos Coelho'})
```

```
{u'_id': u'29089281024316211250b344a0e9e57312c88b8bfc',  
  u'mention_confidence': 1,  
  u'mention_size': 3,  
  u'polarity': {u'PopstarOpinionizer_2': 0},  
  u'random': 0.5188359553449241,  
  u'target_id': ObjectId('50b344a0e9e57312c88b8bfc'),  
  u'target_mention': u'Pedro Passos Coelho',  
  u'target_name': u'Pedro Passos Coelho',  
  u'tweet_date': datetime.datetime(2013, 1, 14, 18, 47, 5),  
  u'tweet_id': u'290892810243162112',  
  u'tweet_user': u'18186509'}
```

```
db.news_mentions.find_one({'target_name': 'Marcelo Rebelo de Sousa'})
```

```
{u'_id': u'374619250ec6ca8e9e573741a23a73c',  
  u'news_date': datetime.datetime(2011, 6, 21, 10, 29),  
  u'news_id': u'3746192',  
  u'news_link': u'http://aeiou.expresso.pt/assuncao-esteves-candidata-a-liderar-a-ar=f656920',  
  u'news_source': u'Expresso',  
  u'target_id': ObjectId('50ec6ca8e9e573741a23a73c'),  
  u'target_name': u'Marcelo Rebelo de Sousa'}
```

Tweets & News DB

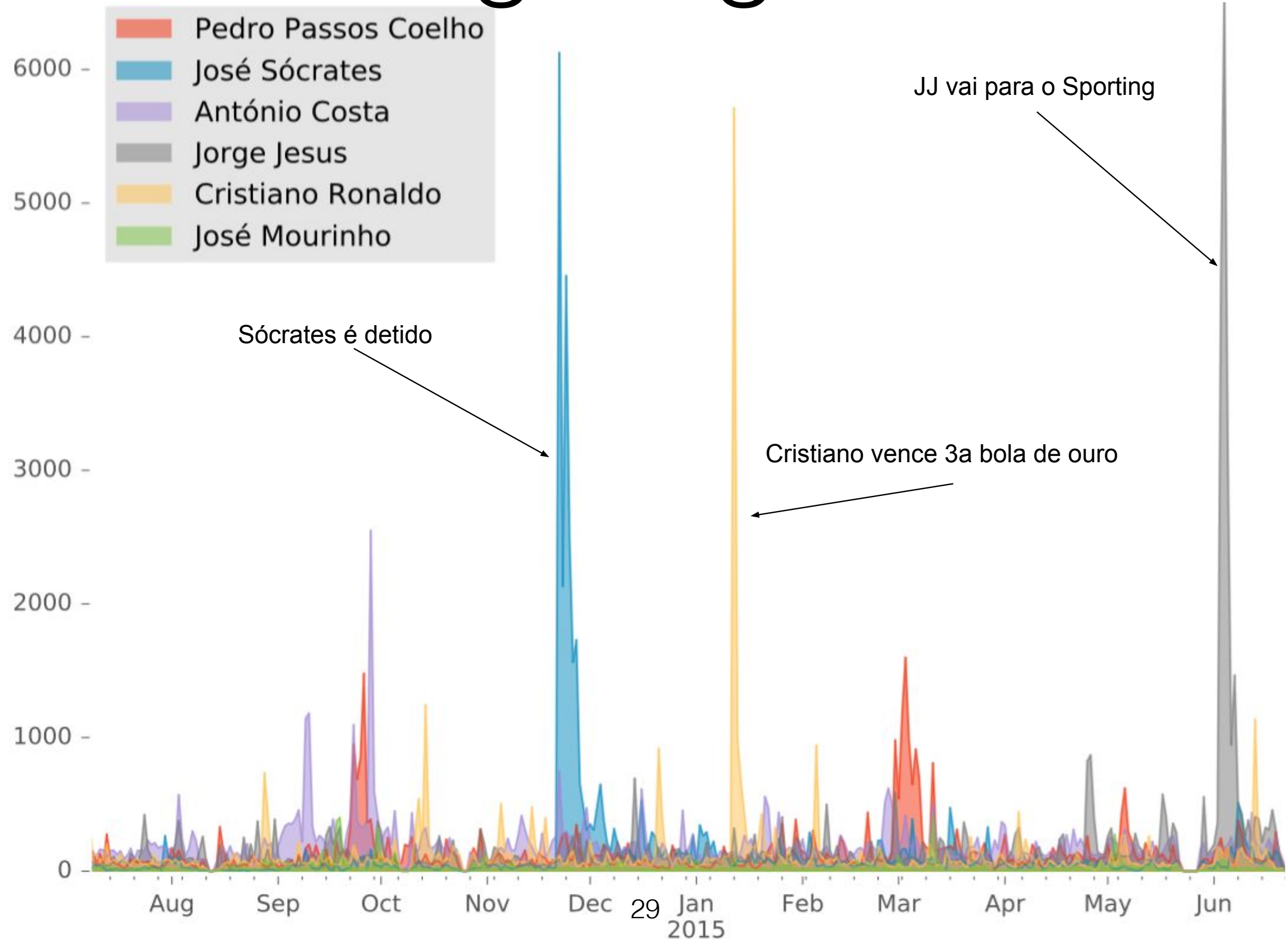
```
db.tweets_hourly_sentiment.find_one({'_id.target_name': 'Jorge Jesus'})
```

```
{u'_id': {u'hour': u'2011-06-21T00', u'target_name': u'Jorge Jesus'},  
u'value': {u'negative_mentions': 3.0,  
u'neutral_mentions': 0.0,  
u'positive_mentions': 0.0,  
u'total_mentions': 3.0}}
```

