Department of Statistics University of California, Los Angeles

#### The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Ryan Robert Rosario

March 14, 2017

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# Outline





- 2 Existing Methods
- Semantic Bootstrap Proposal
- Experiments 4









Rvan Robert Rosario

#### Introduction

- Existing Methods
- 3 Semantic Bootstrap Proposal

#### 4 Experiments

- 5 Data
- 6 Models
- Results
- 8 Conclusion

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Short Texts

Short texts are ubiquitous across the World Wide Web.

Short texts allow the Web to be more accessible to the world as users can communicate thoughts and desires and ingest new information in a very quick manner.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# Example: Google Suggested Search

Suggested Queries for query "ucla."



ucla		Advanced Search
ucla basketball	1,020,000 results	Language Tools
ucla extension	286,000 results	
ucla medical center	721,000 results	
ucla library	480,000 results	
ucla athletics	582,000 results	
ucla map	4,310,000 results	
ucla football	970,000 results	
ucla law	562,000 results	
ucla admissions	544,000 results	
ucla basketball schedule	972,000 results	
	close	

A classifier takes what a user types, as well as their location, user profile and other information and is able to suggest neighboring queries based on some cluster or classification. We extract meaning using external data.

Ryan Robert Rosario

## Example: Google Search Results

Customized Search Results: Another example of using external data to customize search results.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# Example: Facebook Ad Creative



#### Sign Up Today! 24HOURFITNESS.COM

A 24 Hour Fitness membership gets you tons of studio classes for the price of none.

Open Your Online Store www.volusion.com Get everything you need to succeed with Volusion's all-in-one ecommerce platform. Start vo...

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

### Example: Tweet

Twitter still rules the short text world.



2+ Follow	~
-----------	---

Openwashing, where a company secretly plans to restrict access to "open source", is one of most toxic ideas in CS. Don't do it!



But there a whole series of issues regarding their use in machine learning.

Ryan Robert Rosario

## Example: Title of a Blog Post



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## Problems with Short Text Classification

In various works including [7][14][9][6], the main problems associated with *classifying* short texts are as follows

severe data/feature sparsity;

- words have less discriminative power since less relational information is available;
- Iimited context;

Our goal is to create a method that improves short text classification despite these limitations via text augmentation.

Ryan Robert Rosario

#### Introduction

#### 2 Existing Methods

3 Semantic Bootstrap Proposal

#### 4 Experiments

#### 5 Data

#### 6 Models



#### 8 Conclusion

#### Ryan Robert Rosario

# Existing Methods

The difficulties of working with short texts date back several decades, gaining mass interest in the 1980s for database search. Text augmentation methods (also called query expansion) can be classified into three categories

- Relevance Feedback
- Obmain-Specific
- Lexical Modeling

Throughout the decades, the definition of "short" has also changed.

Ryan Robert Rosario

# Existing Methods: Relevance Feedback

A relevance feedback method works as follows

- user issues a short query (short text) q;
- system returns an initial set of (similar) results;
- user marks results as relevant (similar) or not;
- system computes ranking or retrieval criteria for a new query q' (augmented query);
- System displays a revised set of retrieval results.

Ryan Robert Rosario

# Existing Methods: Relevance Feedback, Xu and Croft Example

One key contribution is provided in Xu and Croft[13], where the human step is removed.

- the retrieval engine uses global<sup>1</sup> context to return what it considers relevant results;
- the results are scored and ranked as *local* context;
- Concepts from the top n results are included in augmented text/query.

Their method improved precision by about 10% on collections of long texts, but did not improve precision on smaller texts.

Ryan Robert Rosario

 $<sup>^1\</sup>mbox{The}$  researchers defined global context as using word cooccurrence or other statistics.

## Existing Methods: Relevance Feedback

The relevance feedback class of algorithms has the following drawbacks:

- assumes that a retrieval system actually provides relevant results in the first place, and
- Implementation Detail: the wording strongly suggests there is in fact a separate system, called the retrieval system;
- based on a strict ranking system that may be too arbitrary to be used deterministically;
- the type of data used in augmentation is arbitrary (concepts? phrases?);
- authors of more modern methods (e.g. [10]) have cited these methods as performing poorly.

Ryan Robert Rosario

## Existing Methods: Domain Based

In a domain-based approach, **external datasets** are used to augment or annotate texts for classification. The most common being Wikipedia<sup>2</sup> and WordNet<sup>3</sup>.

In such domain-specific methods, there are some other features in  $\operatorname{common}$ 

- they use some measure of term cooccurrence rather than retrieval and ranking;
- they generate other features from the corpus as a way of transferring knowledge, such as part-of-speech.

<sup>3</sup>https://wordnet.princeton.edu/

Ryan Robert Rosario

<sup>&</sup>lt;sup>2</sup>http://www.wikipedia.org

# Existing Methods: Domain Based, Mandala et. al. Example

In [8], researchers used three different datasets (thesauri) to serve as the body where sampling occurs:

- WordNet to measure distance between groups of terms (external data);
- term cooccurrence metrics to identity synonyms;
- "head-modifier-based" thesaurus which considered language structure such as subject-verb, adjective-noun etc. (feature generation).

Ryan Robert Rosario

# Existing Methods: Domain Based, Mandala et. al. Example

Then, all combinations of thesauri are considered as an ensemble to expand queries to an *arbitrary* length of 20 terms for evaluation.

