

NLP-CIC @ PRELEARN: Mastering Prerequisites Relations, from Handcrafted Features to Embeddings*

Jason Angel
Instituto Politécnico Nacional
Mexico City, Mexico
ajason08@gmail.com

Segun Taofeek Aroyehun
Instituto Politécnico Nacional
Mexico City, Mexico
aroyehun.segun@gmail.com

Alexander Gelbukh
Instituto Politécnico Nacional
Mexico City, Mexico
www.gelbukh.com

Abstract

We present our systems and findings for the prerequisite relation learning task (PRELEARN) at EVALITA 2020. The task aims to classify whether a pair of concepts hold a prerequisite relation or not. We model the problem using handcrafted features and embedding representations for in-domain and cross-domain scenarios. Our submissions ranked first place in both scenarios with average F1 score of 0.887 and 0.690 respectively across domains on the test sets. We made our code freely available¹.

1 Introduction

A prerequisite relation is a pedagogical relation that indicates the order in which concepts can be presented to learners. The relation can be used to guide the presentation sequence of topics and subjects during the design of academic programs, lectures, and curricula or instructional materials.

In this work, we present our systems to automatically detect prerequisite relations for Italian language in the context of the PRELEARN shared task (Alzetta et al., 2020) at EVALITA 2020 (Basile et al., 2020). The evaluation of submissions considers: (1) in-domain and cross-domain scenarios defined by either the inclusion (in-domain) or exclusion (cross-domain) of the target domain in the training set. The four domains are 'data mining' (DM), 'geometry' (Geo), 'pre-calculus' (Prec), and 'physics' (Phy). (2) the type

of resources (features) used to train the model – raw text VS. structured information.

The combination of these settings defined the four PRELEARN subtasks. Formally, a prerequisite relation exists between two concepts if one has to be known beforehand in order to understand the other. For the PRELEARN task, given a pair of concepts, the relation exists only if the latter concept is a prerequisite for the former. Therefore, the task is a binary classification task.

We approach the problem from two perspectives: handcrafted features based on lexical complexity and pre-trained embeddings. We employed static embeddings from Wikipedia and Wikidata, and contextual embeddings from Italian-BERT model.

2 Related works

Prerequisite relation learning has been mostly studied for the English language (Liang et al., 2018; Talukdar and Cohen, 2012). Adorni et al. (2019) performed unsupervised prerequisite relations extraction from textbooks using word co-occurrence and order of words appearance in the text. In the case of Italian language there is *ITA-PREREQ* (Miaschi et al., 2019), the first dataset for prerequisite learning, and actually the one used for the present work. It was automatically built as a projection of *AL-CPL* (Liang et al., 2018) from the English Wikipedia to the Italian Wikipedia. In addition, Miaschi et al. (2019) examines the utility of lexical features for individual concepts and features derived from the concept pairs.

3 Methodology

This section describes the data analysis, the features we used to model the task, and the system we finally submitted to the PRELEARN competition.

*“Copyright © 2020 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).”

¹https://github.com/ajason08/EVALITA2020_PRELEARN

3.1 Dataset

The dataset provided by the organizers includes the concept pairs splitted into the following domains: 'data mining', 'geometry', 'precalculus' and 'physics'. The dataset contains the list of concepts with a link to the corresponding Wikipedia article. The first paragraph of such article is named the concept description. All concept descriptions are cleaned in order to facilitate the extraction of information from the text, e.g. the mathematical expressions are already tagged using this pattern `formula_<number>`.

Table 1 displays the number of samples and the distribution over the prerequisite relations (positive samples) across domains for the training set. The test sets in turn exhibits a 50-50 distribution over positive and negative samples.

The only preprocessing we did was lowercase the concept description and remove line-breaks.