```
db.news_hourly_buzz.find_one({'_id.target_name': 'Luís Filipe Vieira'})
```

```
{u'_id': {u'hour': u'2011-06-22T06', u'target_name': u'Lu\xeds Filipe Vieira'},  
u'value': {u'total_mentions': 1.0}}
```


Generating Target Variable



Generating Target Variable

```
def getSeries(db, entity, start_date, end_date):  
    rng_date = pandas.date_range(start_date, end_date, freq='h')  
    tweets_hour = {}  
    for each in db.tweets_hourly_sentiment.find({'_id.target_name': entity}):  
        tweets_hour[each['_id']['hour']] = each['value']['total_mentions']  
    ts = pandas.Series(tweets_hour)  
    ts = ts.reindex(rng_date)  
    ts.fillna(ts.median(), inplace=True)  
    # or we can interpolate: ts = ts.interpolate()  
    return ts  
  
def reSample(ts, start_horizon, end_horizon):  
    ts_hour = ts.index.hour  
    selector = ((ts_hour >= start_horizon) & (ts_hour <= end_horizon))  
    series = ts[selector].resample('D', how='sum')  
    return series
```

Extracting Features

```
dataset = [{ 'date': ..., 'titles': ..., ... },  
            { 'date': ..., 'titles': ..., ... },  
            ...  
            ]  
dataset[date] ["semantic_entities_num"] = len(set(entities.split()))
```

Extracting Features

```
from sklearn.decomposition import TruncatedSVD
from sklearn.feature_extraction.text import TfidfVectorizer

def getTFIDF(training_titles):
    #training_titles = [list of concatenated titles per day]
    vec_tfidf = TfidfVectorizer(min_df=5, max_df=0.9, ngram_range=(
        1, 2), smooth_idf=True, sublinear_tf=True, use_idf=True, max_features=10000)
    tfidf_matrix = vec_tfidf.fit_transform(training_titles)
    lsa_vec = TruncatedSVD(n_components=10, random_state=0)
    lsa_matrix = lsa_vec.fit_transform(tfidf_matrix)
    return vec_tfidf, tfidf_matrix, lsa_vec, lsa_matrix
```

```
... # TF_IDF TITLES
... if 'titles_tfidf' in features:
...     X_test_titles_tfidf = vec_tfidf.transform(each[date]['titles'])
...     X_test_titles = lsa_vec.transform(X_test_titles_tfidf)
```


Extracting Features

```
import lda
def getLDA(training_titles):
    cvec_titles = CountVectorizer(ngram_range=(1, 1), binary = True, min_df = 2, max_df=0.9)
    cvec_matrix = cvec_titles.fit_transform(training_data)
    lda_model = lda.LDA(n_topics=10, n_iter=500, random_state=1)
    lda_model.fit(cvec_matrix)
    topic_word = lda_model.topic_word_ # model.components_ also works
    doc_topic = lda_model.doc_topic_
    return doc_topic, cvec_titles, lda_model
```