		Expanded with						
Topic Type	Base	WordNet	Head-Mod	Cooccur	WordNet+	WordNet+	Head-Mod+	Combined
		only	only	only	Head-Mod	Cooccur	Cooccur	method
Title	0.1175	0.1299	0.1505	0.1637	0.1611	0.1698	0.1859	0.2337
		(+10.6%)	(+28.1%)	(+39.3%)	(+37.1%)	(+44.5%)	(+58.2%)	(+98.9%)
Description	0.1428	0.1525	0.1705	0.1950	0.1832	0.1973	0.2315	0.2689
		(+6.8%)	(+19.4%)	(+33.4%)	(+28.3%)	(+38.2%)	(+62.1%)	(+88.3%)
All	0.1976	0.2018	0.2249	0.2395	0.2276	0.2423	0.2565	0.2751
		(+2.1%)	(+13.8%)	(+21.2%)	(+15.2%)	(+22.6%)	(+29.8%)	(+39.2%)

Titles are the shortest text available in this corpus, and we see that from a base of  $\approx$  0.12, the combined weighted ensemble of thesauri accomplish a 98.9% improvement.

Ryan Robert Rosario

# Existing Methods: Domain Based

The domain-based class of algorithms has the following drawbacks:

- the corpus likely does not match the experimental dataset (e.g. WordNet, Wikipedia) in style or sophistication;
- due to potential mismatches in the data, distance metric selection becomes too much of an art and/or arbitrary;
- again, sampling is done greedily rather than probabilistically ("always pick the best term");
- requires feature generation likely to be costly (e.g. part-of-speech).

Ryan Robert Rosario

# Existing Methods: Lexical Modeling

Methods in the lexical modeling category focus more on tackling the short text issue head-on within a *model* rather than within the data or the features.

Most methods reviewed in this category use either Bayesian topic models, and neural networks are the next frontier.

## Existing Methods: Lexical Modeling, Yan et. al. Example

One such model by [14] modifies Latent Dirichlet Allocation[1]:

- by modeling *biterms* rather than individual words. Biterms are collections of 2 words that may appear anywhere in a sentence (i.e. "i visit apple store" →
  {"visit apple", "apple store", "visit store"})
- researchers believed that such a method models cooccurrence *patterns* and not just cooccurrences.
- moreover, the model considers biterms, not words in documents, to be generated by a topic.

This model is called BTM, for Biterm Topic Model.

Ryan Robert Rosario

## Existing Methods: Lexical Modeling, Yan et. al. Example

Researchers found that based on *coherence*, BTM performed better than LDA even on longer texts.

Table 6: Average coherence score on the top T words (ordered by P(w|z)) in topics discovered by LDA, LDA-U, mixture of unigrams, and BTM. A larger coherence score means the topics are more coherent. It suggests that BTM outperforms others significantly (P-value < 0.01 by t-test).

Т	5	10	20
LDA	$-55.0\pm0.4$	$-236.4\pm2.0$	$-1015.7\pm5.9$
LDA-U	$-54.2\pm0.8$	$-234.8\pm1.1$	$-1009.4\pm4.4$
Mix	$-53.8\pm0.1$	$-233.0\pm1.4$	$-1007.6\pm6.7$
BTM	$-52.4\pm0.1$	$-227.8\pm0.3$	$-990.2\pm3.8$

#### But...

TABLE 8 Time cost (seconds) per iteration of BTM and LDA on Tweets2011 collection.							
	K	50	100	150	200	250	
	LDA	38.07	74.38	108.13	143.47	178.66	
	BTM	128.64	250.07	362.27	476.19	591.24	
Me	BTM 128.64 250.07 362.27 476.19 591.24 TABLE 9 Memory cost (megabytes) per iteration of BTM and LDA on Tweets2011 collection.						

K	50	100	150	200	250
LDA	3177	5524	7890	10218	12561
BTM	927	946	964	984	1002

Ryan Robert Rosario

# Existing Methods: Lexical Modeling

The lexical modeling class of algorithms has the following drawbacks:

- parameter estimation for such models can be very slow;
- it is the researcher's opinion that building a full-blown model is overkill for this problem;
- Solution biterms seem arbitrary why not triterms? <sup>4</sup>
- native neural networks tend to suffer from overfitting problems[12].
- Sayesian topic models have mixed success in scalability.

Ryan Robert Rosario

<sup>&</sup>lt;sup>4</sup>many methods in this category use *biterms*.

## Motivation: So, Why Yet Another Method?

The previous work raises the following questions:

- Can't we just use the existing terms in the text rather than using other data or feature generation?
  - this is faster and more computationally efficient;
  - we do not have to worry about concept drift, transfer or lexicon mismatch since we use the same data;
  - we can make fewer major arbitrary decisions (i.e. distance metrics);
- Can we avoid building a special model and treat text augmentation as a preprocessing step?
- Scan we propose something more general purpose?

The researcher believes that we already have the statistical tools to do so...

Ryan Robert Rosario

#### Introduction

#### Existing Methods



#### Experiments

#### 5 Data

#### 6 Models

#### Results

#### 8 Conclusion

#### Ryan Robert Rosario

### The Classical Bootstrap

The classical bootstrap[4] is a technique for constructing the sampling distribution of a sample statistic to calculate an estimate or its confidence interval. It consists of

- Draw a sample of size *n* from some population where the observations are
  - with replacement;
  - independent;
  - (3) of a sufficiently large size.
- **2** Compute some function of the data *f* on such sample.
- Solution Repeat steps 1 and 2 for a total of N times, where N is large.

By the Central Limit Theorem, the metrics calculated by f follow a normal distribution.

Ryan Robert Rosario

## The Bootstrap for Small Samples

While the bootstrap is most often used to estimate a sampling distribution of a parameter, it has also been used when a sample is **small**, with mixed results[11][3].



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# The Bootstrap for Small Samples

**But how small is too small?** One author [3] has stated **8** to be an approximate lower bound.

But, many texts are even *shorter* than 8 words. In my experimental data, the average is 4-5.

