

# Joint Sentence–Document Model for Manifesto Text Analysis

Anonymous ALTA submission

## Abstract

Election programs (so-called manifestos) are published verbal declaration of the intentions, motives, or views of the political party. Political scientists use manifestos to understand policy-relevant themes and also quantify a party’s position on the left–right spectrum. Rather than handling the two tasks separately, we propose a joint sentence–document model for sentence-level thematic classification and document-level position quantification using manifestos in different languages. In order to handle text from multiple languages we exploit continuous neural embeddings for semantic text representation. We empirically show that the proposed joint model performs better than state-of-art approaches for the document-level task using manifestos from thirteen countries, written in six different languages.

## 1 Introduction

Among many actors, political parties are at the core of contemporary democratic systems. One of the widely used datasets by political scientists is the Comparative Manifesto Project (CMP) dataset, initiated by [Volkens et al. \(2011\)](#), that collects party election programs (so-called manifestos) from elections in many countries around the world. The goal of the project is to provide a large data collection to support political studies on electoral processes. A sub-part of the manifestos has been manually annotated at the sentence-level with one of over fifty fine-grained political themes, divided into 7 coarse-grained topics (see Table 5). While manual annotations are very useful for political analyses, they come with two major drawbacks. First, it is very time-consuming and labor-

intensive to manually annotate each sentence with the correct category from a complex annotation scheme. Secondly, coder preferences towards particular categories might lead to annotation inconsistencies and affect comparability between manifestos annotated by different coders ([Mikhaylov et al., 2012](#)). In order to overcome these challenges, fine and coarse-level manifesto sentence classification was addressed using supervised machine learning techniques ([Verberne et al., 2014](#); [Zirn et al., 2016](#)). Nonetheless, manually-coded manifestos remain the crucial data source for studies in computational political science ([Lowe et al., 2011](#); [Nanni et al., 2016](#)).

Other than the sentence-level labels, the manifesto text also has document-level signals, which quantify its position on the left–right spectrum ([Slapin and Proksch, 2008](#)). Though sentence-level classification and document-level quantification tasks are inter-dependent, existing work handles them separately. We instead propose a joint sentence–document model to jointly model the two tasks. Overall, the contributions of this work are as follows:

- we empirically study the utility of multi-lingual embeddings for cross-lingual manifesto text analysis — at the sentence and document-levels.
- we evaluate the effectiveness of modelling the sentence- and document-level tasks together.
- we study the value of *country* information used in conjunction with text for the document-level regression task.

## 2 Related Work

The recent adoption of NLP methods has led to significant advances in the field of Compu-

tational Social Science (Lazer et al., 2009), including political science (Grimmer and Stewart, 2013). Some popular tasks addressed with political text include: party position analysis (Biessmann, 2016); political leaning categorization (Akoglu, 2014; Zhou et al., 2011); stance classification (Sridhar et al., 2014); identifying keywords, themes & topics (Karan et al., 2016; Nallapati et al., 2004; Ding et al., 2011); emotion analysis (Rheault, 2016); and sentiment analysis (Bakliwal et al., 2013). The source data includes manifestos, political speeches, news articles, floor debates and social media posts.

With the increasing availability of large-scale datasets and computational resources, large-scale comparative political text analysis has gained the attention of political scientists (Lucas et al., 2015). For example, rather than analyzing the political manifestos of a particular party during an election, mining different manifestos across countries over time can provide deeper comparative insights into political change.

Existing classification models, except (Glavaš et al., 2017), utilize discrete representation of text (i.e., bag of words). Also, most of the work analyzes manifesto text at the country level. Recent work has demonstrated the utility of neural embeddings for multi-lingual coarse-level topic classification (7 major categories) over manifesto text (Glavaš et al., 2017). The authors show that multi-lingual embeddings are more effective in the cross-lingual setting, where labeled data is used from multiple languages. In this work, we focus on cross-lingual fine-grained thematic classification (57 categories in total), where we have labeled data for all the languages.

For the document-level quantification task, much work has used label count aggregation of manually-annotated sentences as features (Lowe et al., 2011; Benoit and Däubler, 2014), while other work has used dictionary-based supervised methods, or unsupervised factor analysis based techniques (Hjorth et al., 2015; Bruinsma and Gemenis, 2017). The latter method uses discrete word representations and deals with mono-lingual text only. In Glavas et al. (2017), the authors leverage neural embeddings for cross-lingual EU parliament speech text quantification with two pivot texts for extreme left and right positions. They represent the documents using word embeddings averaged with TF-IDF scores as weights. All these

approaches model the sentence and document-level tasks separately.

