



Zestaw algorytmów zrównoleglenia i przechowywania
wielojęzycznych zanurzeń słów
na potrzeby obliczania prawdopodobieństwa tłumaczenia

dr Rafał Jaworski
Zakład Sztucznej Inteligencji
Uniwersytet im. Adama Mickiewicza w Poznaniu
Seminarium, IPI PAN, Warszawa



Zrównoleglenie na poziomie słów

- W swoich badaniach koncentruję się algorytmach przetwarzania języka naturalnego
 - Ich implementacje znajdują zastosowanie w systemach wspomagania tłumaczenia
 - Ponadto są wykorzystywane w badaniach lingwistycznych
 - Przedstawię moje główne osiągnięcia badawcze
-



Zrównoleganie na poziomie słów

Potencjalne zastosowania we wspomaganiu tłumaczenia:

- Automatyczne pozycjonowanie elementów nietłumaczonych
 - Znajdowanie błędów w tłumaczeniu (słów nieprzetłumaczonych)
 - Wyszukiwanie tłumaczenia wybranej frazy
-



Zrównoleglanie na poziomie słów

Istniejące rozwiązania:

- Giza++
 - fast_align (tylko model I i II IBM)
 - rozwiązania oparte na sieciach neuronowych:
 - a) Simalign
 - b) Joël Legrand, Michael Auli, Ronan Collobert: *Neural Network-based Word Alignment through Score Aggregation*
-



Inter-language Vector Space

- Inter-language Vector Space to zestaw metod do generowania, zrównoleglania, przechowywania i odczytywania zanurzeń słów w wielu językach.
 - Badania te przełożyły się na wdrożenia w oprogramowaniu XTM Cloud firmy XTM International.
-

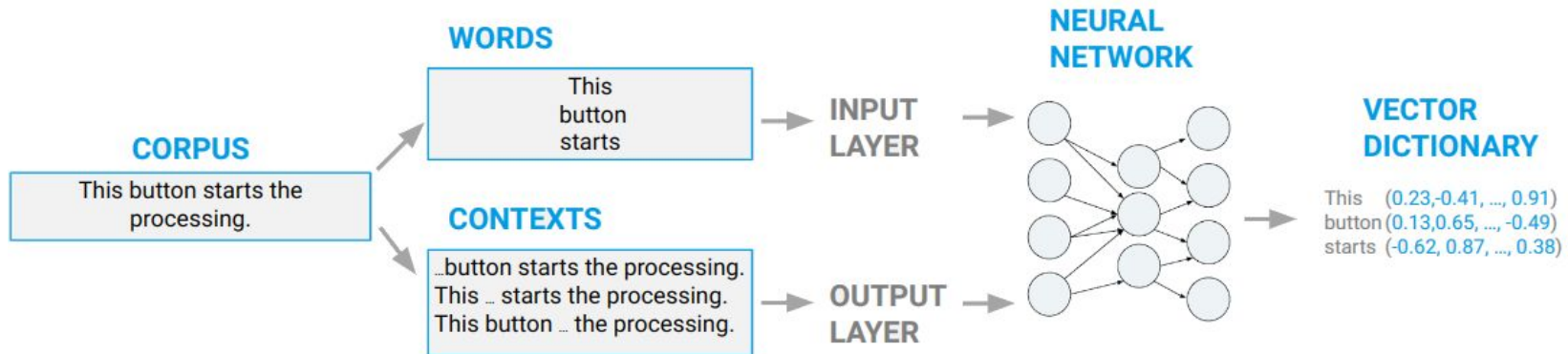


Inter-language Vector Space - word2vec

- Pierwszą funkcją zestawu jest generowanie reprezentacji wektorowych słów metodą word2vec. Są one generowane przy użyciu sieci neuronowej typu auto-encoder na podstawie danych tekstowych - korpusu tekstu.
 - Korpus dzielony jest na poszczególne słowa oraz konteksty, w których te słowa występują. Słowa są podawane warstwie wejściowej sieci neuronowej, natomiast konteksty są oczekiwanym wyjściem.
-