```
if 'titles_lda' in features:
    cvec_matrix = cvec_titles.transform(each[date]['titles'])
    doc_topic = lda_model.transform(cvec_matrix)
```

Features

Number	Feature	Description	Type
Signal			
1	<i>news</i>	number of news mentions of e_i in $[0, t_p]$ in d_i	<i>Int</i>
2	<i>news d_{i-1}</i>	number of news mentions of e_i in $[0, t_p]$ in d_{i-1}	<i>Int</i>
3	<i>news total d_{i-1}</i>	number of news mentions of e_i in $[0, 24[$ in d_{i-1}	<i>Int</i>
4	<i>news titles</i>	number of title mentions in news of e_i in $[0, t_p]$ in d_i	<i>Int</i>
5	<i>avg content</i>	average content length of news of e_i in $[0, t_p]$ in d_i	<i>Float</i>
6	<i>sources</i>	number of different news sources of e_i in $[0, t_p]$ in d_i	<i>Int</i>
7	<i>weekday</i>	day of week	<i>Categ</i>
8	<i>is weekend</i>	true if weekend, false otherwise	<i>Bool</i>

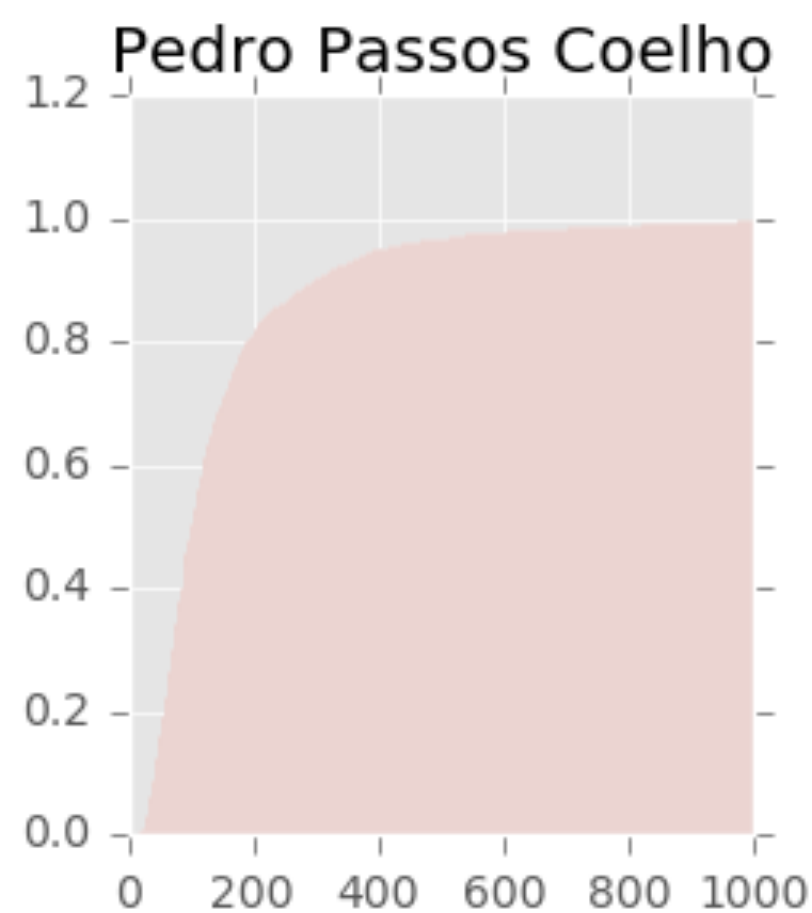
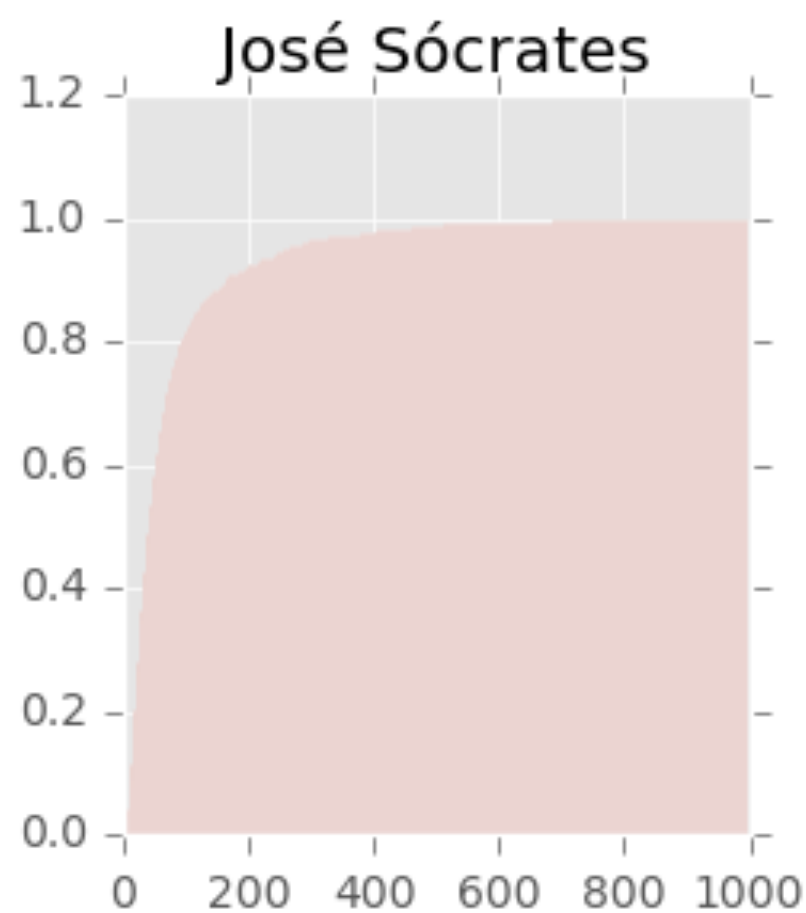
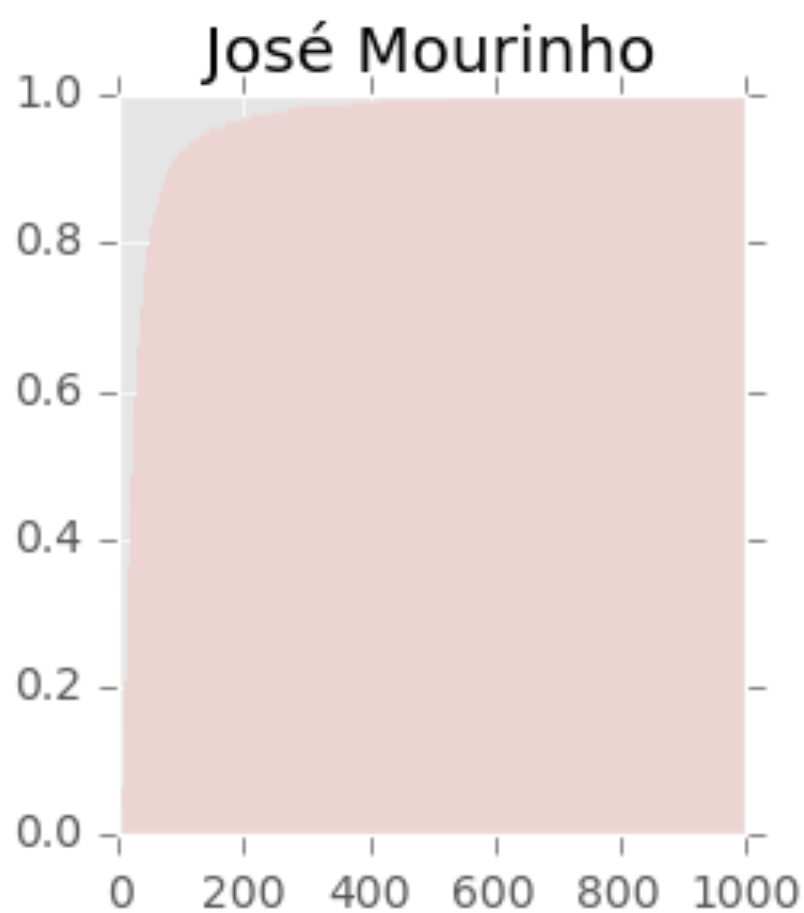
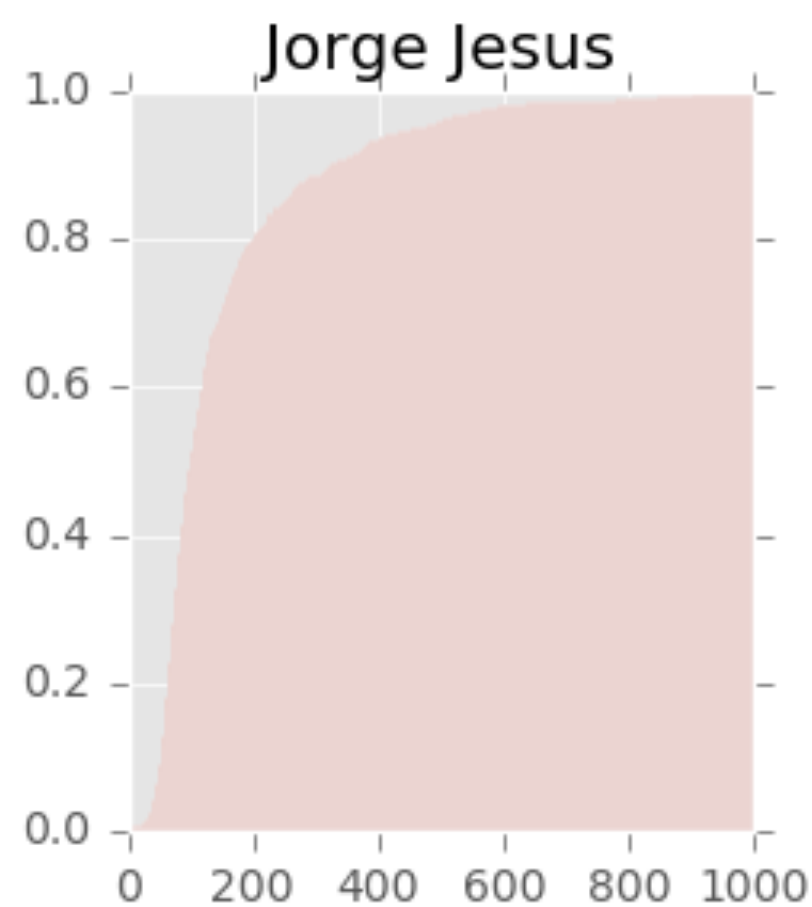
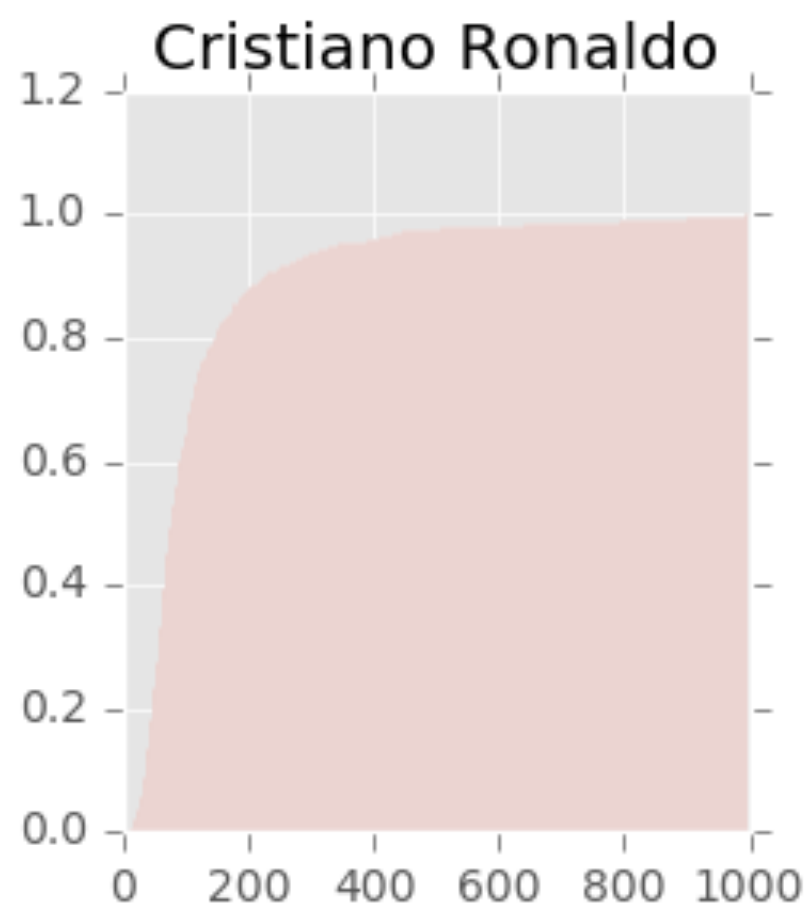
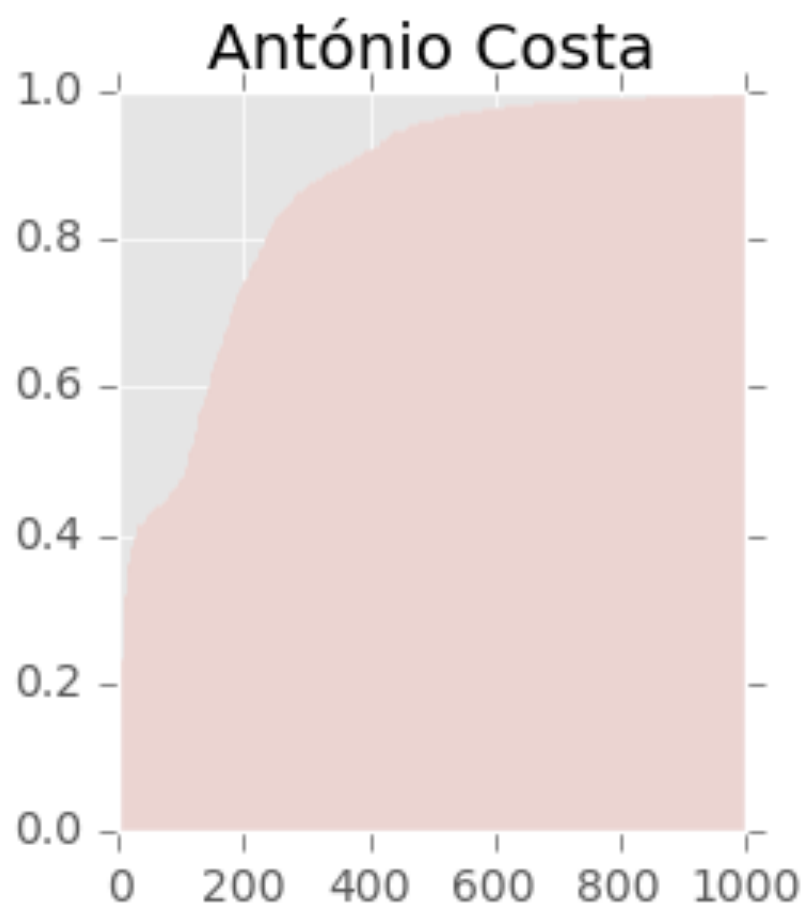
Features

Number	Feature	Description	Type
Textual			
9–18	<i>tfidf titles</i>	TF-IDF of news titles $[0, t_p]$ in d_i	<i>Float</i>
19–28	<i>LDA titles</i>	LDA-10 of news titles $[0, t_p]$ in d_i	<i>Float</i>
Sentiment			
29	<i>pos</i>	number of positive words in news titles $[0, t_p]$ in d_i	<i>Int</i>
30	<i>neg</i>	number of negative words in news titles $[0, t_p]$ in d_i	<i>Int</i>
31	<i>neu</i>	number of neutral words in news titles $[0, t_p]$ in d_i	<i>Int</i>
32	<i>ratio</i>	<i>positive/negative</i>	<i>Float</i>
33	<i>diff</i>	<i>positive – negative</i>	<i>Int</i>
34	<i>subjectivity</i>	$(positive + negative + neutral) / \sum words$	<i>Float</i>
35–44	<i>tfidf subj</i>	TF-IDF of subjective words (pos, neg and neu)	<i>Float</i>
Semantic			
45	<i>entities</i>	number of entities in news $[0, t_p]$ in d_i	<i>Int</i>
46	<i>tags</i>	number of tags in news $[0, t_p]$ in d_i	<i>Int</i>
47–56	<i>tfidf entities</i>	TF-IDF of entities in news $[0, t_p]$ in d_i	<i>Float</i>
57–66	<i>tfidf tags</i>	TF-IDF of news tags $[0, t_p]$ in d_i	<i>Float</i>

Prediction

$$\hat{f}_p = \begin{cases} 0(\text{low}), & \text{if } P(f_p(e_i, N, T) \leq \delta) = k \\ 1(\text{high}), & \text{if } P(f_p(e_i, N, T) > \delta) = 1 - k \end{cases}$$

δ is the inverse of cumulative distribution function at k of $f_p(e_i, N, T)$



Combining stuff

