

Automated Vocabulary Building
for Characterizing and Forecasting Elections
using Social Media Analytics

Aravindan Mahendiran

Thesis submitted to the Faculty of the
Virginia Polytechnic Institute and State University
in partial fulfillment of the requirements for the degree of

Master of Science
in
Computer Science and Applications

Naren Ramakrishnan, Chair
Cal Ribbens
Aditya Prakash

January 10, 2014
Blacksburg, Virginia

Keywords: Election Forecasting, Twitter, Query Expansion, Social Group Modeling, Probabilistic
Soft Logic

Copyright 2014, Aravindan Mahendiran

Automated Vocabulary Building
for Characterizing and Forecasting Elections
using Social Media Analytics

Aravindan Mahendiran

(ABSTRACT)

Twitter has become a popular data source in the recent decade and garnered a significant amount of attention as a surrogate data source for many important forecasting problems. Strong correlations have been observed between Twitter indicators and real-world trends spanning elections, stock markets, book sales, and flu outbreaks. A key ingredient to all methods that use Twitter for forecasting is to agree on a domain-specific vocabulary to track the pertinent tweets, which is typically provided by subject matter experts (SMEs). The language used in Twitter drastically differs from other forms of online discourse, such as news articles and blogs. It constantly evolves over time as users adopt popular hashtags to express their opinions. Thus, the vocabulary used by forecasting algorithms needs to be dynamic in nature and should capture emerging trends over time. This thesis proposes a novel unsupervised learning algorithm that builds a dynamic vocabulary using Probabilistic Soft Logic (PSL), a framework for probabilistic reasoning over relational domains. Using eight presidential elections from Latin America, we show how our query expansion methodology improves the performance of traditional election forecasting algorithms. Through this approach we demonstrate how we can achieve close to a two-fold increase in the number of tweets retrieved for predictions and a 36.90% reduction in prediction error.

Dedication

To my awesomest brother, mother and friends

Acknowledgments

First and foremost, I would like to thank my advisor Dr.Naren Ramakrishnan. He was not only instrumental in kindling my interest in machine learning and data mining but is also my mentor in many ways.

I would like to thank Dr.Cal Ribbens and Dr.Aditya Prakash for their guidance throughout my research.

I would also like to thank Dr.Bert Huang and Dr.Lise Getoor for introducing me to Probabilistic Soft Logic and providing me the PSL engine for my research.

I would like to thank my fellow graduate students Nabeel Mohammed, Rupinder Paul Khandpur and Sathappan Muthiah. Their help and support were invaluable through my graduate studies.

Contents

1	Introduction	1
1.1	Related Work	2
1.1.1	Volume based approaches	3
1.1.2	Profile Modeling	3
1.1.3	Flaws in current state of the art	4
1.2	Motivation	6
1.3	Document Overview	7
2	Dynamic Query Expansion	8
2.1	Query Expansion	8
2.2	Probabilistic Soft Logic	9
2.3	Dynamic Query Expansion using PSL	11
2.4	Results	15
3	Prediction Models	20
3.1	Unique Visitor Model	20
3.2	Regression Model	21
3.3	Performance	23

4	Evaluations and Results	25
5	Conclusions and Future Work	28
	Bibliography	30

List of Figures

2.1	Design of the query expansion pipeline.	12
2.2	Vocabulary growth with each iteration.	14
2.3	Evolution of hashtags for Henrique Capriles	16
2.4	Time series comparison for different hashtags identified for Hugo Chavez.	17
2.5	Hashtags identified for Michelle Bachelet	18
2.6	Vocabulary of hashtags identified for different elections	19
4.1	Recall of seed vocabulary vs PSL vocabulary	26

List of Tables

3.1	Regression coefficients learned for features	23
3.2	Track Record of Prediction Algorithms	23
4.1	Reduction in prediction error for Unique Visitor Model. All values shown are percentages.	26
4.2	Reduction in prediction error for Regression Model. All values shown are percentages.	27

Chapter 1

Introduction

The last decade has seen a massive explosion of online data in multiple forms, e.g., news articles, blogs and social media like Twitter, Facebook and MySpace. Twitter, in particular, is a novel micro-blogging service launched in 2006. Twitter users post messages called *tweets* on a public message board and these tweets are limited to 140 characters. Originally the tweets were meant to be personal status updates but over the years these tweets have evolved to serve multiple purposes. Today, in addition to simple status updates, tweets can be URL references to websites, or even directed messages to specific individuals. Due to the short nature of the messages users often combine multiple words into *hashtags* to convey their views. Therefore, such hashtags become the most important part of a tweet as often the entire essence of the tweet is captured in single hashtag. They evolve over time and gain more traction as users adopt the popular ones. This makes the language used in Twitter very different from other textual web content like blogs and articles.

Today twitter has grown so big¹ that it has come to be looked at as a treasure trove of mine-able data. With official APIs that are open to public, the easy access to large volumes of data has piqued the interest of scientists in the data mining community. Researchers have studied various real world phenomena like book sales [1], box office earnings [2] and even stock prices [3], and have not only demonstrated strong correlations to the chatter on Twitter but were

¹As of May 7, 2013 twitter has 555 million active registered users with 135000 new users signing up every day and approximately 1 billion tweets created every 5 days.

also able to make forecasts about future trends too.

An emerging trend is to study how online chatter can be used to model the social, economic or political landscape of a country. Political leaders have started using Twitter as a channel to mobilize supporters for their ideologies. For instance, during the 2008 US presidential election Barack Obama used social media and specifically Twitter extensively in his campaign. His victory established Twitter as a channel to garner support for a particular ideology be it political or otherwise. Bollen et al. in their research [4] used a version of the well-established psychometric instrument- Profile of Mood States (POMS) to model the mood of Twitter traffic and found correlations to a number of social and economic events that occurred during the same time period. The results from this research encouraged more researchers to study and quantify political sentiment in social media and if possible even forecast elections.

Nevertheless, there is a constant debate among political scientists about whether Twitter can be used as a surrogate for political opinion of the masses. Some believe Twitter indeed is an indicator of political opinions while others question the validity of such results. In this work we aim to answer these questions by modeling online social groups through enhanced vocabulary building and improving the performance of election prediction algorithms that use Twitter. First, we review the current state-of-the-art in election forecasting using Twitter.

1.1 Related Work

We organize the literature review into three parts. First we look at a selection of volume-based approaches to predict elections i.e., models that predict election results by merely counting the number of times a particular candidate is mentioned on Twitter. Then we review more sophisticated approaches that model the demographics of an election to make more informed predictions. Finally, we shall summarize a quite prevalent pessimistic view on such methodologies' capability to predict elections.

