

Towards a Universal Semantic Dictionary

Maria Jose Castro-Bleda ^{1,*} , Eszter Iklódi ²  and Gábor Borbély ²

¹ VRain Valencian Research Institute for Artificial Intelligence, Universitat Politècnica de València, 46022 Valencia, Spain

² Budapest University of Technology and Economics, 1111 Budapest, Hungary; eszter.iklodi@gmail.com (E.I.); recski.gabor@aut.bme.h (G.R.); borbely@math.bme.hu (G.B.)

* Correspondence: mcastro@dsic.upv.es; Tel.: +34-96-387-7351

Received: 24 June 2019; Accepted: 24 September 2019; Published: 28 September 2019

Abstract: A novel method for finding linear mappings among word embeddings for several languages, taking as pivot a shared, multilingual embedding space, is proposed in this paper. Previous approaches learned translation matrices between two specific languages, while this method learns translation matrices between a given language and a shared, multilingual space. The system was first trained on bilingual, and later on multilingual corpora as well. In the first case, two different training data were applied: Dinu’s English–Italian benchmark data, and English–Italian translation pairs extracted from the PanLex database. In the second case, only the PanLex database was used. The system performs on English–Italian languages with the best setting significantly better than the baseline system given by Mikolov, and it provides a comparable performance with more sophisticated systems. Exploiting the richness of the PanLex database, the proposed method makes it possible to learn linear mappings among an arbitrary number of languages.

Keywords: natural language processing; semantics; word embeddings; multilingual embeddings; translation; artificial neural networks

1. Introduction

Computer-driven natural language processing plays an increasingly important role in our everyday life. In the current digital world, using natural language for human–machine communication has become a basic requirement. To meet this requirement, it is inevitable to analyze human languages semantically. Nowadays, state-of-the-art systems represent word meaning with high dimensional vectors, known as word embeddings.

Current embedding models are learned from monolingual corpora, and therefore infer language dependency. However, one might ask if the structure of the different embeddings, i.e., different meaning representations, are universal among all human languages. Youn et al. [1] proposed a procedure for building graphs from concepts of different languages. They found that these graphs reflected a certain structure of meaning with respect to the languages they were built of. They concluded that the structural properties of these graphs are consistent across different language groups, and largely independent of geography, environment, and the presence or absence of literary traditions. Such findings led to a new research direction within the field of computational semantics, which focuses on the construction of universal meaning representations, most of the times in the form of cross-lingual word embedding models [2]. One way to create such models is to find mappings between embeddings of different languages [3–5]. These mappings are usually transformation matrices that serve to map an embedding from one language space to another. In most of the cases, one pivot space is chosen, typically English, where every other language will be mapped into. In our work, we do not chose any language to serve as a pivot space, but instead, we create a language independent space, i.e., a universal space, to where we map all languages.

Our work proposes a novel procedure for learning such mappings in the form of translation matrices that serve to map each language to a shared space. The method was first tested on bilingual, and later on multilingual corpora as well. With the bilingual experiments, we obtained on Dinu's benchmark data [6] a 0.377 precision@1 score for English–Italian and a 0.310 precision@1 score for Italian–English translation. These results, however, are far from the current state-of-the-art result on this dataset [5], but they are in the same order of magnitude or even better than many previous attempts [3,6,7]. For further bilingual and for some multilingual experiments, an own dataset was created from the PanLex database [8]. We published the obtained scores of various experimental settings using this dataset [9]. Generally, bilingual experiments using only the PanLex dataset resulted in worse scores than using only Dinu's dataset, but combining the two showed a slight improvement in the Italian–English direction. Multilingual experiments were carried out using three different languages—English, Italian, and Spanish—at the same time. The obtained pairwise precision values showed worse results than when the system was trained in bilingual mode. However, these results are still promising considering that a completely new approach was implemented, and they showed that the system definitely learned from data which are available for a wide range of languages. To reach state-of-the-art performance, the system has to deal with dimension reduction in the multilingual space and further experimentation in multilingual mode with an extended number of languages could also provide meaningful outputs.

Section 2 summarizes the progress made on learning translation matrices between word embeddings over the last five years. Section 3 discusses the proposed method in detail. Following that, Section 4 describes our experimental setup, and Sections 5 and 6 report and analyze the obtained results. Finally, Section 7 concludes with the advantages and disadvantages of the proposed model, and also discusses some improvements for future work.

2. Related Work

2.1. Word Embeddings

One way to build semantic representations is to use distributional models. The idea is based on the observation that synonyms or words with similar meanings tend to occur in similar contexts, or, as it was phrased by Firth in 1957: "You shall know a word by the company it keeps" [10]. For example, in the following two sentences "The cat is walking in the bedroom" and "A dog was running in a room" words such as "dog" and "cat" have exactly the same semantic and grammatical roles, therefore we could easily imagine the two sentences in the following variations: "The dog is walking in the bedroom" and "A cat was running in a room" [11]. Based on this intuition, what distributional models are aiming to do is to compute the meaning of a word from the distribution of words around it [12]. The obtained meaning representations are usually high dimensional vectors, called word embeddings, which refer to their characteristic feature that they model a world by embedding it into a vector space.