# Research Question: Can we modify the Bootstrap to work with text?

Doing so requires modifications to the Bootstrap procedure.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## Motivation: Why Use the Bootstrap?

Several advantages to the Bootstrap over existing methods:

- very easy to implement;
  - improved scalability and embarrassingly parallel;
- very accessible to practitioners even outside machine learning;
- uses the existing data as a population rather than external data or queries;
- makes no transformations to individual observations (only in their dependence and collection);
- the Bootstrap is one of the most popular simulation methods for working with sample sizes.
  - Regularization is another common method;
  - A general purpose method that does not require special data or special models.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

**First we need to define**, what is an observation, and what is the population?



With a standard dataset, the dataset itself is the population, and each ball is an observation to be sampled. There is a simple nesting of observations within data.

#### Ryan Robert Rosario

The situation is not so simple for text. We have an inherent nesting of *words* within *documents* within a *corpus*.



We have a few choices of how to define an *observation* and *population*.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

**Option 1** If we *directly* apply the resampling from the bootstrap on the corpus level, treating the corpus as the population, and the documents as the observations, we end up just duplicating existing documents.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

**Option 2** If we *directly* apply the resampling from the bootstrap on the document level, treating the document as the population, and the words as the observations, we end up just duplicating existing words.



While duplicating words increases sample size, it does nothing to address sparsity.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

#### **Proposed Variation of Option 2: The Semantic Bootstrap**

Suppose instead we sample *indirectly* from the document by choosing words that are semantically similar to the existing words in each document.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## Semantic Bootstrap as a Variation of Classical Bootstrap

What we just described can be summarized as the classical bootstrap with three variations:

- only the resampling step is of interest;
- Itreat each term in a document as an observation;
- treat a semantic space S|d as the population;
- independence is lost as certain words are more discriminant than others

But what is the semantic space S?

Ryan Robert Rosario

## The Semantic Space S

We must define the population - the semantic space S.

**Belief:** any matrix factorization that maintains distance between terms, documents and term/document pairs.

**For this research**, we use *Latent Semantic Analysis (LSA)*, a very commonly used matrix factorization in text mining and natural language processing.

Ryan Robert Rosario
## The Semantic Space S: Latent Semantic Analysis

We start by representing a *training* corpus as a term-document matrix **X** with |D| rows, one for each document  $d_i$ , and |V| columns, one for each term  $t_j$ .

$$\mathbf{X} = \begin{bmatrix} t_{0,0} & t_{1,0} & t_{2,0} & \dots & t_{|V|,0} \\ t_{0,1} & t_{1,1} & t_{2,1} & \dots & t_{|V|,1} \\ \dots & \dots & \dots & \vdots & \vdots \\ t_{0,|D|} & t_{1,|D|} & t_{2,|D|} & \dots & t_{|V|,|D|} \end{bmatrix}$$

Each entry in the matrix  $X_{ij}$  represents some relationship between each term in each document: presence/absence (0/1), word count, or...

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## The Semantic Space S: Latent Semantic Analysis

A common metric relating terms and documents is TF-IDF<sup>5</sup> and is the score used for  $X_{ij}$ .

$$X_{ij} = \mathsf{TF}\mathsf{-}\mathsf{IDF}_{ij} = \mathsf{TF}_{ij} \times \mathsf{IDF}_i = c_{ij} \times \left[ \log\left(\frac{|D|}{|\{d: t_i \in d\}|}\right) + 1 \right]$$

Properties of TF-IDF

- high TF-IDF is associated with highly discriminative or influential words within a document;
- low TF-IDF is associated with stopwords and other words that do not impart much meaning in the document.

Ryan Robert Rosario

<sup>&</sup>lt;sup>5</sup>where  $c_{ij}$  is the number of times  $t_i$  appears in document  $d_j$ , |D| is the number of documents, and  $\{d : t_i \in d\}$  is the set of documents containing t.

### The Semantic Space S: Latent Semantic Analysis

Then we can decompose  $\mathbf{X}$  as follows:

$$\mathbf{X} \approx \mathbf{D_k} \mathbf{\Sigma_k} \mathbf{W_k}^\mathsf{T}$$

We can then construct term and document similarities as follows

$$\begin{split} \mathbf{S}_{t} &= \mathbf{X}\mathbf{X}^{\mathsf{T}} \\ &= \left(\mathbf{D}_{k}\boldsymbol{\Sigma}_{k}\mathbf{W}_{k}^{\mathsf{T}}\right)\left(\mathbf{D}_{k}\boldsymbol{\Sigma}_{k}\mathbf{W}_{k}^{\mathsf{T}}\right)^{\mathsf{T}} \\ &= \mathbf{D}_{k}\boldsymbol{\Sigma}_{k}\mathbf{W}_{k}^{\mathsf{T}}\mathbf{W}_{k}\boldsymbol{\Sigma}_{k}^{\mathsf{T}}\mathbf{D}_{k}^{\mathsf{T}} \\ &= \mathbf{D}_{k}\boldsymbol{\Sigma}_{k}^{2}\mathbf{D}_{k}^{\mathsf{T}} \end{split}$$

and similarly,

$$\mathbf{S}_{d} = \mathbf{X}^{\mathsf{T}} \mathbf{X} = \mathbf{W}_{k} \mathbf{\Sigma}_{k}^{2} \mathbf{W}_{k}^{\mathsf{T}}$$

Ryan Robert Rosario

### The Semantic Space S: Cosine Similarity

Once we have computed  $\boldsymbol{S}_t$  and  $\boldsymbol{S}_d,$  we can make comparisons across terms and documents using a distance metric. Cosine similarity is perhaps the most common metric used in text mining and NLP.

$$\cos\theta = \frac{u \cdot v}{||u||||v||}$$

For terms,

$$\delta_{t_i,t_j} = \mathsf{sim}(t_i,t_j) = rac{t_i \cdot t_j}{||t_i||||t_j||}$$

For documents,

$$\delta_{d_i,d_j} = \operatorname{sim}(d_i,d_j) = rac{d_i \cdot d_j}{||d_i|||d_j||}$$

#### Ryan Robert Rosario

## A Sampling Scheme

Up to this point, we have decomposed **X** into a semantic space *S* using SVD and computed pairwise similarities using cosine similarity. We also have information about the relationship between a document and its terms in  $X_{ij}$ , the TF-IDF score.

