

Studying Language Properties and Divergences Via Automatically Classified Errors

[Introduction:](#)

[Background:](#)

[Dataset statistics:](#)

[Main:](#)

[Methodology:](#)

[Steps:](#)

[Decisions:](#)

[error profile](#)

[Create an error-profile for every language:](#)

[lang to vec](#)

[Experiments and results:](#)

[Questions](#)

[Graphs](#)

[Classified Syntactic Errors for the Whole DS:](#)

[Classified Syntactic Errors by level, minus the average across levels:](#)

[The error profile for each nationality:](#)

[Prominent Words](#)

[Classified Syntactic Errors by nationality, minus the average across nationalities:](#)

[Classified Syntactic Errors by level:](#)

[All Levels together:](#)

[Distances](#)

[Comparing to other work:](#)

[Results](#)

[Future work:](#)

[Link to github](#)

[References:](#)

Introduction:

Previous studies have examined the usefulness of a generalized unified framework analyzing cross-lingual syntactic divergences. However, there are few empirical studies that have examined the syntactic divergences across language pairs quantitatively.

In a recent study, Nikolaev and colleagues, proposed a framework for extracting divergence patterns across language pairs from a parallel corpus, building on Universal Dependencies. They defined Universal Dependencies (UD) as '*a framework for treebank annotation, whose objectives include satisfactory analyses of individual languages, providing a suitable basis for bringing out cross linguistic parallelism, suitability for rapid consistent annotation and accurate automatic parsing, ease of comprehension by non-linguists, and effective support for downstream tasks*' (Nikolaev et al., 2020).

They used '*The Parallel Universal Dependencies (PUD) corpus consists of 1000 sentences translated into various languages by professional translators*'

In this work we took a different approach, using a parallel corpus, consisting of sentences written in English by kids of various nationalities, and the corrected sentences.

Also worthwhile to mention "Classifying Syntactic Errors in Learner Language" (Choshen et al. 2020.) This work is mentioned in the footnotes as well, when we used their github repository. The abstract:

'We present a method for classifying syntactic errors in learner language, namely errors whose correction alters the morphosyntactic structure of a sentence. The methodology builds on the established Universal Dependencies syntactic representation scheme, and provides complementary information to other error-classification systems. Unlike existing error classification methods, our method is applicable across languages, which we showcase by producing a detailed picture of syntactic errors in learner English and learner Russian. We further demonstrate the utility of the methodology for analyzing the outputs of leading Grammatical Error Correction (GEC) systems.'

The Current Study

The current study examines error patterns made in English by individuals of various nationalities. We used the EFCAMDAT dataset which specifically records nationalities and not languages. *Nationality is used as a proxy for language when comparing to other work on languages.* In this work we took a different approach to language comparisons, using a parallel corpus, consisting of sentences written in English by kids of various nationalities, and the corrected sentence, classifying the errors and analysing the error POS, both in the original and corrected sentence, and comparing to previous work on these Cross-Linguistic Syntactic Divergences.

This work is exploratory in nature.

We classified error-types in EFCAMDAT, and show visualisations of the distribution according to POS. We sliced the data according to nationality and proficiency level as well. We measured distance between languages, computed from their error-profile matrix, and compared that to distances computed from their lan2vec vector. Measuring distances between languages with our error-profile compared to lang2vec, showed high correlation for most languages, with the highest being Spanish (mexico).

Jp (Japanese) and to a lesser extent de showed the least total similarity to other languages in both distance measurements.

We showed some prominent missed POS, ('the', 'a', ',', '.') as well prominent missed POS in the DET POS, with 'the' and 'a' dominating this category.

For brevity sake we used the following names:

POS - part of speech.

DS - Dataset.

DP - data points (usually referring to a single error, and its metadata.)

DF - Dataframe. A python object that holds all the relevant processed data, in a convenient format for analysis.

Background:

Datasets:

The main Dataset we are working with is the EFCAMDAT dataset.
(<https://philarion.mml.cam.ac.uk/>)

From the site:

The EF-Cambridge Open Language Database (EFCAMDAT) is a publicly available resource to facilitate second language research and teaching.

It contains written samples from thousands of adult learners of English as a second language, world wide.

EFCAMDAT currently contains over 83 million words from 1 million assignments written by 174,000 learners, across a wide range of levels (CEFR stages A1-C2). This text corpus includes information on learner errors, part of speech, and grammatical relationships.

Researchers can search for language patterns using a range of criteria, including learner nationality and level.

The resource is actively developed by the Department of Theoretical and Applied Linguistics at the University of Cambridge in partnership with Education First.

Dataset statistics:

Total nationalities: 186

Total proficiency levels: 15. (1-16, no level 5)

Total students: 114,554

(obtained from step 2 described in STEPS):

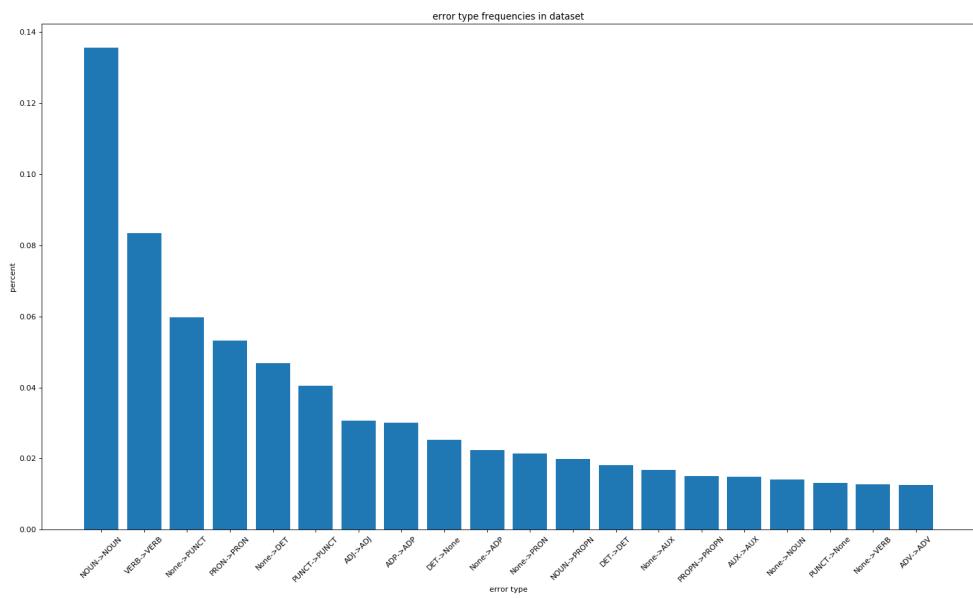
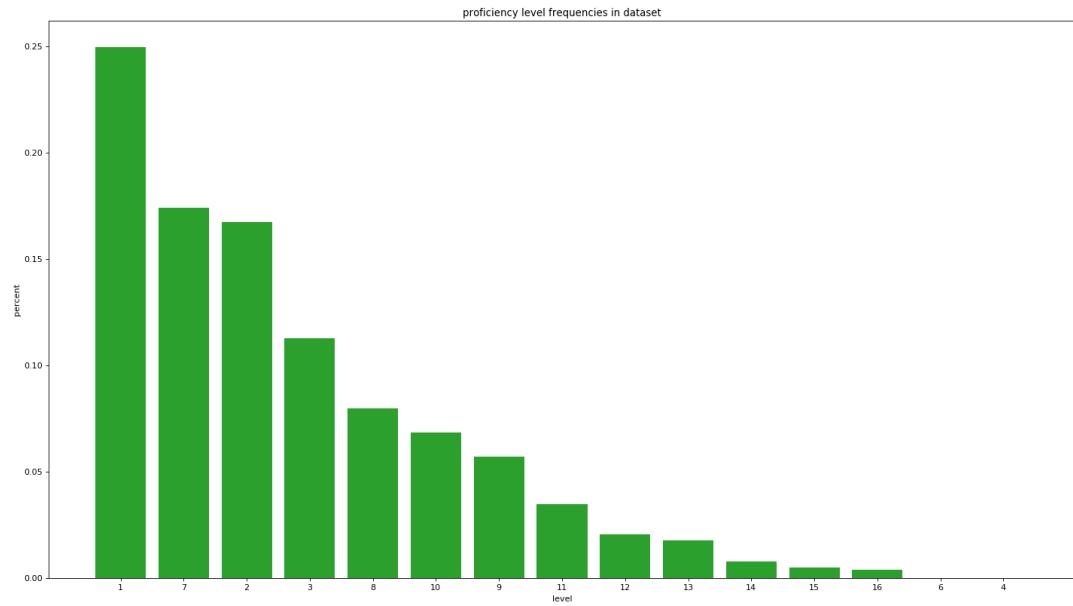
Total amount of classified errors: 4,564,523

The EFCAMDAT dataset provided us with the original and corrected sentences.

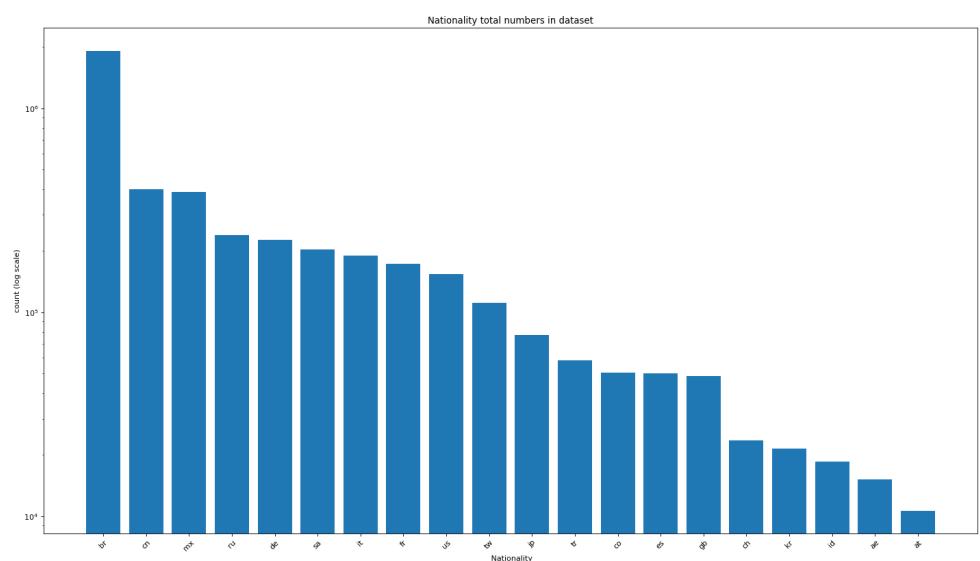
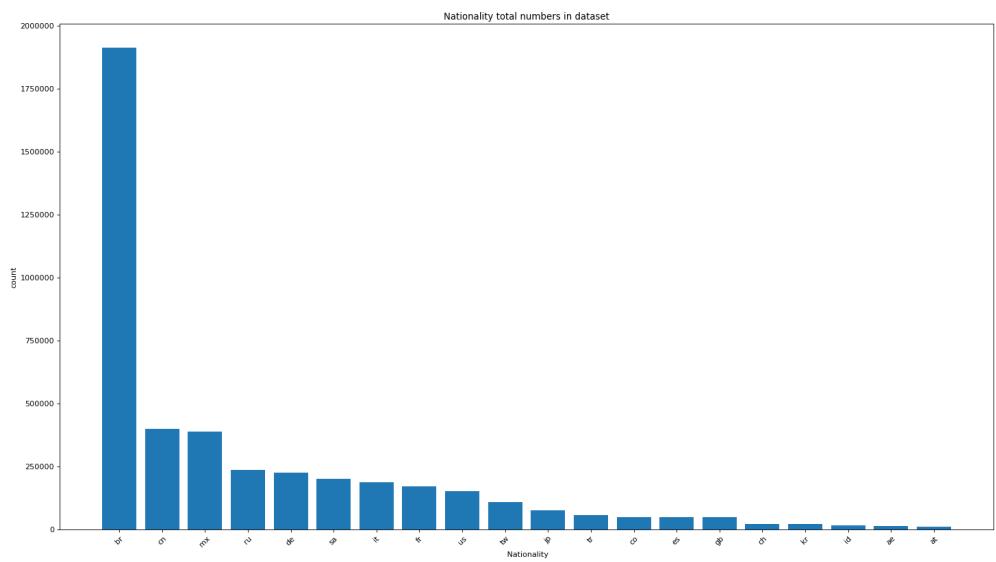
Also student uid, nationality, date, topic, grade, level.

We decided to provide our own error-type classification using the ErrAnt¹ model, for consistency and reproducibility.

Lets see the distribution of various parameters in the dataset, proficiency level, nationality and error-types:



¹ <https://github.com/chrisjbryant/errant> see [references](#).



Main:

Methodology:

The first big effort was made to parse the DS into a DF

Steps:

1. Parsing the dataset XML, extracting the original sentence and the corrected sentence.

The DS included correction type tagging, but for consistency and reproducibility we decided to annotate the correction type ourselves.