2.2. Monolingual Word Embeddings

Mikolov et al. [13] suggested a Bag-of-words Neural Network, more specifically two architectures, for learning monolingual word embeddings. The first one, denoted as the Continuous Bag-of-Words Model (CBOW), tried to predict the current word based on the context, whereas the second one, denoted as the continuous skip-gram model, tried to maximize the classification of a word based on another word in the same sentence. The CBOW turned out to be slightly better on syntactic tasks and the skip-gram on semantic tasks. Mikolov's procedure has become known as the *word2vec* (<http://deeplearning4j.org/word2vec>) procedure.

2.3. Multilingual Word Embeddings

In 2013, Mikolov et al. [3] published a simple two-step procedure for creating universal embeddings. In the first step, they built monolingual models of languages using huge corpora,

and, in the second step, a small bilingual dictionary was used to learn linear projection between the languages. The optimization problem was the following:

$$\min_W \sum_{i=1}^n \|Wx_i - z_i\|^2 \quad (1)$$

where W denotes the transformation matrix, and $\{x_i, z_i\}_{i=1}^n$ are the continuous vector representations of word translation pairs, with x_i being in the source language space and z_i in the target language space.

Faruqui and Dyer [7] proposed a procedure to obtain multilingual word embeddings by concatenating the two word vectors coming from the two languages, applying Canonical Correlation Analysis. Xing et al. [14] found that bilingual translation can be largely improved by normalizing the embeddings and by restricting the transformation matrices into orthogonal ones. Dinu et al. [6] showed that the neighborhoods of the mapped vectors are strongly polluted by hubs, which are vectors that tend to be near a high proportion of items. They proposed a method that computes hubness scores for target space vectors and penalizes those vectors that are close to many words, i.e., hubs are down-ranked in the neighboring lists. Lazaridou et al. [15] studied some theoretical and empirical properties of a general cross-space mapping function, and tested them on cross-linguistic (word translation) and cross-modal (image labeling) tasks. They also introduced the use of negative samples during the learning process. Amar et al. [16] proposed methods for estimating and evaluating embeddings of words in more than fifty languages in a single shared embedding space. Since English usually offers the largest corpora and bilingual dictionaries, they used the English embeddings to serve as the shared embedding space. Artetxe et al. [17] built a generic framework that generalizes previous works made on cross-linguistic embeddings and they concluded that the best systems were the ones with orthogonality constraint and a global pre-processing with length normalization and dimension-wise mean centering. Smith et al. [4] also proved that translation matrices should be orthogonal, for which they applied Singular Value Decomposition (SVD) on the transformation matrices. Besides, they also introduced a novel “inverted softmax” method for identifying translation pairs. All these works listed above applied supervised learning. However, in 2018, Lample et al. [5] introduced an unsupervised way for aligning monolingual word embedding spaces between two languages without using any parallel corpora. Their method consisted of an adversarial approach for aligning two monolingual word embedding spaces, followed by a refinement step using frequent aligned words (according to the adversarial mapping). This unsupervised procedure holds the current state-of-the-art results on Dinu’s benchmark word translation task.

3. Proposed Method

We propose a method that learns linear mappings between word translation pairs in the form of translation matrices. For each language, exactly one translation matrix is learned. During the training, the cosine similarity of word translation pairs is maximized; the similarity is calculated after applying the linear mapping of the translation matrices being learned. Since this space, where the translation matrices map the original languages, is not equal to any of the original language spaces, we refer to it as a multilingual or universal space. This space keeps changing during the training process, as the translation matrices become more and more optimized. At test time, the system is evaluated with the precision metric, principally used for word translation tasks.

3.1. Cosine Similarity and Precision

Cosine similarity (https://en.wikipedia.org/wiki/Cosine_similarity) is a measure of similarity between two non-zero vectors. It is calculated as the normalized dot product of two vectors, as shown in Equation (2). In fact, cosine similarity is a space that measures the cosine of the angle of two vectors. It is important to note that cosine similarity is not a proper distance metric, since the triangle inequality property does not apply. In word similarity tasks, however, this metric is used for measuring the

similarity of two words represented as word vectors. Although cosine similarity values by definition are in range of $[-1, 1]$, in word similarity tasks, it is particularly used in positive space, $[0, 1]$, where parallel vectors are similar and orthogonal vectors are dissimilar.

$$\text{cos_sim} = \cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \cdot \|\vec{b}\|} \quad (2)$$

Precision is a metric used for measuring the performance of translator systems, which intend to learn to translate from a source language into a target language. On the target side, a look-up space is defined, which could, for example, correspond to the most frequent 200K words of the target language, as in our experiments. After translating a word, the N word vectors of the look-up space that are closest to the translated one are regarded. The Precision @ N metric denotes the percentage of how many times the real translation of a word is found among the N closest word vectors in the look-up space. Usual N values are 1, 5, and 10.