We can now describe a sampling scheme:

- **(**) we have a probability distribution over term similarities  $\delta_t$ ,
- **②** we have a probability distribution over document similarities  $\delta_d$ ,
- we can also have a probability distribution over word discrimination X<sub>i</sub>..

### A Conditional Probability Distribution over S

We can now propose a sampling from S:

$$P(t_j|t_i) = \frac{\delta_{t_i,t_j} + |\min \delta_{t_i,\cdot}|}{\sum_m (\delta_{t_i,t_m} + |\min \delta_{t_i,\cdot}|)}$$
$$P(d_j|d_i) = \frac{\delta_{d_i,d_j} + |\min \delta_{d_i,\cdot}|}{\sum_m (\delta_{d_i,d_m} + |\min \delta_{d_i,\cdot}|)}$$

Ryan Robert Rosario

## A Probability Distribution over d

Within each document d, each term has a different level of discriminative power quantified by TF-IDF. So if we want to select a random term based on discriminative power, we use  $X_{ij}$  as follows

$$P(t_i|d_j) = \frac{x_i}{\sum_j x_j}$$

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## The Semantic Bootstrap Algorithm

Now that we have a population S to sample new words from, and a probability distribution induced from it, and we also have a probability distribution over the terms in a document d, we can propose a sampling scheme as follows...

Let  $\varepsilon$  be the *augmentation rate* – a value greater than 1, specifying how much longer the new document  $d^*$  should be relative to the original document d:

$$d^* = \varepsilon |d|$$

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

### The Semantic Bootstrap Algorithm

**Proposal:** We sample terms from S until the new augmented document reaches the desired length. That is, while  $|d^*| < \varepsilon |d|$ :

- Pick a target term t in d according to P(t<sub>i</sub>|d<sub>j</sub>). This yields a term with high discriminating value with highest probability. We want to find a word similar to it given the context of the current document.
- Pick a candidate document d' according to P(d<sub>j</sub>|d<sub>i</sub>). This yields a document d' that is most similar to d according to S, with highest probability.
- Pick a candidate term t' from d' according to P(t<sub>j</sub>|t<sub>i</sub>). This yields the most semantically related term from the most semantically related document, with highest probability.

Using probabilities rather than simply picking the most relevant terms and documents reduces the possibility of bias from a bad selection. It also eliminates the need for arbitrary ranking cutoffs.

## The Semantic Bootstrap Algorithm: A Final Sanity Check

In [13], the researchers did one final check when augmenting terms to protect against *concept drift*, where selecting a bad term causes the results to be irrelevant.

- In this research, all samplings are based on d, and not on iterations of d\*, so concept drift is not possible, BUT
- A poor sampling will still yield bad results.

We want to make sure that any sampled terms contribute as much semantic cohesion as possible. So...

#### The Semantic Bootstrap Algorithm: Mutual Information

We use Mutual Information (MI) as a final check to try to ensure the augmented bag-of-words  $d^*$  is as semantically related to d as possible. The higher the change in MI, the more candidate term t'adds to semantic cohesiveness.

$$I(d) = \sum_{(t_i, t_j, i \neq j) \in d} P(t_i, t_j) \log \left( \frac{P(t_i, t_j)}{P(t_i, )P(t_j)} \right)$$

We accept the term t' into  $d^*$  according to the transition probability

$$P(d 
ightarrow d^*) \propto \min\left\{1, rac{1}{Z} \exp\left(I(d^*) - I(d)
ight)
ight\}$$

Ryan Robert Rosario

## The Semantic Bootstrap Algorithm

# The Semantic Bootstrap algorithm can be illustrated as the following flowchart



ε is referred to in this research as the augmentation extension and measures |d\*| relative to |d| (e.g. ε = 2 means doubling the length).

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Suppose we have the following text we want to augment.

Yom Kippur, Tel Aviv style

The text is typically short and vague that is also not a complete sentence. It could be about the Jewish *religion*, it could be about *politics* in the region, or it could be more *geographical/cultural* in nature. Maybe we can augment it with more terms to assist in classification.

First we preprocess the text and turn it into a bag of words representation:

yom kippur tel aviv style

Next, we *probabilistically* select a document d' from semantic space *S* according to how similar it is to *d*.

Ryan Robert Rosario

The d' selected here had a high cosine similarity to d and was chosen with probability 0.878.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Next, from d, we sample a target word t. We will then sample a word t' from d' according to how similar they are in the semantic space S.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Now we have a candidate document  $d^*$ . We accept the new term t' and thus  $d^*$  according to the change in semantic cohesiveness,  $P(d \rightarrow d^*)$ :

## d\* yom kippur tel aviv style hashanah

- If  $P(d 
  ightarrow d^*)$  is relatively "large", then accept t' into  $d^*$ .
- Otherwise, reject t' and keep d.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

The sampling process continues until the desired number of terms has been added.

Some minutiae:

- It is possible to select the same document d' across iterations;
- It is possible to select the same t and/or t' across iterations. The final term acceptance step had some effect in preventing these words from overwhelming d\*;
- It is possible to get into a situation where d\* is repeatedly rejected. It was found that as many as 5 retries for each iteration were required for as few as 0.1% of cases.