Inter-language Vector Space - word2vec



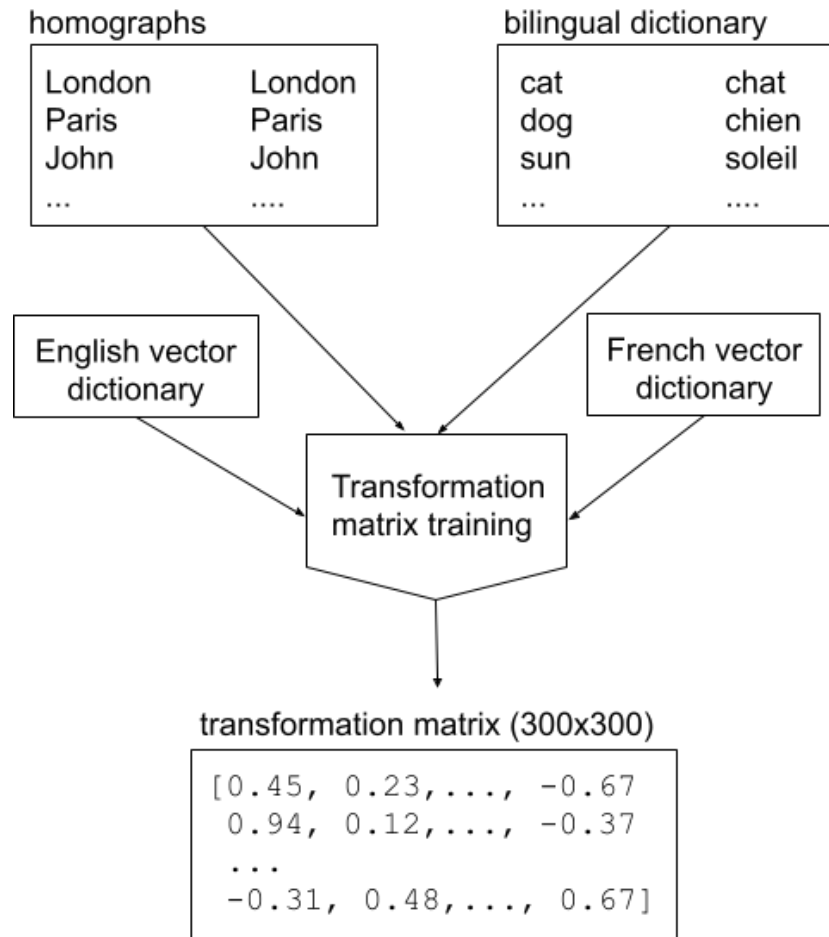


Inter-language Vector Space

- Zanurzenia trenowane są osobno dla wielu języków.
 - Zanurzenia słów nie są porównywalne pomiędzy językami.
 - Aby rozwiązać ten problem, dokonujemy zrównoleglenia przestrzeni wektorowych.
-



Zrównoleglanie na poziomie słów





Inter-language Vector Space

Trenowanie macierzy transformacji na podstawie:

Smith, S.L., Turban, D.H., Hamblin, S., & Hammerla, N.Y. (2017). Offline bilingual word vectors, orthogonal transformations and the inverted softmax. ArXiv, [abs/1702.03859](https://arxiv.org/abs/1702.03859).



