

Cross-Linguistic Neural Phoneme Embeddings for Computational Historical and Typological Linguistics

Marlon Betz

August 1, 2016

Contents

1	Introduction	2
2	Theoretical and Practical Motivation	3
3	Related Research	4
4	Embedding Models	4
4.1	Neural Embedding Models	4
4.1.1	Neural Language Models	4
4.1.2	Word2Vec	5
4.1.3	Polyglot Neural Language Models	10
5	Evaluation	13
5.1	Data	13
5.2	Intrinsic Evaluation	13
5.3	Discussion	15
5.4	Phoneme Inventory Clustering	17
6	Resume	17
7	Appendix	22

1 Introduction

Embeddings nowadays build the backbone of every Deep Learning NLP architecture. Due to their capacity to encode a vast amount of latent semantic and syntactic information without the need of previous manual construction, they gave rise to the contemporary Deep Learning boom, i.e. deep neural networks that can capture hidden features in data sets and by that have allowed for huge performance gains in numerous NLP fields.

While embeddings are currently being used for several linguistic units such as characters [Kim et al.2015, dos Santos and Zadrozny2014, Zhang et al.2015], words [Mikolov et al.2013a, Mikolov et al.2013b, Pennington et al.2014] or entire sentences [Kiros et al.2015], proper phonemes have been largely been excluded from this trend, although there are some data that would indeed allow for it and research areas where they could be of great use, especially in the fields of Computational Typological and Historical Linguistics, which the current deep learning boom has hardly touched yet. There, even though it is clear that some phonemes share more common features and such form natural classes, they are often treated as pure symbols that share a common distance between each other. Even phoneme representation models that do incorporate phonological features are often hand-crafted [Kondrak2000, Rama2016] or include task-specific information that inherently suffer

from restricted generalization abilities when used for other tasks [Jäger2014]. Moreover, those methods usually reduce the number of possible phonemes to a minimum, getting rid of important information such as secondary or co-articulations.

In this paper, I will first discuss a theoretical motivation for using data-driven phoneme embeddings instead of plain symbolic representations or hand-crafted feature encodings. I will then shortly take a look on the related research of character embeddings. A big part of this paper will then review several embedding models. I will then discuss intrinsic evaluation methods for the embeddings and use those to compare the performances of the previously described models. This is then followed by a discussion of several use cases that could be interesting for those interested in data-driven approaches to Typological and Historical Linguistics. Finally, I will recapitulate the benefits and drawbacks of phoneme embeddings in a final resume.

2 Theoretical and Practical Motivation

It is well-established that certain phonemes share phonological features, such as [NASAL], [VOICE] or [LABIAL]. This implies that some phonemes are closer to each other than to other phonemes, given a certain distance measure incorporating such features. However, phonemes are nevertheless handled as pure symbols quite often, where the distance between phonemes does not vary. Even in the field of Historical Linguistics, where words are often given as phoneme sequences, symbolic string comparisons remain the backbone of most models. For instance, given the three hypothetical words /apa/, /aba/ and /aq”a/, purely symbolic string distances would estimate the distances between the words to be the same. However, /b/ and /p/ only differ in voice but are both labials, while /q”/ is an uvular ejective and hence should be less close to the labial phonemes. It is indeed possible to incorporate such features by hand. However, this decision would either yield a tremendously tedious work if done for all phonemes to be worked with, or would be incomplete insofar as features that might be seen to be unnecessary for a certain task are left out and a big number of distinctive phonemes would merge. Moreover, such hand-crafted feature architectures do not directly allow for the conversion of one feature architecture into the other: A model designed for IPA cannot directly work with ASJP data and vice versa, as the symbols and their respective features are different for each model. Embeddings can solve those problems: On the one hand, they allow for fine-grained distance measures between phonemes and thus overcome the symbolic constraint of equidistance. On the other hand, they are trained on data, which means all phonemes found in a corpus can be embedded in the embedding space without the need to create features by hand, given a corpus that is big enough to capture a variety of phoneme inventories to allow for unbiased phoneme comparisons beyond language boundaries. Once having a reliable embedding space, it is further possible to use that space for other tasks without the need to train new task-specific embeddings. Embedding spaces in general allow for one-shot or zero-shot

learning. This means that embeddings trained on some task can either be used for another task directly or by fitting a simple regression from one embedding space into the other, without the need to train new embeddings for a new task. This follows the idea that embeddings should incorporate as much information as possible, while for certain tasks some information might be irrelevant and can be left out.

3 Related Research

Word embeddings have proven to be a reliable tool for several NLP tasks. However, finer grained embedding models that take care of the internal structure of a word have been shown to give similar or even improved performance over traditional word embeddings as they are able to encode morphological [dos Santos and Zadrozny2014] or semantic [Chen et al.2015b] information even for out-of-vocabulary items [Ling et al.2015]. This follows the simple idea that since a word can be split up into its respective morphemes that encode morphosyntactic and semantic information, encoding subparts of a word instead of a the whole word should allow for better generalization for unseen words without the need to compute and store vector representations for every word type.

4 Embedding Models

4.1 Neural Embedding Models

Models that make use of cooccurrence count matrix factorizations usually suffer from the need of storing huge amounts of data into working memory. Neural networks with hidden layers, on the other hand, have been proven to be able to approximate any given hidden function while they still allow for online or mini-batch training, which is a huge improvement over models that only work with whole data batches. Moreover they allow for fine-grained objective functions, where additional data can be introduced into the network to enhance the amount of information stored in the embeddings.

4.1.1 Neural Language Models

[Bengio et al.2003a] pioneered as the first to introduce Neural Language Models in general. Here, a network with two hidden layers is proposed, where a first non-linear hidden layer with locally shared weights gives embeddings of single words in the given context window, while the next layer is a concatenation of the single word embeddings and another non-linear activation function and gives a context embedding. The last layer than is a softmax classifier that predicts the next word (cf. Figure 1). It later became clear that this model had two major drawbacks:

- The softmax classifier needs to compute the probability of all word types in the vocabulary. This is a major flaw also of following models (see below).

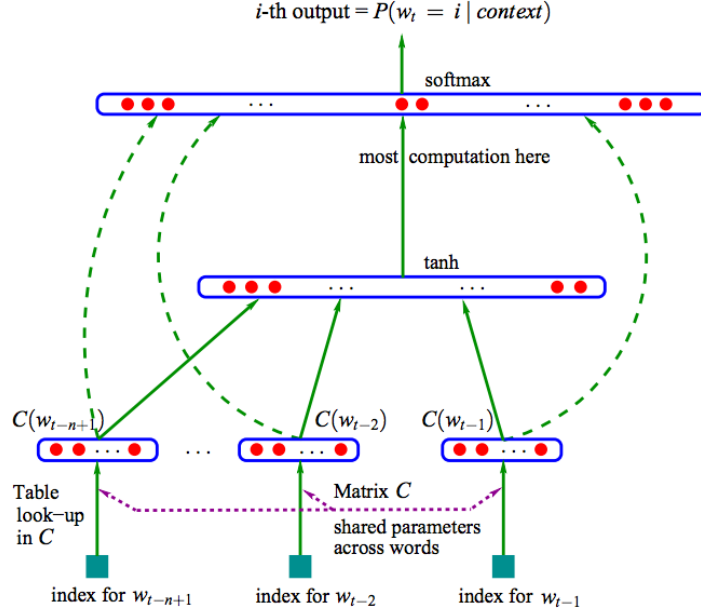


Figure 1: The neural language model as proposed in [Bengio et al.2003a].