1.1.1 Volume based approaches

In one of the most cited papers in this space [5], the authors claim that “*The mere number of tweets reflect voter preferences and comes close to traditional polls..*” while predicting the 2010 German federal election. They go on to strongly conclude that Twitter can indeed be a valid indicator of political opinion. This was followed by [6, 7, 8, 9] all of which use volume based approaches combined with sentiment analysis. Both [6, 8] fit a regression model to opinion polls with volume of mentions and sentiment as independent variables and the opinion polls as the dependent variable. They conclude that sentiment is a weak predictor compared to share of volume. In general, the methodologies described in the above publications count the occurrence of certain handcrafted keywords and classify such tweets as positive or negative using a classifier trained on human annotated lexicons. Some advanced sentiment classifiers also provide the likelihood that the given sample of text belongs to an empirically defined psychological and structural category like anxiety, anger, and sadness.

1.1.2 Profile Modeling

More sophisticated approaches are presented in [10, 11, 12]. The authors either model the candidates or the voters in the elections rather than compute the aggregated sentiment of the mass. In [11] the authors build a support vector machine (SVM) classifier trained on manually labeled tweets and classify users into ‘left’ and ‘right’ aligned. Through latent semantic analysis (LSA), they claim to have identified the hidden structure in the data that is strongly associated with users’ political affiliations. Using this information and how political information diffuses in a network, they show an accuracy of 95% in predicting the political alignment of Twitter users. Livne et al. in [10] analyze the Twitter profiles of candidates who contested in the 2010 mid-term elections in the U.S. They identify topics specific to groups of candidates, split according to their known political orientations and use the features obtained as inputs to a regression model to predict the elections. In a similar technique Diaz-Aviles in [12] model the candidates by building an emotional

vector for each candidate using the mentions of that candidate and sentiments associated with each mention learned using the NRC Emotion Lexicon (EmoLex). They then use such profiles learned to predict the rise and fall of a candidate's popularity.

In another thread of research, Mustafaraj et al. [13] model the distribution of political content among Twitter users. They divide the users into two groups the “vocal minority” and the “silent majority”. They observe that these two groups engage in different ways over social media. The vocal minority aims to broaden the impact of tweets by re-tweeting and linking to other web content whereas the silent majority who tweet significantly lesser are more inclined to share their personal view points. Though Mustafaraj et al. do not make any election forecasts, they make observations such as *“Because of this difference between content generated by different groups, one should be aware of aggregating data and building models upon them, without verifying the underlying model that has generated the data.”*.

1.1.3 Flaws in current state of the art

Of late there have been a lot of studies showing how such models that predict elections using social media feeds are flawed [14, 15, 16, 17]. These publications not only list the obvious issues in using Twitter to predict elections but also detail recommendations on how to make such methodologies better.

Daniel Gayo-Avello surveys almost all the state-of-the-art approaches in predicting elections in his paper [15] most of which have been detailed above. According to him, post-hoc analysis of elections in retrospect must not count as valid predictions and that researchers must be wary of the *file drawer* effect, i.e., the act of filing away negative results and publishing only the positive results. His major points of argument against such models are:

- Lack of explainability: The models are tailor made to fit a particular election and that they need to be generic enough to reproduce similar results when run on other elections. In particular Metaxas et al in [14] state that any method claiming predictive power on the basis of Twitter

data should be a clearly defined algorithm and should be “explainable”, i.e., black box approaches should be avoided.

- **Vote modeling:** There is no predefined notion of “vote” that has been used to predict the elections. Most of the models aim to predict elections merely by counting the tweets related to a candidate.
- **Self-selection bias:** Biases in Twitter are ignored. Twitter is not a representative sample of the electorate demographic as not every age, gender, or social group is represented. Gayo-Avello also notes that since people tweet on a voluntary basis, the data gathered is only by those who are politically active. Another point of contention is the credibility of tweets i.e., whether the tweets are rumors, campaign propaganda, or contain misleading information intended to maliciously attack a candidate’s on-line popularity.
- **Incumbency modeling:** Since in 2008 and 2010 , 91.6% and 84% of elections were won by the the incumbent candidate respectively, Gayo-Avello argues that incumbency should be used as the baseline measure rather than just chance. He also notes that most of the methodologies are only slightly better than chance.
- **Sentiment modeling:** Gayo-Avello states that even though sentiment classifiers are a highly researched subspace of natural language processing (NLP), the accuracy of such methods are only slightly better than those of random classifiers. Further, these classifiers do not detect humor and sarcasm which in his opinion play a major role in political discussions.
- **Absence of political opinions:** Lastly in [16] Gayo-Avello akin to [13] states that abstaining from tweeting about politics can play even more important role than the ones mentioning the candidates and hence researches should also model this lack of chatter about a particular candidate or political party.

1.2 Motivation

Notwithstanding the above drawbacks, a common element among all past (and future) methods for forecasting elections using Twitter will be the aspect of filtering tweets that match a vocabulary to obtain those that are pertinent to a given election. Usually this is done by a process of querying the Twitter API for a set of keywords such as the names of the candidate and names/symbols of political parties contesting the elections. Given that the language used on Twitter is completely different from the language in newspapers and magazines, this process will likely yield a very low recall. For example, for the 2012 presidential election in Venezuela users preferred the hashtags *#elmundocochavez* and *#hayuncamino* to show their support for Hugo Chavez and Henrique Capriles respectively. Such hashtags are not known *a priori* and gain more traction and adoption closer to the election. Querying just for “chavez” or “capriles” would result in discarding a huge bank of tweets that are indicative of a candidate’s popularity. Thus, it becomes vital that any methodology that predicts elections accounts for such memes that become popular during the time period leading up to the election.

In this work we build on earlier work in social group modeling [18] and address this issue by:

- Designing and implementing a new dynamic query expansion algorithm using Probabilistic Soft Logic for vocabulary building and expansion
- Showing how the vocabulary obtained from the Dynamic Query Expansion exercise improves the retrieval of relevant tweets and improving the accuracy of prediction algorithms.
- Conducting an exhaustive evaluation of our methodology across a set of national, state-level, and mayoral elections in multiple countries of Latin America.

1.3 Document Overview

The rest of the document is organized as follows: In the next chapter, we outline the Probabilistic Soft Logic (PSL) framework and how we take advantage of this domain specific language to build our dynamic query expansion algorithm. In the third chapter, we detail two state-of-the-art algorithms used to predict elections and which can benefit from our dynamic query expansion algorithm. We then detail our experiments and present the results confirming our hypothesis in the fourth chapter. In the final chapter, we make our final conclusions and present a road map for future work.

Chapter 2

Dynamic Query Expansion

2.1 Query Expansion

In most document corpora, a single concept can be referred using multiple terms. In information retrieval (IR) this is called *synonymy* and has a huge impact on the recall of documents pertaining to the concept. Researchers address this problem by creating as exhaustive a query as possible. But when exploring the Twitter corpora it becomes almost impossible to handcraft such an expansive query as the meme and hashtag adaptations are not known *a priori*.