### 3 Manifesto Text Analysis

In the CMP, trained annotators manually label manifesto sentences according to the 57 fine-grained political categories (shown in Table 5), which are grouped into seven policy areas: External Relations, Freedom and Democracy, Political System, Economy, Welfare and Quality of Life, Fabric of Society, and Social Groups. Political parties either write their promises as a bulleted list of individual sentences, or structured as paragraphs (an example is given in Figure 4), providing more information on topic coherence. Also the length of documents, measured as the number of sentences, varies greatly between manifestos. The typical length (in sentences) over manifestos (948 in total) from 13 countries — Austria, Australia, Denmark, Finland, France, Germany, Italy, Ireland, New Zealand, South Africa, Switzerland, United Kingdom and United States — is  $516.7 \pm 667$ . Variance in the number of sentences across documents in conjunction with class imbalance makes automated thematic classification a challenging task.

A sentence is split into multiple segments if it discusses unrelated topics or different aspects of a larger policy, e.g. (as indicated by the different colors, and associated integer labels):

We need to address our close ties with our neighbours (107) as well as the unique challenges facing small business owners in this time of economic hardship. (402)

Such examples are not common, however.<sup>1</sup> Also the segmentation was shown to be inconsistent and to have no effect on quantifying the proportion of sentences discussing various topics and document-level regression tasks (Däubler et al., 2012). Hence, consistent with previous work (Biessmann, 2016; Glavaš et al., 2017), we consider the sentence-level classification to be a multi-class single-label problem. We use the segmented text when available (especially for evaluation), and complete sentences otherwise.

A manifesto as a whole can be positioned on the left-right spectrum based on the proportion of

<sup>1</sup>In Däubler et al. (2012), based on a sample of 15 manifestos, the authors noted that around 7.7% of sentences encode multiple topics.

topics discussed. We use the RILE score, which is defined as the difference between the count of sentences discussing left- and right-leaning topics (Budge and Laver, 1992):

$$\text{RILE} = \sum_{r \in R} \text{per}_r - \sum_{l \in L} \text{per}_l \quad (1)$$

where,  $R = \{104, 201, 203, 305, 401, 402, 407, 414, 505, 601, 603, 605, 606\}$  and  $L = \{103, 105, 106, 107, 202, 403, 404, 406, 412, 413, 504, 506, 701\}$ , and  $\text{per}_t$  denotes the share of each topic  $t$  as given in Table 5, per document. Note that the RILE score is provided for almost all the manifestos in CMP dataset, but the sentence-level annotations are provided only for a subset of manifestos. That is, in some cases, the underlying annotations that the RILE score calculation was based on is often not available for a given manifesto.

## 4 Proposed Approach

We propose a joint sentence–document model to classify manifesto sentences into one out of 57 categories and also quantify the document-level RILE score. The joint formulation is not only to capture the task inter-dependencies, but also to use annotations at different levels of granularity (sentence and document) effectively — a RILE score is available for 948 manifestos from 13 countries, whereas sentence-level annotations are available only for 235 manifestos. We use a neural network to model the sentence-level classification and document-level regression tasks. The proposed architecture is given in Figure 1. Since the text across countries is multi-lingual in nature, we use multi-lingual embeddings to represent words ( $e_w$ ). We refer to the total set of manifestos available for training as  $D$ , and the subset which is annotated with sentence-level labels as  $D_s$ . We denote each manifesto as  $d$ , which has  $l_d$  sentences  $s_1, s_2, \dots, s_{l_d}$ .

### 4.1 Sentence-level Model

We represent each sentence using the average embedding of its constituent words,  $s_i = \frac{1}{|s_i|} \sum_{w \in s_i} e_w$ . The average embedding representation is given as input to hidden layer with rectified linear activation units (ReLU) to get the hidden representation. Finally, the predictions are obtained using a softmax layer, which takes the hidden representation as input and gives probability

of 57 classes as output. We use cross-entropy loss function for sentence-level model. For sentences in  $D_s$ , with ground truth labels  $Y_s$  (one-of-K encoding), the loss function is given as follows:

$$L_S(D_s, Y_s) = - \frac{1}{\sum_{i=1}^{D_s} l_{d_i}} \sum_{i=1}^{D_s} \sum_{j=1}^{l_{d_i}} \sum_{k=1}^K Y_{s_{ijk}} \log \hat{Y}_{s_{ijk}} \quad (2)$$