For this we created parallel files the model can use to annotate.

2. Annotating the correction type with ²[ERRANT]

We used a trained model³ to annotate the correction.

3. Creating an m2 file of the corrections

This is needed for step 5

4. Saving the metadata in a Dataframe.

5. creating CoNLL-U files for the original and corrected sentences using udpipe ⁴ library.

6. Convert the m2 and CoNLL-U files to a syntax based (p.o.s.) m2 file using ⁵

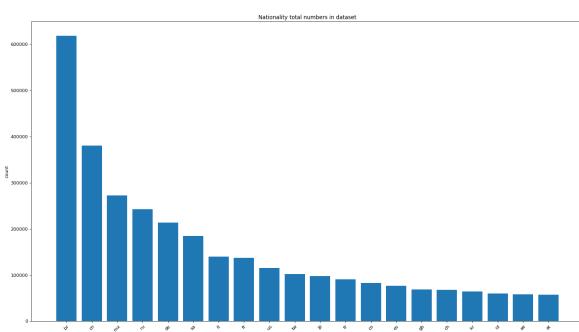
GEC_UD_divergences

7. Add syntax based correction-types to the Dataframe (DF).

Now we have the DF, We can run our analysis. In our DF, a data-point is a single error, classified by the model.

Decisions:

Levels 6, 4 were not included in the analysis, as they had negligible data-points.
(15 and 8 respectively.)



² ERRANT (<https://github.com/chrisjbryant/errant>)

³ <https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-3131/english-ewt-ud-2.5-191206.udpipe>

⁴ <https://github.com/ufal/udpipe>

⁵ https://github.com/borgr/GEC_UD_divergences

Seeing there are 186 nationalities, we decided to showcase mainly the largest 10, making the threshold 100,000 data points.

We decided to remove nationalities with less than 10,000 data-points. So the analysis refers to 20 nationalities.

Since the EFCAMDAT records nationalities and not languages, we use nationality as a proxy for language when comparing to other work on languages.

We only used for the comparison parts nationalities we can map confidently and without duplicates in our dataset.

The mapping is as follows:

Country	Name in EFCAMDAT	lang2vec	Language	Number DP in DS
Brazil	br	pt	portuguese	476,817
chine	cn	zh	chinese	165,162
mexico	mx	es	Spanish	87,259
russia	ru	ru	Russian	70K
Germany	de	de	German	54K
Saudi	sa	ar	Arabic	47,340
Italy	it	it	Italian	4K
France	fr	fr	French	41K
Japan	jp	jp	Japanese	21K

error profile

Create an error-profile for every language:

After creating the DF, I created an error-profile matrix, by using the syntax based (p.o.s.) error classification. Every cell in the matrix is the count of an error type.

Every line in the matrix is divided by its sum, to normalize it.

This was done for the whole dataset, as well as for single nationalities, student proficiency (level) and for a combination of both.

The same is shown after subtracting the average across all languages, per cell.

For a visual representation of the relationship between the proficiency and error profile, We created a matrix, where every cell shows a bar graph representing that cell's value, for all levels, ordered by proficiency.

lang to vec

Next we Downloaded the lang2vec⁶ dataset and extracted the features vectors, for comparison. We Quantified the correlation between distance metrics of languages represented by their lang2vec vector, compared to the same languages represented by their error-profile matrix. Distances were computed with cosine similarity. For lang2vec this was done on the whole vector.

For our error-profile matrices, this was done per line, then averaged.

Quantifying the correlation between the 2 methods for measuring the distances was done by Spearman-Rank correlation.

⁶ <https://github.com/antonisa/lang2vec> see [references](#)

Experiments and results:

Questions

Are there specific features that distinguish specific languages?

Do they correlate with known Cross-Linguistic Syntactic Divergences?

Is there a correlation between proficiency level and some error type?

Graphs

In the following graphs, the rows (marked ‘english’) are the corrected POS, and the columns (marked according to the current DS slice: the whole DS, nationality, Level etc.) are the original POS, as classified by ERRANT⁷.

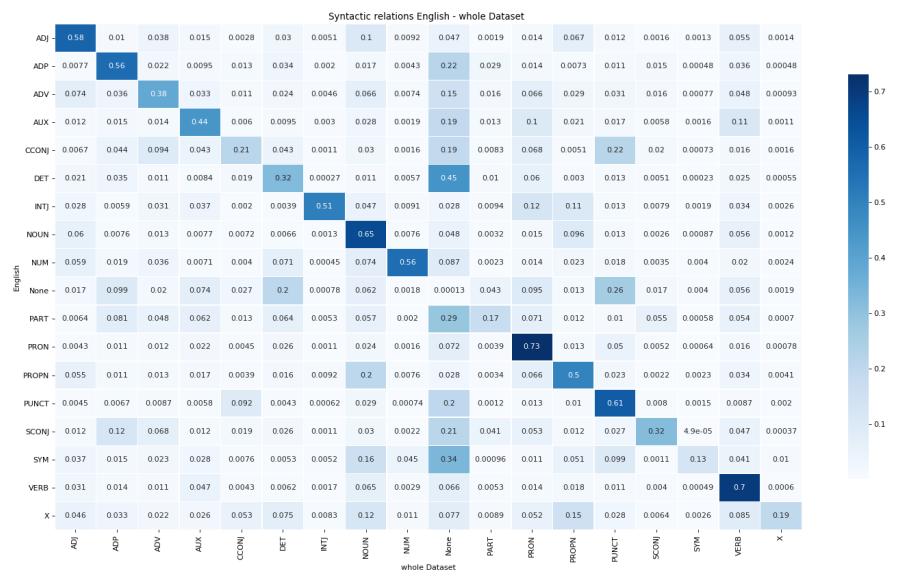
Classified Syntactic Errors for the Whole DS:



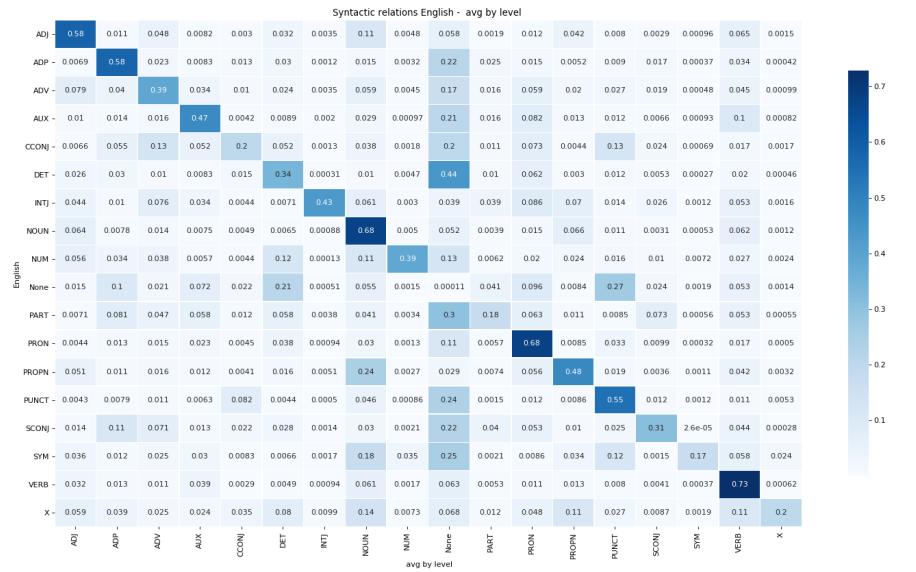
We can see the ‘NOUN->NOUN’ is prominent. Although it is attenuated when we normalize per row.

From here on all graphs are normalized per row. (Every cell is divided by the row total.)

⁷ <https://github.com/chrisibryant/errant>



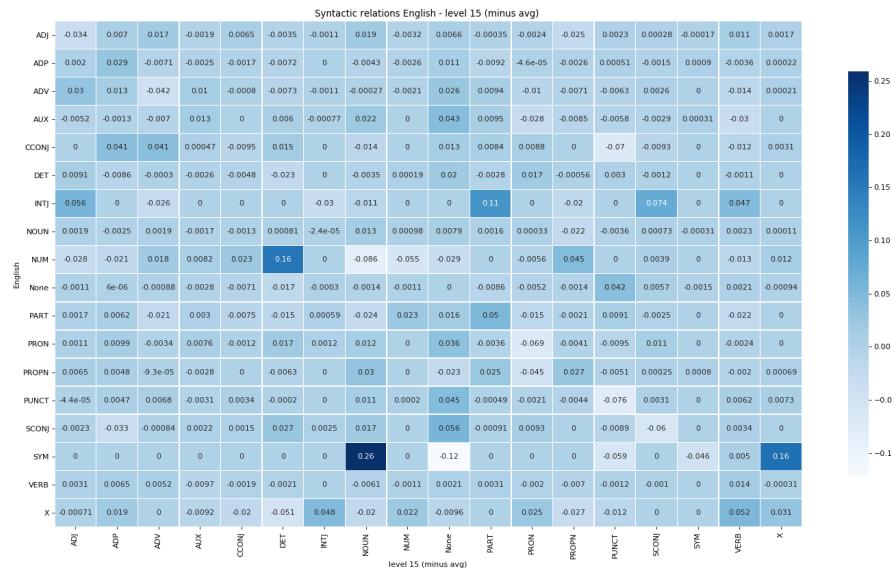
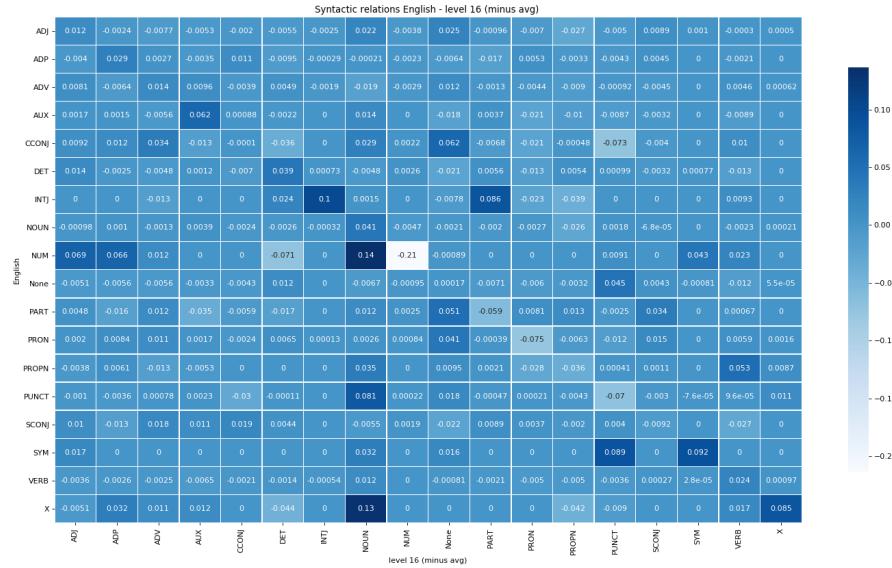
Classified Syntactic Errors for the Whole DS, averaged by level:
(each level gets the same weight, regardless of size.)