3.2. The Objective Function

The objective of the proposed method is to learn linear mappings in the form of translation matrices that are obtained by maximizing the cosine similarity of gold word translation pairs in the above defined multilingual space. Therefore, for each language, one single translation matrix is searched that maps the language from its original vector space to the multilingual one. The method tries to bring the translation pairs as close together as possible in the multilingual space. Therefore, it is not only applicable for language pairs but for any number of languages as well. The main advantage is that by introducing new languages the number of the learned parameters remains linear to the number of languages since instead of learning pair-wise translation matrices, for each language, only one matrix is learned, the one that maps directly to the shared, multilingual space.

Let L be a set of languages, and TP a set of translation pairs where each entry is a tuple of two in the form of (w_1, w_2) where w_1 is a word in language L_1 and w_2 is a word in language L_2 , and both L_1 and L_2 are in L . Then, let us consider the following equation to optimize:

$$\frac{1}{|TP|} \cdot \sum_{\substack{L_1, L_2 \\ \in L}} \sum_{\substack{(w_1, w_2) \\ \in TP}} \text{cos_sim}(w_1 \cdot T_1, w_2 \cdot T_2) \quad (3)$$

where T_1 and T_2 are translation matrices mapping L_1 and L_2 to the multilingual space. Since the equation is normalized with the number of translation pairs in the TP set, the optimal value of this function is 1. Off-the-shelf optimizers are programmed to find local minimum values, thus, during the training process, the loss function is multiplied by -1 . Word vectors are always normalized, thus the cos_sim reduces to a simple dot product.

At test time, first, both source and target language words are mapped into the multilingual space, and from the most frequent 200K mapped target language words a look-up space is defined. Then, the system is evaluated with the Precision metric, more specifically with Precision @1, @5, and @10. The distance assigned to the word vectors when searching in the look-up space is the cos_sim .

Previous works, such as by Mikolov et al. [4] or Lample et al. [5], suggested restricting the transformation matrix to an orthogonal one. From an arbitrary transformation matrix T , an orthogonal T' can be obtained by applying the SVD procedure. Our experiments showed that, by applying SVD on the transformation matrices, the learning is significantly faster. The best results were obtained when applying the SVD only once, at the beginning of the learning process.

4. Experimental Setup

4.1. Pre-Trained Word Embeddings

For pre-trained word embeddings, we took the *fastText* embeddings published by Lample et al. [5]. These embeddings were trained by applying their novel method where words are represented as a bag of character n-grams [18]. This model outperformed Mikolov's [13] CBOW and skipgram baseline systems that did not take any sub-word information into account. The pre-trained word vectors trained on Wikipedia are available for 294 languages (<http://github.com/facebookresearch/fastText>).

Some experiments were also run by using the same embedding that was used by Dinu et al. [6] in their experiments. These word vectors were trained with *word2vec* and then the 200K most common words in both the English and Italian corpora were extracted. The English word vectors were trained on the WackyPedia/ukWaC and BNC corpora, while the Italian word vectors were trained on the WackyPedia/itWaC corpus. This word embedding is referred to as the *WaCky* embedding.

4.2. English–Italian Setup of Dinu

Dinu et al. [6] constructed an English–Italian gold dictionary split into a training and a test set that is now being used as benchmark data for evaluating English–Italian word translation tasks. Both training and test translation pairs were extracted from a dictionary built from Europarl Eng-Ita (<http://opus.lingfil.uu.se/> [19]).

For the test set, they used 1500 English words split into five frequency bins, 300 randomly chosen in each bin. The bins were defined in terms of rank in the frequency-sorted lexicon: [1K–5K], [5K–20K], [20K–50K], [50K–100K], and [100K–200K]. Some of these 1500 English words had multiple Italian translations in the Europarl dictionary, thus the resulting test set contained 1869 word pairs all together, with 1500 different English and 1849 different Italian words (see Table 1).

For the training set, the above-mentioned Europarl dictionary was first sorted by the English frequency. Then, the top 5K entries were extracted and care was taken to avoid any overlap with the test elements on the English side. On the Italian side, however, an overlap of 113 words was still present. In the end, the training set contained 5K word pairs with 3442 different English, and 4549 different Italian words (see Table 1).

Table 1. Statistics of word counts.

Set	Language	No. Words
Train (5000 word pairs)	Eng.	3442
	Ita.	4549
Test (1869 word pairs)	Eng.	1500
	Ita.	1849

4.3. The PanLex Corpus

PanLex [8] is a nonprofit organization that aims to build a multilingual lexical database from available dictionaries in all languages. The name PanLex is coming from the words *panlingual* and *lexical*, which reflect the main objective of this project. They are basically digitizing and centering the content of different, already existing dictionaries made by domain experts. Own translations are not accepted. To each translation pair, a confidence value is assigned, which can be used for filtering the extracted data. These confidence values are in the range of [1, 9], with 9 meaning high and 1 meaning low confidentiality. The main purpose is to preserve the diversity of languages, so the collection of “threatened” or “endangered” languages and dictionaries of rare language combinations are top priority. Some examples of the English–Italian PanLex data can be seen in Table 2.

Table 2. Sample of PanLex entries of the extracted tsv file.

