

# SEQUENCE MIXUP FOR ZERO-SHOT CROSS-LINGUAL PART-OF-SPEECH TAGGING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

There have been efforts in cross-lingual transfer learning for various tasks. We present an approach utilizing an interpolative data augmentation method, Mixup, to improve the generalizability of models for part-of-speech tagging trained on a source language, improving its performance on unseen target languages. Through experiments on ten languages with diverse structures and language roots, we put forward its applicability for downstream zero-shot cross-lingual tasks.

## 1 INTRODUCTION

Recently, neural network models have obtained state-of-the-art results in part-of-speech (POS) tagging tasks across multiple languages. Since numerous languages lack suitable corpora annotated with POS labels, there have been efforts to design models for cross-lingual transfer learning. Cross-lingual learning enables us to utilize the annotated corpora of a source language to train models that are effective over a different target language. Interpolative data augmentation methods have been proposed to mitigate overfitting in models in the absence of enough training data. Sequence-based mixup Chen et al. (2020) is an interpolative data augmentation method for named entity recognition. However, these methods have not been explored for cross-lingual transferability.

Interpolative data augmentation methods are aimed at increasing the diversity of the training distribution and, as a result, improving the generalizability of underlying models. We leverage this capability of sequence Mixup (Seq. Mixup) to capture rich linguistic information for cross-lingual transferability of POS tagging tasks for ten languages with different structures and language roots. To this end, we first measure the dataset level cosine similarity across languages, defined as the average of sentence-level embedding of dataset samples. This gives an overview of the syntactical and semantical relationship among the different languages. We then finetune multilingual models over a source language using sequence Mixup and evaluate it across a target language for varying language similarities, probing sequence-based interpolative data augmentation. We also evaluate sequence Mixup on combinations of similar and dissimilar languages and verify its transferability on various target languages.

## 2 METHODOLOGY AND SETUP

Mixup Zhang et al. (2018) is a data augmentation technique that generates virtual training samples from convex combinations of individual inputs and labels. For a pair of data points  $(x, y)$  and  $(x', y')$ , Mixup creates a new sample  $(\tilde{x}, \tilde{y})$  by interpolating the data points using a ratio  $\lambda$ , sampled from a Beta distribution, where  $\tilde{x} = \lambda \cdot x + (1 - \lambda) \cdot x'$  and corresponding mixed label  $\tilde{y} = \lambda \cdot y + (1 - \lambda) \cdot y'$ .

We perform Mixup over the latent space representations for interpolating sequences. Pair of sentences  $(x, y)$  and  $(x', y')$  are randomly sampled and interpolated in the hidden space using a  $L$ -layer encoder  $f(\cdot, \theta)$ . The hidden layer representations for  $x$  and  $x'$  upto the  $k^{th}$  layer are given as,

$$\begin{aligned} h_l &= f_l(h_{l-1}; \theta), \quad l \in [1, k] \\ h'_l &= f_l(h'_{l-1}; \theta), \quad l \in [1, k] \end{aligned} \tag{1}$$

At the  $k^{th}$  layer, the hidden representations of each token in  $x$  are linearly interpolated with each token in  $x'$ . After this,  $\tilde{h}_k$  is fed to the upper layers,

$$\begin{aligned} \tilde{h}_k &= \lambda h_k + (1 - \lambda) h'_k \\ \tilde{h}_l &= f_l(\tilde{h}_{l-1}; \theta), \quad l \in [k + 1, L] \end{aligned} \tag{2}$$

We evaluate the performance of zero-shot learning on sequence Mixup, where the model is trained on one source language or a set of source languages and evaluated on a target language. To choose the source and target, for each language we average sentence level embeddings of dataset instances and use the average embeddings as language representation. Using these representations, we find the cosine similarity among the languages as shown in Figure 1. This approach can be extended for tasks on under-resourced languages, where models can be trained on a similar high-resourced language or a set of languages.

**Setup** We evaluate our approach on POS tagging with datasets from the Universal Dependencies (UD) dataset <sup>1</sup> for ten different languages - Arabic (ar), Dutch (nl), French (fr), German (de), Hindi (hi), Indonesian (id), Italian (it), Marathi (mr), Vietnamese (vi) and Urdu (ur). For each experiment, we use 800 sentences for each source language for training, 100 sentences for validation and test each from the target language for evaluation.

