

Fundação **Casa de Rui Barbosa**

3º Seminário Tecnologia e Cultura

convergência entre acervos digitais de arquivos, bibliotecas e museus.

Iniciativas em Humanidades Digitais

Renato Rocha Souza
renato.souza@fgv.br

Using artificial intelligence to identify state secrets



A Twofold Mission

- Applied research to help government explore and manage sensitive information
- Tools to help people understand what information is released, and what needs to be protected



HISTORY LAB



History as Data Science

We turn documents into data and develop tools to explore history.

A Coordinated Approach

A Research &
Development Team of
Data Scientists, Social
Scientists, Engineers, and
Web Developers, and
Stakeholders



The Biggest Database of Declassified Documents

- **The Foreign Relations of the United States (1945-1980).** A curated collection of the ~80,000 most important declassified documents selected by State Department historians with access to every government department and agency.
- **The State Department Central Foreign Policy Files (1973-1978).** 1.7 million State Department Cables and metadata from ~500,000 more still classified cables and documents delivered by diplomatic pouch.
- **Henry Kissinger Telephone Conversations (1973-1976).** 4.5 thousand transcripts of Kissinger Telephone Conversations during his tenure as Secretary of State.
- **The Hillary Clinton Emails (2009-2012).** As of November release, 16,246 email chains with a total of 40,737 individual messages
- **Other Collections To Come:**
 - President's Daily Briefs
 - FBI "Vault" of FOIA'ed documents
 - NATO Archives
 - Aramis (UK Foreign Office Cables 1992-2000)

IV: Predicting Classification

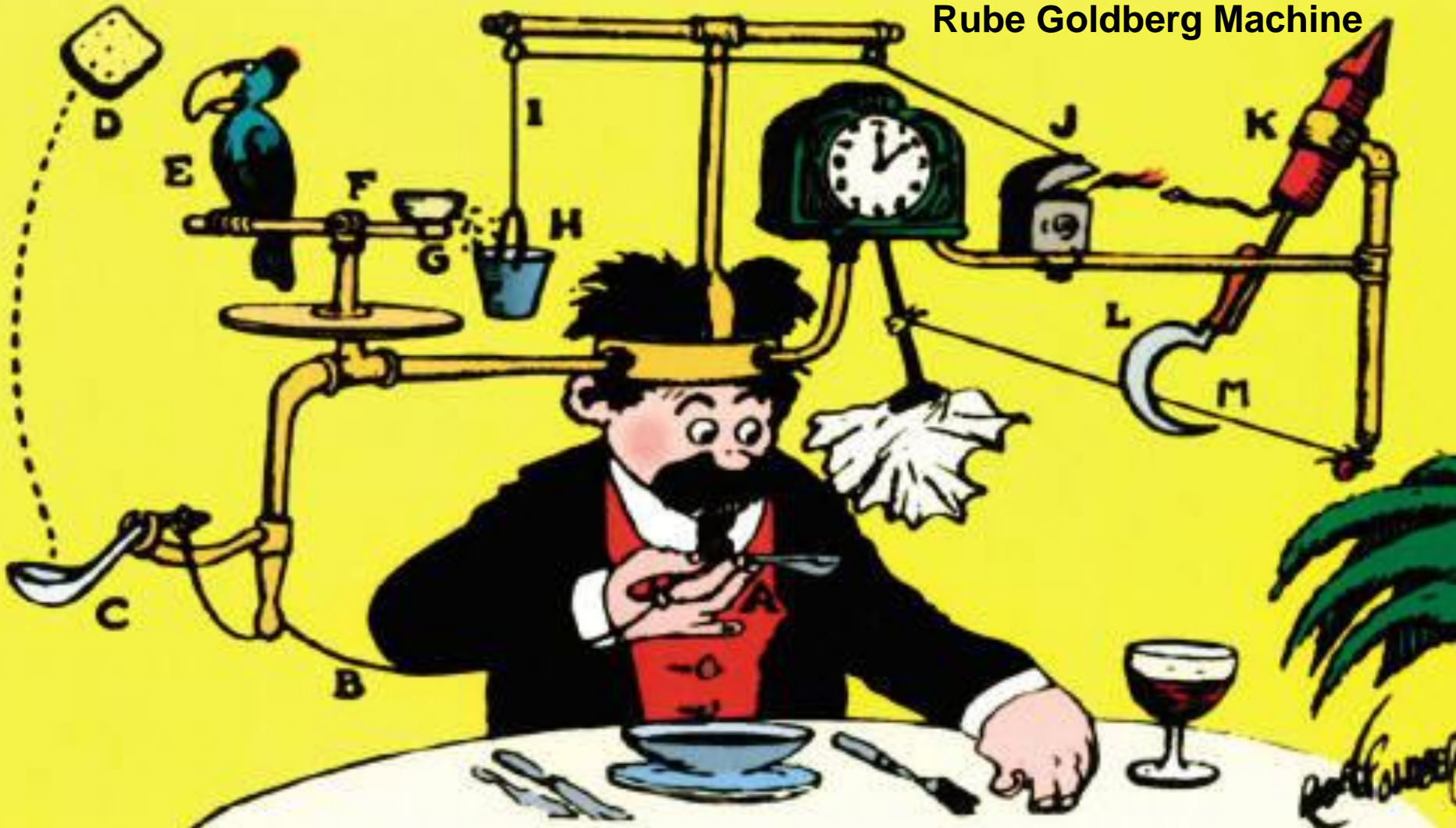


- Data from State Dept. Cables collection, 1973-1978
- Use Machine Learning techniques to predict the classification levels, which are tagged in the documents

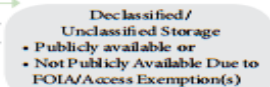
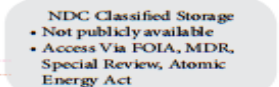
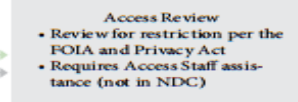
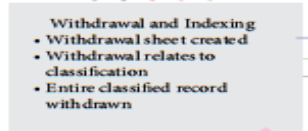
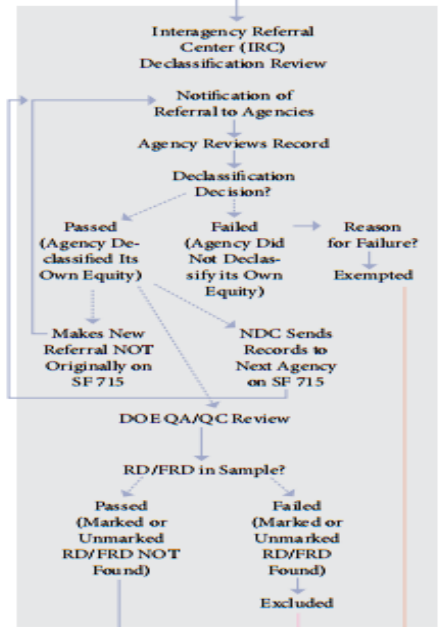
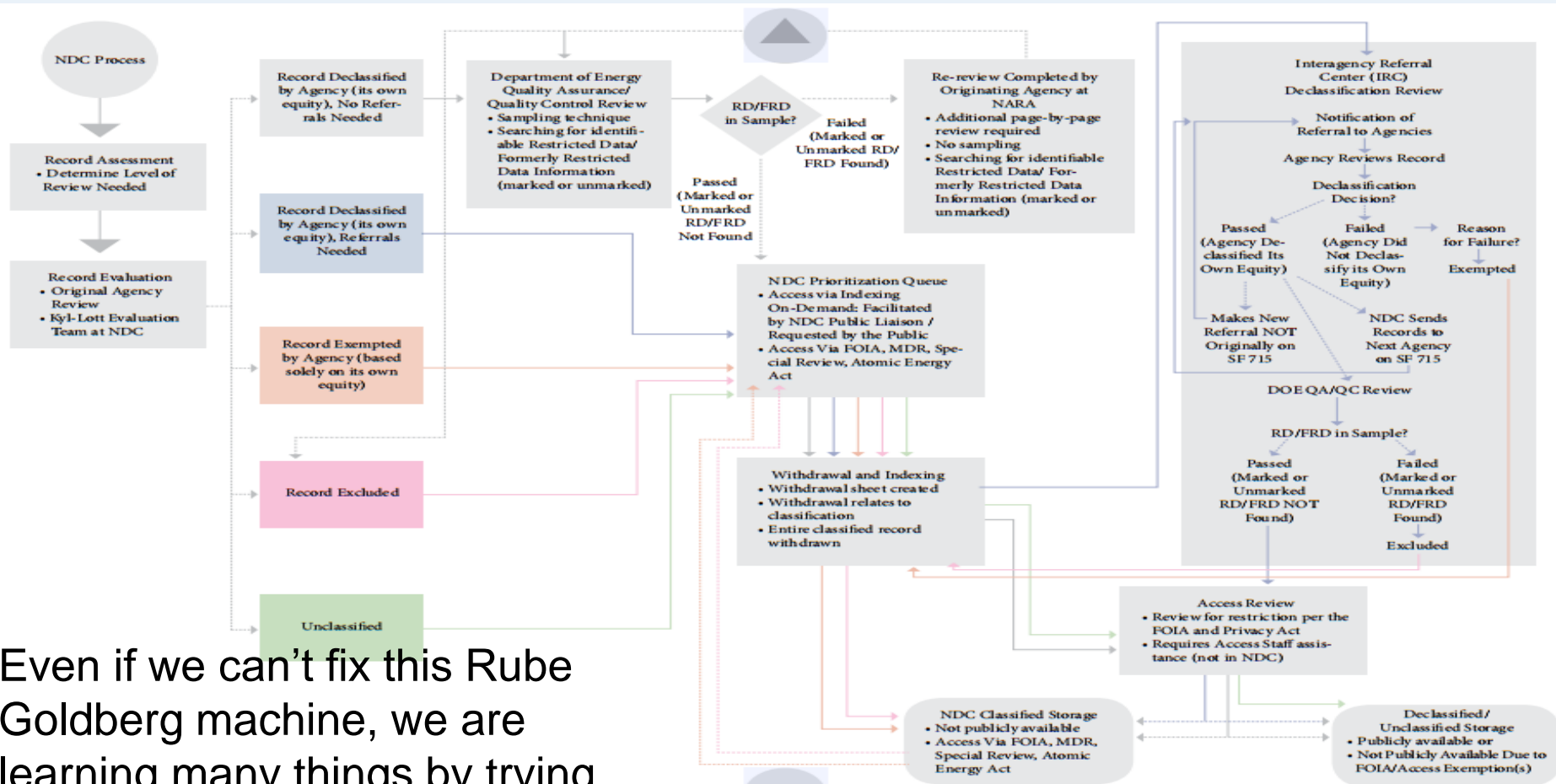