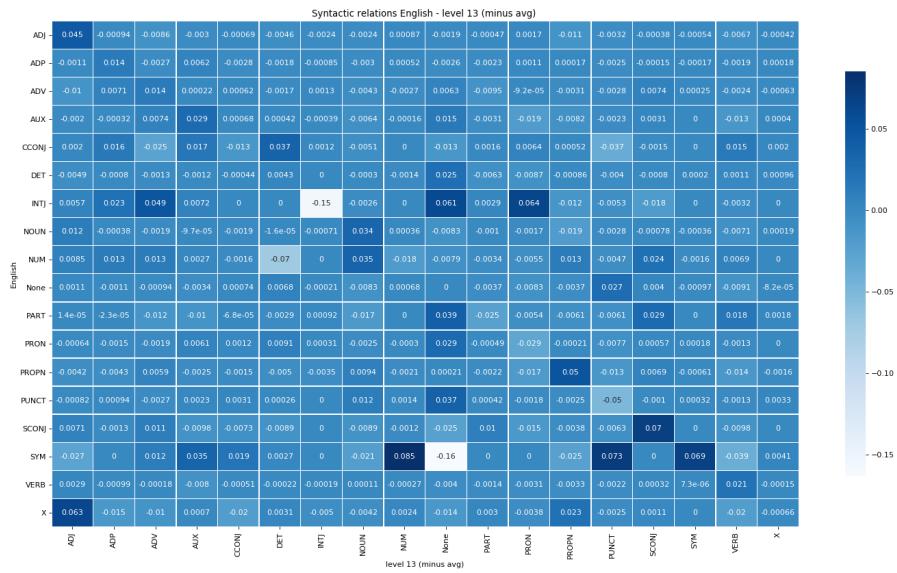
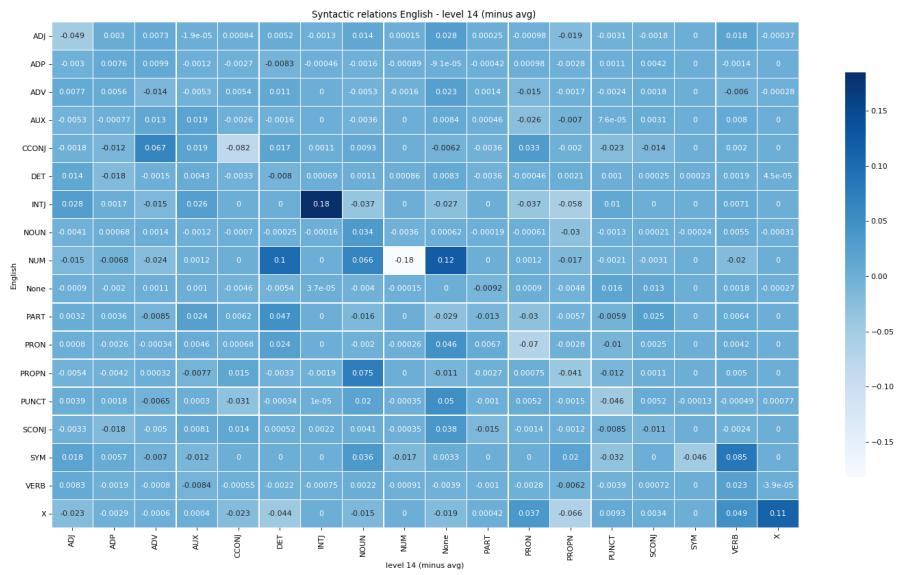


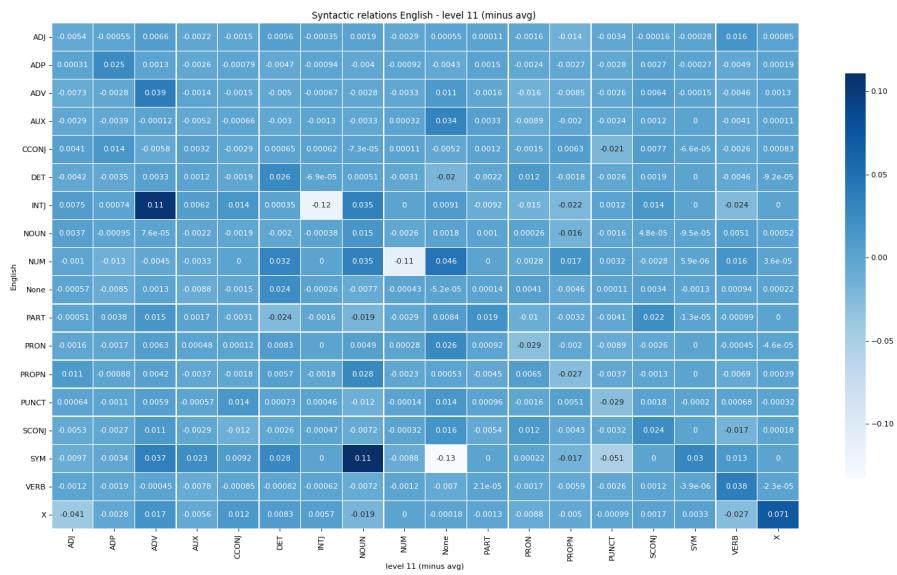
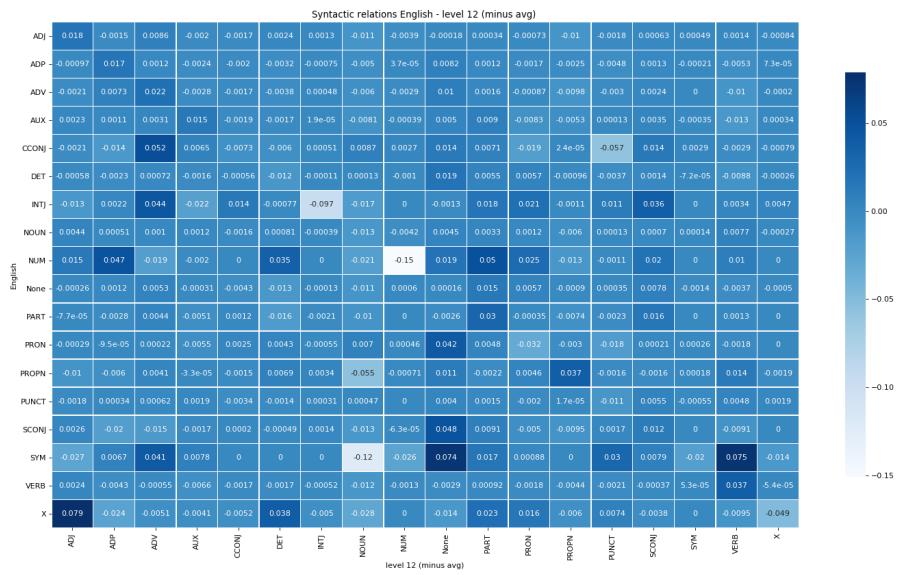
Classified Syntactic Errors by level, minus the average across levels:

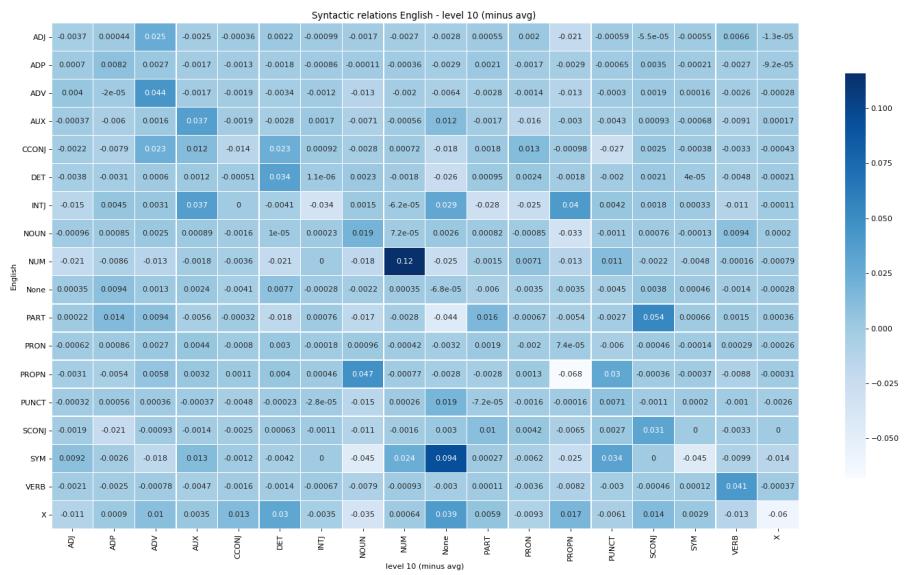
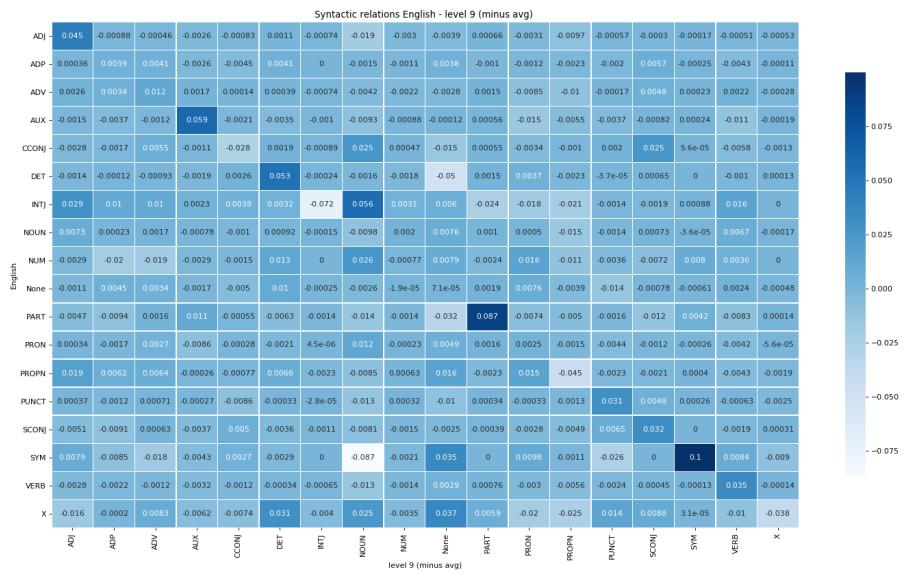
The error profile for the dataset minus the average value for each cell.

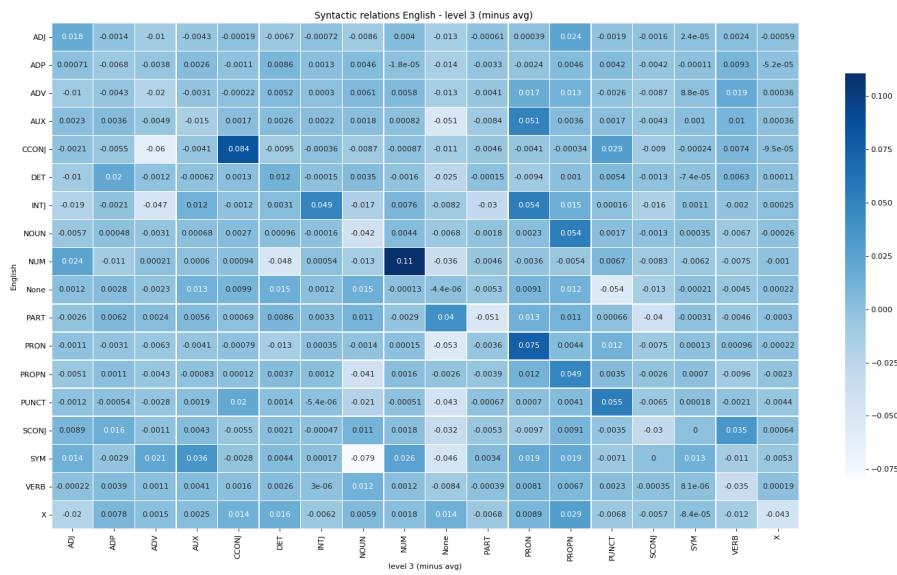
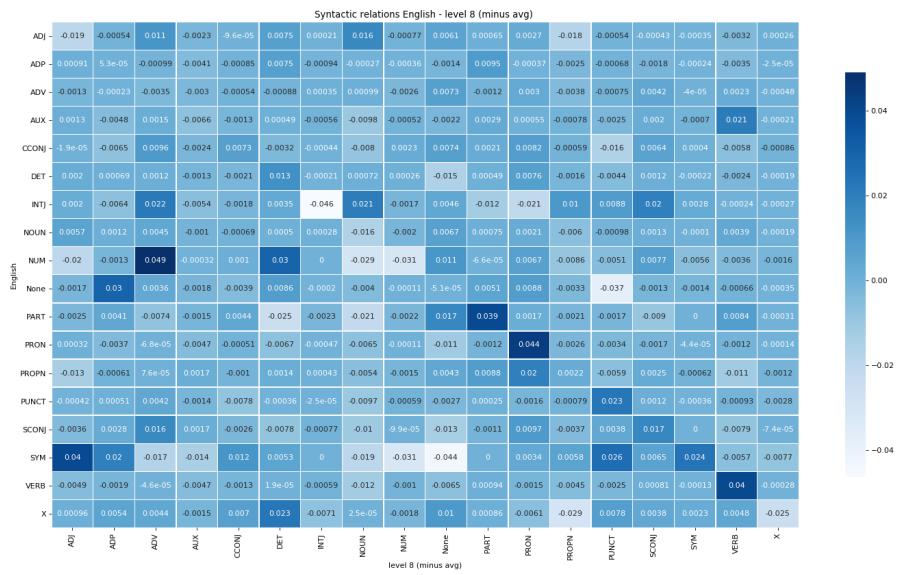
The average is computed across Levels.