### 2.1 SINGLE LANGUAGE TRANSFER

We train the model on a single source language and a different target language dataset to evaluate its performance, as shown in Table 1. As sequence Mixup trains on interpolated sequences, it regularizes the model and prevents overfitting, outperforming mBERT. For the target language Italian, we observe higher F1 when the source language is French and lower scores for source languages Hindi and Arabic. This is in line with the trend observed in Fig 1, where the cosine similarity for the language pair (French, Italian) is highest, lower for (Hindi, Italian) and lowest for (Arabic, Italian). We observe large improvements when sequence Mixup is applied over dissimilar languages, validating that Mixup is able to generate more diverse input samples which intersect with the target language structure and semantics.

Table 1: F1-scores for POS tagging on Seq. Mixup and mBERT (mean of 10 runs). Improvements are shown with **blue** (↑) over mBERT.

Source	Target	mBERT	Seq. Mixup
<b>High Similarity</b>			
French	Italian	94.52	94.75
Indonesian	Vietnamese	56.08	56.34
German	Dutch	85.32	85.48
Hindi	Marathi	64.41	64.97
Urdu	Arabic	44.61	47.38
<b>Low Similarity</b>			
Hindi	Italian	58.63	63.71
French	Arabic	39.63	40.55
Arabic	Italian	25.41	28.42

### 2.2 MULTI LANGUAGE TRANSFER

To extend our experiments, we choose a pair of languages on which the model is trained and present the results in Table 2. This helps to infer in what manner additional language data impacts the performance. Languages Dutch, German and French have high cosine similarity, leading to larger improvement for the Dutch language compared to single language transfer. For target language Italian, F1-score decreases when trained on both Arabic and French data; this can be reasoned by the low cosine similarity of Arabic and Italian language.

Table 2: F1-scores for POS tagging on Seq. Mixup and mBERT (mean of 10 runs) when trained on two source languages (New+Original). Improvements are shown with **blue** (↑) and poorer performance with **red** (↓).

Source	Target	Single Source	Dual Source
de+fr	it	94.75	94.85
fr+de	nl	85.48	86.11
ar+fr	it	94.75	28.00
ar+hi	mr	64.97	30.07

## 3 CONCLUSION

We analyze interpolative regularization-based data augmentation over tokens for zero-shot cross-lingual transfer of part-of-speech tagging across ten languages. Through extensive experiments over languages with varying syntactic and semantic structures on single and pair of languages, we pave the way for using interpolative data augmentation to improve the generalizability of neural networks for zero-shot transfer learning on downstream tasks.

<sup>1</sup><https://universaldependencies.org/>

## URM STATEMENT

The authors acknowledge that at least one key author of this work meets certain URM criteria of ICLR 2024 Tiny Papers Track.

## REFERENCES

Jiaao Chen, Zhenghui Wang, Ran Tian, Zichao Yang, and Diyi Yang. Local additivity based data augmentation for semi-supervised NER. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1241–1251, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.95. URL <https://aclanthology.org/2020.emnlp-main.95>.

Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=r1Ddp1-Rb>.

## A APPENDIX

### A.1 COSINE SIMILARITY OF LANGUAGES

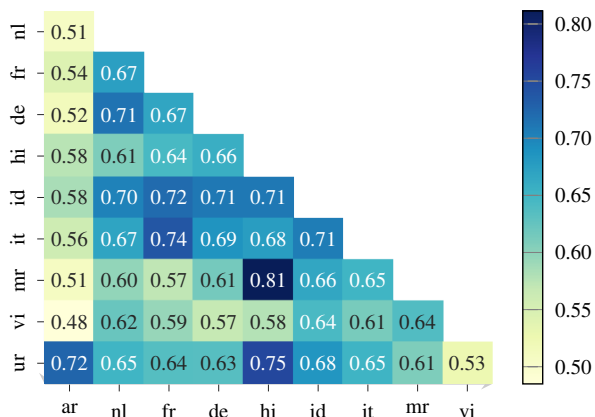


Figure 1: Cosine Similarity of the ten languages.

BERT-base-multilingual-cased (mBERT) has been used as an encoder in sequence Mixup and for obtaining the embeddings to evaluate cosine similarity.

### A.2 TRAINING SETUP

The learning rate is  $5e-5$  with Adam optimizer and batch size 16. All hyperparameters are selected based on validation F1-score.