English	Italian	Confidence Values
Sarajevo	Sarajevo	9
euro	euro	9
simple	semplice	8
difficult	difficile	8
college	università	7
plausible	verisimile	7
sea	mare	6
sky	cielo	6
better	migliore	5
inform	informare	5
combustible	combustibile	4
office	ufficio	4
sorcerer	conscitore	3
it	ella	3
Great Wall of China	Grande muraglia cinese	2
factory workers	lavoratori dell'industria	2
stay	restare	1
sometimes	qualche volta	1

PanLex also exhibits different *language varieties* that include, among others, regional variations and different writing systems. A *language variety* is denoted with a three-letter *language code*, e.g., eng for English, and with a three-digit *variety code*, e.g., 000. To the most widely spoken variety of a language, usually the 000 *variety code* is assigned. When extracting data from the PanLex database, in all cases, the *language variety* with the smallest *variety code* was taken.

4.4. Dataset Creation from PanLex

The procedure applied for extracting a proper data from the PanLex database for training multilingual embedding models roughly follows the same steps as in [6]. After extracting the raw translation pairs from the PanLex database, a filtered version of entries was formed by dropping translations with a confidence value below 7 and those for which no word vector was found in the *fastText* embedding. This resulted in an English–Italian word translation set containing 69,623 entries.

For the test set, 1500 English words were taken and split into five frequency bins, 300 randomly assigned to each bin. The bins were defined as in [6], i.e., in terms of rank in the frequency-sorted lexicon: [1K–5K], [5K–20K], [20K–50K], [50K–100K], and [100K–200K]. In [5], the word vectors sorted by their frequency in descending order were published, and this order was used as the source of English word frequency data. In the PanLex database it is a common issue that one English word has sometimes as many as 10 different Italian translations. Therefore, to avoid having an undesirably huge test set with many Italian synonyms only those English words were selected, for which in the corresponding bin only one Italian translation was present. This way the obtained test set contains exactly 1500 word pairs, which are made up of 1500 different English words and their Italian translations.

For the training set, the 69,623 entries were first sorted by their English frequency, then the top 5K entries were extracted and, as in [6], care was taken to avoid any overlap with test elements on the English side. Then, the top 5K entries were selected in three different ways:

1. The first 5K entries were taken.
2. The first 5K different English words were taken with the most frequent Italian translation.
3. Only those English words were taken for which only one Italian translation was present.

4.5. Baseline Experimental Setting

For the baseline system, the *fastText* embedding was used as a pre-trained embedding and the system was trained on Dinu’s English–Italian data. For parameter adjustment, Dinu’s training data were split into train and validation sets such that no overlap was present on the English side, i.e., no word appeared in both sets; this follows Dinu’s procedure of constructing their original training and test sets. It should be noted that this does not apply for Italian words. For the word count and overlap statistics of Dinu’s original training and test sets, see Table 3, and for the same statistics of the newly produced training and validation sets, see Table 4.

Table 3. Statistics of the original train and test split of Dinu’s data.

Number of English words	train	3442
Number of Italian words		4549
Number of English words	test	1500
Number of Italian words		1849
Overlap English		0
Overlap Italian		113

Table 4. Statistics of the new train and validation split of Dinu’s data.

Number of English words	train	3098
Number of Italian words		4129
Number of English words	valid	344
Number of Italian words		499
Overlap English		0
Overlap Italian		80

The system was adjusted on the previously described training and validation split. For the optimizer, the tensorflow implementation (https://www.tensorflow.org/api_docs/python/tf/train/AdagradOptimizer) of the Adagrad algorithm [20] was used. For evaluation the most frequent 200K words of the target space embedding were used as look-up space for calculating Precision @1, @5, and @10. In all cases, both English–Italian and Italian–English precision scores were observed. In addition, the average cosine similarity value of the validation set was also checked. During training and validation the precision and similarity values were also all calculated in the multilingual space. Gold dictionaries were constructed from the input data files themselves. Following Dinu, any word appearing in the dictionary was considered a valid translation. Various translations may come from synonyms or different male–female forms on the Italian side.

5. Experimental Results

5.1. Parameter Adjustment Using Dinu’s Data

First, parameter adjustment was performed using Dinu’s data, which gave 0.1 as the best learning rate and 64 as the best batch size, where batch size is equal to the number of translation pairs used in one iteration. With applying SVD only once at the beginning the obtained results of our best system are significantly worse than state-of-the-art results on this benchmark data, but they are comparable with or even better than some of the previous models discussed in Section 2.

5.2. Experimenting with SVD

Previous works, (e.g., [4,5]) suggested restricting the transformation matrix to an orthogonal one. Based on these findings, this system also features a configuration option of applying an SVD. Three different SVD modes were studied:

- 0: Not using SVD at all
- 1: Using SVD after every n th epoch
- 2: Using SVD only once, at the beginning

In the following experiments, the same datasets were used as for parameter adjustment. Learning rate was set to 0.1 and batch size to 64, as found the best setup before. Altogether 200 epochs were done and evaluation was performed on every 10th epoch.

5.2.1. Not Using SVD

This experiment was carried out without applying any SVD. Translation matrices were initialized with random numbers. Figure 1 shows that similarity values are monotone increasing, meaning that the system is learning. However, the learning process is relatively slow since even after 200 epochs the similarity score is still quite low, bearing in mind that the optimal value is 1.0.