Rube Goldberg Machine



Our Focus: Tools to Manage Classification and Declassification Risks



PROCESS RESULTS: 60% DECLASSIFICATION RATE

Even if we can't fix this Rube Goldberg machine, we are learning many things by trying

Identifying features

- ① File Designation
- ② Action Distribution
- ③ Information Distribution
- ④ Security Classification
- ⑤ Airgram Number
- ⑥ Addressees
- ⑦ Originating Post
- ⑧ Preparation Date
- ⑨ Subject
- ⑩ References
- ⑪ Text

DEPARTMENT OF STATE
AIRGRAM *AD(US)15-6 the Congo*

Original to be Filed in _____ Disseminated File _____ FILE DESIGNATION _____

UNCLASSIFIED A-377

TO SECRETARY OF STATE
FOR AGRICULTURE

FROM American Embassy, Kinshasa DATE September 14, 1966

SUBJECT PL AEO Usual Marketing Requirement - Dairy Products

REF (1) A-18, September 6, 1966 (2) AGR-4, August 22, 1966

Import statistics are not available here which show breakdown between U.S. noncommercial sales and commercial shipments. However our check of figures submitted in ref. 2 indicates that if we subtract all U.S. shipments of any nature, the Congo still imported 1,933.1 metric tons of dairy products in the first half of FY 66. All non-U.S. shipments were commercial. Therefore Congo has completed its 1800 M.T. usual marketing requirement for dairy products. We assume no further information on this point is required.

FOR THE AMBASSADOR
John F. McMillan
Acting Agricultural Attaché

UNCLASSIFIED

Original to be Filed in _____ Disseminated File _____

Approved by: _____ Date: 9/14/66

Checked by: _____ Date: _____

01000

① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫

Used Fields	Description
<i>origclass</i>	The original class of the cable (the classification target)
<i>body</i>	Full text of the cable
<i>subject</i>	Keywords of subjects dealt with in the document.
<i>concepts</i>	Concepts attributed to the document
<i>TAGS</i>	Traffic Analysis by Geography and Subject
<i>from</i>	Who/where sent the document.
<i>to</i>	Who/where received the document.
<i>office</i>	Which State Department office or bureau sent the document.
<i>date</i>	Document creation date

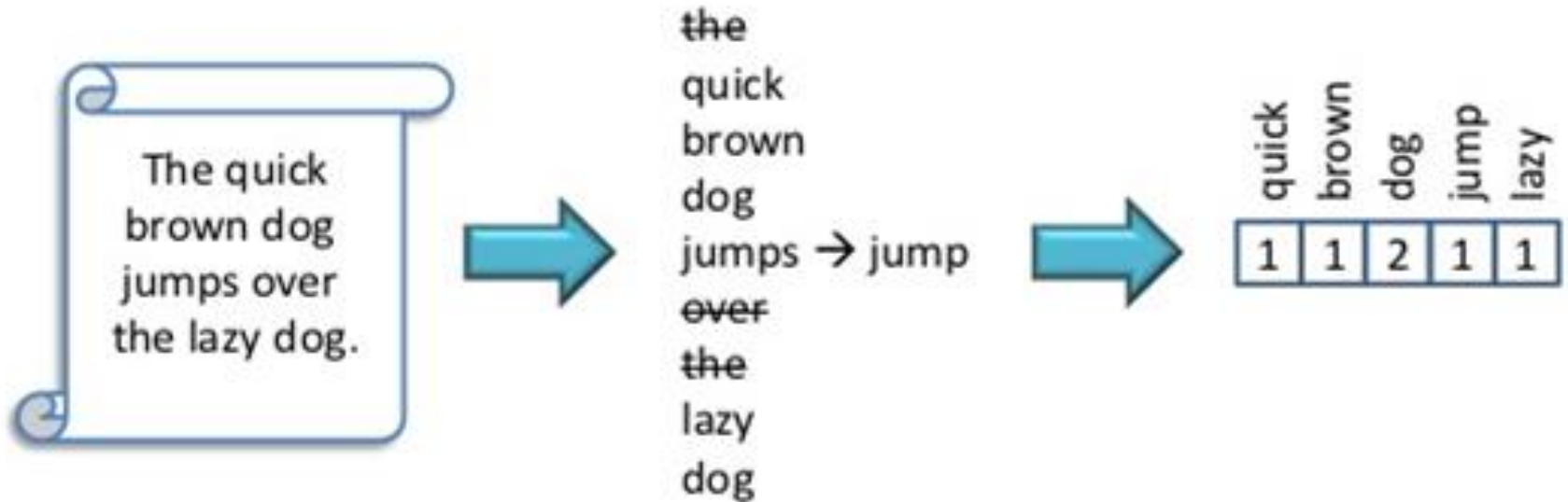
Situation	Total in Database	Unclassified	Limited Official Use	Confidential	Secret
declassified cables	1.758.279	876.797	411.973	375.690	93.635
Error messages for <i>body</i>	119.744	53.935	21.744	25.233	18.832
blank <i>body</i>	8282	2.726	1.645	1.924	1.987
blank or n/a <i>concepts</i>	634.967	445.300	114.507	65.502	9.658
blank or n/a <i>subject</i>	26.109	16.490	5.820	2.914	885
blank or n/a <i>from</i>	17	7	6	3	1
blank or n/a <i>to</i>	9.740	6.027	1.572	1.698	443
Used for classifier	981.083	368.043	280.251	270.477	62.312

Feature Engineering

- Hyphenation was eliminated from textual fields, as there were garbage from the original printed versions that were scraped from the web;
- Compound names of places in textual fields were aggregated, enabling them to be treated as a single token (i. e. NEW YORK was transformed to NEWYORK). They were present in all textual features, but that step was specially important in the case of the *from* and *to* fields, which represent the names of the embassies. These fields were aggregated under a new field *embassy*, for the vectorization process;
- Tokenization was made and all the trailing punctuation and words with length of 1 were eliminated. Underscores and hyphens in the middle of words were maintained;
- Stopwords were removed, using NLTK english stopwords list;
- Tests were made using stemmed forms of words, but it didn't enhance the performance and the stemming was discarded.
- The field *date* was transformed in a boolean field *weekday* - indicating whether the date fell in a weekend or not; and another field *year+month*, used to test hypothesis on the temporal series of cables regarding to classification windows. That doesn't prove useful, though, for the whole span, although could be promising for small periods of time.
- The fields *body*, *subject*, *concepts*, *tags*, *embassy* and *office* were added altogether in a new field/feature *all_text*, which was tested as an alternative to combining all the other features by concatenating those vectors, with very similar results.