```
... X_train = np.hstack([X_train_signal, X_train_textual, X_train_sentiment, X_train_semantic])  
...  
... y_train = np.array(training_tweets[:, vec_feat.index("tweets_series")])  
... threshold_65 = np.percentile(y_train, 65.0)  
... y_train[y_train <= threshold_65] = 0  
... y_train[y_train > threshold_65] = 1
```


Train - Predict

```
def fitModel(X_train, y_train, classifier):  
    classifier.fit(X_train, y_train)  
    model = pickle.dumps(classifier)  
    return classifier  
  
def predictModel(classifier, X_test):  
    #classifier = pickle.loads(model_path)  
    y_predicted = classifier.predict(X_test)  
    return y_predicted
```

Predictors

`linear_model.ARDRegression` ([n_iter, tol, ...])

`linear_model.BayesianRidge` ([n_iter, tol, ...])

`linear_model.ElasticNet` ([alpha, l1_ratio, ...])

`linear_model.ElasticNetCV` ([l1_ratio, eps, ...])

`linear_model.HuberRegressor` ([epsilon, ...])

`linear_model.Lars` ([fit_intercept, verbose, ...])

`linear_model.LarsCV` ([fit_intercept, ...])

`linear_model.Lasso` ([alpha, fit_intercept, ...])

`linear_model.LassoCV` ([eps, n_alphas, ...])

`linear_model.LassoLars` ([alpha, ...])

`linear_model.LassoLarsCV` ([fit_intercept, ...])

`linear_model.LassoLarsIC` ([criterion, ...])

`linear_model.LinearRegression` ([...])

`linear_model.LogisticRegression` ([penalty, ...])