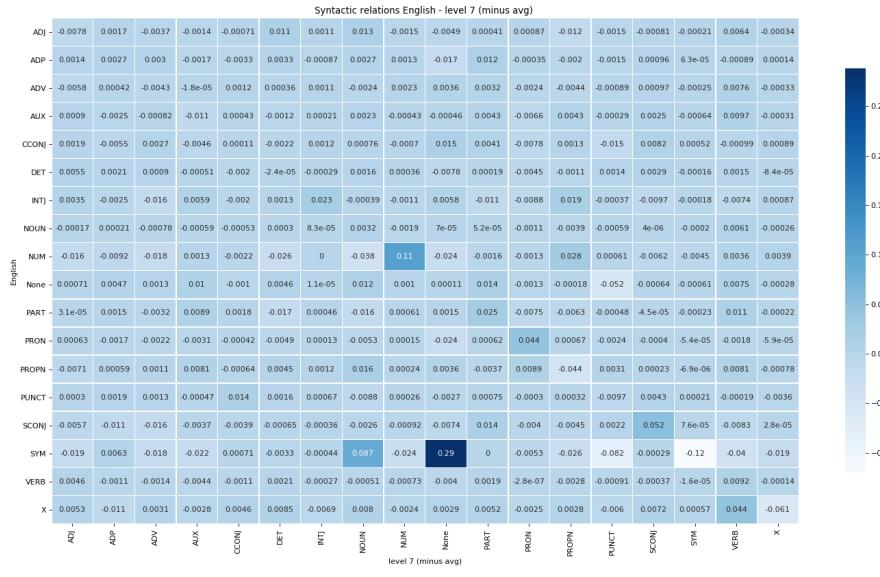
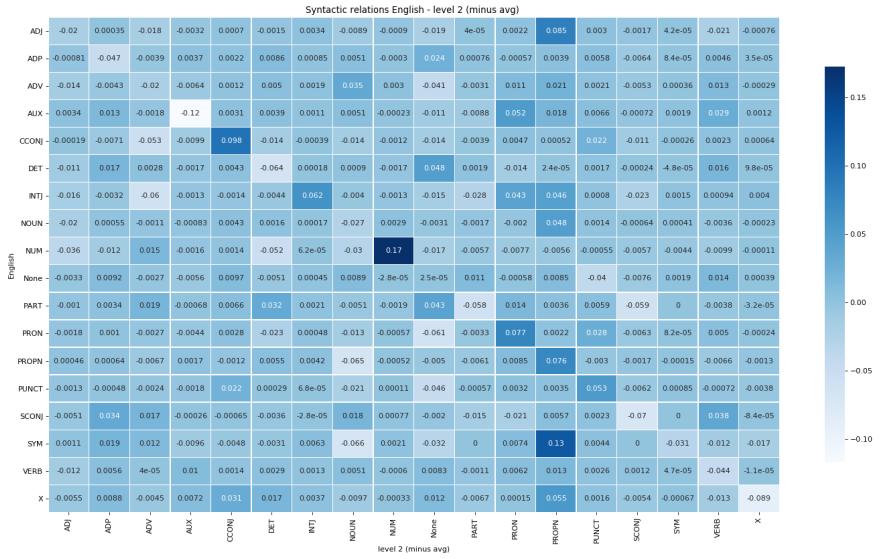


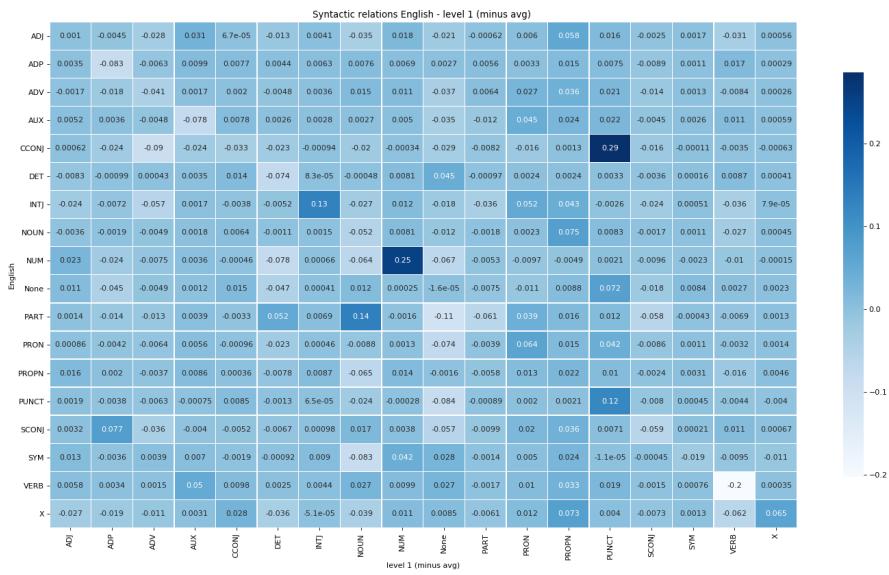










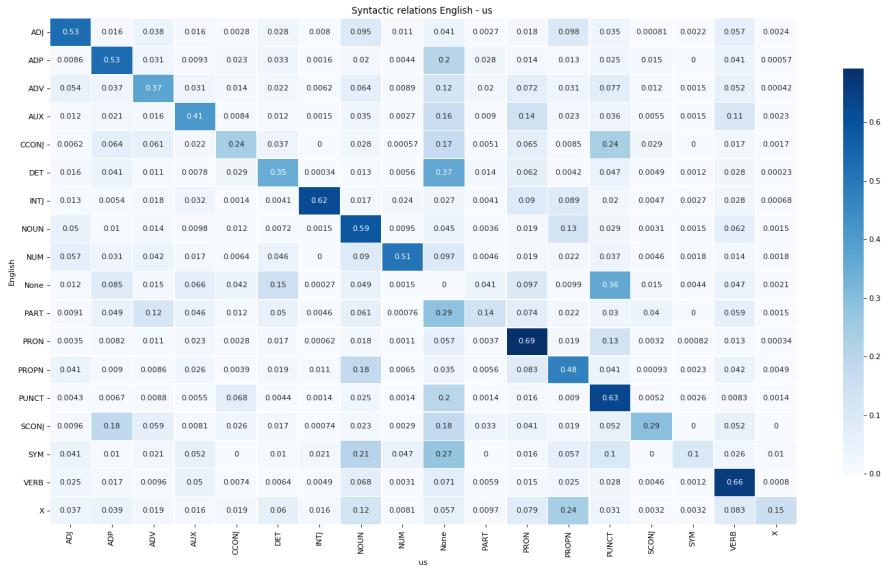
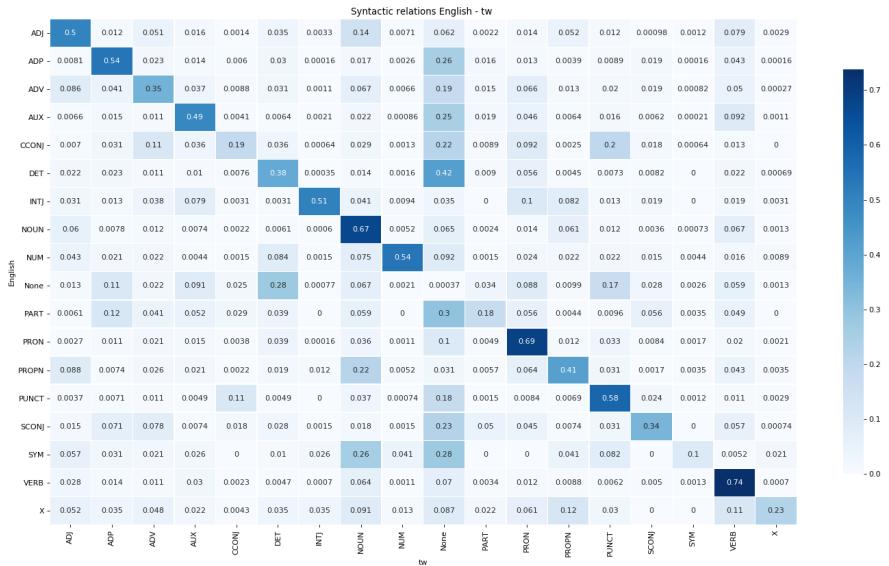


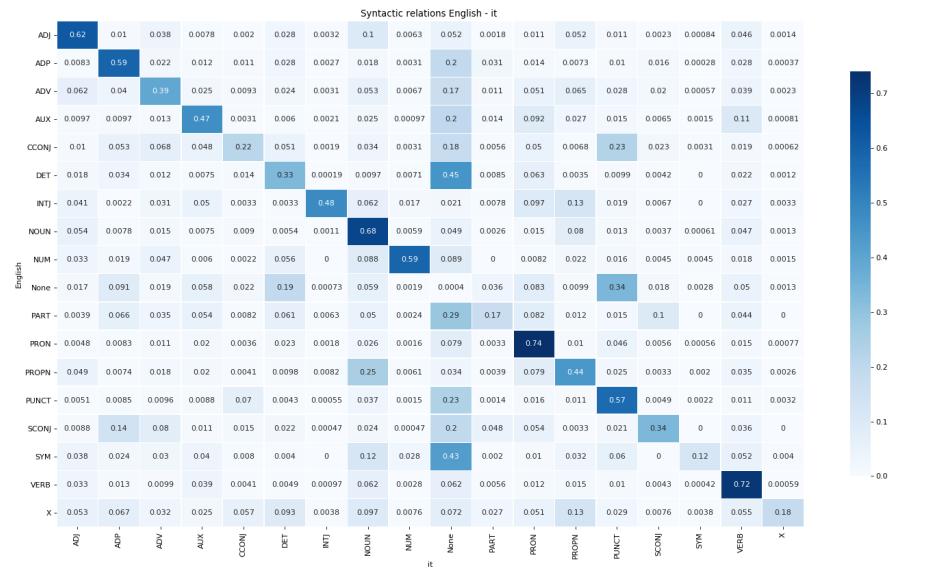
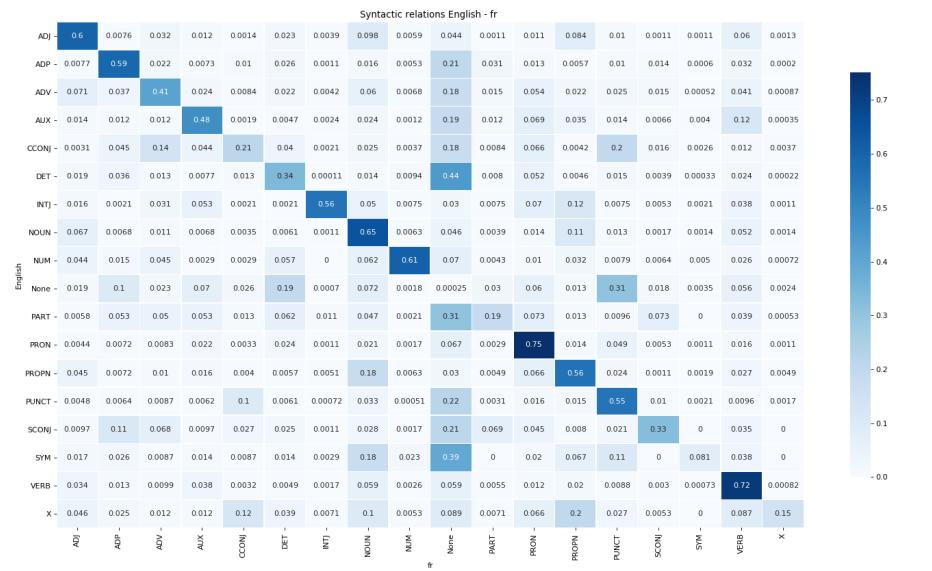
The error profile for each nationality:

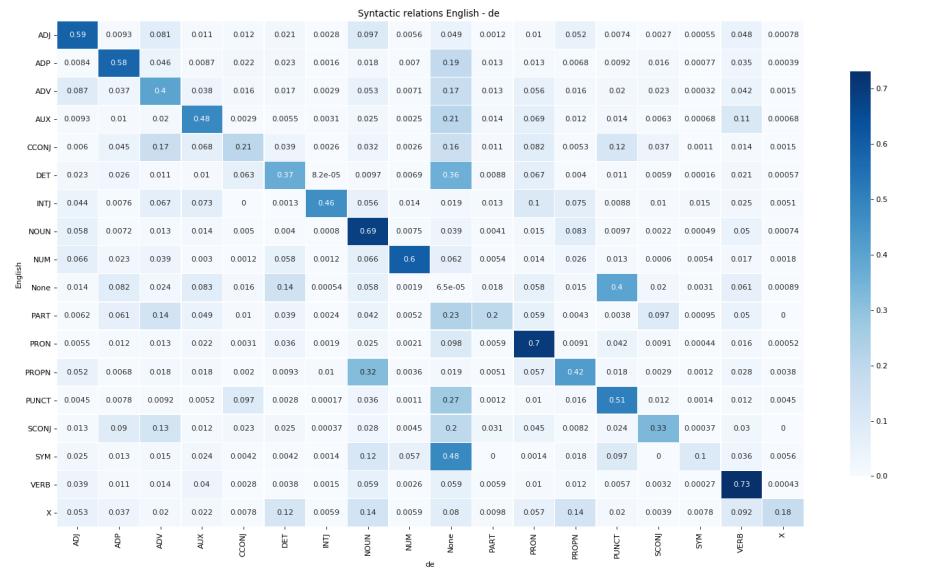
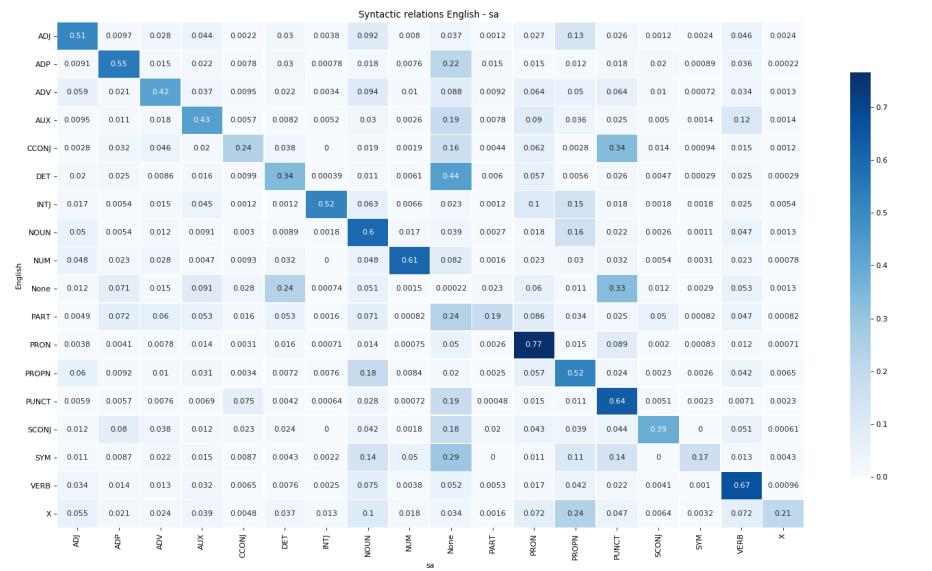
Classified Syntactic Errors by nationality:

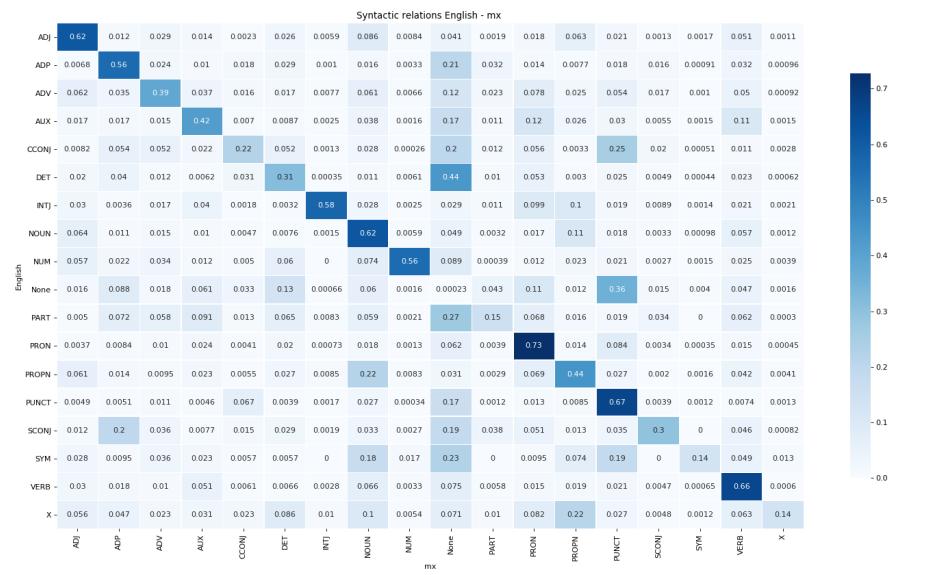
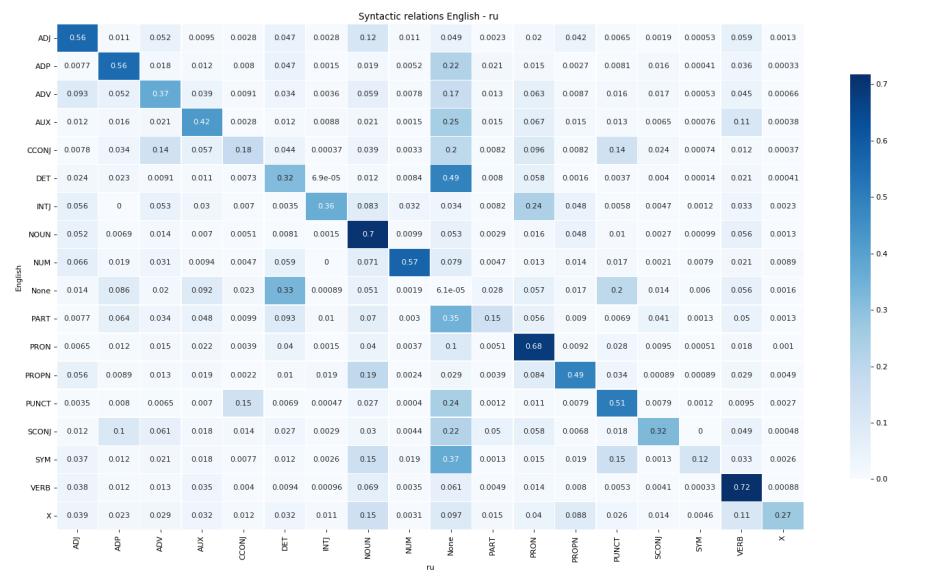
The error profile for the dataset minus the average value for each cell.

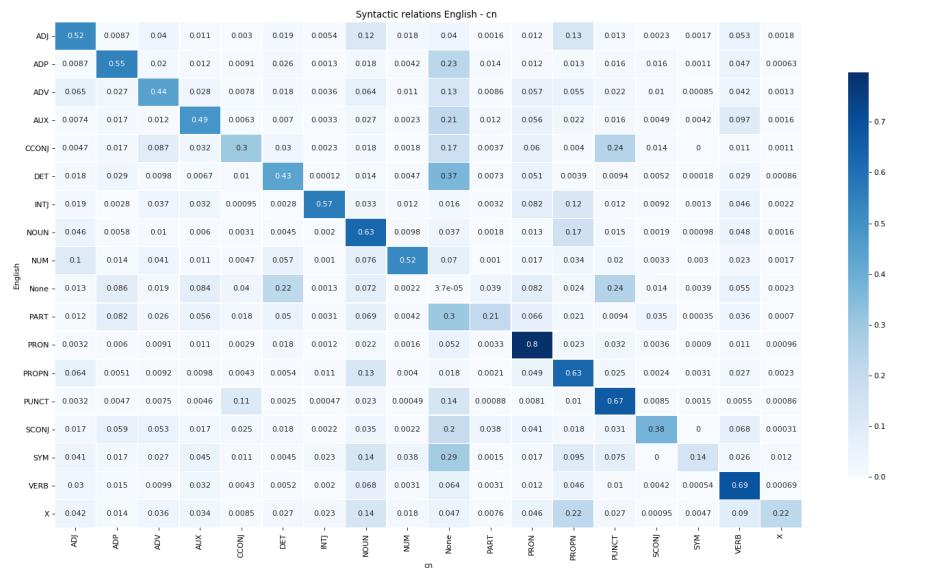
The average is computed across languages.

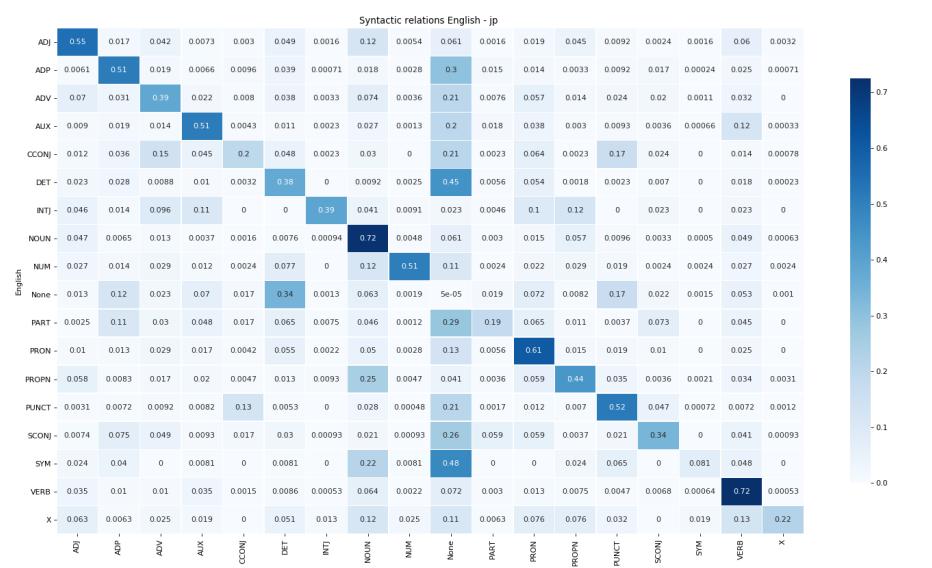
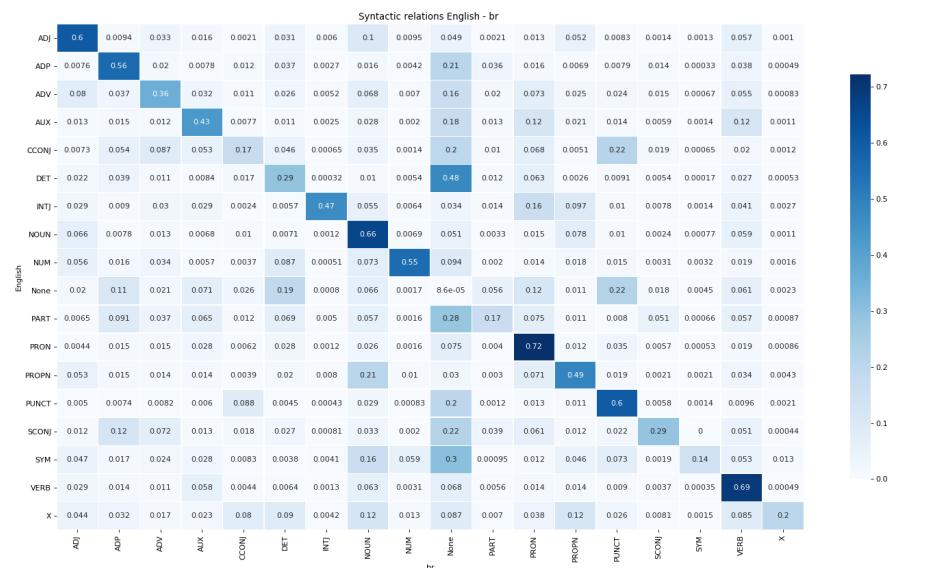






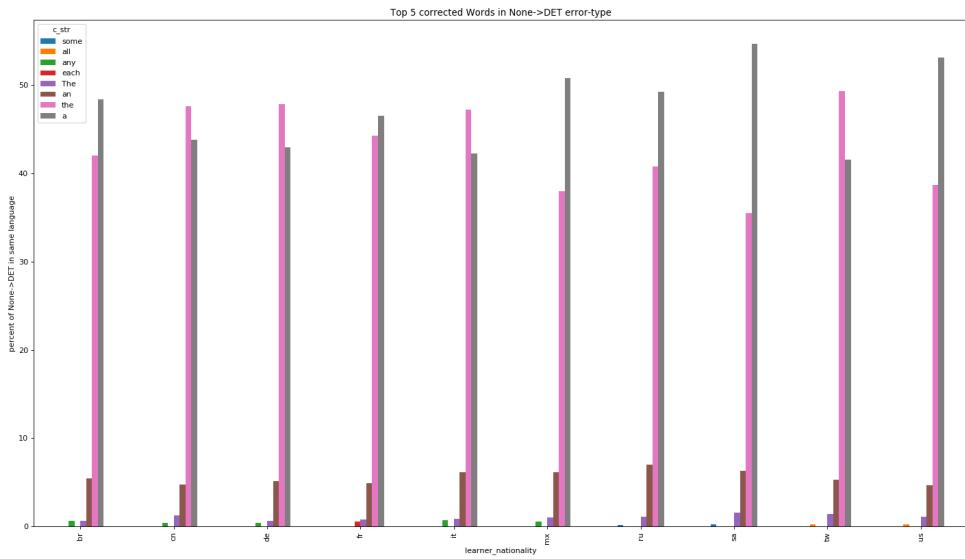






Prominent Words

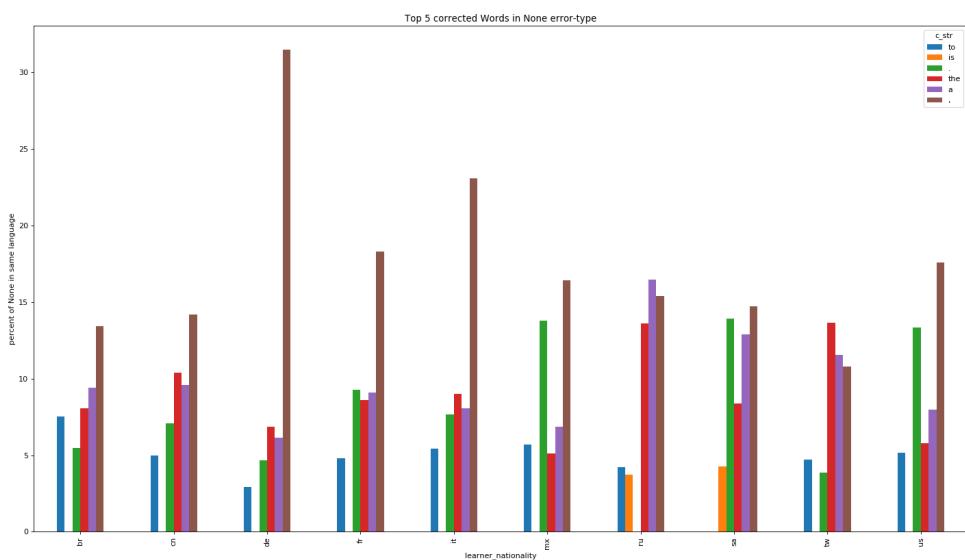
We can see **None->DET** is standing out, So let's see what words are prominent in this error-type.