To address this issue IR experts use *query expansion* as a key strategy. These are iterative algorithms that are initialized with a small set of query terms. When the documents matching the query terms are returned, after basic processing (e.g., tokenization, stop-word removal and stemming), a richer vocabulary is obtained by ranking the terms in these documents by their frequency counts. The top words from this list are then used to query for documents again. This process is continued until no new terms are added to the vocabulary.

Most query expansion strategies are deterministic. We design and implement a dynamic query expansion pipeline using probabilistic soft logic (PSL) so that we build a vocabulary for a given election. We first review the PSL framework before proceeding to detail our methodology.

2.2 Probabilistic Soft Logic

Probabilistic Soft Logic [19] is a framework for collective probabilistic reasoning on relational domains. PSL models have been developed in various domains, including collective classification [20], ontology alignment [21], personalized medicine [22], opinion diffusion [23], trust in social networks [24], and graph summarization [25]. PSL represents the domain of interest as logical atoms. It uses first order logic rules to capture the dependency structure of the domain, based on which it builds a joint probabilistic model over all atoms. Instead of hard truth values of 0 (false) and 1 (true), PSL uses soft truth values relaxing the truth values to the interval $[0, 1]$. The logical connectives are adapted accordingly. This makes it easy to incorporate similarity or distance functions.

User defined *predicates* are used to encode the relationships and attributes and *rules* capture the dependencies and constraints. Each rule's antecedent is a conjunction of atoms and its consequent is a dis-junction. The rules can also be labeled with non-negative weights which are used during the inference process. The set of predicates and weighted rules thus make up a PSL program where known truth values of ground atoms are set from observed data and unknown truth values for the remaining atoms are learned using the PSL inference mechanism.

Given a set of atoms $\ell = \{\ell_1, \dots, \ell_n\}$, an interpretation defined as $I : \ell \rightarrow [0, 1]^n$ is a mapping from atoms to soft truth values. PSL defines a probability distribution over all such interpretations such that those that satisfy more ground rules are more probable. *Lukasiewicz t-norm* and its corresponding co-norm are used for defining relaxations of the logical AND and OR respectively to determine the degree to which a ground rule is satisfied. Given an interpretation I , PSL defines the formulas for the relaxation of the logical conjunction (\wedge), dis-junction (\vee), and negation (\neg) as follows:

$$\begin{aligned}\ell_1 \tilde{\wedge} \ell_2 &= \max\{0, I(\ell_1) + I(\ell_2) - 1\}, \\ \ell_1 \tilde{\vee} \ell_2 &= \min\{I(\ell_1) + I(\ell_2), 1\}, \\ \neg \ell_1 &= 1 - I(\ell_1),\end{aligned}$$

where we use \sim to indicate the relaxation of the Boolean domain. The interpretation I determines whether the rule is satisfied, if not, the *distance to satisfaction*. A rule $r \equiv r_{body} \rightarrow r_{head}$ is satisfied if and only if the truth value of head is atleast that of the body. The rule's distance to satisfaction measures the degree to which this condition is violated.

$$d_r(I) = \max\{0, I(r_{body}) - I(r_{head})\}$$

PSL then induces a probability distribution over possible interpretations I over the given set of ground atoms l in the domain. If R is the set of all ground rules that are instances of a rule from the system and uses only the atoms in I then, the probability density function f over I is defined as

$$f(I) = \frac{1}{Z} \exp\left[-\sum_{r \in R} \lambda_r (d_r(I))^p\right] \quad (2.1)$$

$$Z = \int_I \exp\left[-\sum_{r \in R} \lambda_r (d_r(I))^p\right] \quad (2.2)$$

where λ_r is the weight of the rule r , Z is the continuous version of the normalization constant used in discrete Markov random fields, and $p \in \{1, 2\}$ provides a choice between two different loss functions, linear and quadratic. The values of the atoms can be further restricted by providing linear equality and inequality constraints allowing one to encode functional constraints from the domain.

PSL provides for two kinds of inferences (a) most probable explanation (MPE) and (b) calculation of the marginal distributions. In MPE inference, given a partial interpretation with grounded atoms based on observed evidence, the PSL program infers the truth values for the unobserved atoms satisfying the most likely interpretation. In the second setting, given ground truth data for all atoms we can learn the weights for the rules in our PSL program. In our work we leverage the MPE inference.

2.3 Dynamic Query Expansion using PSL

In [18], as part of the same larger project, PSL was used to model user affiliations within groups. Specifically the authors built a PSL program for a social network where they use a set of users, their posts, messages to other users and the various groups they intend to model. The rules defined helped capture the dynamics of group affiliations through the various interactions. Through the MPE inference users were classified into different groups based on their hashtag usage and their interactions with other users.

In this thesis, we extend this earlier work to achieve what we refer to as dynamic query expansion through PSL. Similar to the query expansion methodology described earlier, we begin with an initial set of hashtags which we believe are indicative of the affinity of a particular user to a candidate contesting in the election. We refer to these hashtags as seed words. Instead of a single inference, we iteratively perform the inference over successive time windows such that the inference from window w_t is used as a prior to window w_{t+1} and the inference from that is used for window w_{t+2} , and so on. Figure 2.1 illustrates the design of the iterative algorithm for dynamic query expansion. The initial pre-processing starts with the tweet input stream which is filtered by the date range specified by the window size. For each election, tweets from a month leading up to the election were used. After extensive analysis it was determined that the most optimal window size was three days. Smaller window sizes resulted in sub-optimal inferences as there were not enough data points feeding into the PSL stage. Larger window sizes lead to memory issues as the PSL optimization procedure generates rules by substituting groundings of all possible combinations for the random variables in the rules defined. Therefore, with large number of tweets feeding into the inference process the number of rules generated explodes and causing memory issues during optimization. The tweets passing the date filter are then geo-coded using a geo-location algorithm that infers the location of a tweet. This ensures that only tweets originating from the country of interest are used. The geo-location algorithm tags the tweets with a location by looking at the GPS coordinates of the tweet if available or landmarks and locations