Bags of words

- Tokenize
- Remove stop words
- Lemmatize
- Compute weights

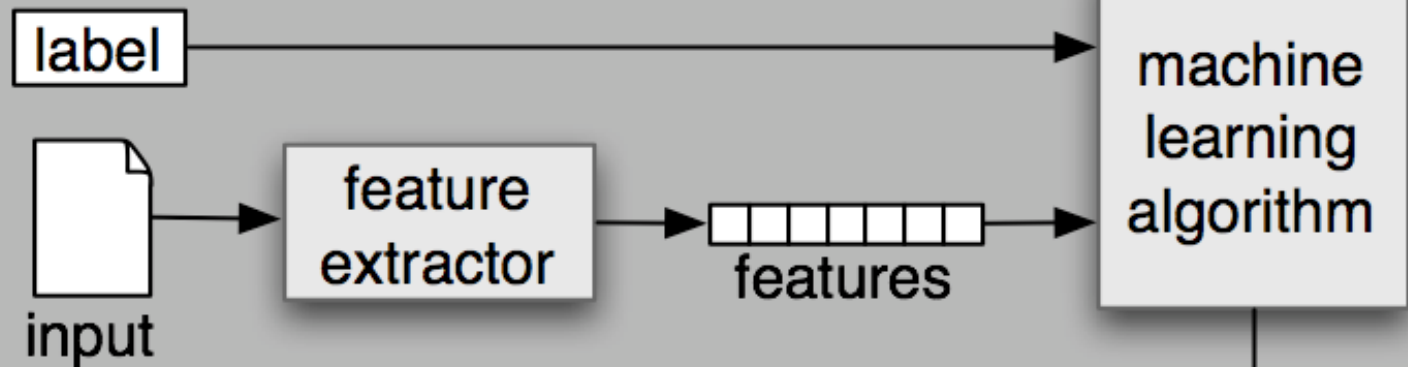


We have tested many alternate weighting schemes, as TfIDf

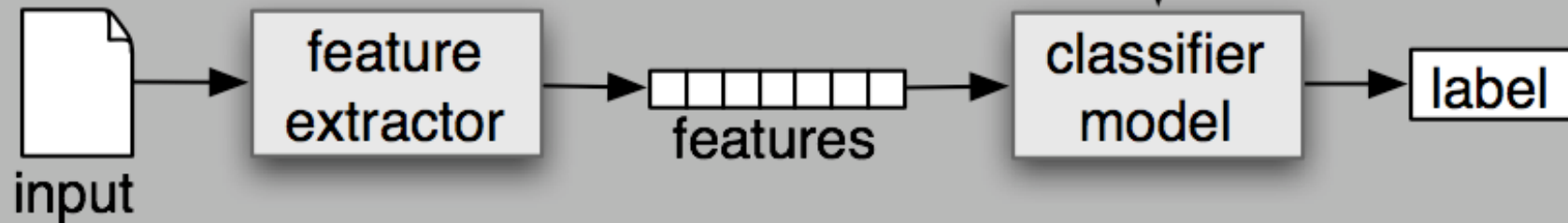
charles bailey WAS indicted for feloniously stealing on **the** 29th of december two dressed deer skins value 20 s **the** property of samuel savage **and** richard savage richard savage **i** am a leather seller 63 chinwell street my partner S name is samuel savage a few days previous to **the** 29th of december **i** looked out seventy skins for an order these skins being of a bad colour **i** directed them to be brimstoned to make them of equal colour pale on **the** 29th in **the** afternoon **i** saw them all smooth on a horse a few hours afterwards they appeared very much tumbled **and** one WAS thrown into **the** yard **and** dirtied **i** caused them to be brought in **the** warehouse **and** counted there WAS two gone our foreman went to workshop street **and** brought armstrong **and** vickrey they searched **and** found this skin in **the** prisoner S breeches **and** **the** other skin WAS found in **the** workshop carter **i** am foreman to samuel **and** richard savage **the** seventy skins **i** WAS with mr SAVAGE looking them out **i** took them out of **the** stove **and** counted them on **the** horse **and** on friday **i** counted them three times over there were no more than sixty eight instead of seventy **i** went to workshop street brought mr armstrong **and** vickrey with me they waited till **the** men left work **and** when they came down they were searched **and** on **the** prisoner one skin WAS found john armstrong **i** went to this gentlemans S house after **the** men came down vickrey **and** **i** were searching in one minute vickrey called me **i** received this skin from him it WAS taken out of **the** prisoner S breeches **i** have had it ever since john vickrey q you were with armstrong