Domain	Samples	Prerequisites rel.
Data mining	424	0.257
Geometry	1548	0.214
Precalculus	2220	0.142
Physics	1716	0.238

Table 1: Training set number of samples and distribution of prerequisite relations (positive samples) across domain

3.2 Features

The following are the set of features we experiment with:

Complexity-based: a set of handcrafted features intended to measure how complex a concept is. The rationale is that less complex concepts are prerequisites for the more complex ones. We used some features that have been found effective for the task of complex word identification (Aroyehun et al., 2018), specifically they are:

- Age of acquisition of concept: we use *ItAoA* (Montefinese et al., 2019), a dataset of age of acquisition norms (we average the values for the different entries per word), to derive the age of acquisition for each concept we compute the geometric mean of values from *ItAoA* for words which occur in the concept description after replacing outliers (by the closest permitted value). In addition, we use the number of matches as a feature.

- Age of acquisition of related concepts: We derived a list of concepts related to each concept by matching which of them appears in the concept description. Then, we average the age of acquisition of those concepts. We also took the count of the related concepts.
- Description length: we count the number of words in the concept description.
- Number of mathematical expressions: we count the occurrence of mathematical expressions. We assume that more complex concepts will have a higher occurrence of mathematical expressions in their descriptions.
- Concept view frequency: the average of the daily unique visits by Wikipedia users (including editors, anonymous editors, and readers) over the last year. We think that the number of visitors will be correlated with the degree of complexity of a concept. To gather this information we used the Pageviews Analysis of Wikipedia ².

Concept-to-Concept features: they aim to model the relation between the concept pairs, specifically we evaluate whether a concept appears as a sub-string in the title or description of the other concept. We did this in both directions resulting in two features. We also represent the domain they belong to as a one-hot vector.

Wiki-embeddings: We map each concept identifier to their corresponding Wikipedia title and Wikidata identifier using the Wikidata Query Service³. Then, we obtain the 100 dimensional vector for each Wikipedia title from a pre-trained Wikipedia embedding⁴ (Yamada et al., 2020). Similarly, we use the Wikidata embedding⁵ (Lerer et al., 2019) to represent the Wikidata identifiers as 200 dimensional vectors.

Italian-BERT features: We used a pre-trained uncased version of Italian BERT (base model)⁶ provided by the MDZ Digital Library team (dbmdz) trained on 13GB of text mainly from

²<https://pageviews.toolforge.org>

³query.wikidata.org

⁴http://wikipedia2vec.s3.amazonaws.com/models/it/2018-04-20/itwiki_20180420_100d.pkl.bz2

⁵https://dl.fbaipublicfiles.com/torchbiggraph/wikidata_translation_v1.tsv.gz

⁶<https://huggingface.co/dbmdz/bert-base-italian-uncased>

Scenario	Resources	System	DM	Geo	Phy	Prec	AVG
in-domain	raw-text	Italian-BERT	0.838	0.925	0.855	0.930	0.887
in-domain	structured	Complex+wd	0.808	0.905	0.795	0.915	0.856
in-domain	structured	Complex	0.828	0.895	0.785	0.885	0.848
cross-domain	raw-text	Italian-BERT	0.565	0.785	0.635	0.775	0.690
cross-domain	structured	Complex+wd	0.535	0.775	0.600	0.760	0.668
cross-domain	structured	Complex	0.494	0.735	0.595	0.730	0.639

Table 2: Test set results for the four PRELEARN subtasks using F1-score

Settings	In-domain	Cross-domain
raw-text	+2.1%	+4.2%
structured	+15.6%	+4.8%

Table 3: Performance advantage over the next best participant on average across domains

Wikipedia and other text sources. With this model, we get the 768 dimensional vector representation for a sequence corresponding to the [CLS] token as in the original implementation of BERT (Devlin et al., 2019). The sequence consists of the combination of the concept and its Wikipedia description.

3.3 Systems

Considering the proposed features and our experimental results at Section 5, we proposed the following three systems to address both, in-domain and cross-domain scenarios. For the in-domain scenario we trained with a combination of all the training samples per domain. In the same way, we combined the remaining three domains for each cross-domain experiment (i.e. excluding samples from the target domain).

Complex: a completely handcrafted machine learning system, it uses all the complexity-based and Concept-to-Concept features (except the domain vector for cross-domain scenario), and we normalize the features using Z-score normalization. This system uses a tree-ensemble learner as classifier⁷ with the default parameters provided by Breiman (2001)⁸. This system participated under the structured resource setting because the “concept view frequency” feature is structured information.