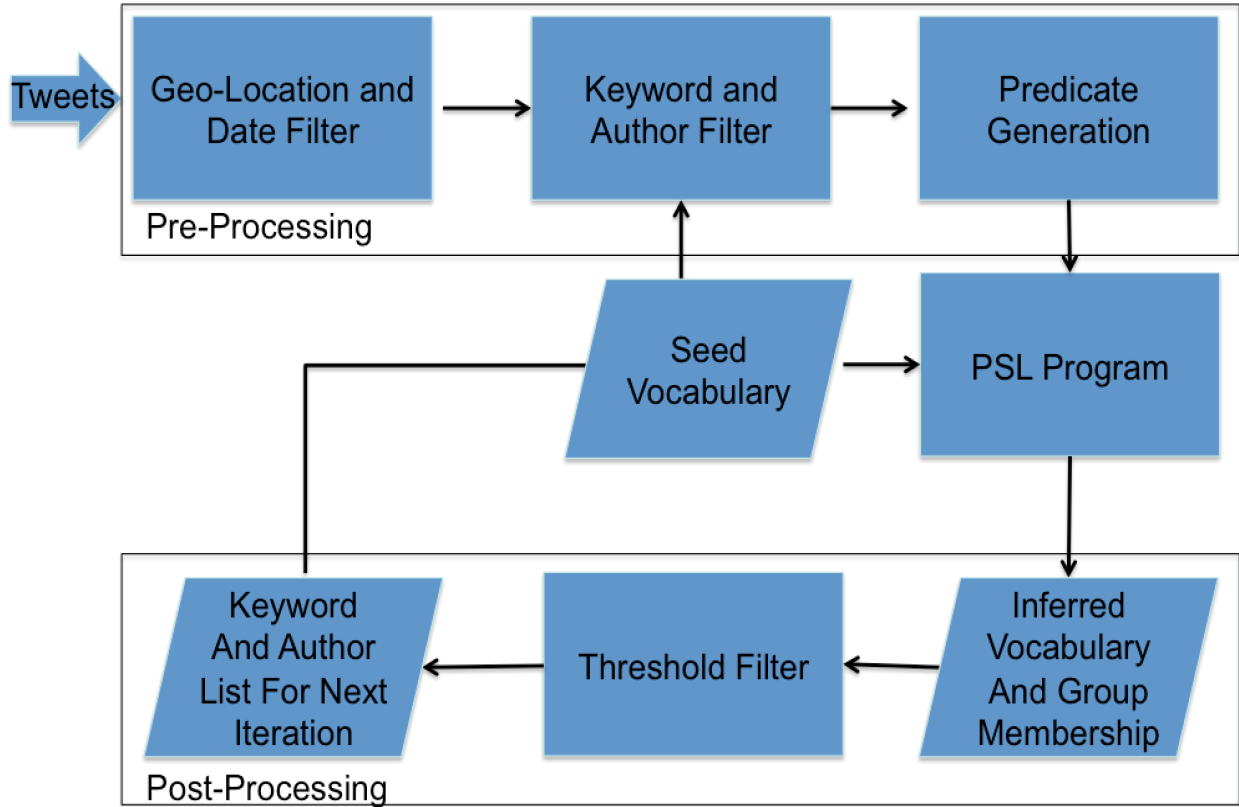


Figure 2.1: Design of the query expansion pipeline.

mentioned in the tweet or author's profile. For tweets that do not have either of these information it uses a label propagation algorithm to infer the author's location through his/her network.

The geo-tagged tweets are then tracked for the presence of a hashtag from the vocabulary for that particular iteration. In addition to filtering tweets using the vocabulary the authors whose affiliations are already inferred by the system are also used as a filtering criteria. The tweets are then converted into PSL predicates and fed into the inference process. The PSL program infers the hashtags and tweeters that are mostly associated with a particular candidate. Each author and hashtag's association with a candidate is measured using the truth value of the predicate grounding. In the post-processing step, these truth values are filtered by a threshold value to identify the hashtags and authors strongly associated to a candidate. These hashtags become a part of the vocabulary of the candidate and along with the users identified

are used as a filter criterion for the next iteration. This iterative process proceeds until the day before the election when we obtain the final vocabulary which are strongly associated with a candidate.

Within the PSL program we define predicates to encode the network. The predicates $Tweeted(U, T)$ and $Contains(T, W)$ capture the fact that a user U tweeted a tweet T and tweet T contains hashtag W respectively. Similarly, the belief that an user U or hashtag W is affiliated/associated to the group G is encoded as $Is_Member(U, G)$ and $Belongs(W, G)$ respectively. In order to capture the temporal connectivity between the iterations, in addition to the initiating the inference process with the rule

$$Seed_Word(W, G) \Rightarrow Belongs(W, G)$$

we define additional rules such as

$$Was_Member(A, G) \Rightarrow Is_Memeber(A, G)$$

$$Belonged(W, G) \Rightarrow Belongs(W, G)$$

where the predicates Was_Member and $Belonged$ are inferences from the previous time window and are loaded in as priors for the current iteration. These rules are weighted slightly lower than the recursive rules below so that the system overcomes the bias it had learned in light of new, more convincing evidence. This way hashtags that are more indicative of a user's affiliation are assigned stronger truth values or weights for every successive iteration and the truth values of hashtags that aren't are reduced. The same reasoning applies to the user-candidate affiliations(memberships) too. Below we outline the recursive PSL rules that grows the hashtag preferences and the user affiliations.

$$\begin{aligned} Tweeted(A, T) \tilde{\wedge} Contains(T, W) \tilde{\wedge} Belongs(W, G) \\ \tilde{\wedge} Positive(T) \Rightarrow Is_Member(A, G) \end{aligned}$$

$$\begin{aligned} Tweeted(A, T) \tilde{\wedge} Contains(T, W) \tilde{\wedge} Belongs(W, G) \\ \tilde{\wedge} Negative(T) \Rightarrow \sim Is_Member(A, G) \end{aligned}$$

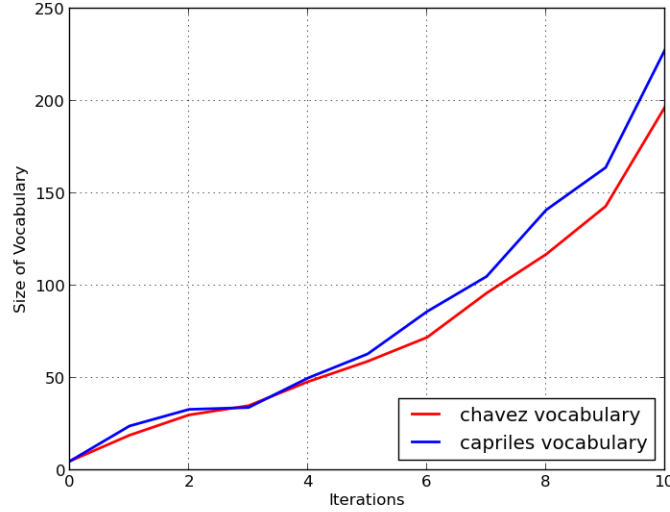


Figure 2.2: Vocabulary growth with each iteration.

$$\begin{aligned}
 &Is_Member(A, G) \tilde{\wedge} Tweeted(A, T) \tilde{\wedge} Contains(T, W) \\
 &\quad \tilde{\wedge} Positive(T) \Rightarrow Belongs(W, G)
 \end{aligned}$$

$$\begin{aligned}
 &Is_Member(A, G) \tilde{\wedge} Tweeted(A, T) \tilde{\wedge} Contains(T, W) \\
 &\quad \tilde{\wedge} Negative(T) \Rightarrow \sim Belongs(W, G)
 \end{aligned}$$

Here *Positive* and *Negative* are predicates whose truth values are calculated from the sentiment of the tweet such that the highly positive tweets get a truth value closer to 1.0 for the predicate *Positive* and highly negative tweets are assigned a truth value of 1.0 for the predicate *Negative*. Since PSL works under the close world assumption, we do not need to specify the groundings that are false i.e., positive tweets are not assigned 0.0 for the predicate *Negative* and vice-versa.