`linear_model.LogisticRegressionCV` ([Cs, ...])

`linear_model.MultiTaskLasso` ([alpha, ...])

`ensemble.AdaBoostClassifier` ([...])

`ensemble.AdaBoostRegressor` ([base_estimator, ...])

`ensemble.BaggingClassifier` ([base_estimator, ...])

`ensemble.BaggingRegressor` ([base_estimator, ...])

`ensemble.ExtraTreesClassifier` ([...])

`ensemble.ExtraTreesRegressor` ([n_estimators, ...])

`ensemble.GradientBoostingClassifier` ([loss, ...])

`ensemble.GradientBoostingRegressor` ([loss, ...])

`ensemble.IsolationForest` ([n_estimators, ...])

`ensemble.RandomForestClassifier` ([...])

`ensemble.RandomTreesEmbedding` ([...])

`ensemble.RandomForestRegressor` ([...])

`ensemble.VotingClassifier` (estimators[, ...])

`svm.SVC` ([C, kernel, degree, gamma, coef0, .

`svm.LinearSVC` ([penalty, loss, dual, tol, C, ...

`svm.NuSVC` ([nu, kernel, degree, gamma, ...])

`svm.SVR` ([kernel, degree, gamma, coef0, tol,

`svm.LinearSVR` ([epsilon, tol, C, loss, ...])

`svm.NuSVR` ([nu, C, kernel, degree, gamma, ..

`svm.OneClassSVM` ([kernel, degree, gamma, .

`svm.l1_min_c` (X, y[, loss, fit_intercept, ...])

Experimental Setup

Entity-Specific Models

2 years training set (+- 720 examples)

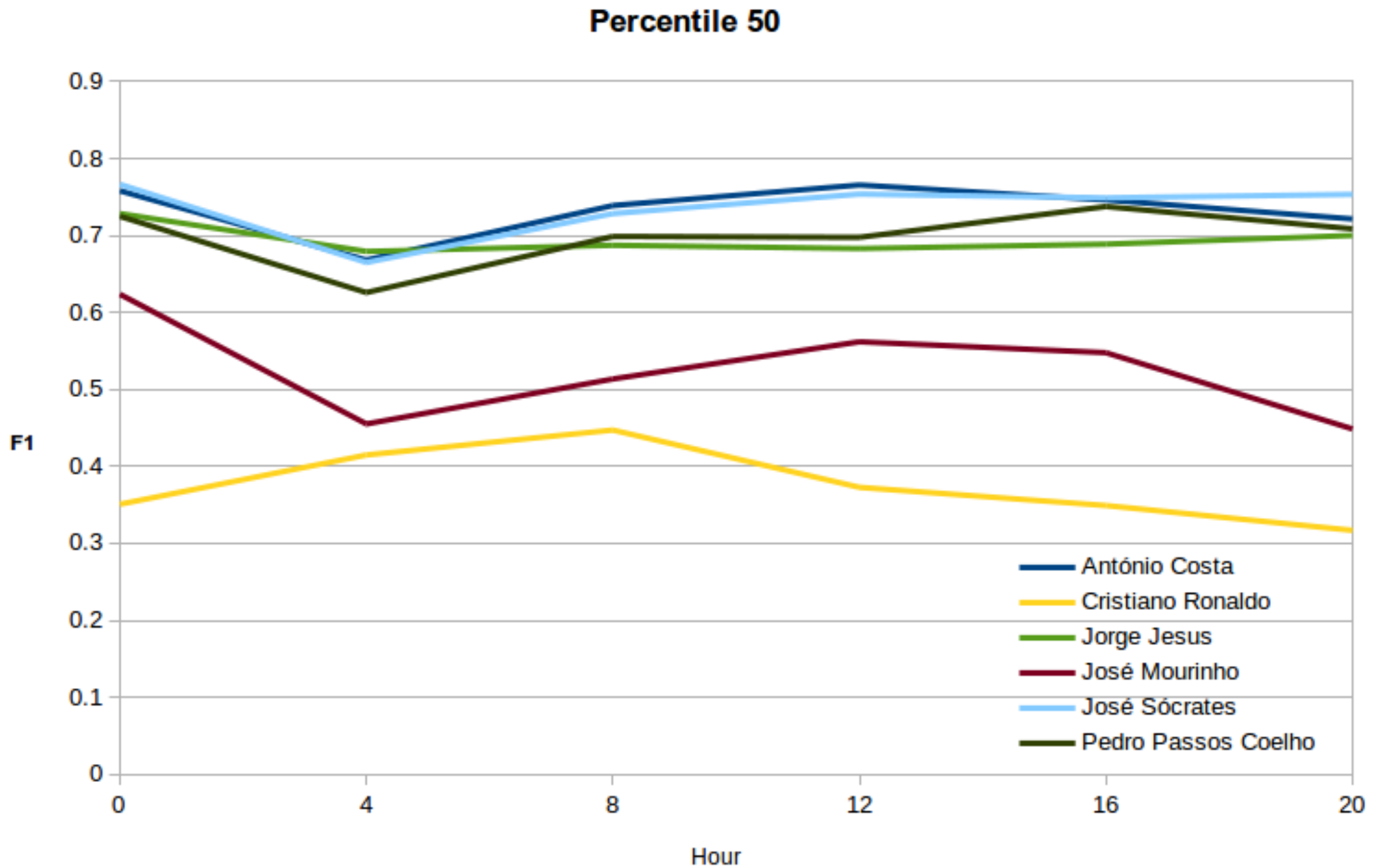
Monthly sliding window

Iteration 1	Training				Test	
	Jan 2013	Feb 2013	...	Dec 2014	Jan 2015	Feb 2015
Iteration 2		Training				Test