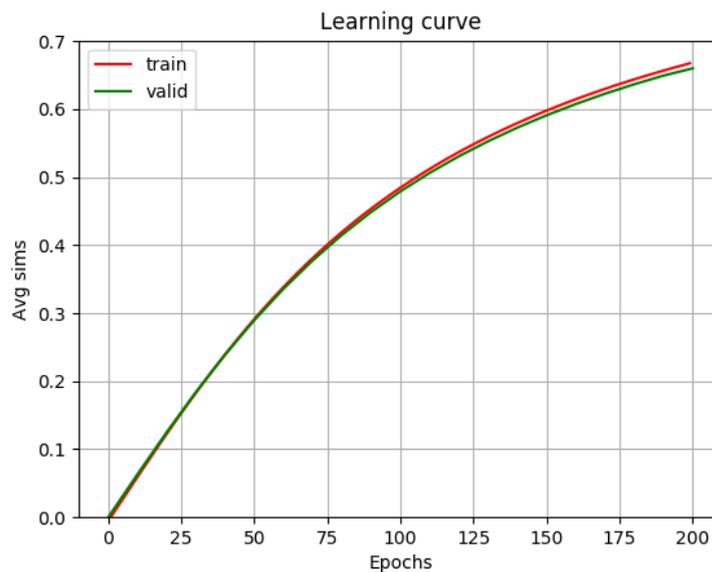


Figure 1. Learning curve of experimenting without using SVD (`svd_mode = 0`).

5.2.2. SVD after Every n th Epoch

This experiment was carried out applying SVD several times over the whole learning process. SVD was made on every 50th epoch, i.e., four times altogether. Figure 2 shows how the learning curve breaks down every time after applying an SVD on the translation matrices, and, also, how fast it is back once again to the previous high similarity values. Besides, this time the average cosine similarity score was higher even at the beginning than it was after 200 epochs with the previous setting, where no SVD was done. Applying SVD on the transformation matrices seems to accelerate the learning process significantly. The learning curve also shows that SVD-to-SVD fractions have exactly the same trajectory regardless of the number of previous epochs done.

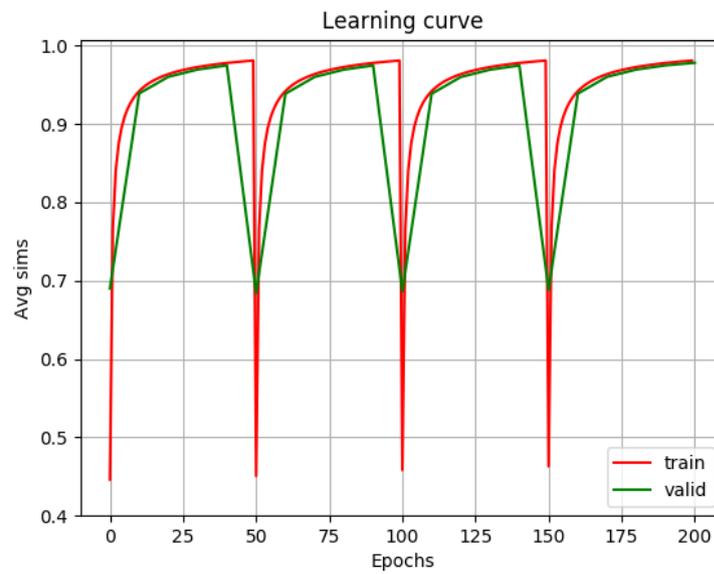


Figure 2. Learning curve of experimenting with SVD after every n th epoch ($svd_mode = 1$).

5.2.3. SVD at the Beginning

This experiment was carried out applying SVD only once, at the very beginning. This means, in simple terms, that instead of a random initial transformation matrix, the system tried to adjust an orthogonal one. Figure 3 shows that the learning curve is monotone increasing, and owing to the initial SVD it gets fairly high right at the beginning.

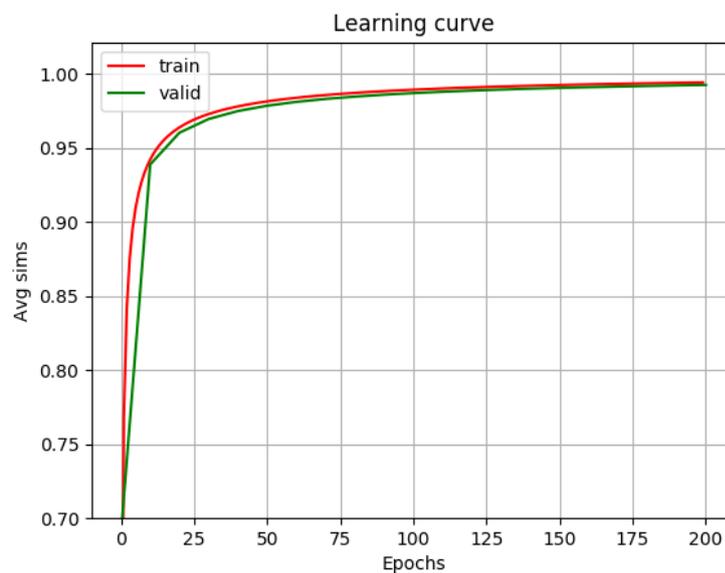


Figure 3. Learning curve of experimenting with SVD at the beginning ($svd_mode = 2$).