- The non-linear hidden layer severely worsened the convergence rate of the network.

[Collobert and Weston2008] then improved the model tremendously by getting rid of the expensive softmax classifier by introducing a score based loss function:

$$J_{\theta} = \sum_{x \in X} \sum_{w \in V} \max(0, 1 - f_{\theta}(x) + f_{\theta}(x^{(w)})) \quad (1)$$

Here, the the correct windows x containing n words are sampled from the set of all possible windows X in the corpus, while for every window x , a corrupted version $x^{(w)}$ is produced by replacing the center word of x by the another word of the vocabulary V . The objective then becomes maximizing the distance between the scores output by the model fore correct and the incorrect window with a margin of 1. The actual model then was still the same as in [Bengio et al.2003a]. This model could already embed semantically similar words close to each other, but still suffered from the non-linear hidden layer, which slowed down learning enormously ([Bengio et al.2003a] report a training time of 7 weeks for a vocabulary of 130000 word types)

4.1.2 Word2Vec

Among general-purpose embedding models in general and among neural embedding models in particular, the word2vec algorithm [Mikolov et al.2013a,Mikolov et al.2013b] stands out

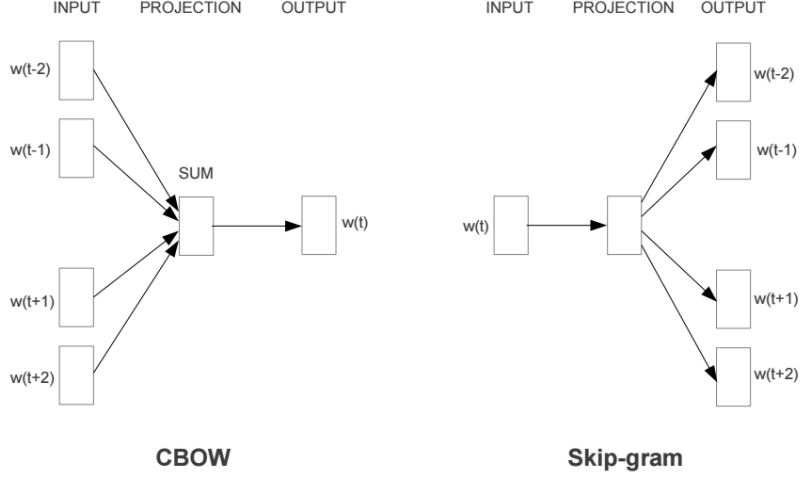


Figure 2: The two common architectures of word2vec. The Continuous Bag-of-Words (CBOW) model predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word. From [Mikolov et al.2013a].

as major breakthrough, as it offers good performance while the parameters to train are less than with the previous models. Different to the previous models in [Bengio et al.2003a, Collobert and Weston2008], which employ non-linear hidden layers, here the model only consists of a linear encoder layer that computes the actual embeddings. It can be seen as a full departure from traditional embedding models as a subgroup of language models towards a new kind of model family that in the first place tries to encode semantic and morphosyntactic information of a word rather than to predict the next word itself. Instead, the context window usually encompasses all surrounding words of a given target word instead of just the previous words. This makes it rather unhandy for true language modeling, but allows for much better incorporation of latent information. Moreover, the model can be trained in two ways: It either tries to predict the target word given its context (Continuous Bag-of-Words; CBOW) or tries to predict the context of a word given the target word itself (Skip-Gram, cf. Figure 2). The corresponding loss function for the CBOW model is given as

$$J_{\theta} = \frac{1}{T} \sum_{t=1}^T \log p(w_t \mid w_{t-n}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+n}) \quad (2)$$

where $p(w_t \mid w_{t-1}, \dots, w_{t-n+1})$ is given as a softmax

$$p(w_t \mid w_{t-1}, \dots, w_{t-n+1}) = \frac{\exp(h^{\top} v'_{w_t})}{\sum_{w_i \in V} \exp(h^{\top} v'_{w_i})} \quad (3)$$

while the loss for the Skip-Gram model is given as

$$J_\theta = \frac{1}{T} \sum_{t=1}^T \sum_{-n \leq j \leq n, j \neq 0} \log p(w_{t+j} \mid w_t) \quad (4)$$

where $p(w_{t+j} \mid w_t)$ is again defined as a softmax

$$\log p(w_{t+j} \mid w_t) = \frac{\exp(v_{w_t}^\top v'_{w_{t+j}})}{\sum_{w_i \in V} \exp(h^\top v'_{w_i})} \quad (5)$$

where v_{w_i} is an input vector and v'_{w_i} an output vector representation of a word w_i and h^\top the transpose of the sum of the output vectors of the words in the context. Again, the calculation of the softmax here includes summing over all contexts in the corpus, which computationally is expensive. Further research largely concentrated on finding a way to define p as a memory- and time-friendly approximation of the true softmax probability. A first solution was given by using a hierarchical softmax classifier [Morin and Bengio2005] instead. Here, the classifier layer is a binary tree with the actual words as leaves. At every node, the network learns to either follow the left or the right branch with a certain probability as in Eq. 5 that is equal to the sum of the child elements of the respective branch, but is actually computed as the product of h^\top and the output vector of the word at node n pulled through a logistic sigmoid:

$$p(right|n, c) = \sigma(h^\top v'_n) \quad (6)$$

This means in order to calculate the softmax probability of word, we only have to follow the path down to the word leaf instead of summing over all vocabulary entries. For binary trees, this means we only have to pass at most $\log_2(|V|)$ nodes to calculate the probability of a word, which is a huge performance boost over the traditional softmax classifier.

Another way of approximating the softmax is using Differentiated Softmax (cf. [Chen et al.2015a]), where the idea is that rare words should help to adjust less parameters than more frequent words. The fully connected softmax layer becomes a sparse matrix here, where certain blocks inside of that matrix are used to predict words with a certain frequency. [Chen et al.2015a] mention their model is fast and accurate, but has less predictive power with rare words, as they are given less parameters to be predicted.