#### Introduction

#### Existing Methods

3 Semantic Bootstrap Proposal

#### 4 Experiments

- 5 Data
- 6 Models





#### Ryan Robert Rosario

### Experimental Questions

The most important question we will review is obviously

#### Does the semantic bootstrap even work?

And more specifically, the following

- O How many terms should be sampled to yield the best results? In other words, what are appropriate values for augmentation rate ε?
- When should the Semantic Bootstrap algorithm be applied? In the training data? In the testing data? Or both, as is typical?
- Is the final probabilistic term acceptance step truly necessary?
- Which of the experimental classification models works best?
- I How can we control for sampling variation?

## A Note on Unmatched Training/Testing Sets

It is standard practice to apply the same transformations to both the training and testing sets.

There are cases where it may be more convenient from a computational standpoint to use a unmatched training or testing set.

Some research suggests that unmatched datasets can sometimes outperform matched datasets[5].

Ryan Robert Rosario

## A Note on Unmatched Training/Testing Sets: Example I

#### Augmenting only the Training Data

Suppose I stand in a noisy room and there is a microphone at the far end of the room that is to record only my voice.

One change we can make is to the *microphone* and its logic, making it more sensitive to the frequencies in my voice. This is akin to training a classifier on an augmented corpus and evaluating on a the standard dataset.

Ryan Robert Rosario

## A Note on Unmatched Training/Testing Sets: Example II

#### Augmenting only the Unseen Testing Data

Suppose I stand in a noisy room and there is a microphone at the far end of the room that is to record only my voice.

Another change we can make is to leave the microphone as it is, and for me to yell loudly over the noise in the room. This is akin to augmenting the unseen signal being classified while leaving the training set the same.

Ryan Robert Rosario

## A Note on Unmatched Training/Testing Sets

Throughout this research, these configurations are called variations:

Variation	Description
SB1	Augment only the training data,
	use original testing data.
SB2	Use original training data and classifier,
	augment only testing data.
SB12	Matched case; augment both sets.
SB0	Raw data; neither is augmented.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## Controlling for Sampling Variation

The Semantic Bootstrap, and the Classical Bootstrap are both probabilistic. Each iteration yields a different sampling. Each iteration of the Semantic Bootstrap yields a new corpus containing augmented documents.

Because of this, the Semantic Bootstrap is applied for 100 iterations and performance metrics averaged over all iterations.

The iterations are for this research only, and not intended for actual use.  $^{\rm 6}$ 

Ryan Robert Rosario

<sup>&</sup>lt;sup>6</sup>Picking an appropriate number of iterations should be future work.

#### Introduction

#### Existing Methods

3 Semantic Bootstrap Proposal

#### Experiments



#### 6 Models



#### 8 Conclusion

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

### Technorati Data



The experimental data from a 2014 widespread crawl of the Technorati blog aggregator and search engine. The titles of each blog post were used as the dataset for this research.

Each blog belongs to one<sup>7</sup> of ten categories.

% of Titles	Category	% of Titles
2.1	politics	7.5
7.4	science	3.0
20.8	sports	14.3
2.1	technology	18.8
21.6	overall	2.4
	% of Titles 2.1 7.4 20.8 2.1 21.6	% of TitlesCategory2.1politics7.4science20.8sports2.1technology21.6overall

<sup>7</sup>Blogs belonging to more than one category were removed from consideration.

Ryan Robert Rosario

#### Technorati Data

On average, blog post titles contain between 4 and 5 words, and the length is log-normal distributed.



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

#### Technorati Data

The text was cleaned using typical text mining procedures

- discard titles not written in English;
- converting foreign characters to ASCII where possible;
- removing known stopwords (i.e. a, the etc.);
- removing words that appear too often or too seldom to be useful;
- removing duplicates and empty documents caused by the above.

The dataset consisted of

- Corpus consisted of 131,519 documents.
- Lexicon/vocabulary size consisted of **10,889** words.
- Dimension of term-document matrix reduced to K = 500 from scree analysis.

- Introduction
  - Existing Methods
- 3 Semantic Bootstrap Proposal

#### Experiments

#### 5 Data







Ryan Robert Rosario

### Experimental Models

To test the efficacy of the Semantic Bootstrap as a preprocessing step for classification, the researcher used two models:

- Linear Support Vector Machine (SVM) [Classical Approach]
  - via libshorttext, a variation designed specifically for short texts, and one state-of-the-art for this research
  - linear SVM with L2 penalty and word count as features
- Latent Dirichlet Allocation Variation (LDA) [Topic Model Approach]
  - Supervised Latent Dirichlet Allocation (sLDA),

 $K = 50, \alpha = 1, \eta = 1$ 

Individual classifiers were constructed for each category and evaluated. The resolution of such a one vs. rest classifier is saved as future work.

## Supervised Latent Dirichlet Allocation (sLDA)

A supervised variant of LDA exists that can remedy the problem of inconsistent labeling by using category labels as a dependent variable.

We can represent a Generalized Linear Model (GLM) as a linear combination of coefficients  $\beta_i$  and K topics  $X_i$ , where each topic consists of a set of terms. In theory, many different link functions can be used. For this research, we used the implemented **logit** link.

#### libshorttext: Linear SVM for Short Texts

libshorttext is a package and variation of SVM developed specifically for use with small texts. It supports

- stemming
- stopword removal
- construction of *n*-grams
- feature normalization/standardization
- L1 and L2 regularization
- several feature types: binary, word count, term frequency (TF), TF-IDF.