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

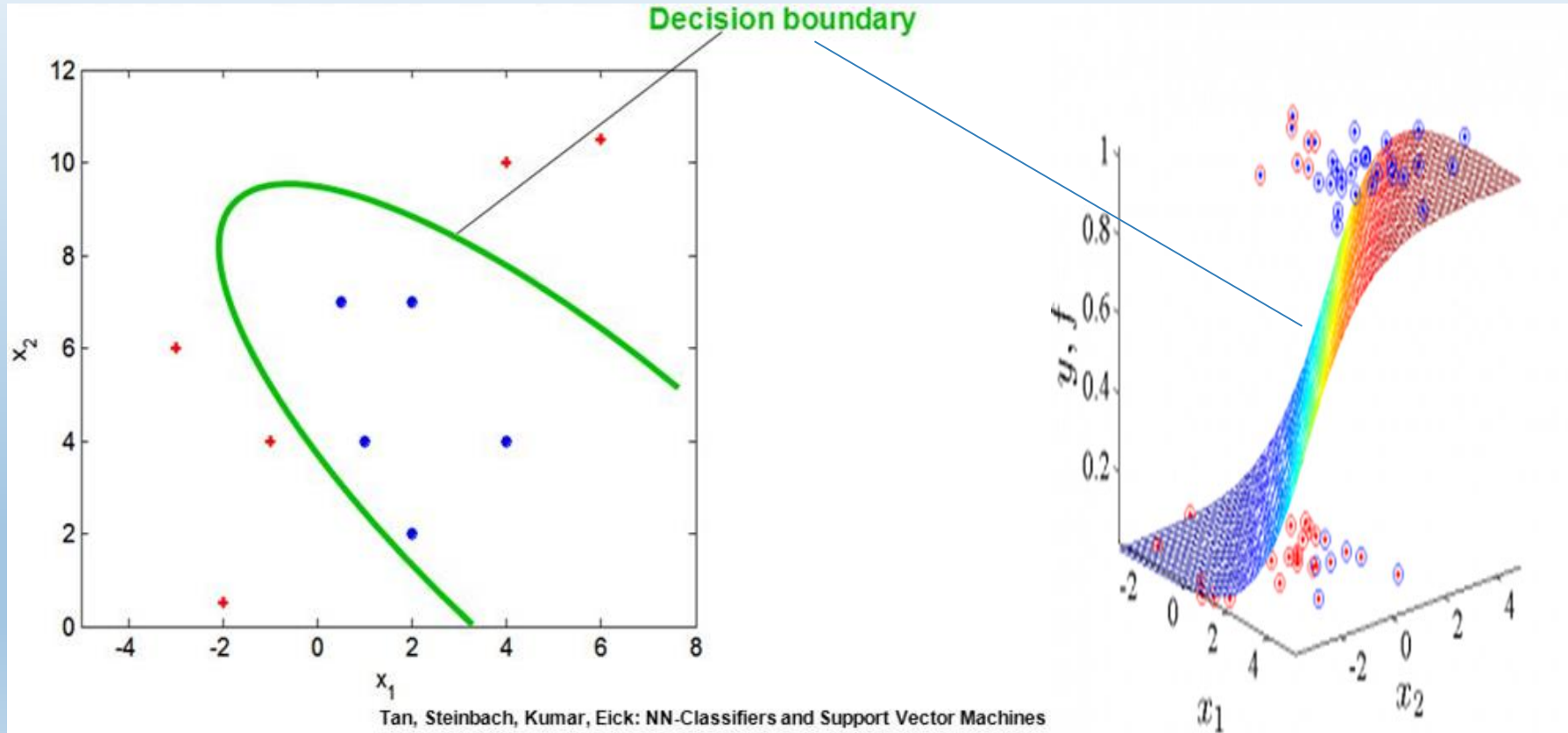
(a) Training



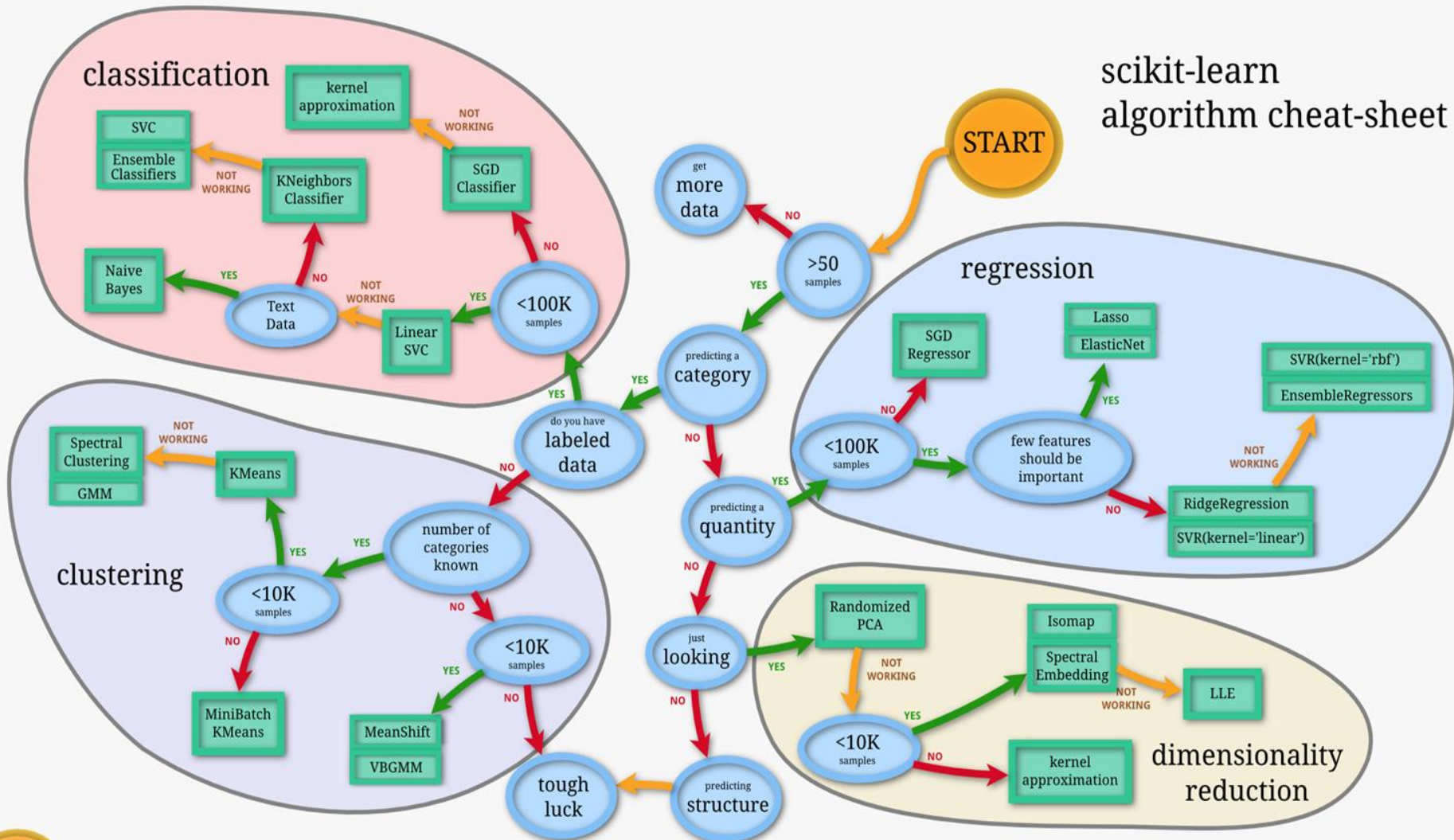
(b) Prediction



Classification in AI can be seen as analogous to “learning good decision boundaries” that separate the examples belonging to different classes in the data set

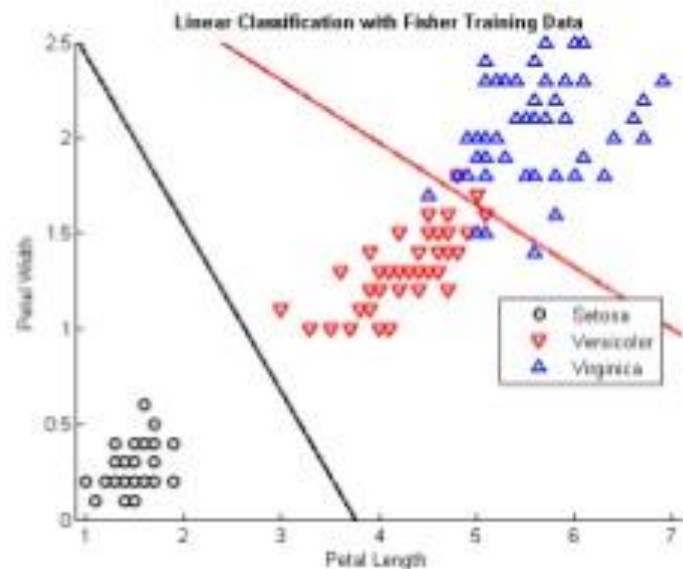


scikit-learn algorithm cheat-sheet

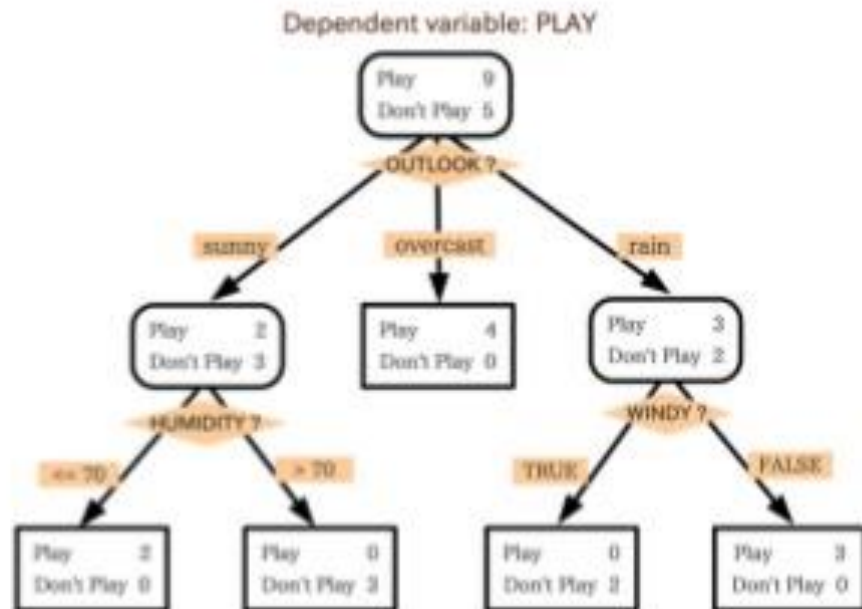


Algorithms

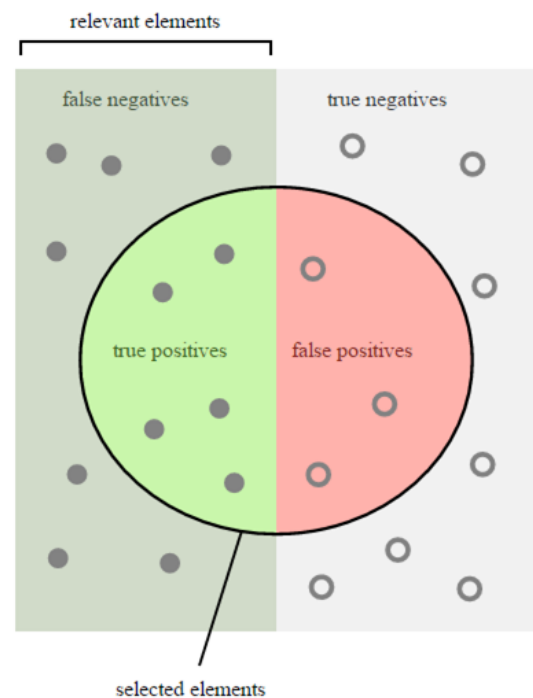
Linear Models



Decision Trees



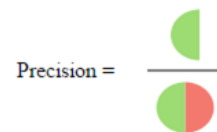
		Predicted condition	
		Predicted Condition positive	Predicted Condition negative
Total population		Predicted Condition positive	Predicted Condition negative
True condition	condition positive	True positive	False Negative (Type II error)
	condition negative	False Positive (Type I error)	True negative
Accuracy (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$		Positive predictive value (PPV), Precision $= \frac{\Sigma \text{ True positive}}{\Sigma \text{ Test outcome positive}}$	False omission rate (FOR) $= \frac{\Sigma \text{ False negative}}{\Sigma \text{ Test outcome negative}}$
		False discovery rate (FDR) $= \frac{\Sigma \text{ False positive}}{\Sigma \text{ Test outcome positive}}$	Negative predictive value (NPV) $= \frac{\Sigma \text{ True negative}}{\Sigma \text{ Test outcome negative}}$