Complex+wd: an improved version of the *Complex* system by only concatenating the Wiki-

data embedding of each concept in the concept pair to the feature set. This system participated under the structured resource setting as well. We decided to not include the Wikipedia embeddings considering the ablation analysis we present in Table 4.

Italian-BERT: a single layer neural network mapping the 768 features from the [CLS] to the output space of dimension 2 as a sequence pair classification task. In addition, the pre-trained weights of the base model are fine-tuned on the training dataset. We fine-tune the base model using the huggingface transformers library (version 3.1) for Pytorch (Wolf et al., 2019). In the in-domain scenario, we use the following training parameters: the number epochs is 10, learning rate is $5e-5$, weight decay is 0.01, batch size is 32, warm up steps is 100, optimizer is AdamW with a linear schedule after a period of warm up steps. We find that the model exhibits high variance across runs in our cross-domain experiments. Hence, in addition to the parameter settings for the in-domain experiments, we choose the number of training steps using a validation set for the unseen target domain. Accordingly, we set the maximum training step to 400 and the warm up steps to 100, 200, 150, and 200 for data mining, geometry, physics, and pre-calculus cross-domain scenarios respectively.

4 Results

Table 2 shows our per-domain results for our systems indicating the kind of scenario and resources they used. We observe the clear superiority of Italian-BERT which only relies on raw-text resources. This suggest that just fine-tuning BERT

⁷Other classifiers were tested and obtained lower performance

⁸<https://cran.r-project.org/web/packages/randomForest/index.html>

Scenario	Resources	Feature set	DM	Geo	Phy	Prec	AVG
in-domain	raw	complexity	0.646	0.817	0.622	0.792	0.720
in-domain	raw	wp_embedding	0.705	0.818	0.670	0.827	0.755
in-domain	raw	Italian-BERT	0.947	0.746	0.829	0.842	0.841
in-domain	structured	complexity +page_view	0.648	0.805	0.629	0.804	0.721
in-domain	structured	wd_embedding	0.660	0.814	0.674	0.838	0.746
in-domain	structured	wd+wp_embedding	0.694	0.824	0.672	0.831	0.755
in-domain	structured	complexity +page_view +wd_embedding	0.697	0.823	0.686	0.845	0.763
cross-domain	raw	complexity	0.072	0.592	0.258	0.586	0.377
cross-domain	raw	wp_embedding	0.000	0.622	0.079	0.344	0.261
cross-domain	raw	Italian-BERT	0.145	0.646	0.460	0.570	0.455
cross-domain	structured	complexity +page_view	0.107	0.588	0.297	0.577	0.392
cross-domain	structured	wd_embedding	0.000	0.661	0.355	0.608	0.406
cross-domain	structured	wd+wp_embedding	0.000	0.660	0.332	0.605	0.399
cross-domain	structured	complexity +page_view +wd_embedding	0.064	0.645	0.366	0.630	0.426

Table 4: Ablation analysis results using F1-score (validation set for Italian-BERT and 10-fold for the others)

is enough for gaining a notion of prerequisite relations on concepts. Still, the systems based on handcrafted features and non-contextual embedding exhibit competitive results, with a good enough performance to rank first in the structured resource setting, while being faster, more interpretable and simpler than the Italian-BERT counterpart.

The results showed that there is a huge performance reduction for the cross-domain scenario. The largest performance drop is on the “data mining” domain. Given that we train our models on the combination of examples from all other domains, it is likely that the probable cause is the domain mismatch. Yet, the reduction on the test sets are smaller than what we observe in our K-fold experiments and validation sets.

In addition, we show in Table 3 the performance advantage we obtained over the next best participant based on the ranking released by the organizers.

One can see that the greater performance advantage is from the structured resource setting. This suggests that the “Concept view frequency” and the Wikidata embedding features are effective.