And supports related models using similar optimization problems

- Iogistic regression
- multiclass SVM

#### Ryan Robert Rosario

#### libshorttext: Linear SVM for Short Texts

While there are several different models implemented in libshorttext, the Linear SVM with L2 penalty with Word Count features performed near the best on precision/recall/F1 score<sup>8</sup> and *also uses the same feature type as the sLDA model.* 

Model / Metric	Accuracy Macro	Accuracy Micro	Precision Macro	Recall Macro	F1 Score Macro
Linear SVM L2, TF-IDF	0.911	0.696	0.720	0.624	0.668
Linear SVM L2, Binary	0.889	0.645	0.697	0.628	0.644
Linear SVM L2, Word Count	0.832	0.62	0.725	0.638	0.679
Linear SVM L2, Term Frequency	0.899	0.587	0.701	0.628	0.663
Linear SVM L1, TF-IDF	0.894	0.628	0.703	0.619	0.638
Linear SVM L1, Binary	0.933	0.614	0.691	0.625	0.640
Linear SVM L1, Word Count	0.885	0.622	0.691	0.625	0.639
Linear SVM L1, Term Frequency	0.900	0.601	0.691	0.626	0.641
Logistic Regression, TF-IDF	0.931	0.623	0.700	0.624	0.642
Logistic Regression, Binary	0.840	0.675	0.700	0.619	0.639
Logistic Regression, Word Count	0.898	0.622	0.700	0.620	0.639
Logistic Regression, Term Frequency	0.849	0.567	0.700	0.620	0.639

 $^{8}\mbox{the researcher used precision/recall and F1 score to assist in making the decision because the dataset is imbalanced.$ 

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

- Introduction
  - Existing Methods
- 3 Semantic Bootstrap Proposal

#### 4 Experiments

- 5 Data
- 6 Models
- 7 Results
  - Effect of  $\varepsilon$
  - Effect of Time of Sampling
  - Effect of the Probabilistic Term Acceptance Step
  - Effect of Model Choice



#### Ryan Robert Rosario

### Results

In this section, we present the results from applying the Semantic Bootstrap to the Technorai data. As a reminder, we studied the following questions

- O How much should each document be augmented by? What are appropriate values for ε?
- When should the Semantic Bootstrap be applied: to both the training and testing sets? Only the training set? Or only in the unseen set?
- **③** Is the final probabilistic term acceptance step necessary?
- Which experimental model works better under the Semantic Bootstrap? SVM or sLDA?
# Results

As a result, we also implicitly answer the following important questions via the above questions

- Does the Semantic Bootstrap even work?
- Is a matrix factorization (LSA) an appropriate semantic space S?

# A Note on Performance Metrics

- sLDA natively supports posterior probabilities, so AUC is used as the performance metric for comparing sLDA models.
- **SVM** does not, so **F1** score is used instead for comparing SVM models.
- F1 score is used when comparing SVM to sLDA

This is fine because we are mostly interested in comparing performance with and without the Semantic Bootstrap.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Effect of  $\varepsilon$ 

# Experiment 1: Effect of $\varepsilon$

### Hypothesis:

- a value such as 1.5 should ensure that a minimum of 1 term is chosen as a candidate for sampling.
- a value higher than 2.0 adds more words than there existed in original *d*, introducing noise and bias.

• 
$$1.5 \le \varepsilon^* < 2$$

In the graphs that follow,  $\varepsilon = 1.0$  represents the raw, un-augmented data (SB0).

Ryan Robert Rosario





Augmentation Rate  $\epsilon$ 

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification





Augmentation Rate  $\epsilon$ 

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

#### Effect of $\varepsilon$

# Experiment 1 Conclusion

The hypothesis of  $1.5 \le \varepsilon^* < 2$  was disproven:

- **①** The optimal value of  $\varepsilon$  is approximately 1.3.
- **2**  $1.5 \le \varepsilon \le 1.7$  performance is approximately equal to baseline.
- **③**  $\varepsilon > 1.7$  performance begins to drop substantially.

Since the average document in the corpus had 4 to 5 terms,  $\epsilon = 1.3$  corresponds to sampling approximately **one** term. This does validate that sampling from *S* yields highly discriminative terms that help performance. Introduction Existing Methods Semantic Bootstrap Proposal Experiments Data Models Results Conclusion References

Effect of  $\varepsilon$ 

# Effect of $\varepsilon^* = 1.3$ on Classifier Performance



#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Introduction Existing Methods Semantic Bootstrap Proposal Experiments Data Models Results Conclusion References

Effect of  $\varepsilon$ 

# Effect of $\varepsilon^* = 1.3$ on Classifier Performance

### SVM Classifier Performance for E=1.3

Augmentation Original Data, No Augmentation Without PMI Step With PMI Step Strategy

SB12

SB12













Experimental Variation

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

0.00 SBO SB1 SB2

Original Data

Introduction Existing Methods Semantic Bootstrap Proposal Experiments Data Models Results Conclusion References

Effect of Time of Sampling

# Experiment 2: When Should the Semantic Bootstrap be Applied?

## Hypothesis

- Matched training and testing sets are the norm.
- There is some research suggesting unmatched sets may be better.
- In industry, we have used unmatched sets with success.
- No strong hypothesis; but interesting to study.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



Ryan Robert Rosario The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

#### Effect of Variation and Augmentation Rate on sLDA Classifier Performance



Variation - SB1 - SB12 - SB2

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

#### Effect of Variation and Augmentation Rate on SVM Classifier Performance



Variation - SB1 - SB12 - SB2

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

### sLDA Classifier Performance across Experimental Configurations at $\epsilon$ =1.3



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

### SVM Classifier Performance across Experimental Configurations at $\varepsilon$ =1.3



Experimental Variation and Configuration

#### SB0 SB1 SB2 SB12 Original Data science

green



#### Rvan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Augmentation

Effect of Time of Sampling

# Experiment 2 Conclusion

Results among the SBx variations were mixed:

- At small  $\varepsilon$  and  $\varepsilon^*$ , differences were negligible.
- **②** For sLDA, SB2 consistently performed better as  $\varepsilon$  increased.
- For SVM, SB2 performed best at large  $\varepsilon \ge 1.8$ .
- **③** SB2 performed the best at higher values for  $\varepsilon$ .
- Solution best when performed only on the testing set.
- **Remarkably:** The matched training/testing set scenario, the status quo, consistently performed the worst.