We can see there are 2 prominent words in every language. All the rest are less frequent.

'The' and 'a'. For cn, de, it and tw(taiwan, also chinese) 'the' is the most missed word.
For br, fr, mx, ru, and us 'a' is the most missed word.

If we look at all missed words(**None->***):

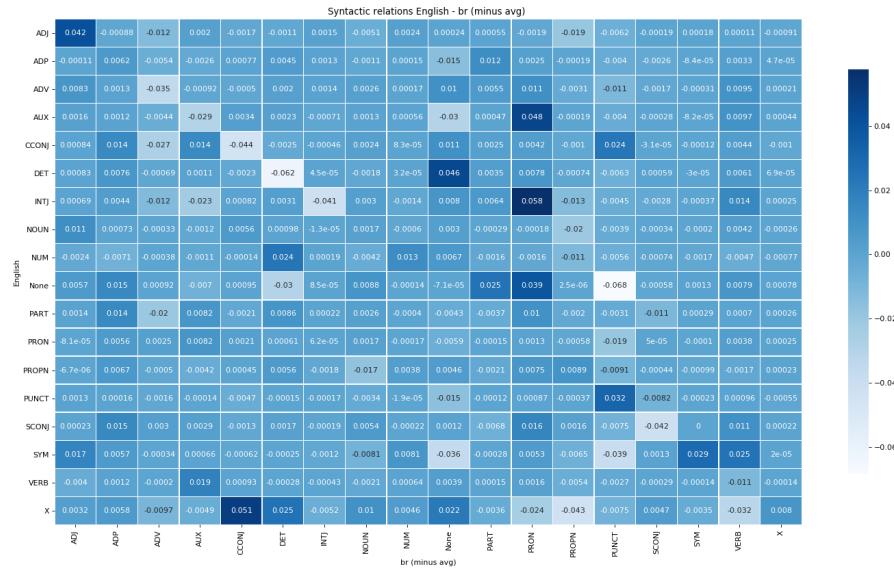
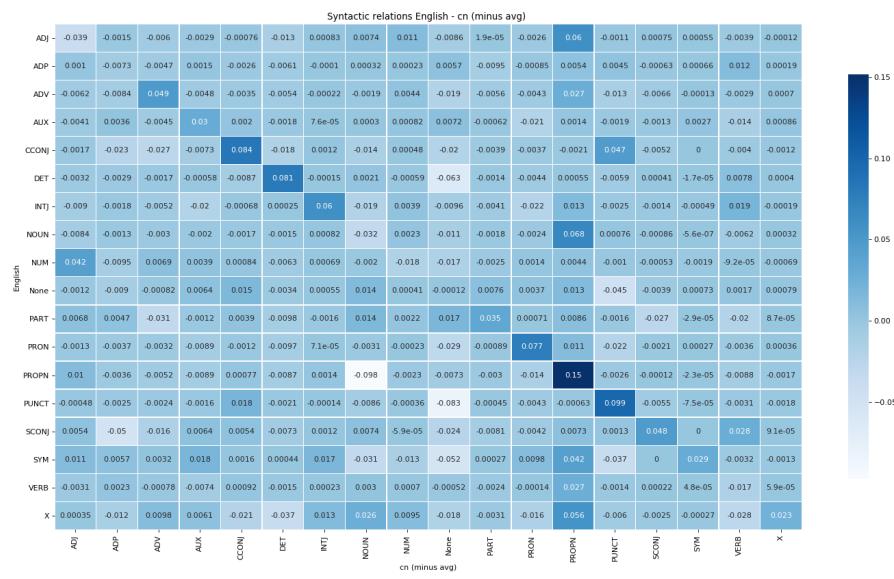


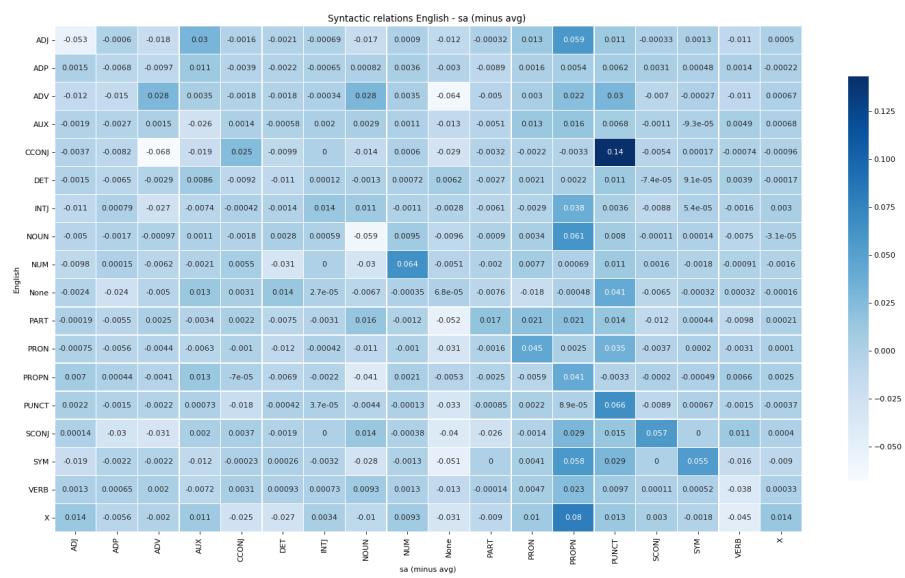
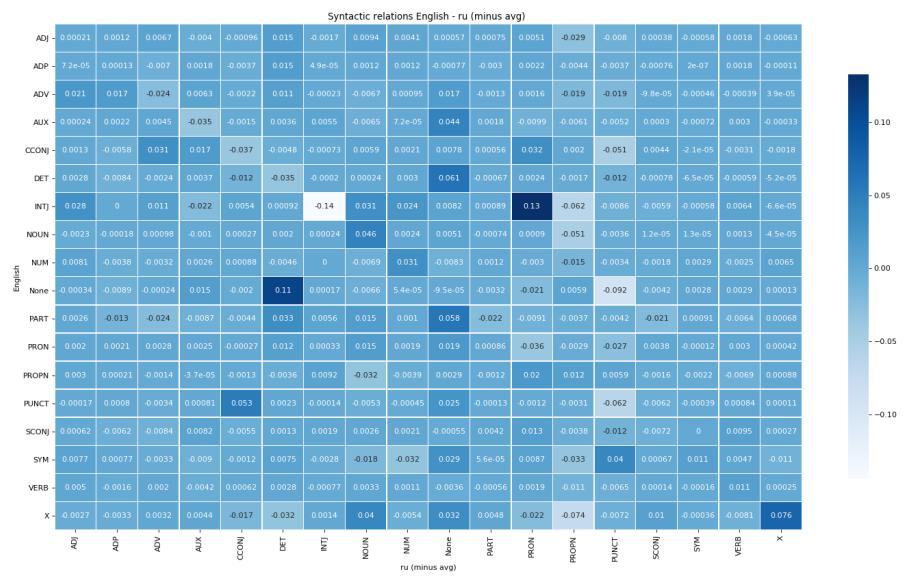
We can see that comma(,) is prominent, except for ru with 'a' and tw with 'the'.

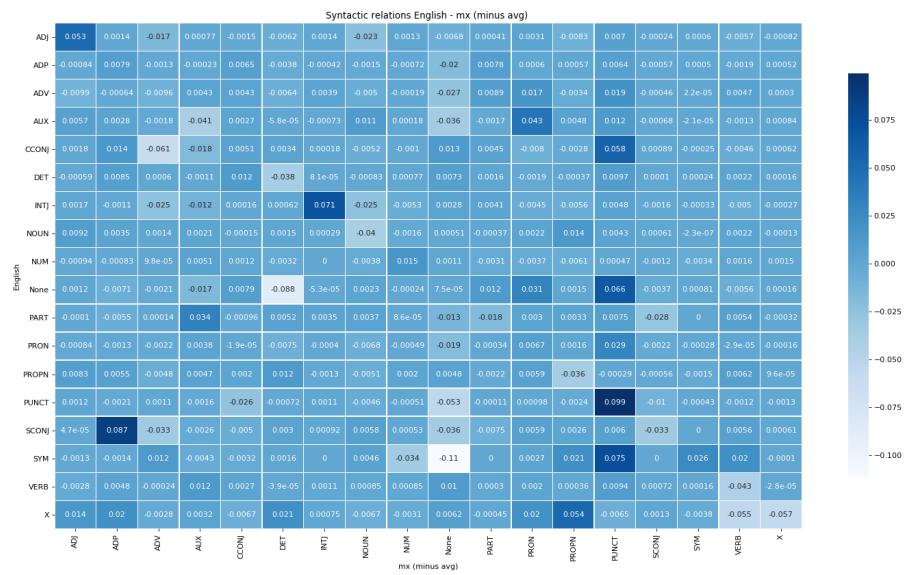
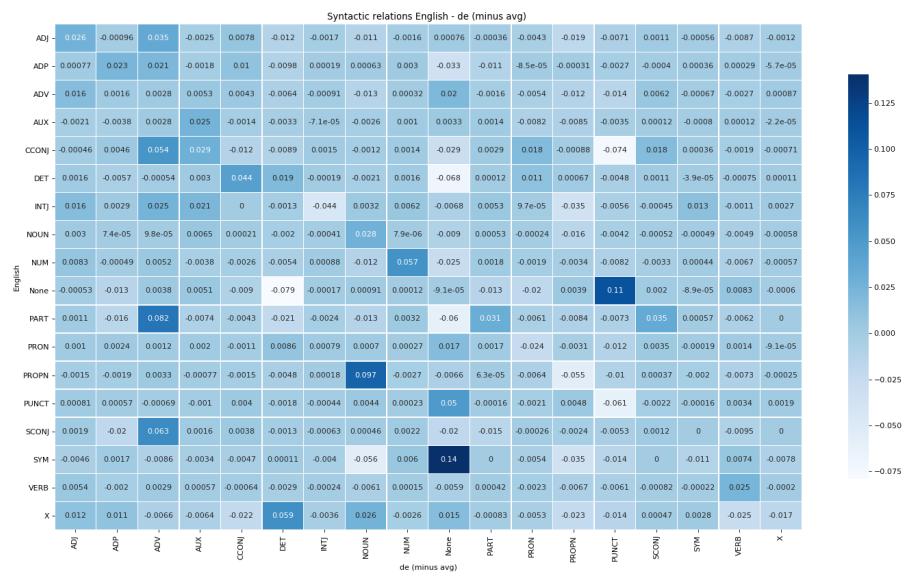
Interestingly, 'to' joined the top 5 missed words, except for sa(saudi - arabic) where 'is' joined instead. 'Is' joined for ru as well.

Interestingly the comma is over-represented in de. (Is comma less prevalent in German?)