5 Discussion: ablation analysis

During the creation of our systems we performed several experiments over the possible features to use. We did 10-fold cross validation for the in-domain experiments except with the Italian-BERT⁹, for which we used a stratified split of 30% for validation set. Table 4 shows the experimental results over the training (validation) set for both, in-domain and cross-domain scenarios. The “Resources” column serves to identify the type of resources used for the current feature.

We observe that the “data mining” domain appears to be difficult in the cross-domain scenario, models based on the non-contextual embedding features obtain results of zero. We suspect that this difficulty is due to the domain mismatch.

Based on these results, we select the Italian-BERT for the raw-text setting, and the “complexity +page_view” and the addition of Wikidata embeddings (“wd_embedding”) for the structured resource setting for our submissions.

⁹Due to its high computational requirements

6 Conclusion

We tackle the task of prerequisite relation learning using a variety of systems that explore three set of features: handcrafted features based on complexity intuitions, embedding models from Wikipedia and Wikidata, and contextual embedding from Italian-BERT model. We examine the capabilities of our models in in-domain and cross-domain scenarios. Our models ranked first in all the sub-task of the PRELEARN competition at EVALITA 2020. We found that although our Italian-BERT model outperformed the others, the simpler models show competitive results.

We plan to further examine the impact of using a combination of all possible domains as training set on the performance of our models.

Acknowledgments

The authors thank CONACYT for the computer resources provided through the INAOE Supercomputing Laboratory’s Deep Learning Platform for Language Technologies.

References

- Giovanni Adorni, Chiara Alzetta, Frosina Koceva, Samuele Passalacqua, and Iliaria Torre. 2019. Towards the identification of propaedeutic relations in textbooks. In *International Conference on Artificial Intelligence in Education*, pages 1–13. Springer.
- Chiara Alzetta, Alessio Miaschi, Felice Dell’Orletta, Frosina Koceva, and Iliaria Torre. 2020. Prelearn @ evalita 2020: Overview of the prerequisite relation learning task for italian. In Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro, editors, *Proceedings of Seventh Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2020)*, Online. CEUR.org.
- Segun Taofeek Aroyehun, Jason Angel, Daniel Alejandro Pérez Alvarez, and Alexander Gelbukh. 2018. Complex word identification: Convolutional neural network vs. feature engineering. In *Proceedings of the thirteenth workshop on innovative use of NLP for building educational applications*, pages 322–327.
- Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro. 2020. Evalita 2020: Overview of the 7th evaluation campaign of natural language processing and speech tools for italian. In Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro, editors, *Proceedings of Seventh Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2020)*, Online. CEUR.org.
- Leo Breiman. 2001. Random forests. *Machine learning*, 45(1):5–32.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and Alex Peysakhovich. 2019. PyTorch-BigGraph: A Large-scale Graph Embedding System. In *Proceedings of the 2nd SysML Conference*, Palo Alto, CA, USA.
- Chen Liang, Jianbo Ye, Han Zhao, Bart Pursel, and C Lee Giles. 2018. Active learning of strict partial orders: A case study on concept prerequisite relations. *arXiv preprint arXiv:1801.06481*.
- Alessio Miaschi, Chiara Alzetta, Franco Alberto Cardillo, and Felice Dell’Orletta. 2019. Linguistically-driven strategy for concept prerequisites learning on italian. In *Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 285–295.
- Maria Montefinese, David Vinson, Gabriella Vigliocco, and Ettore Ambrosini. 2019. Italian age of acquisition norms for a large set of words (itaoa). *Frontiers in psychology*, 10:278.
- Partha Talukdar and William Cohen. 2012. Crowdsourced comprehension: predicting prerequisite structure in wikipedia. In *Proceedings of the Seventh Workshop on Building Educational Applications Using NLP*, pages 307–315.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Ikuya Yamada, Akari Asai, Jin Sakuma, Hiroyuki Shindo, Hideaki Takeda, Yoshiyasu Takefuji, and Yuji Matsumoto. 2020. Wikipedia2Vec: an efficient toolkit for learning and visualizing the embeddings of words and entities from wikipedia. *arXiv preprint 1812.06280v3*.