	Jan 2013	Feb 2013	Mar 2013	...	Jan 2015	Feb 2015

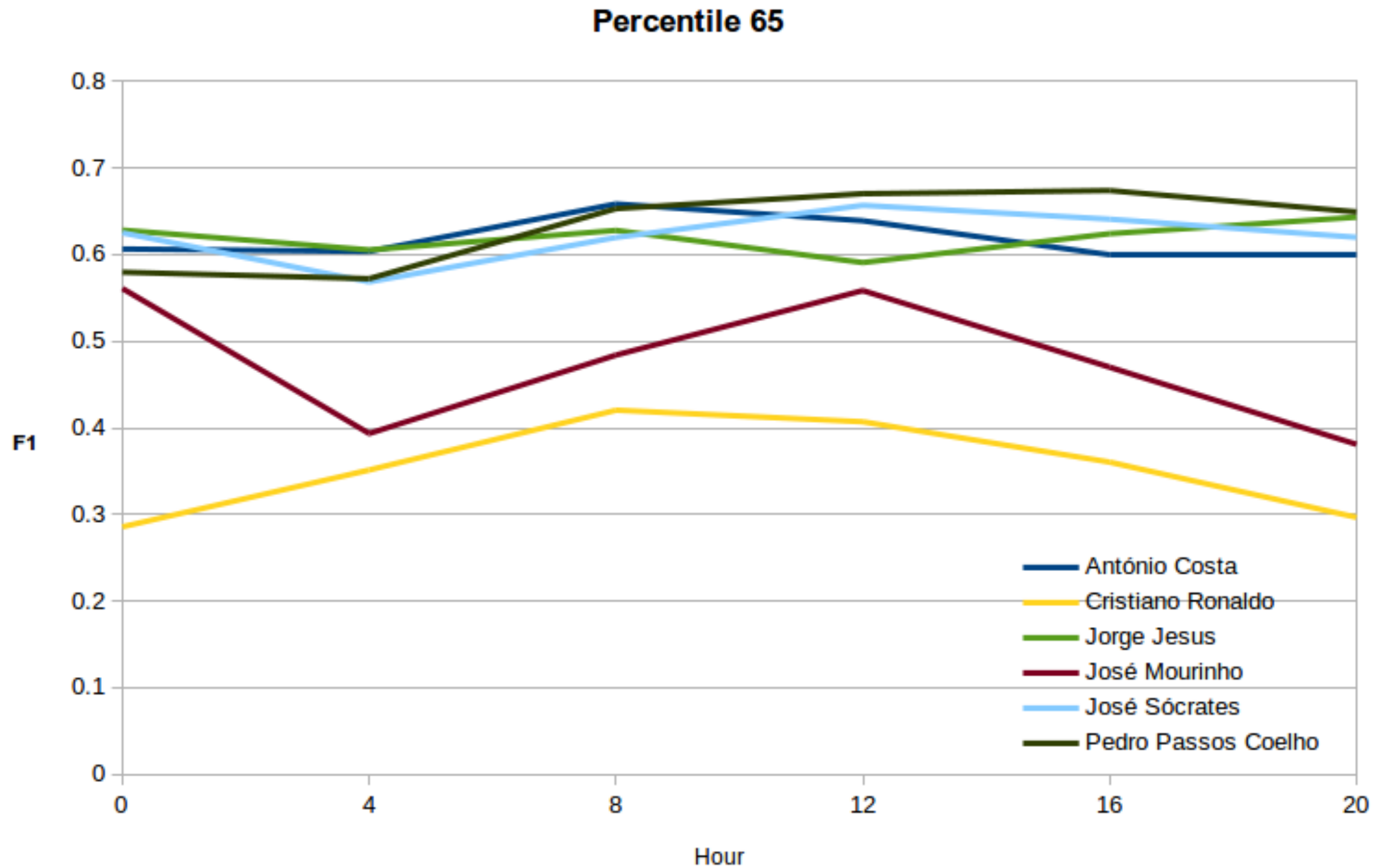
Evaluation

Scoring	Function
Classification	
'accuracy'	<code>metrics.accuracy_score</code>
'average_precision'	<code>metrics.average_precision_score</code>
'f1'	<code>metrics.f1_score</code>
'f1_micro'	<code>metrics.f1_score</code>
'f1_macro'	<code>metrics.f1_score</code>
'f1_weighted'	<code>metrics.f1_score</code>
'f1_samples'	<code>metrics.f1_score</code>
'neg_log_loss'	<code>metrics.log_loss</code>
'precision' etc.	<code>metrics.precision_score</code>
'recall' etc.	<code>metrics.recall_score</code>
'roc_auc'	<code>metrics.roc_auc_score</code>
Clustering	
'adjusted_rand_score'	<code>metrics.adjusted_rand_score</code>
Regression	
'neg_mean_absolute_error'	<code>metrics.mean_absolute_error</code>
'neg_mean_squared_error'	<code>metrics.mean_squared_error</code>
'neg_median_absolute_error'	<code>metrics.median_absolute_error</code>
'r2'	<code>metrics.r2_score</code>

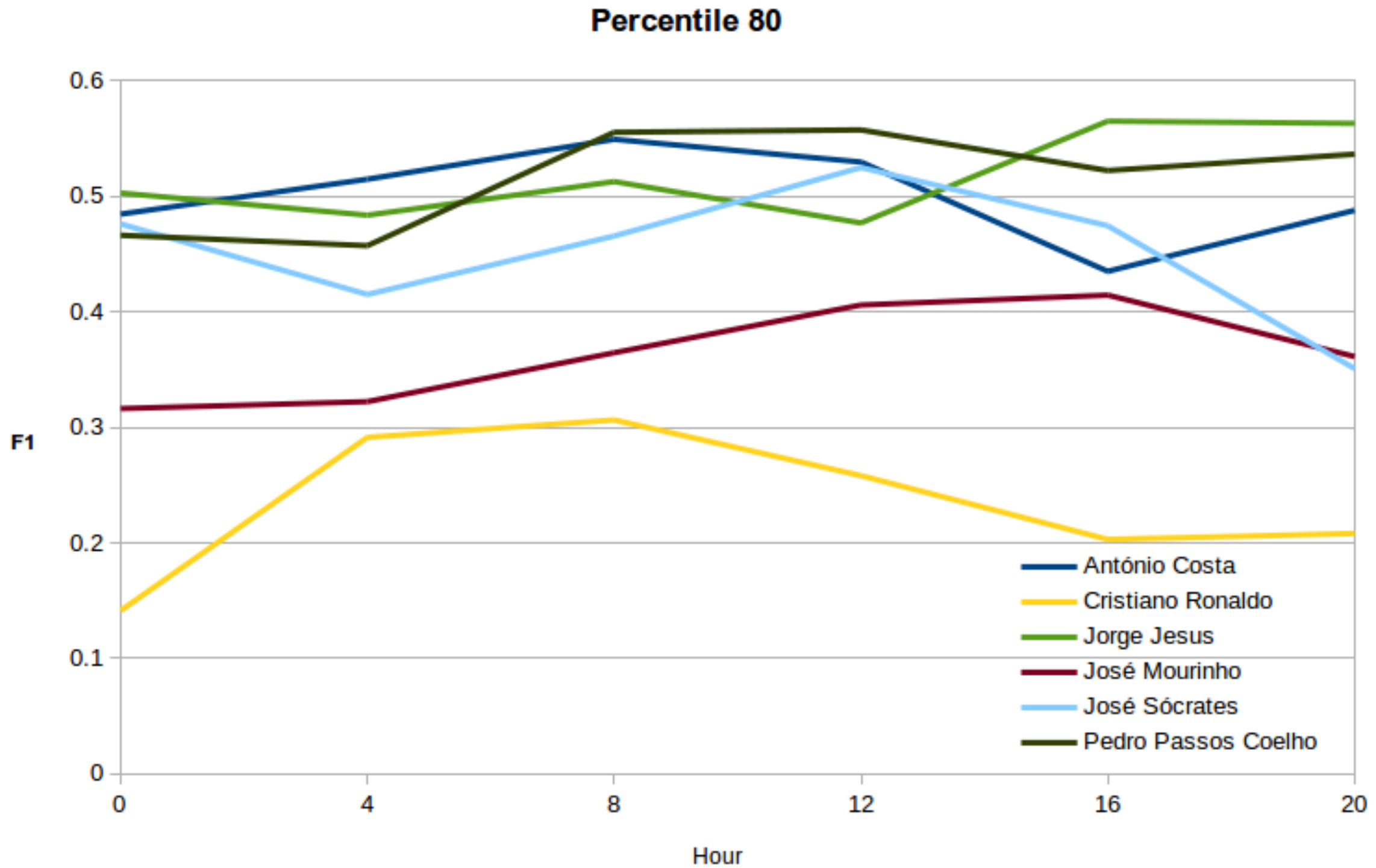
Results



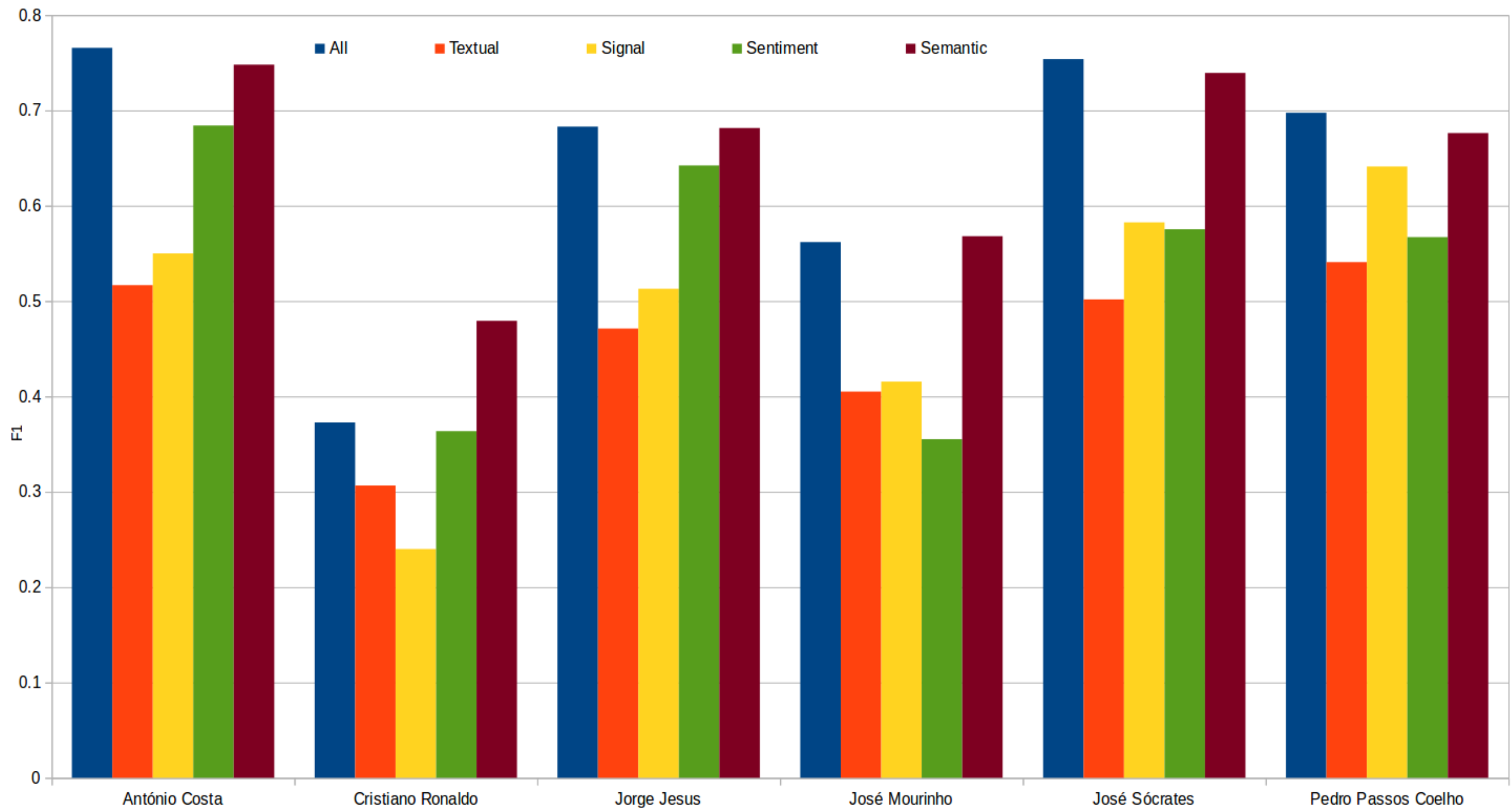
Results



Results



Results $k=0.5$, $tp=12:00$



Results

F1 scores above 0.7 for politicians and balanced dataset ($k=0.5$)

Moment of prediction influences performance

Semantic and Sentiment are the most informative type of features

Not so good results for live events, e.g. TV debates

Questions?

pssc[at]fe.up.pt