Another family of models approximate the true softmax via sampling. Here, it seems to be common to reformulate the softmax definitions from above. For instance, if we use negative log probabilities instead of raw probabilities and if we further use a more general score function $\mathcal{E}(w)$ instead of the more specific dot products of the context and word vectors and start with Eq. 3, the loss function J_θ becomes

$$J_\theta = \mathcal{E}(w) + \log \sum_{w_i \in V} \exp(-\mathcal{E}(w_i)) \quad (7)$$

For back-propagating the error, we calculate the gradients of this term:

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) + \nabla_{\theta} \log \sum_{w_i \in V} \exp(-\mathcal{E}(w_i)) \quad (8)$$

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) + \frac{1}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \nabla_{\theta} \sum_{w_i \in V} \exp(-\mathcal{E}(w_i)) \quad (9)$$

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) + \frac{1}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \sum_{w_i \in V} \nabla_{\theta} \exp(-\mathcal{E}(w_i)) \quad (10)$$

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) + \frac{1}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \sum_{w_i \in V} \exp(-\mathcal{E}(w_i)) \nabla_{\theta} (-\mathcal{E}(w_i)) \quad (11)$$

which can be rewritten as

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) + \sum_{w_i \in V} \frac{\exp(-\mathcal{E}(w_i))}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \nabla_{\theta} (-\mathcal{E}(w_i)) \quad (12)$$

Here, we can replace the softmax in the second term by $P(w_i)$ and reposition the negative coefficient:

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) - \sum_{w_i \in V} P(w_i) \nabla_{\theta} (\mathcal{E}(w_i)) \quad (13)$$

As one can see, we basically end up with a positive and a negative reinforcement term. As the negative reinforcement term is basically the gradient of $\mathbb{E}_{w_i \sim P}$, the final reformulation of the softmax is then given as in Eq. 14.

$$\nabla_{\theta} J_{\theta} = \nabla_{\theta} \mathcal{E}(w) - \mathbb{E}_{w_i \sim P} [\nabla_{\theta} \mathcal{E}(w_i)] \quad (14)$$

The goal of sampling-based approaches is now to approximate the negative reinforcement term $\mathbb{E}_{w_i \sim P} [\nabla_{\theta} \mathcal{E}(w_i)]$ without going over the whole vocabulary. Earlier models used some forms of Importance Sampling [Liu2008, Bengio et al.2003b, Bengio and Sen  cal2008, Cho et al.2015], where a proposal distribution Q is used to approximate P . Although Importance sampling seems to provide a speed-up over the standard softmax, the model suffers from divergence issues. Another group of sampling models does not try to estimate the probability of a word directly but uses an auxiliary loss that as a by-product produces probability estimates of correct words. Noise Contrastive Estimation (NCE, [Gutmann and Hyv  rinen2010, Mnih and Teh2012]) trains a model to differentiate between a true word w_i and noise data \tilde{w}_{ik} , where k is the number of noisy data points. As a start point, one can use a logistic regression loss that minimized the negative log-likelihood as in Eq 15, where Q is a noise distribution:

$$J_\theta = - \sum_{w_i \in V} [\log P(y = 1|w_i, c_i) + k \mathbb{E}_{\tilde{w}_{ik} \sim Q} [\log P(y = 0|\tilde{w}_{ik}, c_i)]] \quad (15)$$

Since we want to avoid summing over all words in the vocabulary to estimate $\mathbb{E}_{\tilde{w}_{ik} \sim Q}$, we substitute it with the mean as a Monte Carlo estimate, which leaves us with the sum over all log probabilities of the k noisy words:

$$J_\theta = - \sum_{w_i \in V} [\log P(y = 1|w_i, c_i) + \sum_{j=1}^k \log P(y = 0|\tilde{w}_{ik}, c_i)] \quad (16)$$

As w_i is sampled from the true Distribution P while \tilde{w}_{ik} comes from the noise distribution Q , this leaves us with a mixture model, where the probability of sampling a specific word with the corresponding label (that is, whether it comes from the true distribution in the training set or is noise) can be given as

$$P(y, w|c) = \frac{1}{k+1} P(w|c) + \frac{k}{k+1} Q(w) \quad (17)$$

This can then be used to estimate the probability of the word to come from the true distribution:

$$P(y = 1|w, c) = \frac{\frac{1}{k+1} P(w|c)}{\frac{1}{k+1} P(w|c) + \frac{k}{k+1} Q(w)} \quad (18)$$

which can be simplified to

$$P(y = 1|w, c) = \frac{P(w|c)}{kP(w|c) + Q(w)} \quad (19)$$

Here, instead of taking the traditional softmax for $P(w|c)$, NCE allows for a reparametrization of it: instead of summing over all probabilities in the lexicon to normalize the data, the denominator of the softmax is substituted by a parameter $Z(c)$ that is estimated by the model itself. In fact, and report equal performance while setting $Z(c)$ to 1, that is, the model learns to normalize the probabilities completely on its own. Coming from the definition of this probability in the CBOW architecture from Eq. 3, this leaves us with the probability of a word given a context as

$$P(w|c) = \exp(h^\top v'_w) \quad (20)$$

Inserting this into Eq. 19, we have

$$P(y = 1|w, c) = \frac{\exp(h^\top v'_w)}{\exp(h^\top v'_w) + kQ(w)} \quad (21)$$

which can be inserted into the logistic loss from Eq. 16 to yield the final NCE loss:

$$J_\theta = - \sum_{w_i \in V} [\log \frac{\exp(h^\top v'_{w_i})}{\exp(h^\top v'_{w_i}) + kQ(w_i)} + \sum_{j=1}^k \log(1 - \frac{\exp(h^\top v'_{\tilde{w}_{ij}})}{\exp(h^\top v'_{\tilde{w}_{ij}}) + kQ(\tilde{w}_{ij})})] \quad (22)$$

The probably most prominent sampling model, however, is the Negative Sampling Model (NEG). Building on NCE, it diverges from it insofar as it naively assumes that $kQ(w) = 1$. This yields Eq. 21 as

$$P(y = 1|w, c) = \frac{\exp(h^\top v'_w)}{\exp(h^\top v'_w) + 1} \quad (23)$$

which is equal to the NCE loss as long as $k = |V|$ and Q is uniform. Eq. 23 can then be reformulated as a logistic sigmoid:

$$P(y = 1|w, c) = \frac{1}{1 + \exp(-h^\top v'_w)} \quad (24)$$

Inserted into the logistic regression loss, this yields

$$J_\theta = - \sum_{w_i \in V} [\log \frac{1}{1 + \exp(-h^\top v'_w)} + \sum_{j=1}^k \log(1 - \frac{1}{1 + \exp(-h^\top v'_{\tilde{w}_{ij}})})] \quad (25)$$

By simplification and substitution of the sigmoid term with $\sigma = \frac{1}{1 + \exp(-x)}$, this derives the NEG loss:

$$J_\theta = - \sum_{w_i \in V} [\log \sigma(h^\top v'_w) + \sum_{j=1}^k \log \sigma(-h^\top v'_{\tilde{w}_{ij}})] \quad (26)$$

For a more in-depth analysis of NEG, see [Goldberg and Levy2014].

As NEG is only equivalent to NCE when $k = |V|$ holds and Q is uniform, this has a large impact on the performance of NEG when used as a proper language model. As it cannot approximate the true softmax in practice, it gives poor performance on typical intrinsic language model evaluation tasks such as perplexity estimation, while it shows superior performance on embedding-specific evaluations such as similarity measures.