### 4.2 Joint Sentence–Document Model

Using the neural network, we model the sentence-level classification and document-level regression tasks together. In the joint model, we use an unrolled (time-distributed) neural network model for the sentences in a manifesto. Here, the model minimizes cross-entropy loss for sentences over each temporal layer (1 to  $n$ ). We use average-pooling with the concatenated hidden representations ( $H_s$ ) and predicted output distributions ( $\hat{Y}_s$ ) of individual sentences, to represent a document.<sup>2</sup>

$$S_r = \frac{1}{|l_{d_i}|} \sum_{s \in d_i} [\hat{Y}_s, H_s] \quad (3)$$

Since the range of RILE is  $[-100, 100]$ , we scale it to the range  $[-1, 1]$  based on a final tanh layer, with  $z = W_r^\top H_d$ , where  $H_d = \text{ReLU}(W_d^\top S_r)$ , as input. Since it is a regression task, we minimize the mean-squared error loss function between the predicted  $\hat{y}_d$  and actual RILE score  $y_d$ . With the given RILE scores for training documents ( $Y_d$ ) and estimated scores ( $\hat{Y}_d$ ), the loss function is given as follows:

$$L_d(\hat{Y}_d, Y_d) = \frac{1}{|D|} (\hat{Y}_d - Y_d)^\top (\hat{Y}_d - Y_d). \quad (4)$$

Overall, the loss function for the joint model, combining Equations 2 and 4, is:

$$\alpha L_S(D_s, Y_s) + (1 - \alpha) L_d(\hat{Y}_d, Y_d) \quad (5)$$

where  $0 \leq \alpha \leq 1$ . We evaluate both cascaded and joint training for this objective function:

**Cascaded Training:** The sentence-level model is trained using  $D_s$ , to minimize  $L_S(D_s, Y_s)$  in Equation 2, and the pre-trained model is used to obtain document-level representation  $S_r$  for all the manifestos in training set  $D$ .

<sup>2</sup>We observed that the concatenated representation performed better than using either hidden representation or output distribution.

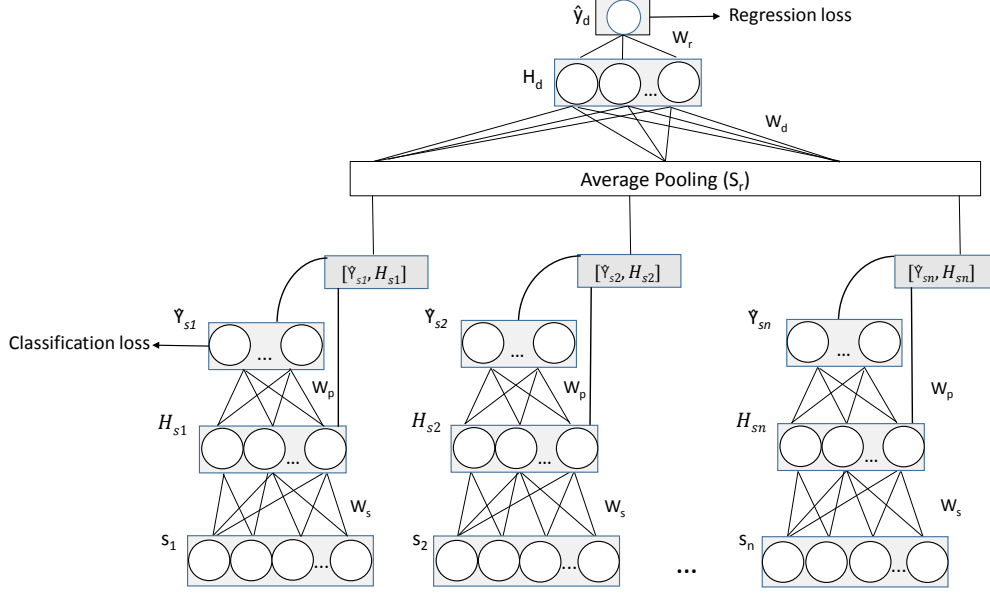


Figure 1: Hierarchical Neural Network for Joint Sentence–Document Analysis.  $s_1, s_2, \dots, s_n$  are input sentences,  $W_s$  and  $W_p$  are shared across unrolled sentences.  $\hat{Y}_{s_i}$  denotes 57 classes and  $\hat{y}_d$  denotes the estimated RILE score

Then the document-level regression task is trained to minimize  $L_d(\hat{Y}_d, Y_d)$  from Equation 4. Here, the sentence-level model parameters are fixed when the document-level regression model is trained using  $S_r$ .

**Joint Training:** As in cascaded training, the sentence-level model is pre-trained using labeled sentences. Then the entire network is updated by minimizing the joint loss function from Equation 5. Here the sentence-level model uses both labeled and unlabeled data.

We use the Adam optimizer (Kingma and Ba, 2014) for parameter estimation.

The proposed architecture evaluates the effectiveness of posing sentence-level topic classification as a precursor to perform document-level RILE prediction, rather than learning a model directly. We also study the effect of the quantity of annotated text at both the sentence- and document-level for the RILE prediction task.