5.3. Experiments with the PanLex Data

5.3.1. Comparing Different Dataset Construction Methods

Tables 5 and 6 compare the results of different dataset construction methods. It is important to note that in the first case the English word of the 5000th translation pair is only the 845th most frequent English word, meaning that there are only 845 different English words in the training set and that,

on average, there are 5–6 different Italian translations for each of them. In the second case, where every English word is kept but only with the most frequent Italian translation, this number is 9007. In the last case, however, the 5000th entry is made up of the 39,426th most frequent English and the 31,543rd most frequent Italian words. However, this last training set provides the best results, thus for further experiment this construction method was applied.

Table 5. English–Italian precision values with the different training sets.

Precision	@1	@5	@10
First 5K entry	0.0093	0.0253	0.0367
First 5K English words with retaining one translation	0.1120	0.2073	0.2427
First 5K English words with one translation	0.1960	0.3087	0.3440

Table 6. Italian–English precision values with the different training sets.

Precision	@1	@5	@10
First 5K entry	0.0000	0.0007	0.0007
First 5K English words with retaining one translation	0.1114	0.2052	0.2440
First 5K English words with one translation	0.1838	0.3059	0.3443

5.3.2. Experimenting with Different Training Set Sizes

Table 7 summarizes the results of experiments with different training set sizes. The 3K dataset proved to be the best on the English–Italian translation, but on the Italian–English it is only slightly better, than the 5K dataset. This behavior of performing better on the smaller training sets is fairly understandable since, as a consequence of the way the training set was constructed, as we are taking in more and more entries, we are actually taking in less and less frequent English words and their Italian translations, for which words neither the embeddings nor the translations are precise enough. Since Dinu’s benchmark data contain 5K entries in the training set, despite the slightly worse performance, we kept using the 5K dataset for the sake of comparability with other results.

Table 7. Experiments with different training set sizes.

Prec.	Eng-Ita			Ita-Eng		
	@1	@5	@10	@1	@5	@10
1K	0.1500	0.2847	0.3340	0.1391	0.2761	0.3256
3K	0.2127	0.3473	0.3933	0.2232	0.3650	0.4152
5K	0.1980	0.3193	0.3620	0.2212	0.3555	0.4030
10K	0.1613	0.2807	0.3227	0.1879	0.3012	0.3372

5.4. Comparison of Systems Trained on Dinu’s and PanLex Data

In the next step, some experiments were made to determine which data are more apt for learning linear mappings between embeddings. To compare all experiments objectively, subsets of the original test sets were created. These subsets do not contain any English word present either in the Dinu training set or in the PanLex training set. Table 8 summarizes the number of word pairs in the old and new test sets. It should be noted that by this reduction mainly the most common English words are affected, and therefore worse scores are expected compared to the previous train-on-Dinu-test-on-Dinu, or train-on-PanLex-test-on-PanLex top results. Scores on Dinu’s test set are shown in Table 9 and on the PanLex data in Table 10. The obtained results show that training on the PanLex data cannot beat the system trained on Dinu’s data, which performs better both on Dinu’s and on the PanLex test sets. Not even combining the two training sets succeeds in achieving significantly better results, although on the PanLex test set it does improve the scores in the Italian–English direction.

Table 8. Word reduction of the new test sets.

Test set	No. Word Pairs in Old	No. Word Pairs in New
Dinu	1869	1455
PanLex	1500	1242

Table 9. Comparing Dinu’s and PanLex data on Dinu’s test set.

Precision	Eng-Ita			Ita-Eng		
	@1	@5	@10	@1	@5	@10
Train:Dinu-Test:old	0.3770	0.5647	0.6245	0.3103	0.5018	0.5474
Train:Dinu-Test:new	0.3560	0.5407	0.5978	0.2917	0.4792	0.5215
Train:PanLex-Test:new	0.1360	0.2309	0.2594	0.1361	0.2556	0.2965
Train:Dinu+PanLex-Test:new	0.2930	0.4349	0.4861	0.2910	0.4556	0.5090

Table 10. Comparing Dinu’s and PanLex data on the PanLex test set.

Precision	Eng-Ita			Ita-Eng		
	@1	@5	@10	@1	@5	@10
Train:PanLex-Test:old	0.1960	0.3087	0.3440	0.1838	0.3059	0.3443
Train:PanLex-Test:new	0.1812	0.2858	0.3196	0.1668	0.2835	0.3213
Train:Dinu-Test:new	0.2295	0.4171	0.4839	0.2227	0.3763	0.4199
Train:Dinu+PanLex-Test:new	0.2295	0.3712	0.4275	0.2498	0.4026	0.4495

5.5. Continuing the Training with PanLex Data

Another experiment was conducted to continue the baseline system trained on Dinu’s data with the PanLex data. In other words, it is the same as initializing the translation matrices of the PanLex training process with previously learned ones. The baseline system reaches its best performance between 2000 and 4000 epochs, depending on which precision value is regarded. Table 11 shows that on the English–Italian task there is no improvement at all, while on the Italian–English task with the best setting slightly better scores are achieved on precision @1 and @10 values.