Ryan Robert Rosario



### sLDA Classifier Performance for Experimental Variation SB2 at $\varepsilon$ =1.3

Augmentation Original Data, No Augmentation Without PMI Step With PMI Step Strategy

#### Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



### SVM Classifier Performance for Experimental Variation SB2 at $\varepsilon$ =1.3

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# Question 3: Does the Final Term Acceptance Step Help?

Although we sample words from S according to semantic similarity, we want to make sure that the terms we sample improve semantic cohesiveness.

### Hypothesis

The hypothesis is that this final term acceptance step should improve performance because only words that improve cohesiveness are accepted whereas words that do not will be accepted with a lower probability.

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



sLDA Classifier Performance vs. Augmentation Strategy at  $\varepsilon$ =1.3

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



### SVM Classifier Performance vs. Augmentation Strategy at $\varepsilon$ =1.3

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

# Experiment 3 Conclusion

The final probabilistic term acceptance yielded no practical difference.

- For sLDA, straight-up sampling from S performed only 1% better on AUC than with the final check for  $\varepsilon \ge 1.8$ .
- For SVM, the final check again only yielded a 1% improvement in F1-score, but across all ε.

**Likely Reason:** The Semantic Bootstrap sampling step already does a good enough job ensuring semantic cohesiveness since both terms are sampled conditioned on a similar document.

Ryan Robert Rosario

# Experiment 4: Which Model Works Best Under the Semantic Bootstrap?

Since two models were chosen for experimentation, it makes sense to compare their performance. For this experiment, we compare the models using a common metric, F1 score.

### Hypothesis

LDA has seen a lot of success in clustering words and texts into latent concepts called topics. SVM is a classic, but uses only symbolic representations, and can be seen as more rigid. The hypothesis is that sLDA will perform better than SVM with the Semantic Bootstrap and will see a bigger boost in performance.

Ryan Robert Rosario



### Effect of Model Type on Classifier Performance vs. Variation

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



### Improvement of SVM over sLDA vs. Experimental Configuration

#### Rvan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



### Improvement of SVM over sLDA vs. Category: Variation SB2, ε=1.3

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



#### Comparing SVM and sLDA against Variation and $\epsilon$ By Category



Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification



The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

Effect of Model Choice

# Experiment 4 Conclusion

- Just like with SBO, SVM (libshorttext) performed superior to sLDA, except in the overall category.
- At ε > 2, the difference between best performing sLDA models and worst performing SVM models becomes negligible.
- If we must use matched training/testing sets (SB12), use SVM with the final term acceptance step.

**Likely Reason:** The SVM *C* parameter has been cited in [2] as being important when data is noisy, and SVM pays more attention to data points closer to the decision boundary, whereas sLDA focus on all data points as a generative algorithm. Adding additional terms may add noise that weakens term-document co-occurrences.

Ryan Robert Rosario

- Introduction
  - Existing Methods
- 3 Semantic Bootstrap Proposal
- 4 Experiments
- 5 Data
- 6 Models
- Results
- 8 Conclusion

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

## Potential I

These results show a lot of potential for the Semantic Bootstrap:

- allows augmenting of small texts to improve classification, a big problem in industry;
- does not require any data external to the corpus of interest, reducing querying times and computation power;
- does not require any niche or clever processing (e.g. NER) of the text, saving computational power;

Ryan Robert Rosario

# Potential II

- eliminates the need to combine small texts into one large document for standard processing;
- the documented improvements apply to *unigrams* and it is hypothesized to be even better with n-grams;
- is fast and embarrassingly parallel, including in the construction of S;
- due to easy implementation and general purpose nature, there are tons of avenues for future research.

Ryan Robert Rosario
#### Future Work

This research was designed to take a big idea from statistics, the Bootstrap, and try to apply it to a text. This work was very broad and there are several things worth pursuing:

- using (s)LDA clusters as features in SVM, or some other machine learning model may provide better results; this is common in the literature;
- evaluating performance on other standard machine learning classifiers such as random forests, Naive Bayes etc.;
- **③** using individually tuned asymmetric priors  $\alpha$ ,  $\beta$  for sLDA;
- determining if ε is specific to the data, or specific to the distribution of document lengths, or a function of some other variable;

Ryan Robert Rosario

#### Future Work

- for evaluation purposes only, we used 100 iterations of the Semantic Bootstrap on each corpus, and it is implied that 1 may be enough on average. How many iterations are best?
- investigating if combining the results of the iterations using boosting or bagging improve performance;
- investigating other matrix factorizations such as NMF, as well as non-matrix factorization methods to construct S.

### Conclusion

In this work, we attempted to adapt the classical bootstrap from statistics for use with text, short text in particular. We call this method the Semantic Bootstrap. We

- introduced the concept of a semantic space S to serve as the population from which we sample new terms;
- made extensive use of term discrimination (TF-IDF) features and semantic similarity to choose additional terms to add to the short text;
- investigated how many terms must be added for best performance and found that adding as few as one term significantly improves classification performance;

Ryan Robert Rosario

## Conclusion

- investigated whether or not using matched training and testing sets performed best and found that only under higher levels of augmentation, it is actually better to apply the method on only the testing/unseen data;
- investigated if the Semantic Bootstrap yields semantically cohesive augmented texts and found no improvement when performing a final check that new terms improved cohesiveness;
- used a variation of SVM for short texts and compared it to a contemporary and supervised topic model and found that the tried and true SVM performed better.