## 5 Experiments

### 5.1 Setting

As mentioned earlier, we use manifestos collected and annotated by political scientists as part of CMP. In this work, we used 948 manifestos from

13 countries, which are written in 6 different languages — Danish (Denmark), English (Australia, Ireland, New Zealand, South Africa, United Kingdom, United States), Finnish (Finland), French (France), German (Austria, Germany, Switzerland), and Italian (Italy). Out of the 948 manifestos, 235 are annotated with sentence level labels (from Table 5). We have RILE scores for all the 948 manifestos. Statistics about number of annotated documents and sentences across languages are given in Table 1.

| Lang. | # Docs (Ann.) | # Sents (Ann.)  |
|-------|---------------|-----------------|
| da    | 175 (36)      | 32161 (8762)    |
| en    | 312 (94)      | 227769 (73682)  |
| fi    | 97 (16)       | 18717 (8503)    |
| fr    | 53 (10)       | 24596 (5559)    |
| de    | 216 (65)      | 146605 (79507)  |
| it    | 95 (14)       | 40010 (4918)    |
| Total | 948 (235)     | 489858 (180931) |

Table 1: Statistics of dataset, ‘Ann.’ refers to annotated. *da* - Danish, *en* - English, *fi* - Finnish, *fr* - French, *de* - German, *it* - Italian

We use off-the-shelf multi-lingual word embeddings<sup>3</sup> to represent words. We chose em-

<sup>3</sup><http://128.2.220.95/multilingual>



beddings trained using translation invariance approach (Ammar et al., 2016), with size 512 for our work, since we found it empirically better compared to other approaches. The Neural Network model has a single hidden layer for all the sentence and document-level approaches.

## 5.2 Sentence-Level Classification

We first compare traditional bag-of-words discrete representation with distributed neural representation for words for fine-grained thematic classification, under mono-lingual training setting (*Mono-lingual*). Hence we compare the following approaches.

Bag of words (*BoW-LR*, *BoW-NN*): We use TF-IDF representation for sentences and build a model for each language separately. We use Logistic Regression classifier, which is the state-of-art approach for fine-grained manifesto sentence classification (Biessmann, 2016). We refer to this approach as *BoW-LR*. We also use Neural Network classifier, which we refer as *BoW-NN*.

Language-wise average embedding (*AE-NN<sub>m</sub>*): We build a Neural Network classifier per language, with average neural embedding as sentence representation.

Since distributed representation allows to leverage text across languages, we evaluate the following approaches with combined training sentences across languages (*Cross-lingual*).

Convolutional Neural Network (*CNN*): CNN was shown to be effective for cross-lingual manifesto text coarse-level (7 major domains as shown in Table 5) topic classification (Glavaš et al., 2017). So, we evaluate CNN with a similar architecture — single convolution layer (32 filters with window size 3), followed by single max pooling layer and finally a softmax layer. We use neural embeddings to represent words.

Combined average embedding (*AE-NN<sub>c</sub>*): We build a Neural Network classifier with training instances combined across languages, with average neural embedding as sentence representation. This is our proposed approach for sentence-level model.

==== HEAD Commonly for all empirical evaluation, we compute micro-averaged performance with 80-20% train-test ratio across 10 runs with random split (at document-level), where the 80% split also contains sentence-level annotated documents proportionally. Optimal model parameters we found for the proposed model (Figure 1) are  $|H_s| = 300$ ,  $|H_d| = 10$ . We compute F-score<sup>4</sup> to evaluate sentence classification performance. Sentence classification performance is given in Table 1. In the mono-lingual setting (Table 1), using word embeddings did not provide better performance compared to bag-of-words. (except for *en* which was the source language for obtaining multi-lingual embeddings). ===== Commonly for all empirical evaluation, we compute micro-averaged performance with 80-20% train-test ratio across 10 runs with random split (at document level), where the 80% split also contains sentence level annotated documents proportionally. Optimal model parameters we found for the proposed model (Figure 1) are  $|H_s| = 300$ ,  $|H_d| = 10$ . We compute F-score<sup>5</sup> to evaluate sentence classification performance. Sentence classification performance is given in Table 2. In the mono-lingual setting (Table 2), using word embeddings did not provide better performance compared to bag-of-words. (except for *en* which was the source language for obtaining multi-lingual embeddings).   
bd7e96b892db911843b42cf19e12d0a930ed902a