Table 11. Continuing the baseline system with the PanLex data.

Precision	Eng-Ita			Ita-Eng		
	@1	@5	@10	@1	@5	@10
Original	0.3770	0.5647	0.6245	0.3103	0.5018	0.5474
Cont. from 2000	0.3426	0.5256	0.5802	0.3229	0.4882	0.5535
Cont. from 3000	0.3535	0.5416	0.5970	0.3229	0.4840	0.5465
Cont. from 4000	0.3510	0.5273	0.5911	0.3118	0.4701	0.5243

5.6. Experiments Using Three Languages

Finally, a multilingual experiment was carried out where the system was trained on three languages—English, Italian, and Spanish—at the same time. During training, the system learns three different translation matrices, one for English–multilingual, one for Italian–multilingual, and one for Spanish–multilingual space mapping. For example, to learn the English–multilingual translation matrix, both the English–Italian and the English–Spanish dictionaries are used, according to Equation (3). Batches are homogeneous, but two following batches are always different in terms of the language origins of the contained data. That is, first an English–Italian batch is fed to the system, then an English–Spanish batch, after that an Italian–Spanish batch, and so on. First, bilingual models were trained to compare them later with the multilingual system. The results of the bilingual models are summarized in Table 12. The results are best on the Italian–Spanish task. Next, the system was trained

using all three languages at the same time. During the training process, the model was evaluated on the bilingual test datasets, for which the results are shown in Table 13. The obtained results show that no advantage was achieved by extending the number of languages, since the multilingual model performs worse than any of the pairwise bilingual models.

Table 12. Results of bilingual models trained pairwise on the three different languages.

Precision	L1-L2			L2-L1		
	@1	@5	@10	@1	@5	@10
Eng-Ita	0.2080	0.3280	0.3687	0.2082	0.3386	0.3904
Eng-Spa	0.2840	0.4320	0.4800	0.2883	0.4331	0.4836
Spa-Ita	0.3920	0.5340	0.5813	0.3655	0.5291	0.5750

Table 13. Bilingual results of the multilingual model trained using three different languages at the same time.

Precision	L1-L2			L2-L1		
	@1	@5	@10	@1	@5	@10
Eng-Ita	0.1573	0.2667	0.3127	0.1638	0.2942	0.3386
Eng-Spa	0.1947	0.2973	0.3447	0.2350	0.3538	0.4064
Spa-Ita	0.2520	0.3640	0.4160	0.2568	0.3723	0.4162

6. Comparison of the Experiments

Tables 14 and 15 show our results on Dinu’s dataset compared to other published works. Our results are worse than those of the current state-of-the-art, but they are still comparable or even better than several previous attempts. The advantage of the proposed method compared to other procedures is that it is applicable for an arbitrary number of languages at the same time. Although the multilingual experiments on the PanLex dataset showed worse results than the bilingual ones, they are still showing convergence and can serve as a baseline for future multilingual experiments.

Table 14. Comparing English–Italian results on Dinu’s data.

Eng-Ita	@1	@5	@10
Mikolov et al. (2013) [3]	0.338	0.483	0.539
Faruqui et al. (2014) [7]	0.361	0.527	0.581
Dinu et al. (2014) [6]	0.385	0.564	0.639
Smith et al. (2017) [4]	0.431	0.607	0.651
Lample et al. (2018) [5]	0.662	0.804	0.834
Proposed method	0.377	0.565	0.625

Table 15. Comparing Italian–English results on Dinu’s data.

Ita-Eng	@1	@5	@10
Mikolov et al. (2013) [3]	0.249	0.410	0.474
Faruqui et al. (2014) [7]	0.310	0.499	0.570
Dinu et al. (2014) [6]	0.246	0.454	0.541
Smith et al. (2017) [4]	0.380	0.585	0.636
Lample et al. (2018) [5]	0.587	0.765	0.809
Proposed method	0.310	0.502	0.547

7. Conclusions and Future Work

This paper proposes a novel method for finding linear mappings between word embeddings in different languages. As a proof of concept, a framework was developed which enabled basic parameter adjustments and flexible configuration for initial experimentation.

An interesting finding was that the system learned much more quickly when an initial SVD was applied on the translation matrices. The results obtained with these settings on Dinu's data show that the proposed model did learn from the data. The obtained precision scores, although far from current state-of-the-art results on this benchmark data, were comparable with the results of previous attempts. The proposed model performed much better using the *fastText* embeddings [5], than using Dinu's WaCky embeddings [6].

Thereafter, an English–Italian dataset was extracted from the PanLex database, from which training and test datasets were constructed roughly following the same steps that Dinu et al. [6] took. The system was trained and tested on both Dinu's and PanLex test sets, and in both cases the matrices trained on Dinu's data were the ones reaching higher scores. On the PanLex data experiments with different training set sizes were executed, out of which the 3K training set gave the best results. Continuing the training of the matrices obtained by using Dinu's data with the PanLex dataset brought a slight improvement on the Italian–English scores, but English–Italian scores only got worse.