4.1.3 Polyglot Neural Language Models

One of the reasons why word2vec has become the most popular neural embedding model (and probably the most popular embedding model in general) is that it does not involve any hidden layers that would affect the training speed performance. However, with the rise of Long Short-term Memory Networks (LSTM, [Hochreiter and Schmidhuber1997]), there is an increase of research on embeddings trained jointly with LSTMs to predict

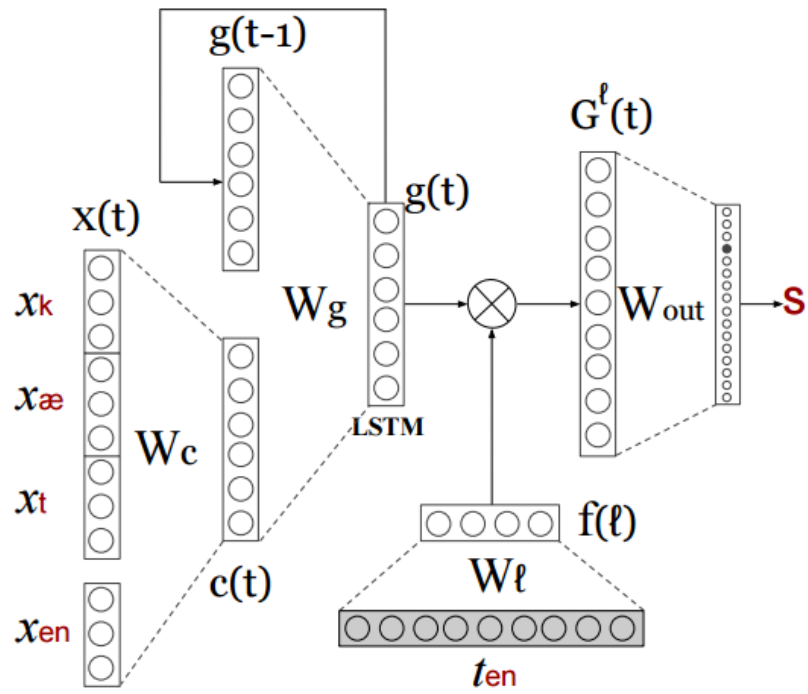


Figure 3: The Polyglot Language Model. From [Tsvetkov et al.2016]

tokens or whole sequences [Chen et al.2015b, Ling et al.2015, Kim et al.2015, dos Santos and Zadrozny2014], while some learn them jointly with convolutional models [Zhang et al.2015]. One family of models important for our case are Polyglot Neural Language Models (PLM) [Tsvetkov et al.2016]. Here, a deep model tries to predict the next phoneme given the previous phonemes, the language of the whole sequence as well as typological information from WALS [Dryer and Haspelmath2013], PHOIBLE [Moran et al.2014] and Ethnologue [Lewis et al.2015] (see Fig. 3). The model first embeds the whole context into a hidden local context representation c_t

$$c_t = W_{xx}x_t + W_{clang}x_{lang} + b_c \quad (27)$$

where x_t is a character embedding for character c_t at position t , x_{lang} is a one-hot vector for the respective language of the sequence and b_c a bias for the character c . This then is fed into an LSTM to yield a global context layer g_t

$$g_t = LSTM(c_t, g_{t-1}) \quad (28)$$

Typological language information (such as if a single segment is represented in the phoneme inventory of a given language, or if there is a minimal contrast between the segment and another one in the phoneme inventory) is first mapped into a low-dimensional space to obtain a dense representation f_ℓ

$$f_\ell = \tan(W_{\ell t_\ell} + b_\ell) \quad (29)$$

and then multiplied with the global context vector using Hadamard multiplication, yielding a global context-language matrix G_t^ℓ

$$G_t^\ell = g_t \otimes f_\ell^\top \quad (30)$$

The resulting matrix is then vectorized into a column vector and an output sequence is computed:

$$p(\phi_t = i | \phi_i, \dots, \phi_{t-1}, \ell) = softmax(W_{out}vec(G_t^\ell) + b_{out})_i \quad (31)$$

Accordingly, the loss function here is a categorical cross-entropy $H(\phi, \hat{\phi})$ over all phonemes:

$$H(\phi, \hat{\phi}) = - \sum_i \hat{\phi}_i \log \phi_i \quad (32)$$

The authors report superior perplexity over a baseline LSTM model, even with the typological information left out. Also downstream experiments on loan word detection and speech synthesis leave the PLM model superior over the baseline LSTM model.

5 Evaluation

5.1 Data

For this paper, the data provided by the Automatic Similarity Judgement Program (ASJP,) will serve as a training corpus. ASJP provides wordlists of roughly 7200 languages, where each word is given as a string of phonemes. This vast amount of different phoneme inventories will allow for the unbiased comparison of phonemes over language boundaries. However, ASJP does not make use of IPA, but a simplified phonemic alphabet that is based on ASCII-characters. This makes it possible to incorporate data where no detailed IPA description is given, but might also reduce the quality of the resulting embedding space, as latent features that might be encoded in more specific IPA descriptions are merged. An overview of the ASJP phonemes and their IPA counterparts are given in the appendix.

5.2 Intrinsic Evaluation

Embeddings are commonly evaluated using *intrinsic* and *extrinsic* tests. Intrinsic evaluation methods allow for the direct comparison of different models of embedding algorithms. Here the underlying assumption is that the embedding space itself should encode certain information. For word embeddings, this usually accumulates in relatedness, analogy, categorization or selectional preference tests [Schnabel et al.2015]. They all have in common that the predicted outcome given a certain model is compared to a gold standard. The problem in our case is that for phoneme embeddings, there is hardly any pre-defined gold standard available, especially not for ASJP data. For this reason I will restrict the intrinsic evaluation part to analogy tests of the sort of those made popular by [Mikolov et al.2013b]. I hand-crafted an analogy table of 108 analogy tasks where each line is of the form

$$v(\textit{phoneme}_1) + v(\textit{phoneme}_2) - v(\textit{phoneme}_3) \approx v(\textit{phoneme}_4) \quad (33)$$

where $v(\textit{phoneme}_1)$ is the vector corresponding to $\textit{phoneme}_1$ and $v(\textit{phoneme}_2) - v(\textit{phoneme}_3)$ can be seen as the latent phonological information available in $\textit{phoneme}_2$ but not in $\textit{phoneme}_3$, while $v(\textit{phoneme}_1) + v(\textit{phoneme}_2) - v(\textit{phoneme}_3)$ is then this latent phonological information added to $\textit{phoneme}_1$, which should be close to the vector corresponding to $\textit{phoneme}_4$. For example, in

$$v(/i/) + v(/u^*/) - v(/u/) \approx v(/i^*/) \quad (34)$$

$v(/u^*/) - v(/u/)$ describes a latent feature that should correspond to [+NASALIZED], which is then added to the phoneme $/i/$ to yield a nasalized $/i^*/$. The latent features should correspond to accepted phonological distinctive features such as [VOICE], [LABIAL] or [ASPIRATED]. For this task it is important that the target phoneme should only differ in one phonological feature from the source phoneme to allow for a detailed evaluation of the generalization abilities of the model in the latent feature space. Most tests concentrate

on plosives and vowels, as they should appear most frequently in the data. The full list of analogy tasks can be found in the appendix. Here, accuracy serves as a test metric as the target phoneme embedding should be the at most n -closest embedding to the predicted point in the embedding space. The parameter n here serves as a threshold to allow for the fact that a prediction should still be good if the target phoneme is at least close to the predicted phoneme. To compare the performance of the different embedding models, an extensive grid search over several hyperparameters was performed:

- Model architecture: CBOW, Skip-gram
- Classifier: Softmax, negative sampling
- Embedding dimensions: [2, 5, 10, 20, 50, 100, 150, 200, 400]
- Context window size: [1, 2, 3, 4, 5, 8]
- Number of negative samples: [1, ..., 21]

An analogy task is thought to be For each hyperparameter setting, an ensemble of three models was trained simultaneously. Each model iterated five times over the corpus, which consists of all words from the ASJP corpus given the language they belong to is neither artificial nor reconstructed. Every word is moreover bracketed by word boundary tokens.