Under cross-lingual setting, *AE-NN<sub>c</sub>* is the sentence-level Neural Network model. We use *AE-NN<sub>c</sub>* in the *cascaded training* for obtaining document-level RILE prediction. Note that in *cascaded training*, sentence and document-level models are trained separately in a cascaded fashion. Joint-training results where the sentence model is trained in a semi-supervised way together with document-level regression task is referred to as *JT<sub>s</sub>*. We set  $\alpha=0.4$  (in equation (4)) empirically which gave the best score for both sentence and document-level tasks. We observed a trade-off in performance with different  $\alpha$ , with lesser  $\alpha$  (0.1), document-level correlation increases (to 0.52) while sentence-level F-score decreases (to 0.33). Higher value of  $\alpha$  (0.9) gives performance closer to cascaded training. *JT<sub>s</sub>* has a comparable

<sup>4</sup>Harmonic mean of precision and recall, [https://en.wikipedia.org/wiki/F1\\_score](https://en.wikipedia.org/wiki/F1_score)

<sup>5</sup>Harmonic mean of precision and recall, [https://en.wikipedia.org/wiki/F1\\_score](https://en.wikipedia.org/wiki/F1_score)

performance with  $AE-NN_c$ . It must be considered that the (sentence and document-level) model built using cascaded training is a special case of jointly trained model, with suitable values of  $\alpha$  (closer to 1). Also  $CNN$ ,  $AE-NN_c$  and  $JT_s$  use combined training instances across languages compared to other approaches. Except for *da* and *it* cross-lingual training performs better than mono-lingual setting. Overall, the proposed approach shows slight improvement over baseline approaches for sentence classification task.

### 5.3 Document-Level Regression

For document-level regression task, following are the baseline approaches. Note that we use tanh output for all the models, since the range of re-scaled RILE is from -1 to +1.

Bag of words ( $BoW-NN_d$ ): We use TF-IDF representation for documents and build a Neural Network model for each language.

Average embedding ( $AE-NN_d$ ): We use average embedding of words in the document to build a Neural Network model across languages.

Bag-of-Centroids ( $BoC$ ): Here the word embeddings are clustered into  $K$  different clusters using K-Means clustering algorithm, and words in each document are assigned to clusters based on its euclidean-distance ( $dist$ ) to cluster-centroids ( $C$ ).

$$cluster(word) = \underset{k}{\operatorname{argmin}} dist(C_k, word)$$

Finally, each document is represented by the distribution of words mapped to different clusters ( $1 \times K$  vector). We use a Neural Network regression model with bag-of-centroids representation. Results with  $K=1000$ , which performed best is given in Table 3.

sentence-level model and RILE formulation ( $AE-NN_c^{rile}$ ): Here the predictions of sentence-level model ( $AE-NN_c$ ) are used directly with RILE formulation (equation (1)) to derive RILE score for manifestos.

Cross-lingual scaling ( $CLS$ ): This is a recent unsupervised approach for cross-lingual political speech text positioning task (Glavas et al., 2017). Authors use average word-embeddings weighed by TF-IDF score to rep-

resent documents.<sup>6</sup> Then a graph is constructed using pair-wise distance of documents. Given two pivots texts for extreme left and right positions  $[-1, +1]$ , label propagation approach is used to quantify other documents in the graph.

RILE score regression performance results are given in Table 3. Other than  $BoW-NN_d$  all other approaches are cross-lingual. We evaluate document-level performance using mean-squared-error (MSE) and Pearson correlation ( $r$ ). The proposed approach's performance, using cascaded training is referred to as  $Cas_d$  and jointly trained model is referred to as  $JT_d$ . Overall the jointly trained model performs best for document-level task, with a comparable performance at sentence-level task.

### 5.4 Quantity of Annotation

We measure the importance of annotated text at sentence and document-level for RILE score regression task. We vary the percentage of labeled data, while keeping the test sample size at 20% as before. In the first setting, we keep the training ratio of documents at 80%, within that 80% we increase the proportion of documents with sentence-level annotations — from 0 (document average embedding setting,  $AE-NN_d$ ) to 80%. Results are given in Figure 2. Similarly, in the other setting, we keep the training set with 80% sentence-level annotated documents (which is  $\sim 20\%$  of the total data), and add documents (with only RILE score), increasing the training set from 20 to 80%. Results of this study are given in Figure 3. We observed that, jointly-trained model uses sentence-level annotations more effectively than cascaded approach (Figure 2) — even with less sentence-level annotations. Also, with less document-level signal (up to 40%) for training, both the approaches perform similarly ( $r$ ). As the training ratio increases, joint-training leverages both sentence and document-level signals effectively.

### 5.5 Use of Country Information

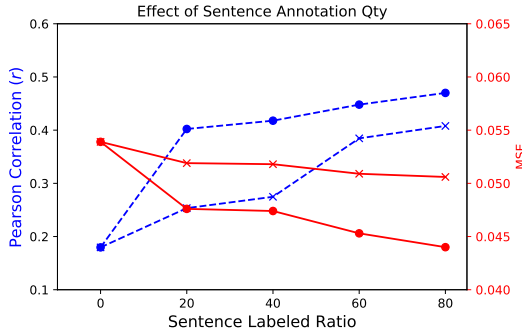
Since the definition of left-right varies between countries, we study the influence of *country* information in the proposed model with *joint-training*.