For tweets that do not have a positive or negative orientation we assign a truth value of 0.5 for both the *Positive* and *Negative* predicates.

The last two rules defined below encode the assumption that when two hashtags co-occur and one is a name of a candidate then the other hashtag is bound to be about the candidate too.

$$\begin{aligned} & \text{Contains}(T, W1) \tilde{\wedge} \text{Contains}(T, W2) \tilde{\wedge} \text{Seed_Word}(W1, G) \\ & \tilde{\wedge} \text{Positive}(T) \Rightarrow \text{Belongs}(W2, G) \end{aligned}$$

$$\begin{aligned} & \text{Contains}(T, W1) \tilde{\wedge} \text{Contains}(T, W2) \tilde{\wedge} \text{Seed_Word}(W1, G) \\ & \tilde{\wedge} \text{Negative}(T) \Rightarrow \sim \text{Belongs}(W2, G) \end{aligned}$$

Once all the tweets are loaded into the PSL program as predicates, we start the inference process by closing all the predicates except *Is_Member* and *Belongs*. This way, only their truth values of these two predicates are inferred and the other groundings of the closed predicates are regarded as facts.

2.4 Results

Figure 2.2 shows how the vocabulary grows with each iteration for the two candidates who contested the Venezuelan presidential election on October 7th 2012. Figure 2.3 shows how the hashtags for Henrique Capriles evolved during the month leading up to the election. Initially in Figure 2.3a the system begins with only a few hand picked hashtags that constitute the seed vocabulary. After a few iterations Figure 2.3b shows how the vocabulary has grown. However, not all the words identified until now remain in the final vocabulary as the system drops certain words in successive iterations. At the same time it is also noticed that hashtags like “capriles” and “hayuncamino” which are very strongly associated with Capriles consistently remain as the top ranked hashtags even after ten iterations (Figure 2.3d). It is also interesting to note that the algorithm identified hashtags like “nochavez” (Figure 2.3c) and attributed it rightly to Hugo Chavez’s primary contender, i.e., Capriles.

In Figure 2.4, the first plot elucidates how hashtags like “*elmnduconchavez*” and “*univistaconchavez*” remain highly associated with Hugo Chavez for the October 7th Presidential election. These hashtags remain indicative of a user’s affiliation throughout the month leading up to the election. Meanwhile hashtags such as “*beatles*” and “*facebook*” (in second plot) show spikes



(b) Day 6



(d) Day 30

Figure 2.3: Evolution of hashtags for Henrique Capriles

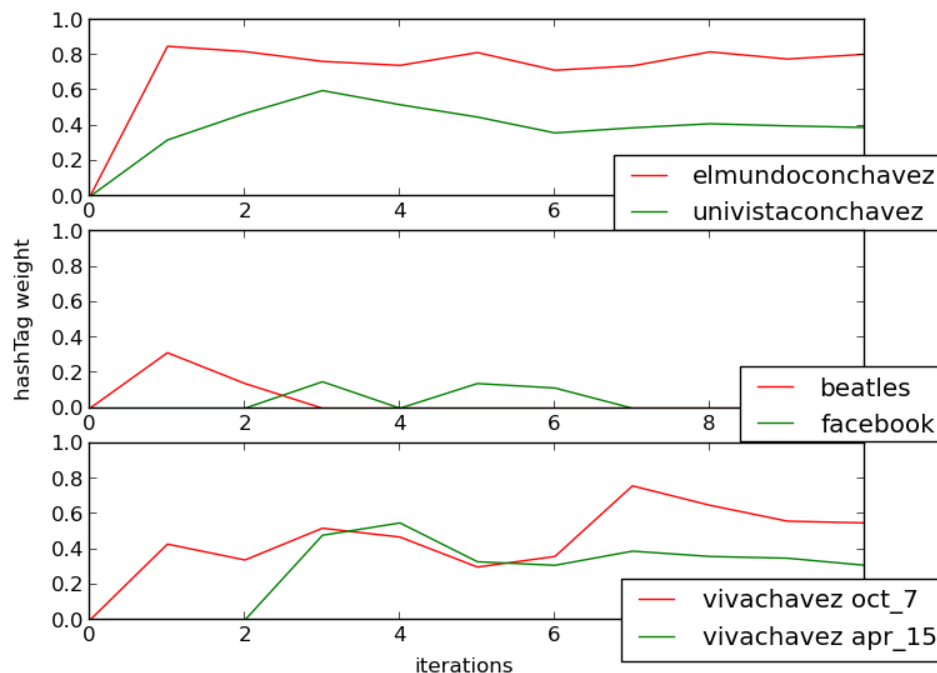


Figure 2.4: Time series comparison for different hashtags identified for Hugo Chavez.

in their time series primarily because users affiliated with Chavez used them during that time window. But as the iterative process continues, the system drops these non-informative words. The third plot presents another interesting observation. Hugo Chavez who had won the election on October 7, 2012 was diagnosed with cancer and passed away before being sworn in as the President. This triggered a re-election on April 15, 2013 where Nicolas Maduro, who had assumed the role of acting president til then, contested against Henrique Capriles in the Presidential race. The hashtag “*vivachavez*” is part of both the elections, despite the fact that Hugo Chavez did not contest the second election. It is picked up as a phrase commonly used by supporters of Nicholas Maduro whose election campaign was strategized around the death of Hugo Chavez to garner sympathy and mobilize support. Similarly variations of the hashtags “*hayuncamino*” and “*unidadvenuzela*” were returned for Henrique Capriles for both these elections.

Figure 2.5 shows the hashtags identified for Michelle Bachelet who won the Chilean Presidential elections that was decided over two rounds. The first



Figure 2.5: Hashtags identified for Michelle Bachelet

round was conducted on the 24th of November, 2013 and the 2nd round was conducted on December 15, 2013. The query expansion pipeline for Bachelet’s group in both these elections were initialized only with three seed words: *Bachelet*, *CommandoMichelle* and *PS*. The first name of the candidate was not used as it introduced a lot of noise because “Michelle” is a very common name. The first figure shows the hashtags identified for Bachelet during the first round and the second figure for the second round. It can be seen that there is a lot of overlap in the vocabulary which is the expected outcome. Similarly there was a lot of common hashtags between the two rounds of election for the other candidate, Evelyn Matthei, too.

Figure 2.6 shows the hashtags identified for the elections from Mexico, Paraguay, Honduras and Ecuador. It can be noticed that the vocabulary for the Honduran and Ecuadorean elections are quite noisy. This is primary because Twitter is not as popular in these two countries as in Venezuela or Chile and therefore the number of tweets used for the inference was significantly lesser. This in turn affected the PSL inference mechanism as the lack of evidence pushed the weights for the hashtags down. Therefore, the threshold value in the post processing step was set lower than normal to avoid losing the informative words. But this in turn also introduced a bit of noise to the vocabulary.