Finally, the system was trained on three different languages at the same time. The obtained pairwise precision values proved to be worse than the results obtained when the system was trained in bilingual mode. However, these results are still promising considering that a completely new approach was implemented, and they show that the system definitely learned from data which are available for a wide range of languages.

The approach is quite promising, but, to reach state-of-the-art performance, the system has to deal with some mathematical issues, for example dimension reduction in the multilingual space. Further experimentation in multilingual mode with an extended number of languages could also provide meaningful outputs. By involving expert linguistic knowledge various sets of languages could be constructed using either only very close languages, or, on the contrary, using very distant languages. Thanks to the PanLex database, bilingual dictionaries can easily be extracted, which can, then, be directly used for multilingual experiments.

Author Contributions: Conceptualization and methodology, M.J.C.-B., G.R., G.B.; experiments, E.I.; investigation, M.J.C.-B., E.I., G.R., G.B.; writing—original draft preparation, E.I.; writing—review and editing, M.J.C.-B. and E.I.; supervision, M.J.C.-B., G.R., G.B.; funding acquisition, M.J.C.-B.

Funding: This research was funded by Spanish MINECO and FEDER grant number TIN2017-85854-C4-2-R.

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

SVD Singular Value Decomposition

References

1. Youn, H.; Sutton, L.; Smith, E.; Moore, C.; Wilkins, J.F.; Maddieson, I.; Croft, W.; Bhattacharya, T. On the universal structure of human lexical semantics. *Proc. Natl. Acad. Sci. USA* **2016**, *113*, 1766–1771. [[CrossRef](#)] [[PubMed](#)]
2. Ruder, S.; Vulić, I.; Søgaard, A. A survey of cross-lingual word embedding models. *arXiv* **2017**, arXiv:1706.04902.
3. Mikolov, T.; Le, Q.V.; Sutskever, I. Exploiting similarities among languages for machine translation. *arXiv* **2013**, arXiv:1309.4168.

4. Smith, S.L.; Turban, D.H.; Hamblin, S.; Hammerla, N.Y. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. In Proceedings of the 5th International Conference on Learning Representations (ICLR), Toulon, France, 24–26 April 2017.
5. Lample, G.; Conneau, A.; Ranzato, M.; Denoyer, L.; Jégou, H. Word Translation Without Parallel Data. In Proceedings of the Sixth International Conference on Learning Representations (ICLR), Vancouver, BC, Canada, 30 April–3 May 2018.
6. Dinu, G.; Lazaridou, A.; Baroni, M. Improving zero-shot learning by mitigating the hubness problem. *arXiv* **2014**, arXiv:1412.6568.
7. Faruqui, M.; Dyer, C. Improving vector space word representations using multilingual correlation. In Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, Gothenburg, Sweden, 26–30 April 2014; pp. 462–471.
8. Kamholz, D.; Pool, J.; Colowick, S.M. PanLex: Building a Resource for Panlingual Lexical Translation. In Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC), Reykjavik, Iceland, 26–31 May 2014; pp. 3145–3150.
9. Eszter, I.; Recski, G.; Borbély, G.; Castro-Bleda, M.J. Building a global dictionary for semantic technologies. *Iberspeech* **2018**, 286–290. [[CrossRef](#)]
10. Firth, J.R. A synopsis of linguistic theory. *Stud. Linguist. Anal.* **1957**, 1930–1955.
11. Bengio, Y.; Ducharme, R.; Vincent, P.; Jauvin, C. A neural probabilistic language model. *J. Mach. Learn. Res.* **2003**, *3*, 1137–1155.
12. Jurafsky, D.; Martin, J.H. *Speech and Language Processing*; Pearson: London, UK, 2017.
13. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient estimation of word representations in vector space. *arXiv* **2013**, arXiv:1301.3781.
14. Xing, C.; Wang, D.; Liu, C.; Lin, Y. Normalized word embedding and orthogonal transform for bilingual word translation. In Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, CO, USA, 31 May–5 June 2015; pp. 1006–1011.
15. Lazaridou, A.; Dinu, G.; Baroni, M. Hubness and pollution: Delving into cross-space mapping for zero-shot learning. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Beijing, China, 26–31 July 2015; Volume 1, pp. 270–280.
16. Ammar, W.; Mulcaire, G.; Tsvetkov, Y.; Lample, G.; Dyer, C.; Smith, N.A. Massively multilingual word embeddings. *arXiv* **2016**, arXiv:1602.01925.
17. Artetxe, M.; Labaka, G.; Agirre, E. Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), Austin, TX, USA, 1–5 November 2016; pp. 2289–2294.
18. Bojanowski, P.; Grave, E.; Joulin, A.; Mikolov, T. Enriching word vectors with subword information. *arXiv* **2016**, arXiv:1607.04606.
19. Tiedemann, J. Parallel Data, Tools and Interfaces in OPUS. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC), Istanbul, Turkey, 21–27 May 2012; pp. 2214–2218.
20. Duchi, J.; Hazan, E.; Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.* **2011**, *12*, 2121–2159.