<sup>6</sup>We use this aggregate representation since it was shown to be better than word alignment and scoring approach (Glavas et al., 2017)

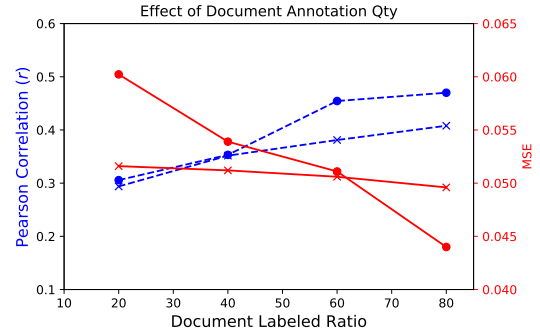
| Lang. | Mono-lingual |             |                    | Cross-lingual |                    |                 |
|-------|--------------|-------------|--------------------|---------------|--------------------|-----------------|
|       | BoW-LR       | BoW-NN      | AE-NN <sub>m</sub> | CNN           | AE-NN <sub>c</sub> | JT <sub>s</sub> |
| da    | 0.29         | <b>0.35</b> | 0.24               | 0.30          | 0.28               | 0.30            |
| en    | 0.36         | 0.38        | <b>0.42</b>        | 0.40          | <b>0.42</b>        | 0.41            |
| fi    | 0.21         | 0.29        | 0.26               | <b>0.30</b>   | 0.27               | 0.26            |
| fr    | 0.28         | 0.36        | 0.24               | 0.36          | 0.37               | <b>0.38</b>     |
| de    | 0.30         | 0.31        | 0.31               | 0.31          | 0.31               | <b>0.33</b>     |
| it    | 0.32         | <b>0.33</b> | 0.25               | 0.30          | 0.32               | 0.26            |
| Avg.  | 0.32         | 0.34        | 0.35               | 0.34          | <b>0.36</b>        | 0.35            |

Table 2: Micro-Averaged F-measure. Best scores are given in bold under each language setting

| Approach                           | MSE(↓)       | $r(\uparrow)$ |
|------------------------------------|--------------|---------------|
| BoW-NN <sub>d</sub>                | 0.054        | 0.23          |
| AE-NN <sub>d</sub>                 | 0.057        | 0.14          |
| BoC                                | 0.052        | 0.33          |
| AE-NN <sub>c</sub> <sup>rile</sup> | 0.060        | 0.35          |
| CLS                                | —            | 0.24          |
| Cas <sub>d</sub>                   | 0.050        | 0.41          |
| JT <sub>d</sub>                    | <b>0.044</b> | <b>0.47</b>   |

Table 3: RILE score prediction performance. Best scores are given in bold (higher is better for  $r$ , and lower is better for MSE).Figure 2: Fixing 80% training documents with RILE score, ratio of documents with sentence-level annotations is varied.  $\times$  denotes Cas<sub>d</sub> and  $\circ$  denotes JT<sub>d</sub>

We use two ways to incorporate country information (Hoang et al., 2016): (a) *stack* — one-hot encoding (13 countries, 1 X 13 vector) of each manifesto’s *country* is concatenated with hidden representation of the document ( $S_r$  in Figure 1) (b) *non-linear stack* — one-hot-encoded country vector is passed through a hidden layer with tanh non-linear activation and concatenated with  $S_r$ . With

Figure 3: Fixing 80% training documents with sentence-level annotations, ratio of documents with RILE score is varied.  $\times$  denotes Cas<sub>d</sub> and  $\circ$  denotes JT<sub>d</sub>

| Approach         | MSE             | $r$           |
|------------------|-----------------|---------------|
| stack            | 0.045 (0.001 ↓) | 0.49 (0.02 ↑) |
| non-linear stack | 0.048 (0.004 ↓) | 0.48 (0.01 ↑) |

Table 4: RILE score prediction performance with *country* information. Difference compared to JT<sub>d</sub> is given within paranthesis.  $\uparrow$  – improvement,  $\downarrow$  – decrease in performance

both the models we observed mild improvement in correlation (given in Table 4).