Figure 4 and 5 show that the models have severe problems to capture the compositional structure of the phonemes. However, in order to check if the model at least learns to cluster related phonemes together, another grid search was performed with an analogy threshold $n = 20$ that allows for less accurate predictions still to be valid. The model indeed seems to learn to capture latent phonological features, even though the compositional quality of the learned features seems to be quite bad.

For a more detailed evaluation, an ensemble of 10 models (CBOW with hierarchical softmax, 100 embedding dimensions, context window of 2 phonemes) was trained on the corpus (Figure 8). Analogy tasks with a threshold of 1 yield reasonable results for vowels, while for consonants the model’s accuracy is far worse. The performance for complex phonemes is even worse. This might be connected to an underrepresentation of those complex phonemes in the corpus, which does not let the model generalize well. In fact, the quality of the phoneme embeddings is correlated with their frequency in the corpus (Pearson correlation 0.9, Spearman correlation 0.9, Figure 10). This leads to the question if one should not indeed follow the example of others [Rama2016, Jäger2014] and pool coarticulations and non-pulmonic phonemes with the plain pulmonic ones to allow for at least reasonable analogy performance over the whole corpus. Indeed, pooling complex articulations together with plain consonants produces better results (Figure 9), although they are still quite improvable.

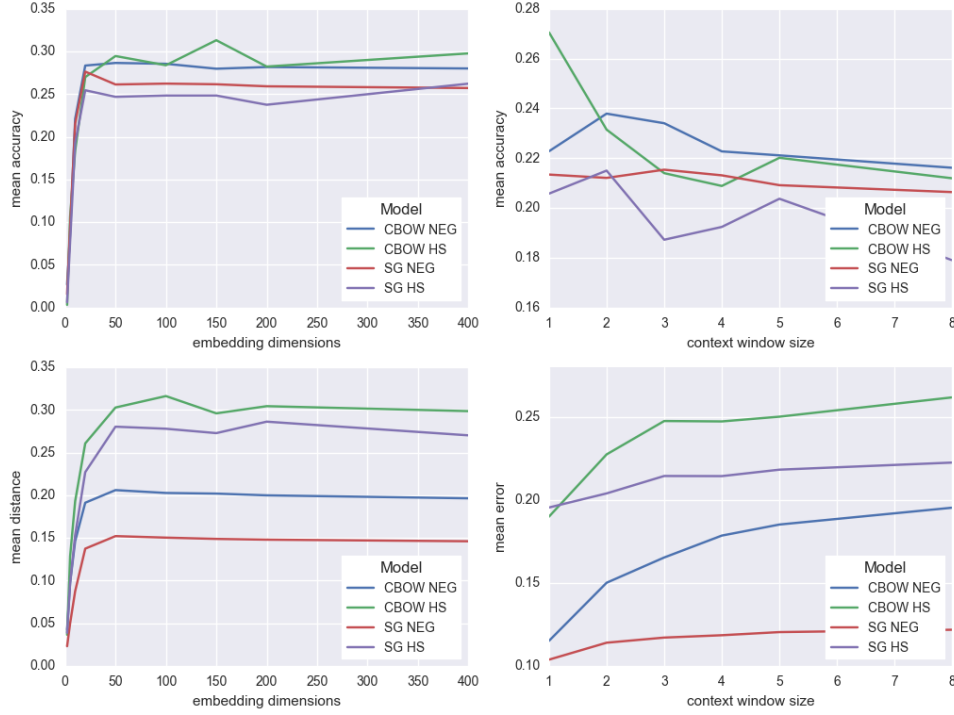


Figure 4: Comparison of the four word2vec models evaluated. (top left) Both CBOW Models perform better than the skip-gram models over all numbers of embedding dimension. (bottom left) Mean distance between the predicted vector and the target with regard to embedding dimensions. Here, negative sampling yields less error than hierarchical softmax. (top right) Mean accuracy for the models with regard to the context window size. The models perform worse the bigger the context is. (bottom right) Mean distance between the predicted vector and the target with regard to context window size. Again, bigger contexts lead to worse predictions.

5.3 Discussion

Simple visual inspection shows that the model learn to cluster phonemes for certain phonological features. However, the analogy tests have shown that the models have problems with encoding the compositional information expected for phonemes. Predictions are often near to the expected phoneme, though, as Figure ?? might suggest. This is probably caused by the size of the training data that is too small, as the correlation between the frequency of a phoneme class and mean accuracy of analogy tests in its respective domain

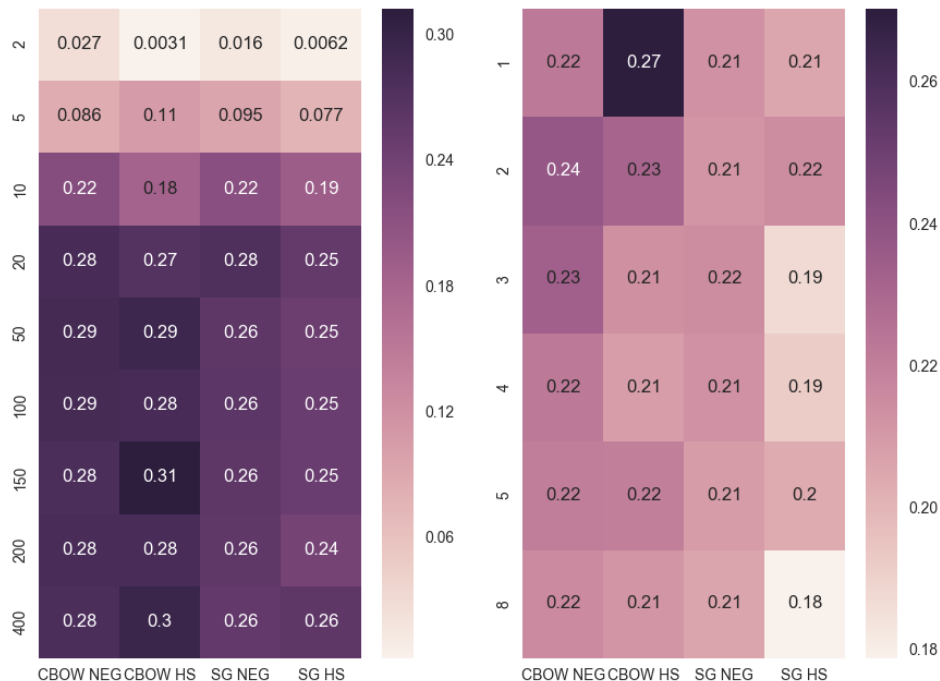


Figure 5: Mean accuracies of the four models tested, over the number of embedding dimensions (left) and the context window size (right)

might suggest. Getting rid of complex articulations and pooling the respective consonants with plain ones indeed improves the analogy test performance. Vowel phonemes, given their relative small type size over the whole corpus, perform much better than consonants, which could be a further indicator that the number of data is simply not enough.