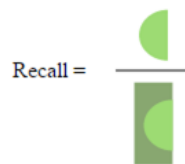
$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$$

$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

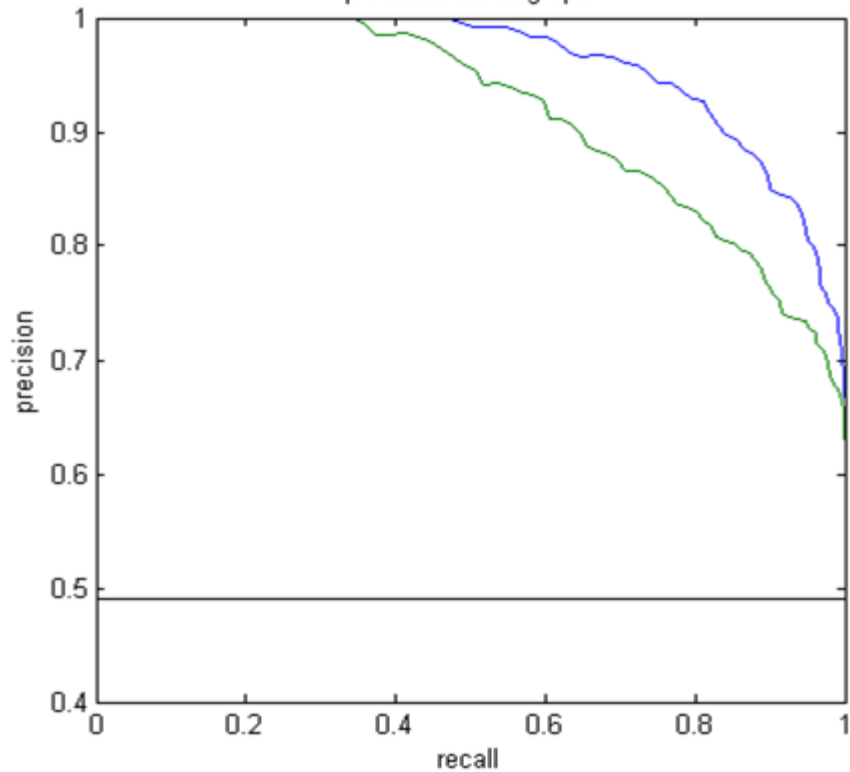
How many selected items are relevant?



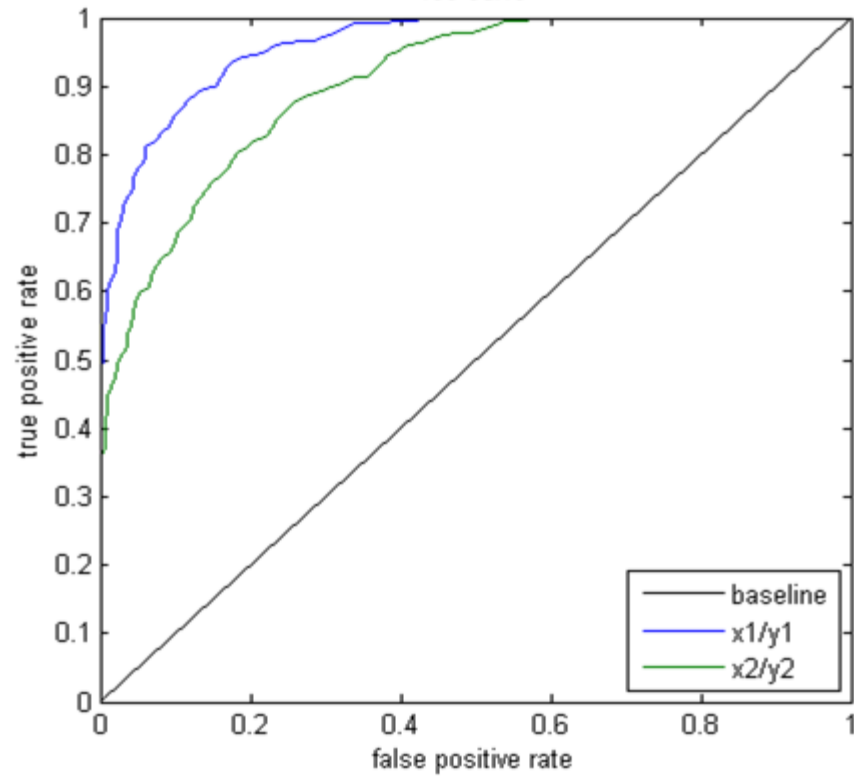
How many relevant items are selected?