## 6 Conclusion and Future Work

In this work we evaluated the utility of a joint sentence–document model for sentence-level thematic classification and document-level RILE score regression tasks. Our observations are as follows: (a) joint model performs better than state-of-art approaches, where cascaded training can be seen as a special case of joint training. (b) joint-training leverages sentence-level annotations more

effectively than cascaded approach for RILE score regression task. There are many extensions possible to the current work. First is to handle class imbalance in the dataset with a cost-sensitive objective function. Secondly, CNN gave a comparable performance with Multi-layer Perceptron (NN), which motivates the need to evaluate an end-end sequential architecture. Off-the-shelf embeddings leads to out-of-vocabulary scenarios. It could be beneficial to adapt word-embeddings with manifesto corpus. Finally, background information such as country can be leveraged more effectively.

## References

- Leman Akoglu. 2014. Quantifying political polarity based on bipartite opinion networks. In *AAAI*.
- Waleed Ammar, George Mulcaire, Yulia Tsvetkov, Guillaume Lample, Chris Dyer, and Noah A Smith. 2016. Massively multilingual word embeddings. *arXiv preprint arXiv:1602.01925*.
- Akshat Bakliwal, Jennifer Foster, Jennifer van der Puil, Ron O'Brien, Lamia Tounsi, and Mark Hughes. 2013. Sentiment analysis of political tweets: Towards an accurate classifier. In *ACL*.
- Kenneth Benoit and Thomas Däubler. 2014. Putting text in context: How to estimate better left-right positions by scaling party manifesto data using item response theory. In *Mapping Policy Preferences from Texts Conference*.
- Felix Biessmann. 2016. Automating political bias prediction.
- B. Bruinsma and K. Gemenis. 2017. Validating Word-scores. *ArXiv e-prints*.
- Ian Budge and Michael Laver. 1992. *Party policy and government coalitions*. St. Martin's Press New York.
- Thomas Däubler, Kenneth Benoit, Slava Mikhaylov, and Michael Laver. 2012. Natural sentences as valid units for coded political texts. *British Journal of Political Science*.
- Zhuoye Ding, Qi Zhang, and Xuanjing Huang. 2011. Keyphrase extraction from online news using binary integer programming.
- Goran Glavaš, Federico Nanni, and Simone Paolo Ponzetto. 2017. Cross-lingual classification of topics in political texts. In *ACL WS*.
- Goran Glavas, Federico Nanni, and Simone Paolo Ponzetto. 2017. Unsupervised cross-lingual scaling of political texts. In *EACL*.
- Justin Grimmer and Brandon M Stewart. 2013. Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political analysis*.
- Frederik Hjorth, Robert Klemmensen, Sara Hobolt, Martin Ejnar Hansen, and Peter Kurrild-Klitgaard. 2015. Computers, coders, and voters: Comparing automated methods for estimating party positions. *Research & Politics*.
- Cong Duy Vu Hoang, Gholamreza Haffari, and Trevor Cohn. 2016. Incorporating side information into recurrent neural network language models. In *NAACL-HLT*.
- Mladen Karan, Jan Šnajder, Daniela Širinic, and Goran Glavaš. 2016. Analysis of policy agendas: Lessons learned from automatic topic classification of croatian political texts.
- Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR* abs/1412.6980.
- David Lazer, Alex Sandy Pentland, Lada Adamic, Sinan Aral, Albert Laszlo Barabasi, Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann, et al. 2009. Life in the network: the coming age of computational social science. *Science (New York, NY)* 323(5915):721.
- Will Lowe, Kenneth Benoit, Slava Mikhaylov, and Michael Laver. 2011. Scaling policy preferences from coded political texts. Wiley Online Library, volume 36, pages 123–155.
- Christopher Lucas, Richard A Nielsen, Margaret E Roberts, Brandon M Stewart, Alex Storer, and Dustin Tingley. 2015. Computer-assisted text analysis for comparative politics. *Political Analysis* 23(2):254–277.
- Slava Mikhaylov, Michael Laver, and Kenneth R Benoit. 2012. Coder reliability and misclassification in the human coding of party manifestos. In *Political Analysis*.
- Ramesh Nallapati, James Allan, and Sridhar Mahadevan. 2004. Extraction of key words from news stories.
- Federico Nanni, Căcilia Zirn, Goran Glavaš, Jason Eichorst, and Simone Paolo Ponzetto. 2016. Topfish: topic-based analysis of political position in us electoral campaigns. In *PolText*.
- Ludovic Rheault. 2016. Expressions of anxiety in political texts. In *NLP+ CSS 2016*, page 92.
- Jonathan B Slapin and Sven-Oliver Proksch. 2008. A scaling model for estimating time-series party positions from texts. Wiley Online Library, volume 52, pages 705–722.



