

Grammatical Accents: Using Machine Learning to Quantify Language Transfer

Motivation Goal: To use machine learning to establish a broad-based method to empirically study the effects of first language syntax on second language (L1->L2 transfer). **Q_:** Does NLI work in languages other than english [cf. 1]? **Q**₂: What grammatical features can we train on successfully? Which are the most informative? Q₂: Which are the most accurate* classifiers Q₂: Can machine learning algorithms learn L1->L2 patterns that generalize across L2s? **Q**: Are only certain parts of input (i.e. language) informative? Which ones? [7] Native-Language Identification (NLI): The process of determining an author's native language (L1) based only on their writings in a second language (L2)



Method:

Compare the results of a variety of state-of-the-art machine learning techniques on NLI in two languages: English and Spanish.

Background

Native-language identification has been proven possible when a wide set of features is applied to the task [1]. Further more, languages besides english have been widely ignored (Q_1) . As a first step, we broaden our language set to include Spanish while simultaneously restricting our feature set to exclusively syntactic features as inspired by [4].

- [3] POS n-grams <= 4-grams, dependency labels.
- [4] POS n-grams <= 4-grams. Used SVMs and shallow neural networks, achieving accuracy > 50%.
- [5] POS n-grams <= tri-grams. Used SVMs to achieve accuracy > 50%

Raw Corpora



Feature Representations



Machine Learning



Classification



Tiwalayo Eisape¹, William Merrill², Sven Dietz¹, Joshua K. Hartshorne¹ ¹Department of Psychology Boston College, ²Department of Linguistics Yale University



REFERENCES [1] Tetreault, Joel, Daniel Blanchard, and Aoife Cahill. "A report on the first native language identification shared task." Proceedings of the eighth workshop on innovative use of NLP for building educational applications. 2013.

[2] Joshi, A. K. and Srinivas, B. Disambiguation of Super Parts of Speech (or Supertags): Almost Parsing. Proceedings of the 17th International Conference on Computational Linguistics Kyoto, Japan. 1994. [3] Berzak, Yevgeni, Roi Reichart, and Boris Katz. "Reconstructing native language typology from foreign language usage." arXiv preprint arXiv:1404.6312 (2014). [4] Gebre, Binyam Gebrekidan, et al. "Improving native language identification with tf-idf weighting." the 8th NAACL Workshop on Innovative Use of NLP for Building Educational Applications (BEA8). 2013. [5] Jarvis, Scott, Yves Bestgen, and Steve Pepper. "Maximizing classification accuracy in native language identification." Proceedings of the eighth workshop on innovative use of NLP for building educational

applications. 2013. [6] Kruengkrai, Canasai, et al. "Language identification based on string kernels." Communications and Information Technology, 2005. ISCIT 2005. IEEE International Symposium on. Vol. 2. IEEE, 2005 [7] Johnson, Jacqueline S., and Elissa L. Newport. "Critical period effects on universal properties of language: The status of subjacency in the acquisition of a second language." Cognition 39.3 (1991): 215-258. [8] Petrov, Slav. "Announcing syntaxnet: The world's most accurate parser goes open source." Google Research Blog 12 (2016).

Arabic Chinese French German Japanese

Spanish

Telugu

Special thanks to William Merrill, Clinton Tak, and the rest of the Language Learning lab. TE is supported by the Ronald E. McNair Scholarship (TRIO) and JKH is supported by the Academic Technology Innovation Grant (Boston College)



Insights:

1. Spanish and Italian - same language family: Italo-Western Romance

2. Hindi and Telugu - high proximity and language sharing

Conclusions / Future Directions

By achieving state of the art accuracy, using strictly syntactic features, we show machine learning can pick up on generalizable, grammatical idiosyncrasies associated with (L1 ->L2) language transfer.

Next Steps:

. Expand features to further encapsulate syntax "Super Tagging" [2]

2. Open up the black box.

Reverse engineer our learning algorithms for interpretation

Acknowledgements