Another factor is the quality of the data itself. Sometimes a given phoneme is simply wrong (e.g. all words in the Ukrainian word list have final /o/ where they should have /a/) or only stems are given when whole words are expected, which clearly worsens the models' ability to learn the distinction between phonemes that rather occur initially or finally (e.g. the Italian word list has *bere* for "to drink" but *ved* for "to see", which even stands against the phonotactics of the language). In general, a lot of latent information is lost with the deliberate merger of voiced and unvoiced phonemes for lots of phoneme classes in the data set.

Another final reason for the mediocre performance might finally be that in general phonemes might behave differently than words with regards to their compositionality in

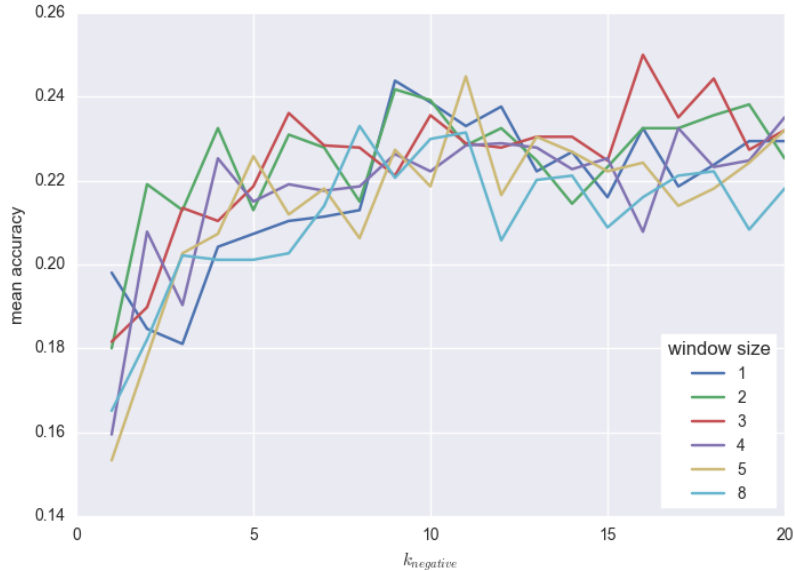


Figure 6: Mean and maximal accuracies of the skip-gram model with regard to the window size and number of negative samples over all other parameters.

the latent feature space. The models indeed learn to cluster phonemes with regard to certain phonological features only using the context of the respective phonemes, but the question leaves open if the context of a phoneme is enough to clearly identify a phoneme beyond language boundaries.

5.4 Phoneme Inventory Clustering

6 Resume

While the theoretical benefits of phoneme embeddings might look appealing, the

References

- [Bengio et al.2003a] Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003a). A neural probabilistic language model. *journal of machine learning research*, 3(Feb):1137–1155.
- [Bengio and Senécal2008] Bengio, Y. and Senécal, J.-S. (2008). Adaptive importance sampling to accelerate training of a neural probabilistic language model. *IEEE Transactions on Neural Networks*, 19(4):713–722.

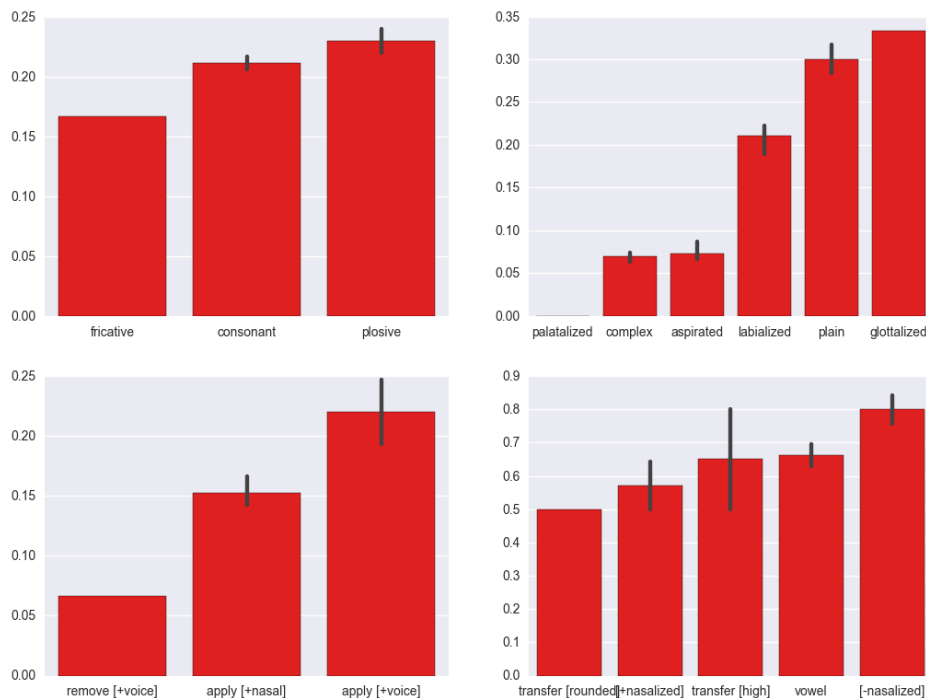


Figure 7: Mean accuracies for analogies with an analogy threshold of 1. The baseline model does not seem to be able to capture the compositionality of the latent features too well. (top left) Accuracy is best for plosives across plain pulmonic as well as more complex articulations, while fricatives perform worse. (top right) Glottalized and Plain pulmonic consonant phonemes yield best performance. Among complex articulations, which all perform bad, labialized phonemes yield best results, while aspirated and palatalized phonemes perform even worse. (bottom left) Adding voice works better than removing it or adding nasality. (bottom right) Vowel analogies work far better than consonant analogies. This should be due to the small number of possible vowel phonemes.

[Bengio et al.2003b] Bengio, Y., Sen  cal, J.-S., et al. (2003b). Quick training of probabilistic neural nets by importance sampling. In *AISTATS*.

[Chen et al.2015a] Chen, W., Grangier, D., and Auli, M. (2015a). Strategies for training large vocabulary neural language models. *arXiv preprint arXiv:1512.04906*.

[Chen et al.2015b] Chen, X., Xu, L., Liu, Z., Sun, M., and Luan, H. (2015b). Joint learning of character and word embeddings. In *Proceedings of IJCAI*, pages 1236–1242.

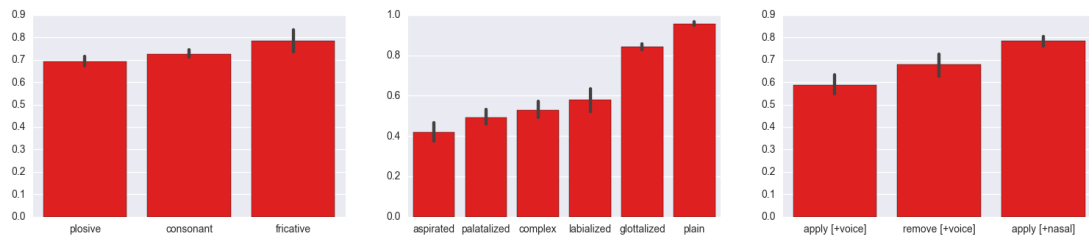


Figure 8: Mean accuracies for consonant analogies with an analogy threshold of 20. Here, the model is expected to yield better results, as wrong predictions are not as penalizing. (left) Accuracy is reasonable for plosives across plain pulmonic as well as more complex articulations. For fricatives, the accuracy is higher, but fricatives were only evaluated for plain pulmonic articulation. (middle) Plain pulmonic consonant phonemes nearly yield optimal accuracy, while glottalized phonemes still yield good results. Among complex articulations, labialized phonemes yield best results, while aspirated and palatalized phonemes perform far under average. (right) Removing voice yields better results than adding it. Applying nasality to plosives even performs better.