precision-recall graph

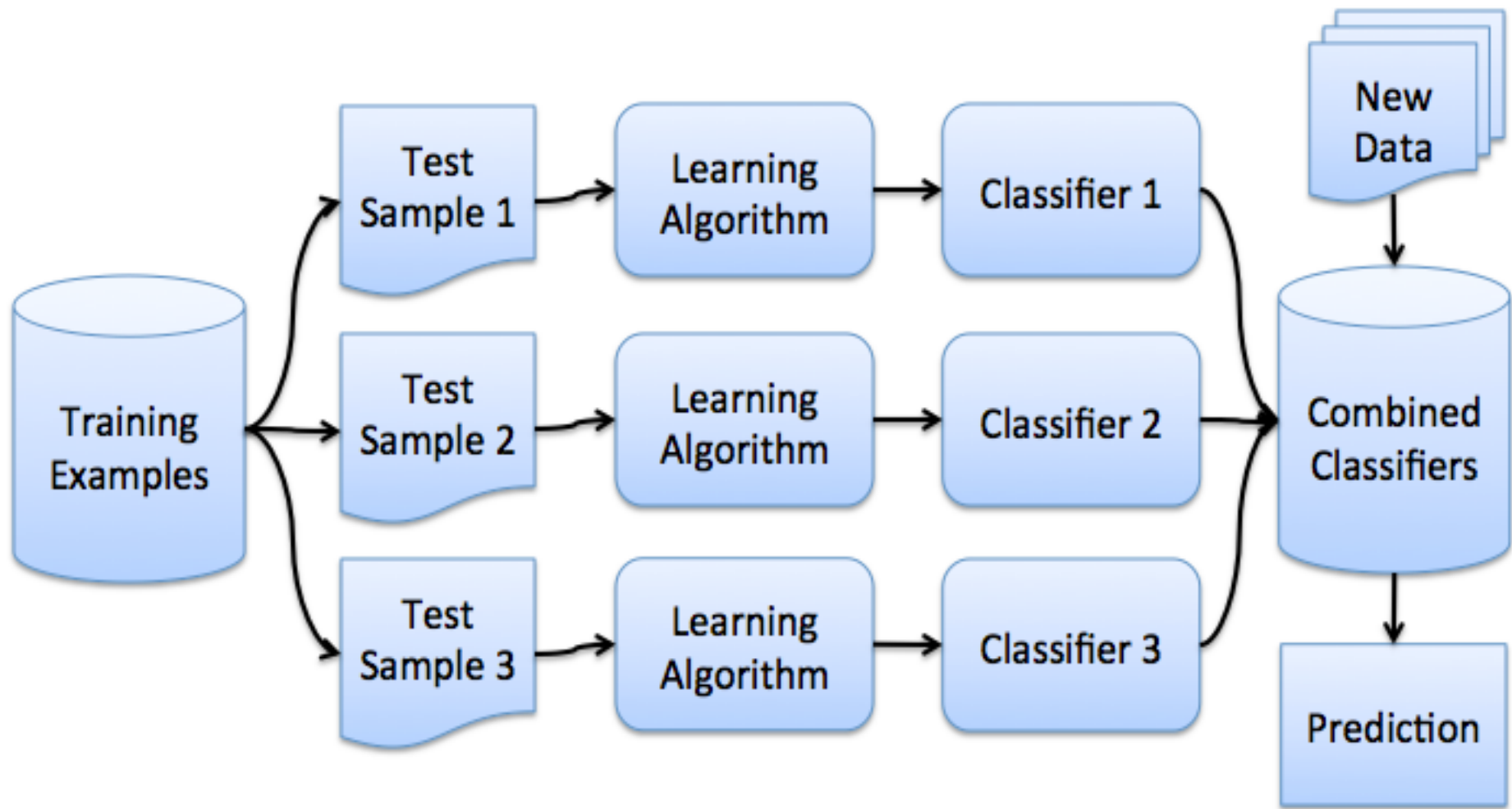


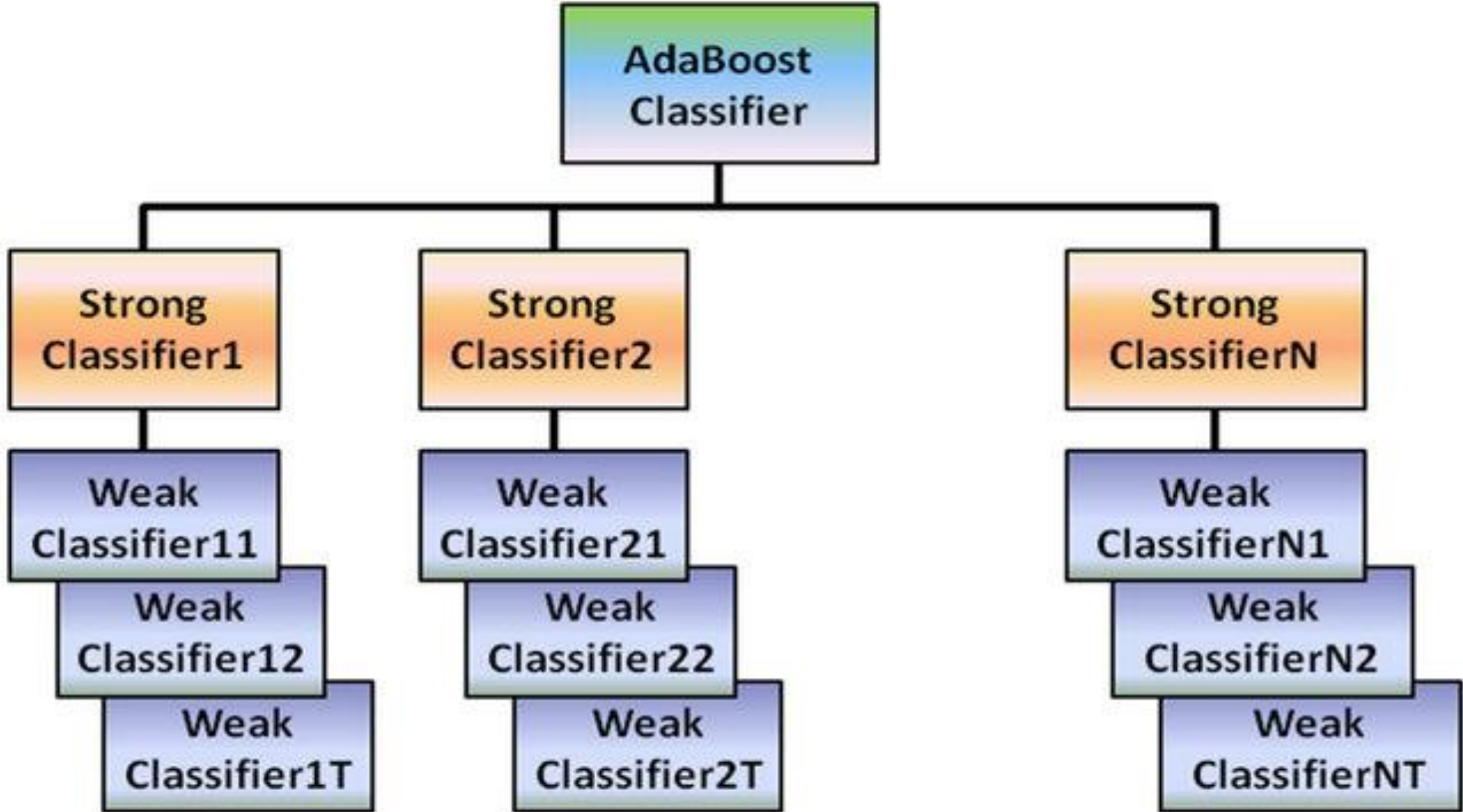
roc curve



Very ordinary results at first... ~0.75

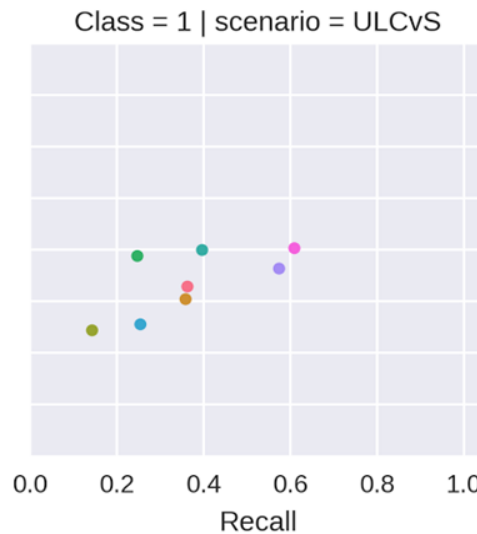
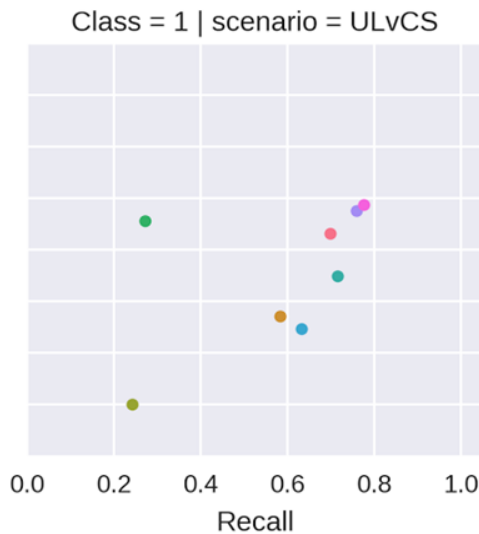
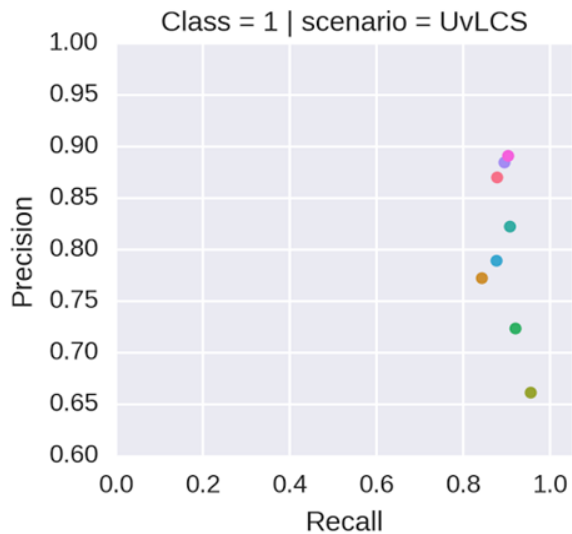
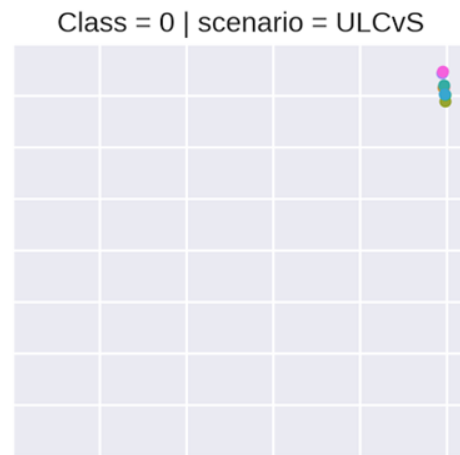
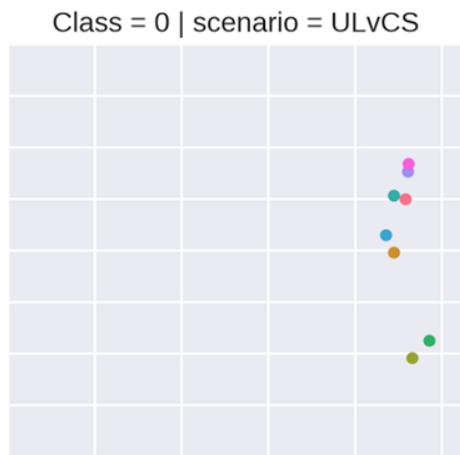
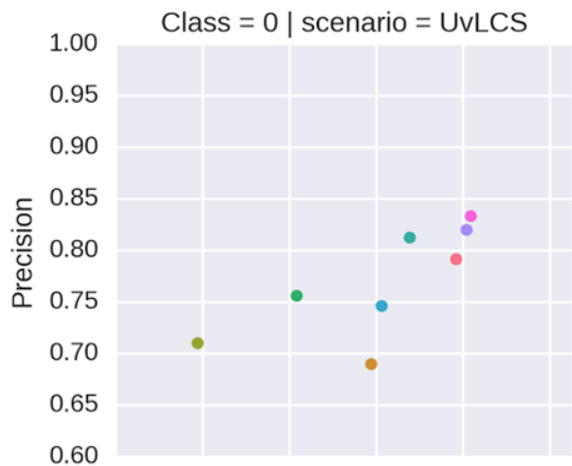
- Simple classifiers
- Scarce feature engineering
- Few data cleansing