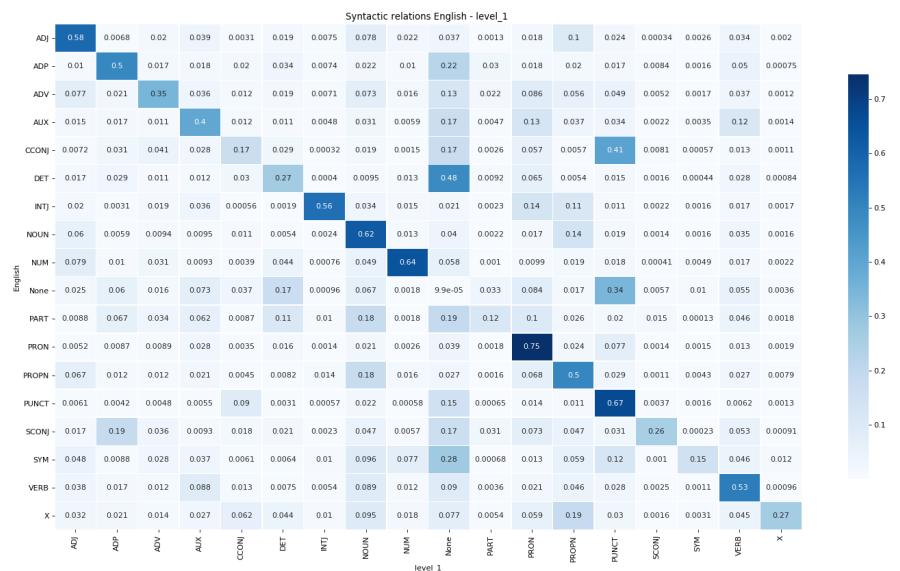
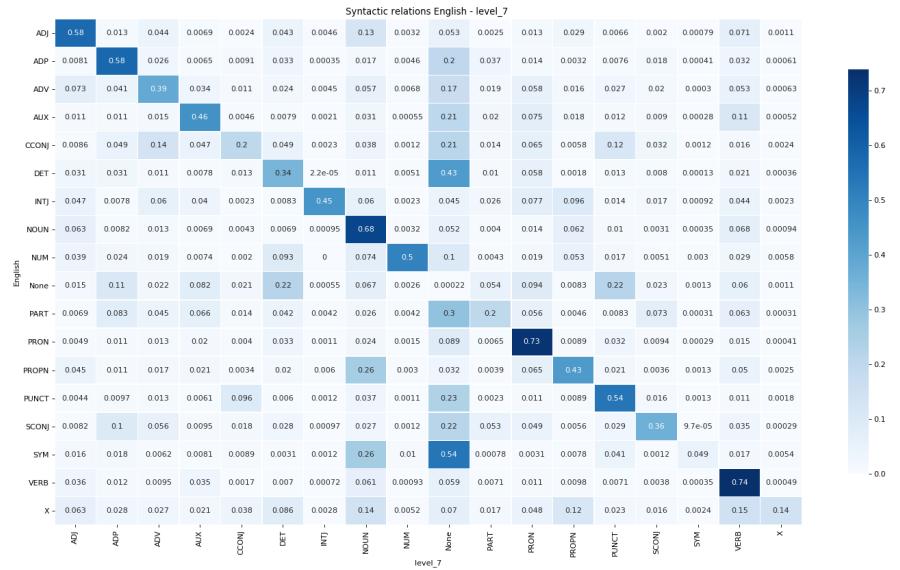
Classified Syntactic Errors by nationality, minus the average across nationalities:

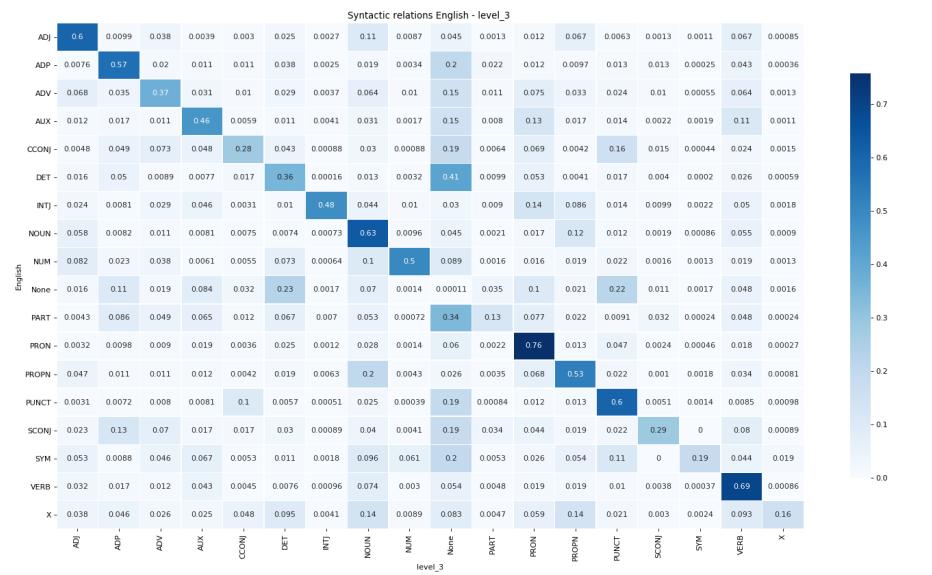
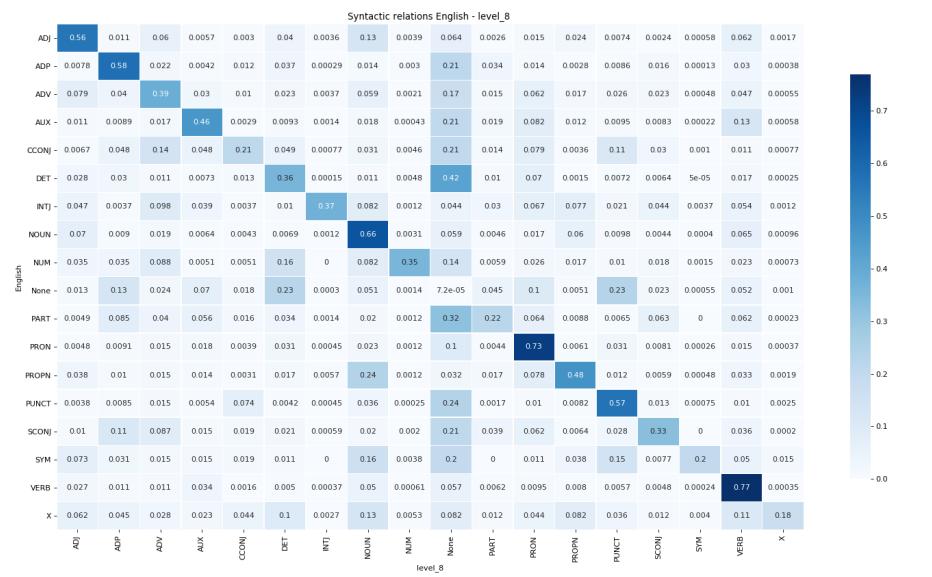


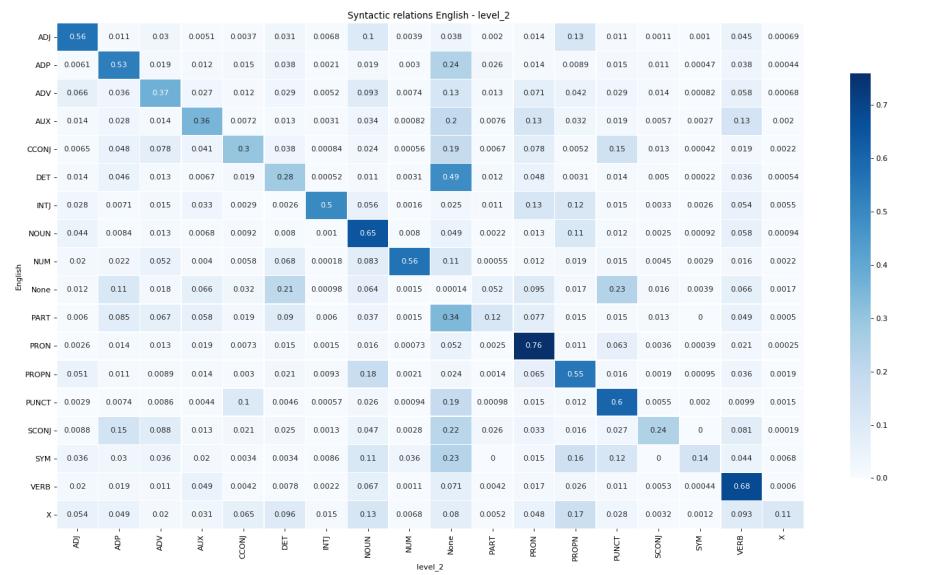
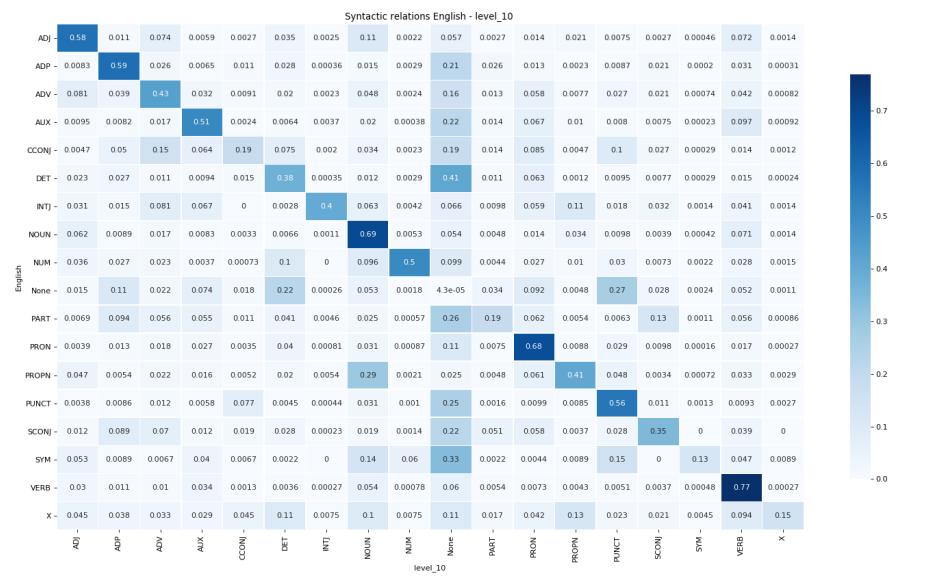




Classified Syntactic Errors by level:

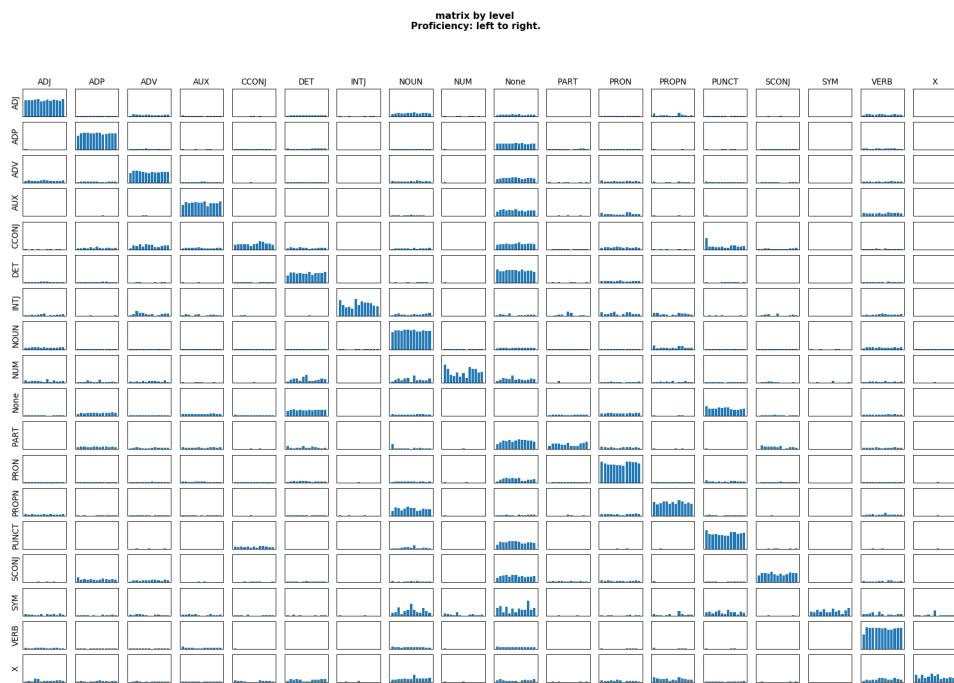






All Levels together:

For visualization purposes, We inserted into each cell all levels together.
In each cell the levels are ordered the same, by proficiency (1,2,3,4...16):

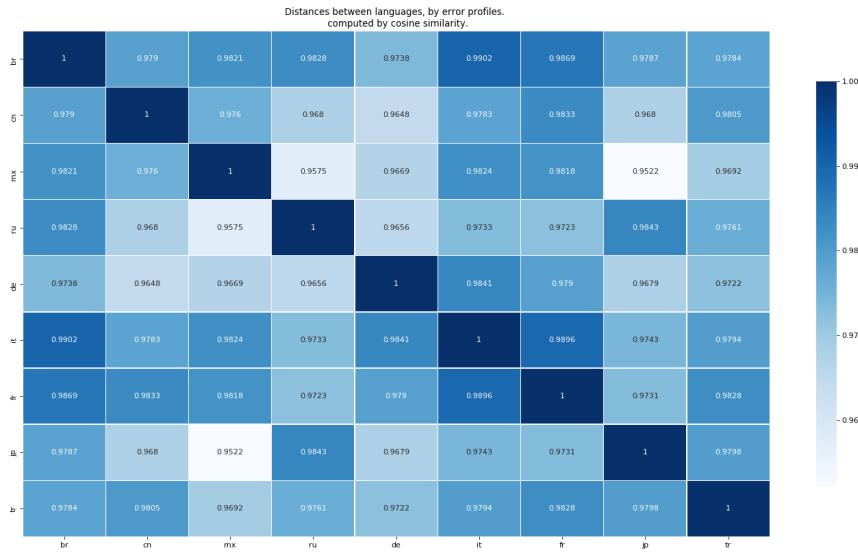


It is hard to see a specific influence or correlation of proficiency level on the error profile, but this gives a first visual guide for anyone wanting to search deeper.