(a) Mexico



(b) Paraguay



(c) Honduras



(d) Ecuador

Figure 2.6: Vocabulary of hashtags identified for different elections

Chapter 3

Prediction Models

In this section we review two prediction models we adapted from current literature to test our hypothesis. Both these models will be used in conjunction with our dynamic query expansion algorithm presented in the earlier chapter. The first model is a naive model that forecasts elections based on the counts of mentions of a candidate. We dub this as “*unique visitor model*” and is adapted from [7] and [5]. The second model uses a regression fit to regress from tweet features to opinion polls and then predicts election. We dub this model as the “*regression model*” and this approach is adapted from [8] and [6].

3.1 Unique Visitor Model

The assumption here is that large parties that are more popular will have a larger social media footprint than smaller and less popular parties. This model takes advantage of this assumption and predicts elections by calculating the relative popularity of candidates contesting the election. We first define a vocabulary for each candidate. This vocabulary is crafted by hand and includes the candidate’s names and aliases, the name and acronyms for his/her political party and the official Twitter handle of the candidate and is the same as the seed vocabulary that was used to initialize the PSL pipeline. For the given time period, the tweets from the country in question are tracked for the occurrence of the terms in the vocabulary. We then build a time series

of sentiment and Klout scores from the tweets returned. The Klout score is a value provided by Klout.com that quantifies the impact each user has on social media. We use the sentiment scores provided as a part of the meta-data of the tweet. The sentiment scores typically fall in the range of $[-15, 15]$ and are provided by Lexalytics. Once a time series of the Klout and sentiment scores are built, we calculate the absolute popularity of a candidate C_d as:

$$C_d = \sum_i K_i * UCS_{id} \quad (3.1)$$

where K_i is the Klout score for user i , and UCS_{id} is User Candidate Score, the average of sentiment scores for all tweets from user i about candidate d . We then normalize the popularity scores across all candidates so that they sum to 1. This gives us the relative popularity of each candidate P_d using which we predict the elections.

$$P_d = \sum_i \frac{C_d}{C_i} \quad (3.2)$$

From the above equations, it is noticeable that each user contributes only once to the popularity score of a candidate. This was preferred to merely counting the mentions of a candidate since we desire to remove the bias of bot-generated tweets from election campaigns that artificially boosted the number of times a candidate is mentioned on Twitter.

3.2 Regression Model

In this model, in addition to Twitter data, we leverage any opinion polls available for the elections to make our predictions. Like the earlier model we track the tweets that contain a word from the vocabulary defined for each candidate. We then define a linear regression fit that uses the opinion polls as dependent variable and features generated from these tweets as independent variable. We reason that by regressing from the Twitter features to the opinion polls the bias due to Twitter being a non-representative sample can be

mitigated. We use a total of six features based on: Klout scores, number of unique users, total number of mentions, sentiment and incumbency. We normalize each of these features across all candidates to obtain the relative share of the volume. For example we define share of positive mentions (*SoPM*) as:

$$SoPM(x) = \frac{\#PositiveMentions(x)}{\sum_i \#PositiveMentions(i)} \quad (3.3)$$

and share of negative users (*SoNU*) as:

$$SoNU(x) = \frac{\sum_j K_j}{\sum_i \sum_j K_j} \quad (3.4)$$

where K_j is the Klout score of user j who tweeted negatively about a candidate. Similarly, we define share of sentiment (*SoS*) as the sum of all sentiment scores normalized across all candidates. We use a binary variable for incumbency. We then build a timeline of opinion polls. For each of the polling dates we calculate these features by using tweets created during the 10 day window leading up to the polling date. When we have more than one polling house publishing its opinion poll (for the same date) we take the average of the polls. Once we create a feature set for all the polling dates, we fit a simple least square regression as:

$$\begin{aligned} Popularity(x) = & \alpha_1 * SoPM(x) + \alpha_2 * SoNM(x) \\ & + \beta_1 * SoPU(x) + \beta_2 * SoNU(x) \\ & + \gamma * SoS(x) + \delta * Incumbency(x) + \epsilon \end{aligned} \quad (3.5)$$

Table 3.1 details the coefficients learned for each feature averaged over all the candidates from all the elections. The values confirm our hypothesis that the number of unique users and sentiment have more predictive power than total number of mentions. Intuitively it is also seen that the coefficients for share of negative users and negative mentions carry a negative weight. Another interesting observation is the fact that the incumbency binary variable is not as predictive which is contradictory to established understanding.

Feature	Coefficient Value
<i>SoPU</i>	0.4622
<i>SoNU</i>	-0.443
<i>SoPM</i>	0.1158
<i>SoNM</i>	-0.065
<i>SoS</i>	0.156
<i>Incumbency</i>	0.0

Table 3.1: Regression coefficients learned for features

Election Type	Number of Elections	Number of Correct Predictions	Accuracy
President/Prime Minister	8	8	100%
Governor	4	3	75%
Mayor	24	12	50%
Overall	36	23	63.88%

Table 3.2: Track Record of Prediction Algorithms

After learning the regression fit, we make a prediction by building such features using the same 10 day window leading up to the prediction date.

3.3 Performance

The Unique Visitor Model and the Regression Model were tested exhaustively on a total of 36 elections from Latin America during 2012 and 2013 ranging from local mayoral elections to presidential elections at the country level. It is important to note that every single election was predicted ahead of time and not in retrospect. The tweets were purchased from DataSift, an infoveillance service that resells Twitter data. On an average we collected close to 2 million unique tweets a day from over 21 countries in Latin America. Then these tweets were geo-coded using a geo-location algorithm we developed to obtain tweets from the country of interest. Only tweets from the locations pertaining to elections were used to make the predictions. For example, for the Rio De Janeiro Mayor elections only tweets from the city of Rio De Janeiro were used and similarly for state level Governor elections only tweets originating from that particular state were used. Once the tweets were filtered by location the time series of Klout and sentiment scores were calculated by tracking the tweets for the mentions of candidate.

Table 3.2 shows the overall performance of the two models. It can be noticed that the accuracy drops as the granularity of the elections increases. This is primarily due to the fact that opinion polls were available only for the country level elections. Therefore, we could not use the Regression Model for the state or city level elections. This increased the error as the predictions were generated only from the naive Unique Visitor Model. Also, from the tweets collected it was noticed that there wasn't much chatter on Twitter about smaller, local, elections. This skewed the results as the model tracking the names of the candidates was not as accurate as desirable. If the city level elections were ignored as outliers the overall accuracy of the models improves to 91.6%.