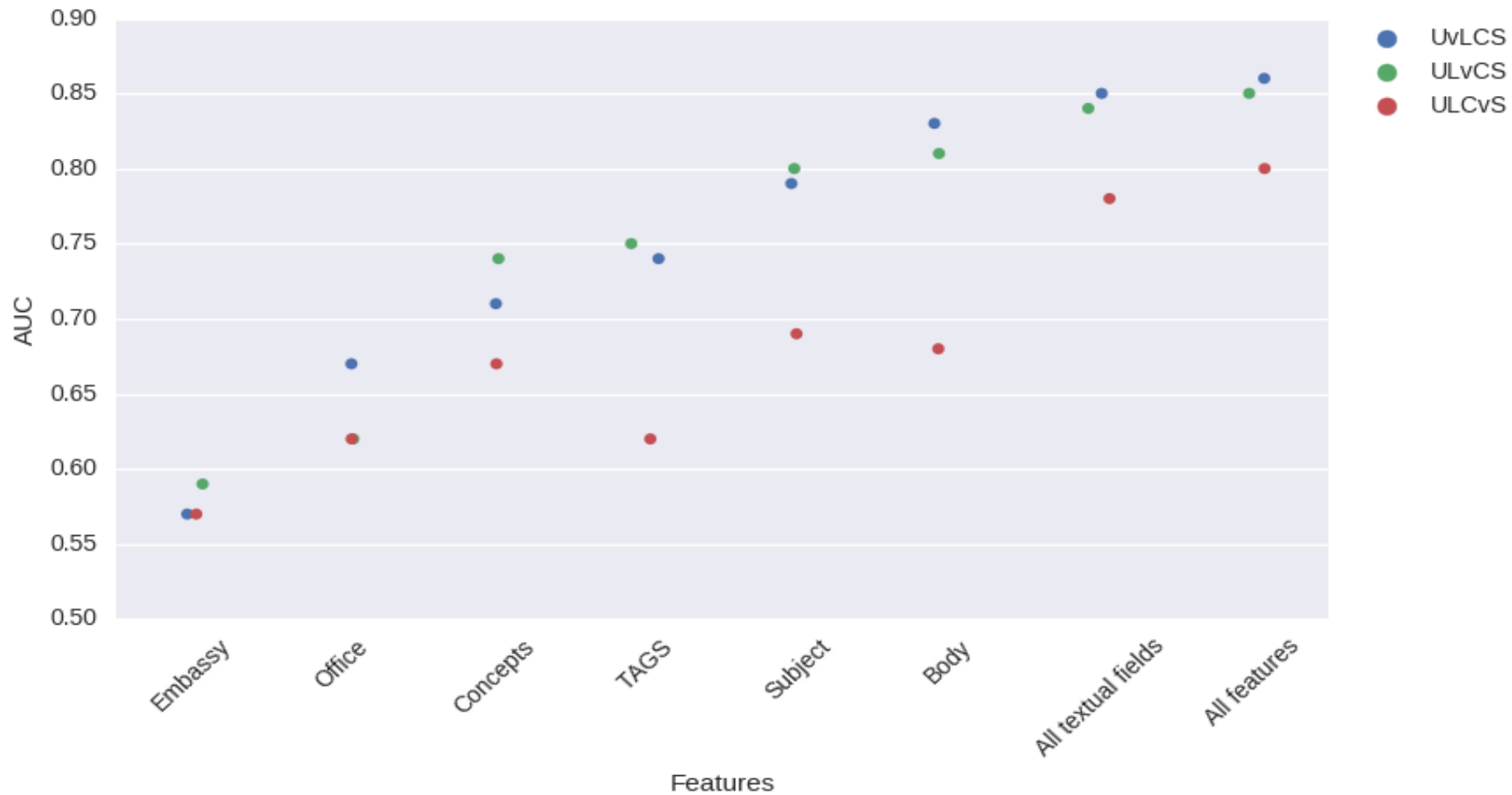


Classifier	ROC/AUC Score	Accuracy Score	Precision (class 0/1)	Recall (class 0/1)	f1-score (class 0/1)
Stochastic Gradient Descent	0.8462/ 0.8338	0.8576/ 0.8512	(0.82/0.88)/ (0.83/0.86)	(0.80/0.89)/ (0.76/0.90)	(0.81/0.89)/ (0.79/0.88)
Logistic Regression	0.8457/ 0.8434	0.8569/ 0.8574	(0.81/0.88)/ (0.82/0.88)	(0.80/0.89)/ (0.79/0.90)	(0.81/0.89)/ (0.81/0.89)
Linear SVM*	0.8454/ 0.8452	0.8563/ 0.8583	(0.81/0.88)/ (0.82/0.88)	(0.80/0.89)/ (0.79/0.90)	(0.81/0.89)/ (0.81/0.89)
Ridge	0.8261/ 0.8373	0.8448/ 0.8546	(0.82/0.86)/ (0.83/0.87)	(0.75/0.90)/ (0.77/0.91)	(0.78/0.88)/ (0.80/0.89)
Bagging (w/ Dec. Tree)	0.8048/ 0.8049	0.8172/ 0.8173	(0.76/0.85)/ (0.76/0.85)	(0.76/0.85)/ (0.76/0.85)	(0.76/0.85)/ (0.76/0.85)
Extremely Randomized Trees	0.8036/ 0.7938	0.8365/ 0.83	(0.86/0.83)/ (0.86/0.82)	(0.67/0.94)/ (0.65/0.94)	(0.76/0.88)/ (0.74/0.87)
AdaBoost (w/ Random F.)	0.8031/ 0.8072	0.8190/ 0.8222	(0.77/0.85)/ (0.77/0.85)	(0.74/0.87)/ (0.75/0.87)	(0.75/0.86)/ (0.76/0.86)
Random Forest	0.7964/ 0.7994	0.8310/ 0.8316	(0.86/0.82)/ (0.85/0.82)	(0.66/0.94)/ (0.67/0.93)	(0.74/0.87)/ (0.75/0.87)
Perceptron*	0.7856/ 0.7963	0.8138/ 0.8112	(0.80/0.82)/ (0.75/0.84)	(0.67/0.90)/ (0.74/0.86)	(0.73/0.86)/ (0.75/0.85)
Passive Aggressive*	0.7745/ 0.8226	0.8095/ 0.837	(0.82/0.81)/ (0.79/0.86)	(0.63/0.91)/ (0.76/0.88)	(0.71/0.86)/ (0.78/0.87)
Multinomial Naive Bayes	0.7735/ 0.7828	0.7614/ 0.7992	(0.64/0.87)/ (0.74/0.83)	(0.82/0.73)/ (0.72/0.85)	(0.72/0.79)/ (0.73/0.84)
Bernoulli Naive Bayes	0.6885/ 0.6885	0.6538/ 0.6538	(0.52/0.84)/ (0.52/0.84)	(0.83/0.55)/ (0.83/0.55)	(0.64/0.66)/ (0.64/0.66)

Feature	Class Combination	ROC/AUC Score	Accuracy Score	Precision (class 0/1)	Recall (class 0/1)	Average f1-score
Subject	(U vs L,C,S)	0.79	0.82	0.81/0.82	0.68/0.91	0.74/0.86
	(U,L vs C,S)	0.80	0.83	0.85/0.77	0.89/0.72	0.87/0.74
	(U,L,C vs S)	0.70	0.96	0.99/0.80	0.99/0.40	0.98/0.53
Concepts	(U vs L,C,S)	0.72	0.75	0.69/0.77	0.59/0.84	0.63/0.81
	(U,L vs C,S)	0.74	0.78	0.80/0.74	0.89/0.58	0.84/0.65
	(U,L,C vs S)	0.68	0.91	0.96/0.75	0.99/0.36	0.97/0.48
Body	(U vs L,C,S)	0.83	0.84	0.79/0.87	0.78/0.88	0.79/0.87
	(U,L vs C,S)	0.81	0.84	0.85/0.82	0.92/0.70	0.88/0.75
	(U,L,C vs S)	0.68	0.95	0.96/0.76	0.99/0.36	0.98/0.49
TAGS	(U vs L,C,S)	0.74	0.78	0.75/0.79	0.61/0.88	0.67/0.83
	(U,L vs C,S)	0.75	0.79	0.82/0.72	0.87/0.63	0.84/0.67
	(U,L,C vs S)	0.62	0.95	0.95/0.73	0.99/0.25	0.97/0.38
Embassies (From/To)	(U vs L,C,S)	0.57	0.67	0.71/0.66	0.19/0.95	0.30/0.78
	(U,L vs C,S)	0.59	0.69	0.70/0.65	0.93/0.24	0.80/0.35
	(U,L,C vs S)	0.57	0.94	0.94/0.72	1.00/0.14	0.97/0.24
Office	(U vs L,C,S)	0.67	0.73	0.76/0.72	0.42/0.92	0.54/0.81
	(U,L vs C,S)	0.62	0.73	0.71/0.83	0.97/0.27	0.82/0.41
	(U,L,C vs S)	0.62	0.95	0.95/0.79	1.00/0.25	0.97/0.38
All_Text	(U vs L,C,S)	0.85	0.86	0.82/0.88	0.81/0.89	0.81/0.89
	(U,L vs C,S)	0.84	0.87	0.88/0.84	0.92/0.76	0.90/0.80
	(U,L,C vs S)	0.78	0.92	0.97/0.78	0.99/0.57	0.98/0.66
All Features (independent vectors)	(U vs L,C,S)	0.86	0.87	0.83/0.89	0.81/0.90	0.82/0.90
	(U,L vs C,S)	0.85	0.87	0.88/0.84	0.92/0.78	0.90/0.81
	(U,L,C vs S)	0.81	0.97	0.97/0.80	0.99/0.61	0.98/0.69
	(U vs C, S)	0.93	0.93	0.93/0.93	0.94/0.92	0.93/0.93