Distances

We computed distances between languages, represented by their error-profile matrix. Distances were computed with cosine similarity

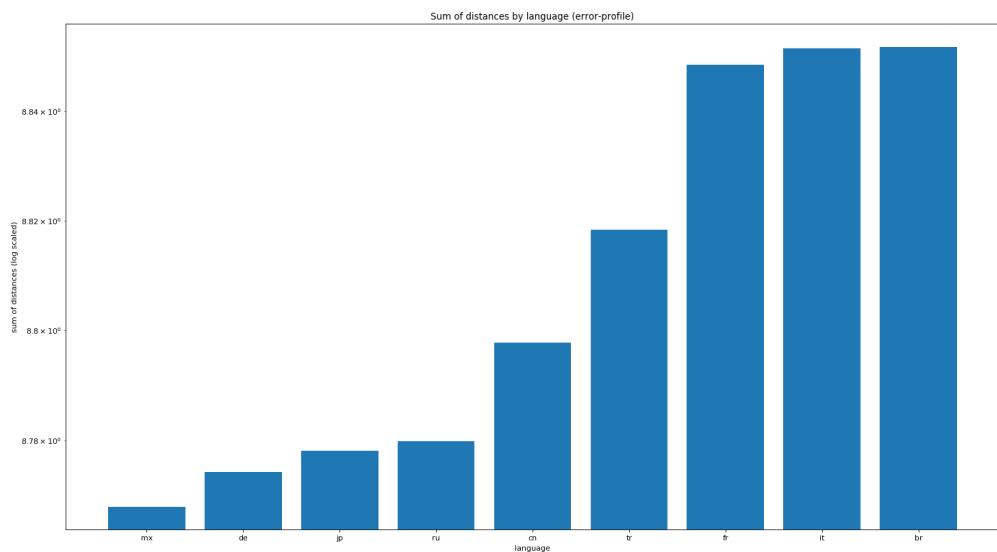
Cosine similarity this was done per line, then averaged:



It seems that Jp and de look the least similar to all other languages.

If we sum the distances:

We used a log scale on the y axis, as the numbers are all close



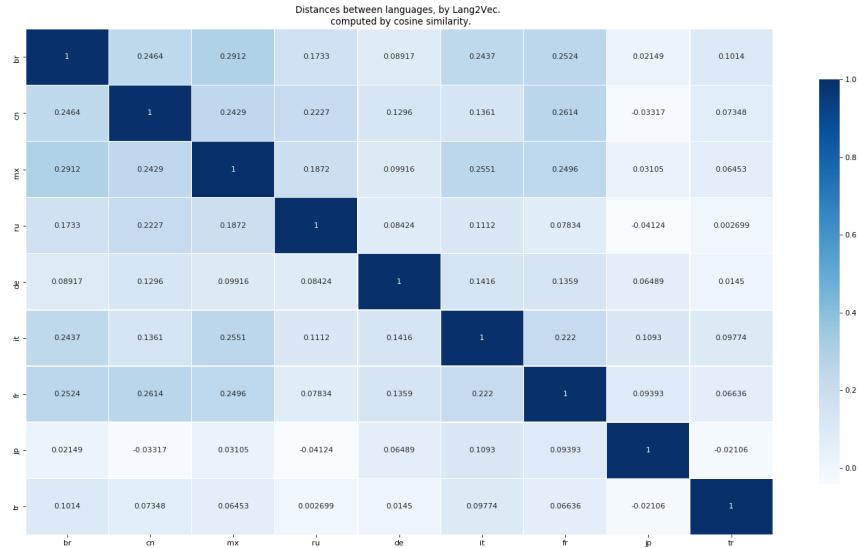
It seems mx, de, jp and ru are dissimilar to other languages, at least in regards to their error-profile.

Next we Downloaded the lang2vec dataset and extracted the features vectors.

We computed distances between languages, represented by their feature vectors.

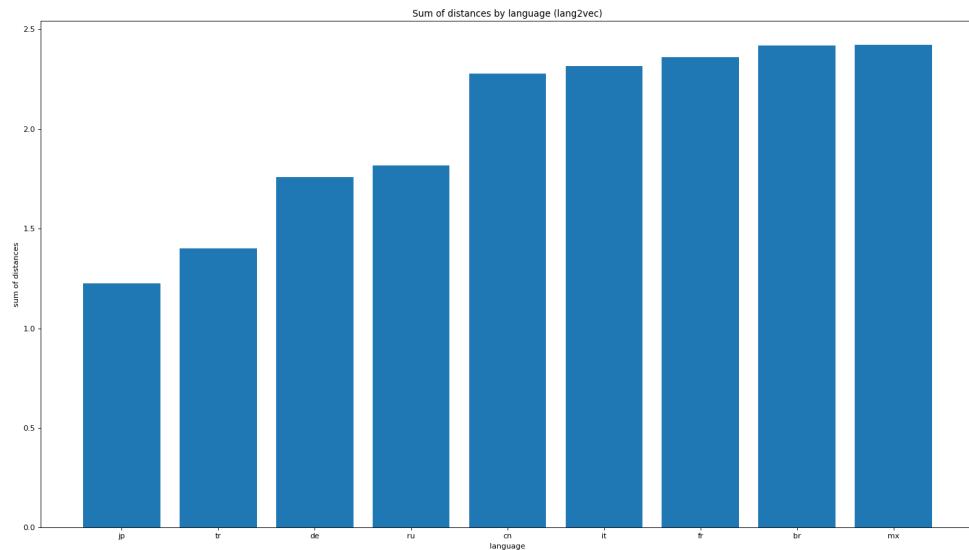
Distances were computed with cosine similarity

Cosine similarity was done on the whole vector:



It seems that jp and tr look the least similar to all other languages.

Let's sum it:



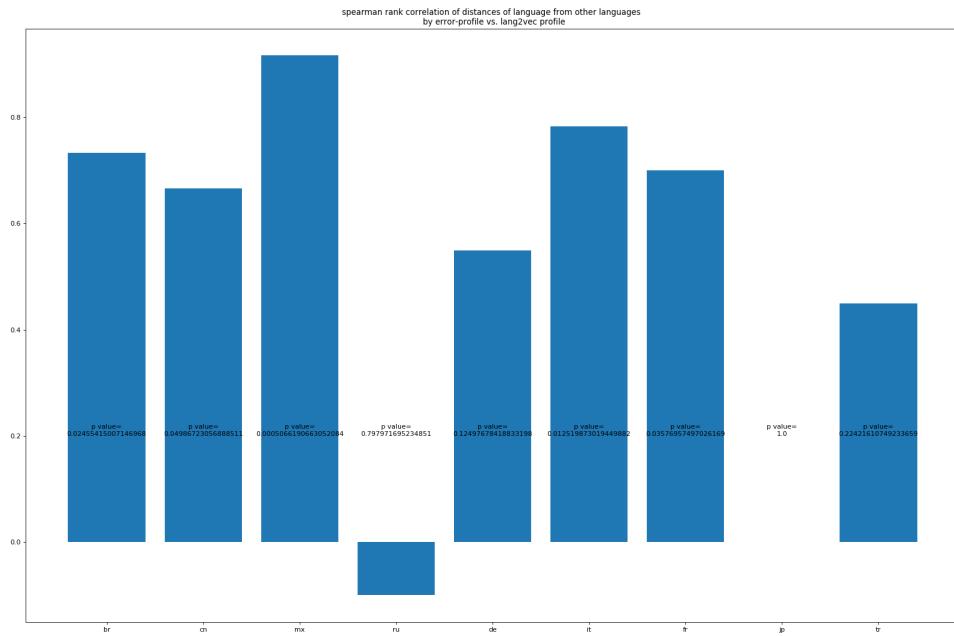
Here we didn't need the log scale..

It seems jp and tr (Turkish) are dissimilar to other languages, at least in regards to their lang2vec vector.

Interestingly, Japanese is dissimilar to other languages in both these comparisons.

Measuring distances between languages with our error-profile compared to lang2vec:
 Quantifying the correlation between the 2 methods for measuring the distances was done by Spearman-Rank correlation.

Measuring distances between languages with our error-profile compared to lang2vec, shows high correlation for most languages, with the highest being spanish (mexico). Interestingly Japanese and Russian showed negative or zero correlation.

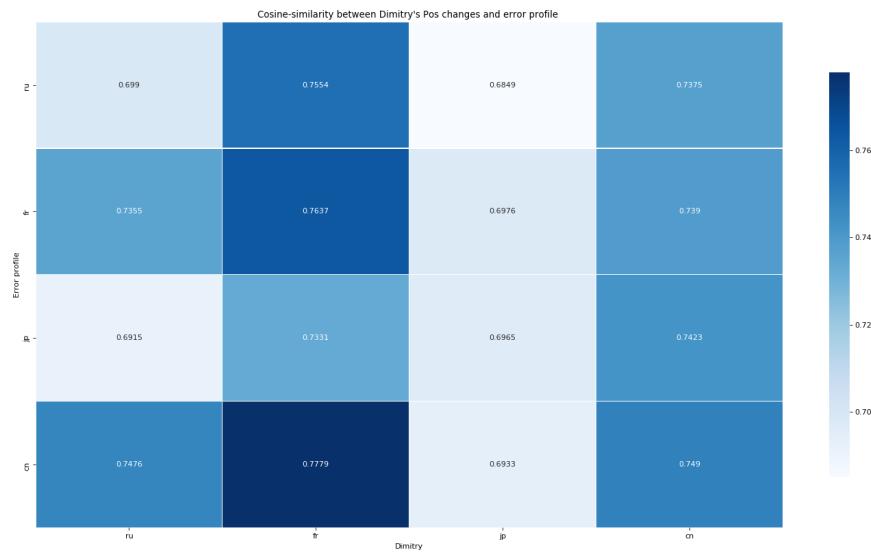


Comparing to other work:

Using data from ‘*Fine-Grained Analysis of Cross-Linguistic Syntactic Divergences*’ (Nikolaev et al. 2020⁸) we compared various languages to our error-profile matrices for the same languages. We measured the distance between our matrices and this data, with cosine-similarity as well.

⁸ <https://arxiv.org/abs/2005.03436>

All distances are similar and we cannot see a preference for the same language across the two datasets. (The diagonal is not emphasized.)



Results

Recap of Some results:

Measuring distances between languages with our error-profile compared to lang2vec, showed high correlation for most languages, with the highest being Spanish (mexico). Interestingly Japanese and Russian showed negative or zero correlation.

It seems jp and tr (Turkish) are dissimilar to other languages, at least in regards to their lang2vec vector. Interestingly, Japanese is dissimilar to other languages in both these comparisons.

It seems mx, de, jp and ru are dissimilar to other languages, at least in regards to their error-profile.

Missed DET POS:

We can see there are 2 prominent words in every language. All the rest are less frequent.

'The' and 'a'. For cn, de, it and tw(taiwan, also chinese) 'the' is the most missed word.

For br, fr, mx, ru, and us 'a' is the most missed word.

All missed POS:

We can see that comma(,) is prominent, except for ru with 'a' and tw with 'the'.

Interestingly, 'to' joined the top 5 missed words, except for sa(saudi - arabic) where 'is' joined instead. 'Is' joined for ru as well. '.' (stop) joined the top 5 for all languages except ru.

Future work:

It can be interesting to look at the distribution of specific words in the errors in different nationalities and correlate that with the respective language syntactic use of that word.

Also studying differences in the above matrices and finding correlation to specific languages, or as a learning process, as correlated to differences in proficiency levels, in the DS as a whole, and in specific languages.

For more matrices see the appendix.

Another thing is parsing another parallel corpus the same way and comparing the results.

Link to github

Github repository of this work :

<https://github.com/amichw/EFCAMDAT>

References:

Works used or cited:

[**https://github.com/borgr/GEC_UD_divergences**](https://github.com/borgr/GEC_UD_divergences):

```
title = "Classifying Syntactic Errors in Learner Language",
author = "Choshen, Leshem and
Nikolaev, Dmitry and
Berzak, Yevgeni and
Abend, Omri",
booktitle = "Proceedings of the 24th Conference on Computational Natural Language Learning",
month = nov,
year = "2020",
address = "Online",
publisher = "Association for Computational Linguistics",
url = "https://www.aclweb.org/anthology/2020.conll-1.7",
pages = "97--107",
```

[**https://github.com/ufal/udpipe**](https://github.com/ufal/udpipe)

[**https://github.com/chrisjbryant/errant**](https://github.com/chrisjbryant/errant)

[**https://github.com/antonisa/lang2vec**](https://github.com/antonisa/lang2vec)

[**https://universaldependencies.org/**](https://universaldependencies.org/)

[**https://arxiv.org/pdf/2005.03436.pdf**](https://arxiv.org/pdf/2005.03436.pdf)

[**https://github.com/macleginn/exploring-clmd-divergences**](https://github.com/macleginn/exploring-clmd-divergences)