Chapter 4

Evaluations and Results

To evaluate our hypothesis about the vocabularies we test our models on eight different presidential elections from Latin America using both the seed vocabulary and the vocabulary generated by the query expansion algorithm. We use the results obtained using the seed vocabulary detailed in the previous section as a baseline score. We then use the same vocabulary to seed our PSL learning algorithm. The prediction algorithms are then run again, now by using the expanded vocabulary obtained through the query expansion. In order to remain consistent with predicting ahead of time, we track only the hashtags identified by the query expansion pipeline until that particular date.

Figure 4.1 shows the increase in the number of documents that were used by the algorithms to make a predictions. It is noticed when averaged across all the eight elections we notice close to a two-fold increase in the number of tweets that were used by these models. This is a substantial increase of relevant tweets for the domain.

To further illustrate the fact that the vocabulary used by such algorithms plays a vital role, we compare the performance of the models using the two different vocabularies. To reduce the effect of outliers we track the popularity of only the top two candidates from each election. Table 4.1 shows the reduction in prediction error for the Unique Visitor model for each candidate if the expanded vocabulary from the PSL approach is used instead of the seed vocabulary. On an average the error was reduced by a 28.60% from the original prediction error obtained by using seed vocabulary. Similarly Ta-

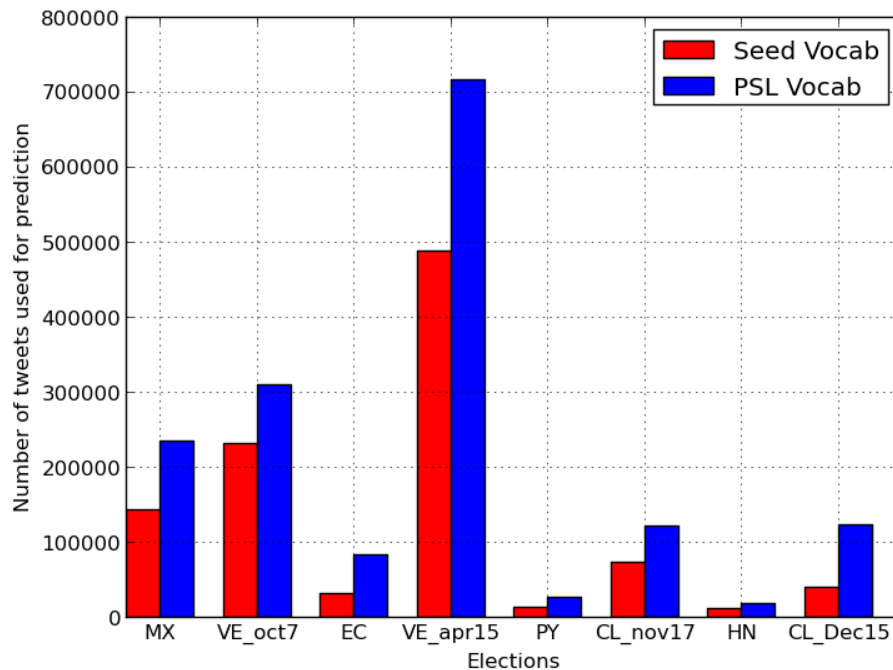


Figure 4.1: Recall of seed vocabulary vs PSL vocabulary

Election	Candidate	Actual Result	Seed Vocab.	Error	PSL Vocab.	Error
Mexico	Pena Nieto	38.1	49.26	11.11	46.43	8.28
	Lopez Obrador	31.64	25.11	6.53	27.65	4.00
Venezuela_Oct7	Hugo Chavez	55.07	63.69	8.62	55.24	0.17
	Henrique Capriles	44.31	36.31	8.00	44.76	0.45
Ecuador	Rafael Correa	57.16	32.36	24.81	32.90	24.27
	Guillermo Lasso	22.68	36.93	14.25	37.88	15.20
Venezuela_Apr15	Nicolas Maduro	50.61	42.08	8.53	44.05	6.56
	Henrique Capriles	49.12	37.98	11.14	37.14	11.98
Paraguay	Horacio Cartes	48.48	29.80	18.68	29.12	19.36
	Efrain Alegre	39.05	27.21	11.84	26.63	12.42
Chile_Nov17	Michelle Bachelet	46.70	26.62	20.08	29.92	16.78
	Evelyn Matthei	25.03	18.76	6.27	19.52	5.51
Honduras	Orlando Hernandez	36.80	28.94	7.86	34.74	2.06
	Xiomara Castro	28.70	9.67	19.03	14.20	14.50
Chile_Dec15	Michelle Bachelet	62.16	57.66	4.50	59.24	2.92
	Evelyn Matthei	37.83	42.34	4.51	40.67	2.84

Table 4.1: Reduction in prediction error for Unique Visitor Model. All values shown are percentages.

Election	Candidate	Actual Result	Seed Vocab.	Error	PSL Vocab.	Error
Mexico	Pena Nieto	38.1	46.80	8.65	39.00	0.85
	Lopez Obrador	31.64	24.67	6.97	28.64	3.00
Venezuela_Oct7	Hugo Chavez	55.07	49.89	5.18	55.89	0.82
	Henrique Capriles	44.31	36.31	8.00	43.91	0.40
Ecuador	Rafael Correa	57.16	53.33	3.84	54.33	2.84
	Guillermo Lasso	22.68	12.27	10.41	12.75	9.93
Venezuela_Apr15	Nicolas Maduro	50.61	51.45	0.84	50.58	0.03
	Henrique Capriles	49.12	35.96	13.16	38.11	11.01
Paraguay	Horacio Cartes	48.48	35.21	13.27	40.63	7.85
	Efrain Alegre	39.05	31.33	7.72	34.44	4.62
Chile_Nov17	Michelle Bachelet	46.70	38.91	7.79	41.80	4.91
	Evelyn Matthei	25.03	19.20	5.83	20.98	4.05
Honduras	Orlando Hernandez	36.80	25.16	11.64	28.30	8.50
	Xiomara Castro	28.70	16.53	12.17	24.90	3.80
Chile_Dec15	Michelle Bachelet	62.16	39.12	23.04	39.80	22.37
	Evelyn Matthei	37.83	20.88	16.95	21.68	16.15

Table 4.2: Reduction in prediction error for Regression Model. All values shown are percentages.

ble 4.2 shows the reduction in error for the Regression Model. Here an even better improvement of 45.19% reduction in error was noted. Averaging the reduction in error for both the models, the query expansion exercise was able to reduce the prediction error by 36.90%. We see greater and more consistent improvement with the regression model as the model weighs each window of tweets differently depending upon the opinion poll time series whereas the unique visitor model values them equally. Therefore, when the algorithm uses the ‘not-so-informative’ hashtags identified during the earlier iterations, the sentiment value and the counts of these mentions bring down the accuracy of the model even though at a later stage hash-tags that are strongly indicative of a user’s preference is picked up. So words such as ”facebook” which occur commonly dominate the counts and therefore skew the results even though they are dropped from the vocabulary at a later point.