- Features
- body
 - Concepts
 - Embassy
 - Office
 - Subject
 - TAGS
 - All textual fields
 - All features



Enriching feature-set with semantic vectors

A extended feature-set including features derived from word2vec Analysis, can improve the performance of the classifier.

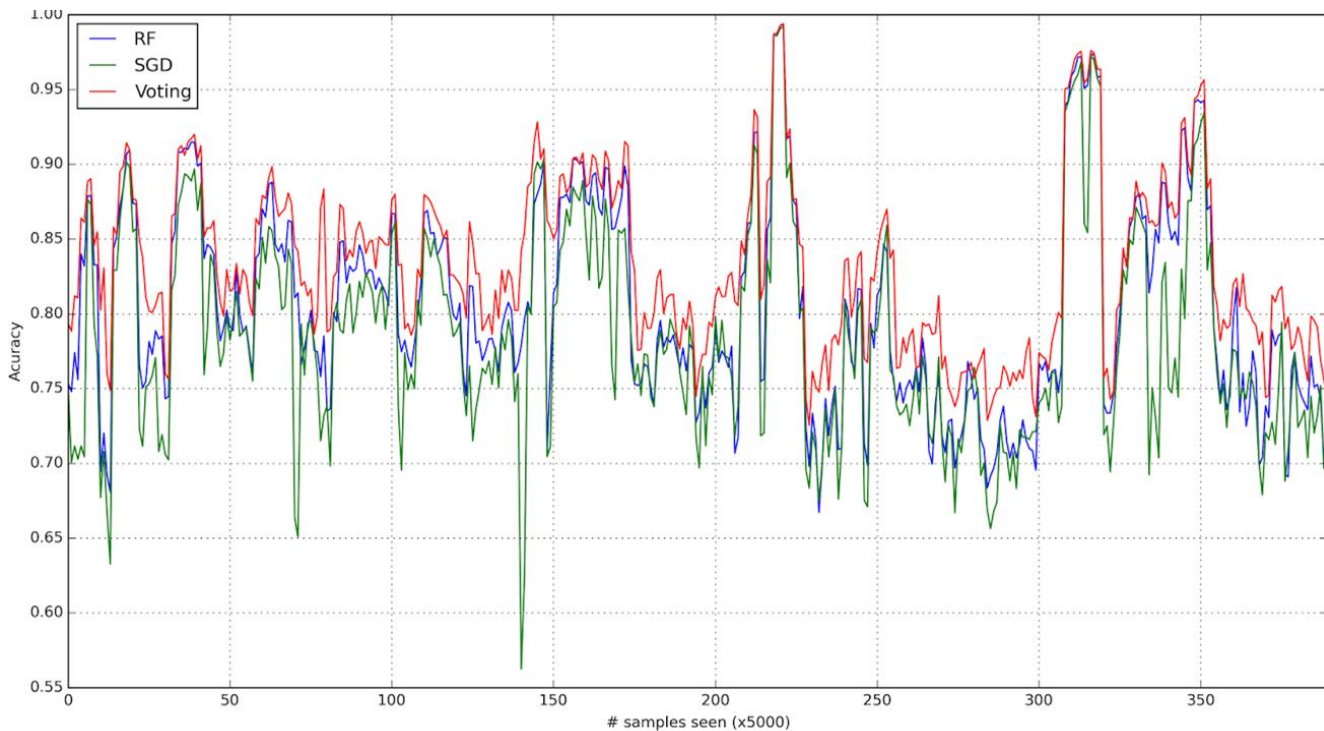
Semantic feature vectors:

$$\vec{d} = \frac{\sum_i tfidf_{w_i,d} \times \vec{w}_i}{\sum_i tfidf_{w_i,d}}$$

Exploring Temporal Evolution

Experiment 1:

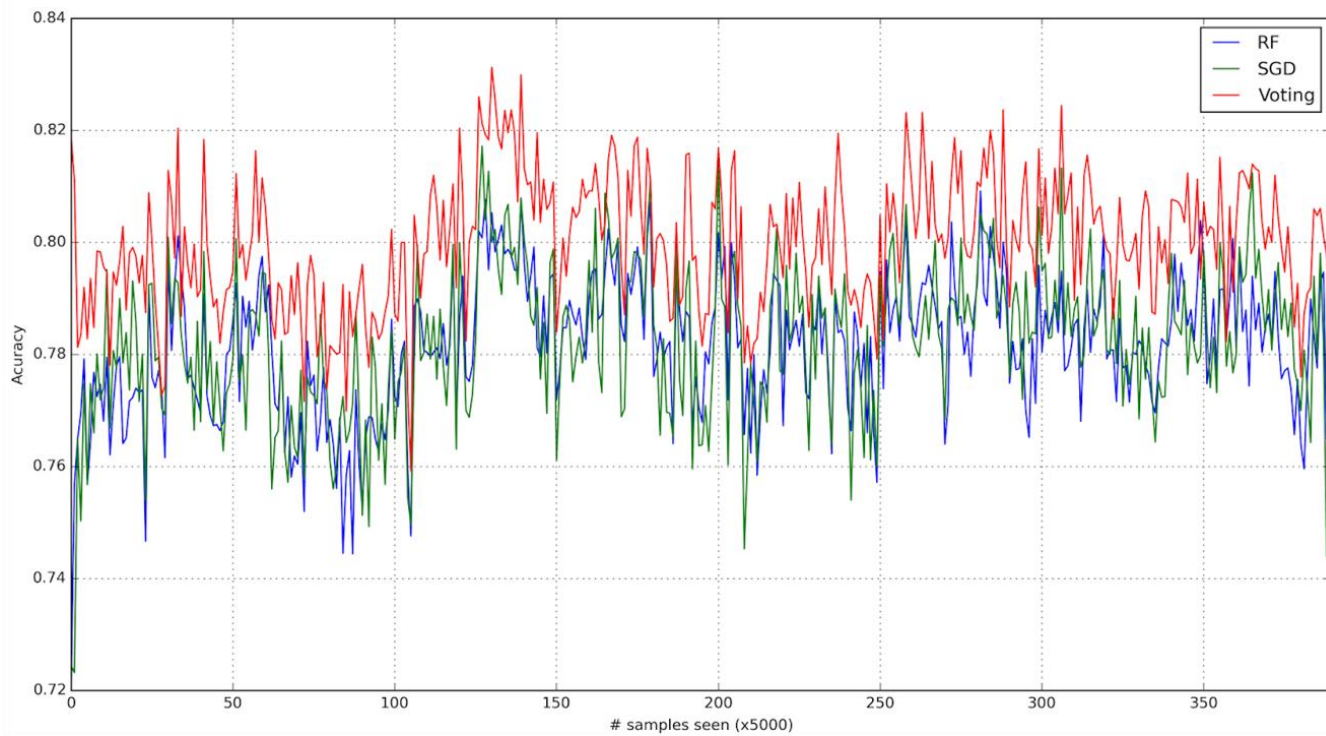
- Split up cable collection if chunks of 5K cables
- Randomize the order
- Run a batch training on this set of batches



Exploring Temporal Evolution

Experiment 2:

- Sort cables in ascending dates
- Split the collection in 5K chunks
- Train the classifier on these chunks from the older ones to the newer.



Next Steps

- Analyze Brazilian and US political reverberations, exploring FGV CPDOC archives diplomatic documents;
- Analyze temporal issues regarding cables
- Dig deeper in the misclassified cables, and identify the misleading features (or the main sources of human errors)
- Identify documents' authorship using ML techniques;
- Build better interfaces for computer aided human classification tasks.