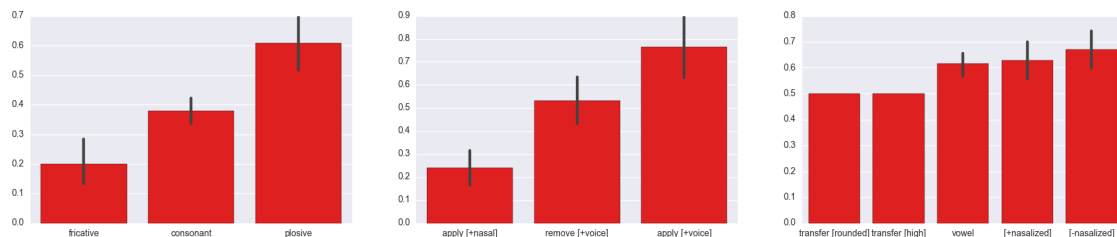


Figure 9: Mean accuracies for analogies with an analogy threshold of 1 on the pooled consonant phonemes. Here, the model is also expected to yield better results, the number of latent features should be reduced. (left) Accuracy is doubled compared with the unpooled phonemes. (middle) Adding voice works quite well, while removing it still yields acceptable performance. Applying nasality again works quite bad. (right) Vowel analogy tasks seem to work reasonably well.

[Cho et al.2015] Cho, S. J. K., Memisevic, R., and Bengio, Y. (2015). On using very large target vocabulary for neural machine translation.

[Collobert and Weston2008] Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM.

[dos Santos and Zadrozny2014] dos Santos, C. N. and Zadrozny, B. (2014). Learning

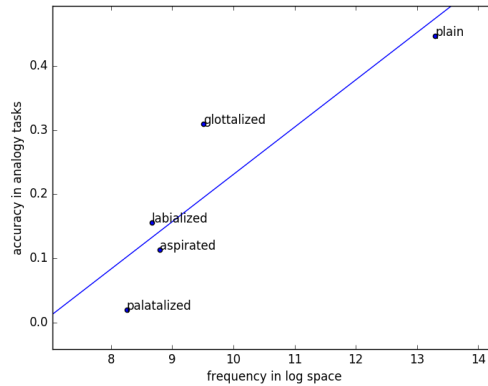


Figure 10: Frequency plotted against mean accuracy in log space. It is clearly visible how the bad performance of the complex phonemes is due to their underrepresentation in the corpus.

character-level representations for part-of-speech tagging. In *ICML*, pages 1818–1826.

[Dryer and Haspelmath2013] Dryer, M. S. and Haspelmath, M., editors (2013). *WALS Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.

[Goldberg and Levy2014] Goldberg, Y. and Levy, O. (2014). word2vec explained: deriving mikolov et al.’s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*.

[Gutmann and Hyvärinen2010] Gutmann, M. and Hyvärinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, volume 1, page 6.

[Hochreiter and Schmidhuber1997] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

[Jäger2014] Jäger, G. (2014). Phylogenetic inference from word lists using weighted alignment with empirically determined weights. In *Quantifying Language Dynamics*, pages 155–204. Brill.

[Kim et al.2015] Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2015). Character-aware neural language models. *arXiv preprint arXiv:1508.06615*.

[Kiros et al.2015] Kiros, R., Zhu, Y., Salakhutdinov, R. R., Zemel, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Skip-thought vectors. In *Advances in neural information processing systems*, pages 3294–3302.

- [Kondrak2000] Kondrak, G. (2000). A new algorithm for the alignment of phonetic sequences. In *Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference*, pages 288–295. Association for Computational Linguistics.
- [Lewis et al.2015] Lewis, M. P., Simons, G. F., and Fennig, C. D. (2015). *Ethnologue: Languages of the world*. dallas, texas: Sil international. online version.
- [Ling et al.2015] Ling, W., Luís, T., Marujo, L., Astudillo, R. F., Amir, S., Dyer, C., Black, A. W., and Trancoso, I. (2015). Finding function in form: Compositional character models for open vocabulary word representation. *arXiv preprint arXiv:1508.02096*.
- [Liu2008] Liu, J. S. (2008). *Monte Carlo strategies in scientific computing*. Springer Science & Business Media.
- [Mikolov et al.2013a] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- [Mikolov et al.2013b] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- [Mnih and Teh2012] Mnih, A. and Teh, Y. W. (2012). A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*.
- [Moran et al.2014] Moran, S., McCloy, D., and Wright, R., editors (2014). *PHOIBLE Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- [Morin and Bengio2005] Morin, F. and Bengio, Y. (2005). Hierarchical probabilistic neural network language model. In *Aistats*, volume 5, pages 246–252. Citeseer.
- [Pennington et al.2014] Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–43.
- [Rama2016] Rama, T. (2016). Siamese convolutional networks based on phonetic features for cognate identification. *arXiv preprint arXiv:1605.05172*.
- [Schnabel et al.2015] Schnabel, T., Labutov, I., Mimno, D., and Joachims, T. (2015). Evaluation methods for unsupervised word embeddings. In *Proc. of EMNLP*.
- [Tsvetkov et al.2016] Tsvetkov, Y., Sitaram, S., Faruqui, M., Lample, G., Littell, P., Mortensen, D., Black, A. W., Levin, L., and Dyer, C. (2016). Polyglot neural language models: A case study in cross-lingual phonetic representation learning. *arXiv preprint arXiv:1605.03832*.

[Zhang et al.2015] Zhang, X., Zhao, J., and LeCun, Y. (2015). Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, pages 649–657.