Chapter 5

Conclusions and Future Work

In this work we built two prediction algorithms to forecast elections and showed how Twitter in addition to standard opinion polls can be used to pulse the political opinion of a country. We also established the reproducibility of election predictions using Twitter by forecasting more than thirty elections from Latin America. We then implemented a novel query expansion methodology using Probabilistic Soft Logic. We showed how vocabulary has a direct impact on the recall of documents and the accuracy of prediction algorithms. Using the query expansion methodology we were able to improve the accuracy of the prediction algorithms by upto 16%. It is also important to note that though we used elections to show performance gains, the query expansion system is generic and can be used to learn a vocabulary for any given domain.

Further, this work was motivated towards a future goal to model the electorate demographics. With more fine grained data about the gender, age and exact location of a user it is possible to infer the preferences at a group level rather than at a user level. This would enable us to study the various interactions between groups and individual users in more detail and thus make more informed election predictions.

While modeling the interactions between users using PSL we used hard coded rule weights. With a labeled data set providing ground truths about user affiliations it is possible to study the various forms of interactions by using the weight learning mechanism of PSL to understand which theories about

social group modeling are more probable and which aren't.

Also, while tracking the candidates using the seed vocabulary we noticed that using the names of candidates introduced a lot of noise as some of the names such as 'Jose' are very common. To avoid this we propose to use more sophisticated named entity recognition to improve the accuracy of the tweets returned about a particular candidate.

In addition to Twitter, we also propose to use other data sources such as Google Search Trends, Facebook and web blogs to track the popularity of candidates online.

Bibliography

- [1] Daniel Gruhl, Ramanathan Guha, Ravi Kumar *et al.*, “The predictive power of online chatter,” in *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. ACM, 2005, pp. 78–87.
- [2] Sitaram Asur and Bernardo A Huberman, “Predicting the future with social media,” in *Web Intelligence and Intelligent Agent Technology (WI-IAT), 2010 IEEE/WIC/ACM International Conference on*, vol. 1. IEEE, 2010, pp. 492–499.
- [3] Johan Bollen, Huina Mao, and Xiaojun Zeng, “Twitter mood predicts the stock market,” *Journal of Computational Science*, vol. 2, no. 1, pp. 1–8, 2011.
- [4] Johan Bollen, Huina Mao, and Alberto Pepe, “Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena.” in *ICWSM*, 2011.
- [5] Andranik Tumasjan, Timm Oliver Sprenger, Philipp G Sandner *et al.*, “Predicting elections with twitter: What 140 characters reveal about political sentiment.” *ICWSM*, vol. 10, pp. 178–185, 2010.
- [6] Brendan O’Connor, Ramnath Balasubramanyan, Bryan R Routledge *et al.*, “From tweets to polls: Linking text sentiment to public opinion time series.” *ICWSM*, vol. 11, pp. 122–129, 2010.
- [7] Diego Saez-Trumper, Wagner Meira, and Virgilio Almeida, “From total hits to unique visitors model for elections forecasting,” in *International Conference on Web Science*, 2011.

- [8] Adam Bermingham and Alan F Smeaton, “On using twitter to monitor political sentiment and predict election results,” 2011.
- [9] Gianluca Demartini, Stefan Siersdorfer, Sergiu Chelaru *et al.*, “Analyzing political trends in the blogosphere.” in *ICWSM*, 2011.
- [10] Avishay Livne, Matthew P Simmons, Eytan Adar *et al.*, “The party is over here: Structure and content in the 2010 election.” in *ICWSM*, 2011.
- [11] Michael D Conover, Bruno Gonçalves, Jacob Ratkiewicz *et al.*, “Predicting the political alignment of twitter users,” in *Privacy, security, risk and trust (passat), 2011 ieee third international conference on and 2011 ieee third international conference on social computing (socialcom)*. IEEE, 2011, pp. 192–199.
- [12] Ernesto Diaz-Aviles, Claudia Orellana-Rodriguez, and Wolfgang Nejdl, “Taking the pulse of political emotions in latin america based on social web streams,” in *Web Congress (LA-WEB), 2012 Eighth Latin American*. IEEE, 2012, pp. 40–47.
- [13] Eni Mustafaraj, Samantha Finn, Carolyn Whitlock *et al.*, “Vocal minority versus silent majority: Discovering the opinions of the long tail,” in *Privacy, security, risk and trust (passat), 2011 ieee third international conference on and 2011 ieee third international conference on social computing (socialcom)*. IEEE, 2011, pp. 103–110.
- [14] Panagiotis Takis Metaxas, Eni Mustafaraj, and Daniel Gayo-Avello, “How (not) to predict elections,” in *Privacy, security, risk and trust (PASSAT), 2011 IEEE third international conference on and 2011 IEEE third international conference on social computing (SocialCom)*. IEEE, 2011, pp. 165–171.
- [15] Daniel Gayo-Avello, “” i wanted to predict elections with twitter and all i got was this lousy paper”—a balanced survey on election prediction using twitter data,” *arXiv preprint arXiv:1204.6441*, 2012.
- [16] —, “Don’t turn social media into another ‘literary digest’ poll,” *Communications of the ACM*, vol. 54, no. 10, pp. 121–128, 2011.

- [17] Daniel Gayo-Avello, Panagiotis Takis Metaxas, and Eni Mustafaraj, “Limits of electoral predictions using twitter.” in *ICWSM*, 2011.
- [18] Bert Huang, Stephen H Bach, Eric Norris *et al.*, “Social group modeling with probabilistic soft logic,” in *NIPS Workshop on Social Network and Social Media Analysis: Methods, Models, and Applications*, 2012.
- [19] Angelika Kimmig, Stephen Bach, Matthias Broecheler *et al.*, “A short introduction to probabilistic soft logic,” in *Proceedings of the NIPS Workshop on Probabilistic Programming: Foundations and Applications*, 2012, pp. 1–4.
- [20] Matthias Broecheler and Lise Getoor, “Computing marginal distributions over continuous markov networks for statistical relational learning,” in *Advances in Neural Information Processing Systems*, 2010, pp. 316–324.
- [21] Matthias Broecheler, Lilyana Mihalkova, and Lise Getoor, “Probabilistic similarity logic,” *arXiv preprint arXiv:1203.3469*, 2012.
- [22] Stephen H Bach, Matthias Broecheler, Stanley Kok *et al.*, “Decision-driven models with probabilistic soft logic,” 2010.
- [23] Stephen Bach, Matthias Broecheler, Lise Getoor *et al.*, “Scaling mpe inference for constrained continuous markov random fields with consensus optimization,” in *Advances in Neural Information Processing Systems*, 2012, pp. 2663–2671.
- [24] Bert Huang, Angelika Kimmig, Lise Getoor *et al.*, “Probabilistic soft logic for trust analysis in social networks,” in *International Workshop on Statistical Relational AI*, 2012.
- [25] Alex Memory, Angelika Kimmig, Stephen Bach *et al.*, “Graph summarization in annotated data using probabilistic soft logic,” *status: accepted*, 2012.