Ryan Robert Rosario

#### References I

- D.M. Blei, A.Y. Ng, and M.I. Jordan. "Latent dirichlet allocation". In: *the Journal of machine Learning research* 3 (2003), pp. 993–1022.
- [2] Vladimir Cherkassky and Yunqian Ma. "Selection of meta-parameters for support vector regression". In: *International Conference on Artificial Neural Networks*. Springer. 2002, pp. 687–693.
- [3] Michael R Chernick. Bootstrap methods: A guide for practitioners and researchers. Vol. 619. John Wiley & Sons, 2011.

#### References II

- [4] Bradley Efron and Robert Tibshirani. An introduction to the bootstrap. Vol. 57. CRC press, 1993.
- [5] Carlos R González and Yaser S Abu-Mostafa. "Mismatched training and test distributions can outperform matched ones". In: *Neural computation* (2015).
- [6] Jian Hu et al. "Enhancing text clustering by leveraging Wikipedia semantics". In: Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval. ACM. 2008, pp. 179–186.

Ryan Robert Rosario

#### References III

- [7] Xia Hu et al. "Exploiting internal and external semantics for the clustering of short texts using world knowledge". In: Proceedings of the 18th ACM conference on Information and knowledge management. ACM. 2009, pp. 919–928.
- [8] Rila Mandala, Takenobu Tokunaga, and Hozumi Tanaka. "Combining multiple evidence from different types of thesaurus for query expansion". In: Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval. ACM. 1999, pp. 191–197.

#### References IV

- [9] Xuan-Hieu Phan, Le-Minh Nguyen, and Susumu Horiguchi. "Learning to classify short and sparse text & web with hidden topics from large-scale data collections". In: *Proceedings of the 17th international conference on World Wide Web.* ACM. 2008, pp. 91–100.
- [10] Yonggang Qiu and Hans-Peter Frei. "Concept based query expansion". In: Proceedings of the 16th annual international ACM SIGIR conference on Research and development in information retrieval. ACM. 1993, pp. 160–169.
- [11] F Scholz. "The bootstrap small sample properties". In: *Tech. Rep.* (2007).

Ryan Robert Rosario

#### References V

- [12] Nitish Srivastava et al. "Dropout: a simple way to prevent neural networks from overfitting." In: Journal of Machine Learning Research 15.1 (2014), pp. 1929–1958.
- [13] Jinxi Xu and W Bruce Croft. "Query expansion using local and global document analysis". In: Proceedings of the 19th annual international ACM SIGIR conference on Research and development in information retrieval. ACM. 1996, pp. 4–11.
- [14] Xiaohui Yan et al. "A biterm topic model for short texts". In: Proceedings of the 22nd international conference on World Wide Web. ACM. 2013, pp. 1445–1456.

Ryan Robert Rosario

## **Thank You**

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

UCLA Department of Statistics

# Backup Slides

Ryan Robert Rosario

The Semantic Bootstrap: Application of the Bootstrap for Small Text Classification

UCLA Department of Statistics

#### Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) has seen much acclaim in text mining and natural language processing.



A generative process that assumes text is generated as follows

- **1** Draw a topic distribution  $\theta_i \sim \text{Dir}(\alpha)$  for each document  $d_i$ .
- **2** Draw a word distribution  $\phi_k \sim \text{Dir}(\beta)$  for each topic k.
- **③** For each word position  $w_{ij}$  in the document,
  - Sample a topic indicator  $z_{ij} \sim \text{Mult}(\theta_j)$ .
  - **2** Sample a word  $w_{ij} \sim \text{Mult}(\phi_{z_{ij}})$ .

Ryan Robert Rosario

#### Latent Dirichlet Allocation

There are several ways to estimate LDA topic models and its many variations

- (Collapsed) Gibbs Sampling
- (Collapsed) Variational Inference
- Section Propagation
- Spectral Decomposition

With the goal begin sampling a correct topic indicator upon convergence

$$P(z_{ij}|z_{-(i,j)}, \mathbf{w}, \alpha, \beta) \propto (n_{w_{i\cdot}} + \alpha_k) \frac{n_{\cdot v} + \beta_v}{\sum_{r=1}^V n_{\cdot r} + \beta_r}$$

Ryan Robert Rosario

#### Latent Dirichlet Allocation

The LDA framework induces some limitations.

- features must be word counts to retain a multinomial likelihood and maintain conjugacy with the Dirichlet priors.
- 2 LDA is unsupervised
- as such, it is unlikely the topics discovered with LDA match the
- topic labeling is inconsistent from run to run
- topic labeling must be performed manually by eye (or with post-hoc analysis)

While researchers have successfully used TF-IDF with LDA, it was not attempted to remove another layer of complexity. Instead we focus on remedying all of the other limitations...

Ryan Robert Rosario

## Supervised Latent Dirichlet Allocation (sLDA)

Graphical model:



The generative process is identical to that of LDA except it adds one final step:

Sample response variable  $y_k | z_{ij}, \eta, \delta \sim \text{GLM}(\bar{z}, \eta, \delta)$ , where

 $ar{z} = rac{1}{N_{w_{d_j}}} \sum_{i=1}^{N_{w_{d_j}}} z_{ij}, \ \eta$  is a vector of regression coefficients and

 $\delta$  is a scale parameter, like  $\sigma^2$  and GLM is a link function.

Ryan Robert Rosario