7 Appendix

Table 1: Consonant analogy tests. All analogy tasks consist of 2 positive values, one negative value and a target value, as described in 5.2.

voiceless plosive → nasal via unvoiced plosive		
positive	negative	target
p, n	t	m
t, m	p	n
k, n	t	N
voiced plosive → nasal via voiced plosive		
b, n	d	m
d, m	b	n
g, n	d	N
plain voiced plosive → plain nasal via plain unvoiced plosive		
b, n	t	m
d, m	p	n
g, n	t	N
plain unvoiced plosive → plain voiced plosive via plain voiced plosive		
p, d	t	b
t, b	p	d
k, d	t	g
plain voiced plosive → plain unvoiced plosive via plain unvoiced plosive		
b, t	d	p
d, p	b	t
g, t	d	k
labialized voiced plosive → labialized nasal via labialized unvoiced plosive		
bw~, nw~	tw~	mw~
dw~, mw~	pw~	nw~
gw~, nw~	tw~	Nw~
labialized unvoiced plosive → labialized voiced plosive via labialized voiced plosive		
pw~, dw~	tw~	bw~
tw~, bw~	pw~	dw~
kw~, dw~	tw~	gw~
labialized voiced plosive → labialized unvoiced plosive via labialized unvoiced plosive		
bw~, tw~	dw~	pw~
dw~, pw~	bw~	tw~
gw~, tw~	dw~	kw~
palatalized voiced plosive → palatalized nasal via palatalized unvoiced plosive		
by~, ny~	ty~	my~
dy~, my~	py~	ny~
gy~, ny~	ty~	Ny~
palatalized unvoiced plosive → palatalized voiced plosive via palatalized voiced plosive		
py~, dy~	ty~	by~
ty~, by~	py~	dy~
ky~, dy~	ty~	gy~
palatalized voiced plosive → palatalized unvoiced plosive via palatalized unvoiced plosive		
by~, ty~	dy~	py~
dy~, py~	by~	ty~
gy~, ty~	dy~	ky~
aspirated voiced plosive → aspirated nasal via aspirated unvoiced plosive		
bh~, nh~	th~	mh~
dh~, mh~	ph~	nh~
gh~, nh~	th~	Nh~
aspirated unvoiced plosive → aspirated voiced plosive via aspirated voiced plosive		
ph~, dh~	th~	bh~
th~, bh~	ph~	dh~
kh~, dh~	th~	gh~
aspirated voiced plosive → aspirated unvoiced plosive via aspirated unvoiced plosive		
bh~, th~	dh~	ph~
dh~, ph~	bh~	th~
gh~, th~	dh~	kh~
glottalized voiced plosive → glottalized nasal via glottalized unvoiced plosive		
b", n"	t"	m"
d", m"	p"	n"
g", n"	t"	N"
glottalized unvoiced plosive → glottalized voiced plosive via glottalized voiced plosive		
p", d"	t"	b"
t", b"	p"	d"
k", d"	t"	g"
glottalized voiced plosive → glottalized unvoiced plosive via glottalized unvoiced plosive		
b", t"	d"	p"
d", p"	b"	t"
g", t"	d"	k"

Table 2: Consonant analogy tests (continued)

unvoiced pulmonic plosive → unvoiced aspirated plosive via unvoiced aspirated plosive		
positive	negative	target
p, th~	t	ph~
t, ph~	p	th~
k, th~	t	kh~
voiced pulmonic plosive → voiced aspirated plosive via voiced aspirated plosive		
b, dh~	d	bh~
d, bh~	b	dh~
g, dh~	d	gh~
unvoiced plain plosive → unvoiced palatalized plosive via unvoiced palatalized plosive		
p, ty~	t	py~
t, py~	p	ty~
k, ty~	t	ky~
voiced plain plosive → voiced palatalized plosive via voiced palatalized plosive		
b, dy~	d	by~
d, by~	b	dy~
g, dy~	d	gy~
unvoiced plain plosive → unvoiced labialized plosive via unvoiced labialized plosive		
p, t"	t	p"
t, p"	p	t"
k, t"	t	k"
voiced plain plosive → voiced labialized plosive via voiced labialized plosive		
b, d"	d	b"
d, b"	b	d"
g, d"	d	g"
unvoiced pulmonic plosive → unvoiced glottalized plosive via unvoiced glottalized plosive		
p, t"	t	p"
t, p"	p	t"
k, t"	t	k"
check voiced pulmonic plosive → voiced glottalized plosive via voiced glottalized plosive		
b, d"	d	b"
d, b"	b	d"
g, d"	d	g"
POA transitions for unvoiced plosives		
p, d	b	t
t, b	d	p
k, d	g	t
POA transitions for voiced plosives		
b, t	p	d
d, p	t	b
g, t	k	d
POA transitions for nasals		
m, d	b	n
n, b	d	m
N, d	g	n
POA transitions for unvoiced fricatives		
f, t	p	s
s, p	t	f
x, t	k	s
X, k	q	x
POA transitions for voiced fricatives		
v, d	b	z
z, b	d	v

Table 3: Vowel analogy tests

apply [+ROUNDED]		
positive	negative	target
i, o	e	u
e, u	i	o
apply [+HIGH]		
i, o	u	e
e, u	o	i
apply [+NASAL]		
i, u*	u	i*
e, o*	o	e*
u, i*	i	u*
o, e*	e	o*
a, o*	o	a*
E, e*	e	E*
3, E*	E	3*
apply [-NASAL]		
i*, u	u*	i
e*, o	o*	e
u*, i	i*	u
o*, e	e*	o
a*, o	o*	a
E*, e	e*	E
3*, E	E*	3

Table 4: ASJP consonants with their IPA equivalents. Glottalization is marked by a following double quote, coarticulation is indicated by a tilde (if one co-articulator) or a dollar sign \$ (if two co-articulators). The tables follow [Jäger2014].

ASJP code symbol	Description	IPA equivalent
p	voiceless bilabial stop and fricative	b, β
b	voiced bilabial stop and fricative	f
f	voiceless labiodental fricative	v
v	voiced labiodental fricative	m
m	bilabial nasal	m
w	voiced bilabial-velar approximant	w
8	voicedless and voiced dental fricative	θ, ð
4	dental nasal	$\underset{~}{n}$
t	voiceless alveolar stop	t
d	voiced alveolar stop	d
s	voiceless alveolar fricative	s
z	voiced alveolar fricative	z
c	voiceless and voiced alveolar affricate	\widehat{ts} , \widehat{dz}
n	alveolar nasal	n
r	rhotics	r, ɾ, ɽ, R
l	voiced alveolar lateral approximant	l
S	voiceless post-alveolar fricative	ʃ
Z	voiced post-alveolar fricative	ʒ
C	voiceless palato-alveolar affricate	$\widehat{tʃ}$
j	voiced palato-alveolar affricate	$\widehat{dʒ}$
T	voiceless and voiced palatal stop	c, ɟ
5	palatal nasal	ɲ
y	palatal approximant	j
k	voiceless velar stop	k
g	voiced velar stop	g
x	voiceless and voiced velar fricative	x, ɣ
N	velar nasal	ŋ
q	voiceless uvular stop	q
G	voiced uvular stop	q
X	voiceless and voiced uvular/pharyngeal fricatives	χ, ʁ, ħ, ʕ
h	voiceless and voiced glottal fricative	h, ħ
7	glottal stop	ʔ
L	all other laterals	L, ɭ, ʎ
!	clicks	ʘ, ɓ, ɗ, ʄ

Table 5: ASJP vowels with their IPA equivalents. Nasality is marked by a following asterisk.

ASJP code symbol	Description	IPA equivalent
i	high front vowel, rounded and unrounded	ɪ, ɪ̃, y, ʏ
e	mid front vowel, rounded and unrounded	e, ø
E	low front vowel, rounded and unrounded	æ, ɛ, œ, œ̃
3	high and mid central vowel, rounded and unrounded	ɨ, ə, ə̃, ɜ, ɵ, ɔ̃, ɵ̃, ɐ
a	low central vowel, unrounded	a, ɐ̃
u	high back vowel, rounded and unrounded	ʊ, u
o	mid and low back vowel, rounded and unrounded	ʊ̃, ʌ, ɑ, ɒ, ɔ̃, ɒ̃