| CMP Coding Scheme   |  |
|---|--|
| <ul style="list-style-type: none"> <li>Domain 1: External Relations               <ul style="list-style-type: none"> <li>101 Foreign Special Relationships: Positive</li> <li>102 Foreign Special Relationships: Negative</li> <li>103 Anti-Imperialism</li> <li>104 Military: Positive</li> <li>105 Military: Negative</li> <li>106 Peace</li> <li>107 Internationalism: Positive</li> <li>108 European Community/Union: Positive</li> <li>109 Internationalism: Negative</li> <li>110 European Community/Union: Negative</li> </ul> </li> <li>Domain 2: Freedom and Democracy               <ul style="list-style-type: none"> <li>201 Freedom and Human Rights</li> <li>202 Democracy</li> <li>203 Constitutionalism: Positive</li> <li>204 Constitutionalism: Negative</li> </ul> </li> <li>Domain 3: Political System               <ul style="list-style-type: none"> <li>301 Decentralisation</li> <li>302 Centralisation</li> <li>303 Governmental and Administrative Efficiency</li> <li>304 Political Corruption</li> <li>305 Political Authority</li> </ul> </li> <li>Domain 4: Economy               <ul style="list-style-type: none"> <li>401 Free Market Economy</li> <li>402 Incentives: Positive</li> <li>403 Market Regulation</li> <li>404 Economic Planning</li> <li>405 Corporatism/Mixed Economy</li> <li>406 Protectionism: Positive</li> <li>407 Protectionism: Negative</li> <li>408 Economic Goals</li> <li>409 Keynesian Demand Management</li> <li>410 Economic Growth: Positive</li> </ul> </li> </ul> | <ul style="list-style-type: none"> <li>411 Technology and Infrastructure: Positive</li> <li>412 Controlled Economy</li> <li>413 Nationalisation</li> <li>414 Economic Orthodoxy</li> <li>415 Marxist Analysis</li> <li>416 Anti-Growth Economy: Positive</li> </ul> <ul style="list-style-type: none"> <li>Domain 5: Welfare and Quality of Life               <ul style="list-style-type: none"> <li>501 Environmental Protection</li> <li>502 Culture: Positive</li> <li>503 Equality: Positive</li> <li>504 Welfare State Expansion</li> <li>505 Welfare State Limitation</li> <li>506 Education Expansion</li> <li>507 Education Limitation</li> </ul> </li> <li>Domain 6: Fabric of Society               <ul style="list-style-type: none"> <li>601 National Way of Life: Positive</li> <li>602 National Way of Life: Negative</li> <li>603 Traditional Morality: Positive</li> <li>604 Traditional Morality: Negative</li> <li>605 Law and Order: Positive</li> <li>606 Civic Mindedness: Positive</li> <li>607 Multiculturalism: Positive</li> <li>608 Multiculturalism: Negative</li> </ul> </li> <li>Domain 7: Social Groups               <ul style="list-style-type: none"> <li>701 Labour Groups: Positive</li> <li>702 Labour Groups: Negative</li> <li>703 Agriculture and Farmers: Positive</li> <li>704 Middle Class and Professional Groups</li> <li>705 Underprivileged Minority Groups</li> <li>706 Non-economic Demographic Groups</li> </ul> </li> </ul> <p>000 No meaningful category applies</p> |

Table 5: *Left* topics are given in *red* and *right* topics are given in *blue*

During our nation's darkest hours, Americans have strived mightily and succeeded in meeting the challenges of their times. The question before us is whether we will do the same during this bright moment; whether we will seize this moment to bring more prosperity and progress to more Americans than ever before; whether, having finally conquered our financial deficits, we will have the courage to conquer the other deficits – in health care, in education, in the environment – that challenge us today.

601  
606  
410  
305

Figure 4: Democratic Party of USA, 2000 —  $\int$  denotes sentence segment. 601 - National Way of Life: Positive, 606 - Civic Mindedness: Positive, 410 - Economic Growth: Positive, 305: Political Authority

Dhanya Sridhar, Lise Getoor, and Marilyn Walker.  
2014. Collective stance classification of posts in on-line debate forums .

Suzan Verberne, Eva D'hondt, Antal van den Bosch, and Maarten Marx. 2014. Automatic thematic classification of election manifestos. In *Information Processing & Management*.

Andrea Volkens, Onawa Lacewell, and Sven Regel Henrike Schultze and Annika Werner. Pola Lehmann. 2011. The manifesto data collection. manifesto project (mrg/cmp/marpor). In *Wissenschaftszentrum Berlin für Sozialforschung (WZB)*.

Daniel Xiaodan Zhou, Paul Resnick, and Qiaozhu Mei. 2011. Classifying the political leaning of news articles and users from user votes.

Ćaćilia Zirn, Goran Glavaš, Federico Nanni, Jason Eichorts, and Heiner Stuckenschmidt. 2016. Classifying topics and detecting topic shifts in political manifestos. University of Zagreb.