Inter-language Vector Space

Rozwiązanie podstawowe (*Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation*) polega na użyciu słownika dwujęzycznego i słownika homografów do wytrenowania macierzy transformacji metodą stochastycznego spadku wzdłuż gradientu:

$$\min_W \sum_{i=1}^n \|y_i - Wx_i\|^2$$



Inter-language Vector Space

Zastosowane rozwiązanie bazuje na idei wytrenowania macierzy transformacji O dokonującej przekształcenia ortogonalnego dwóch przestrzeni wektorowych. Jeśli słownik dwujęzyczny ma postać par wektorów:

$$\{y_i, x_i\}_{i=1}^n$$

wówczas problem znalezienia optymalnego przekształcenia sprowadza się do maksymalizacji podobieństwa cosinusowego na odpowiadających sobie wyrazach ze słownika:

$$\max_O \sum_{i=1}^n y_i^T O x_i$$



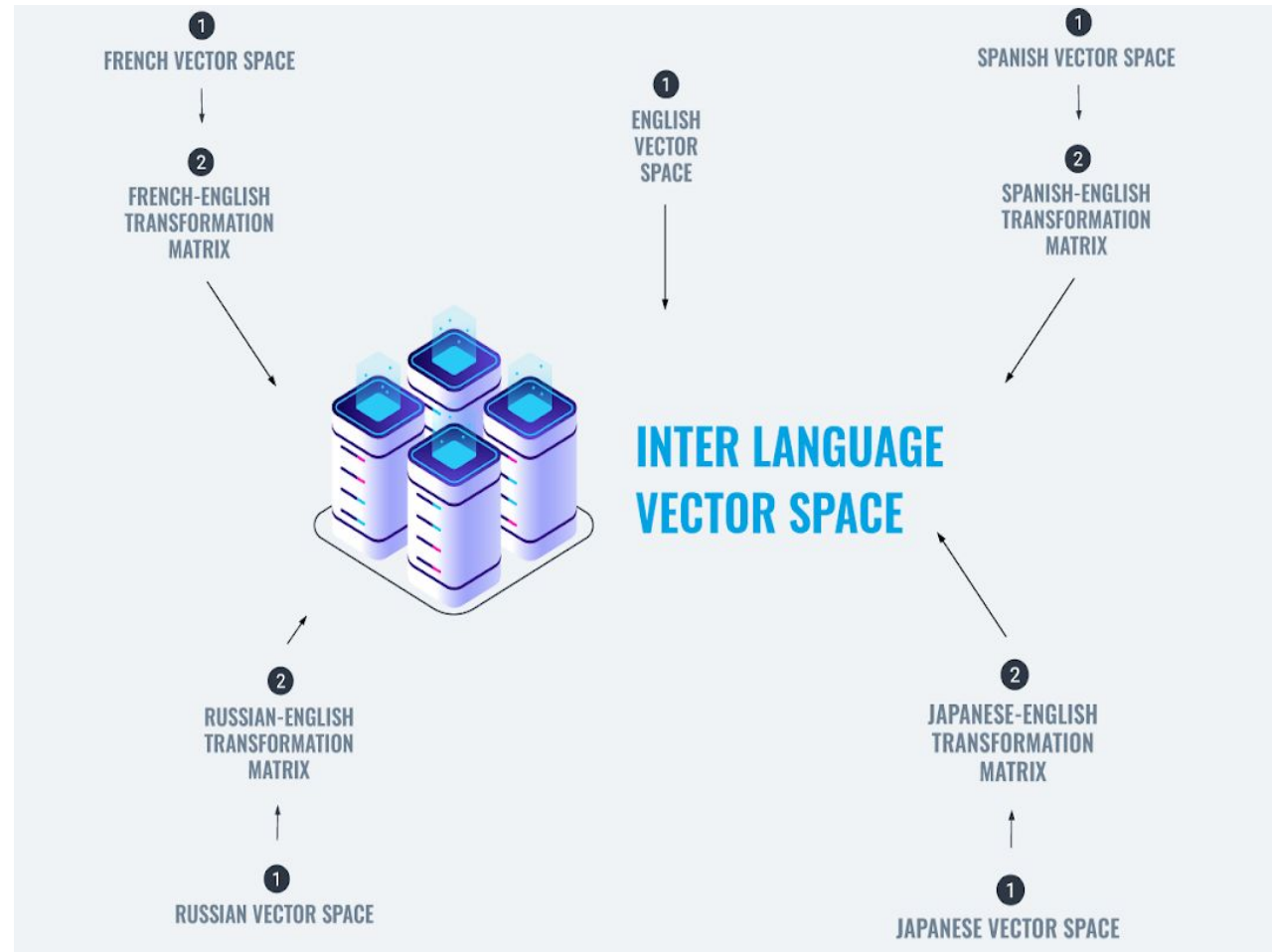
Inter-language Vector Space

Aby znaleźć optymalną macierz O , dokonywane są następujące operacje:

1. Utwórz ze słownika macierze X_D i Y_D tak, aby i -te wiersze zawierały wektory odpowiadających sobie słów.
 2. Wyznacz macierz $M = Y_D^T X_D$
 3. Dokonaj rozkładu według wartości osobliwych M :
$$M = U\Sigma V^T$$
 4. Finalna macierz transformacji jest złożeniem macierzy U^T i V^T
-



Inter-language Vector Space





Inter-language Vector Space

Obliczanie prawdopodobieństwa, że słowa A i B są swoimi tłumaczeniami dokonywane jest na podstawie odległości cosinusowej:

$$P(A, B) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{n=1}^{300} A_i B_i}{\sqrt{\sum_{n=1}^{300} A_i^2} \cdot \sqrt{\sum_{n=1}^{300} B_i^2}}$$



Inter-language Vector Space

Przykładowe wyniki

Słowo angielskie	Słowo włoskie	Podobieństwo
cat	gatto	0.696
cat	gatta	0.552
cat	giorno	0.164
day	giorno	0.692
day	fuoco	0.193
fire	fuoco	0.590



Analiza komponentów zdań

Szkic algorytmu:

1. Wyznacz współczynnik dopasowania dla każdej pary słów źródło-tłumaczenie
 2. Dla każdego słowa źródłowego wyznacz najlepsze dopasowanie wśród słów zdania docelowego.
 3. Dla każdego słowa docelowego wyznacz najlepsze dopasowanie wśród słów zdania źródłowego.
 4. Połącz symetrycznie najlepsze dopasowania.
-



Pozycjonowanie znaczników nietłumaczalnych

- Znaczniki nietłumaczalne znajdujące się w tekstach do tłumaczenia:
 - tagi HTML/XML
 - Znaczniki formatujące
 - Symbole specjalne
 - W podstawowym scenariuszu tłumacz musi te znaczniki przenosić ręcznie ze zdania źródłowego do tłumaczenia.
 - W celu zautomatyzowania tego procesu opracowałem algorytm automatycznego przenoszenia znaczników
-



Pozycjonowanie znaczników nietłumaczalnych

Authentication is required to ` change ` your own user data

Vous devez vous authentifier pour `modifier ` vos propres données utilisateur



Pozycjonowanie znaczników nietłumaczalnych

- Tłumaczenie automatyczne z wykorzystaniem sieci neuronowych osiąga niższą jakość tłumaczenia w przypadku obecności znaczników nietłumaczalnych w zdaniu źródłowym.
- Rozwiązaniem zastosowanym w systemie XTM jest następujący algorytm:
 - usunięcie znaczników ze zdania źródłowego
 - tłumaczenie automatyczne
 - automatyczne pozycjonowanie znaczników



Pozycjonowanie znaczników nietłumaczalnych

- Funkcja auto-inlines systemu XTM jest wywoływana średnio **114** razy w ciągu minuty.
 - Liczba obsługiwanych języków: **56**, par językowych: **1540**
 - Skuteczność: **98%** zdań z automatycznie pozycjonowanymi znacznikami nie wymaga korekty
-



Pozycjonowanie znaczników nietłumaczalnych

- Analiza RMSE średniej odległości znacznika od pozycji prawidłowej:

FROM	TO	TOTAL	PHITS	RMSE
EN	DE	1238	1076	≈2.146
EN	ES	1270	1151	≈1.705
EN	FR	1258	1141	≈2.081
EN	IT	1270	1113	≈1.633
EN	JA	1256	1040	≈2.762
EN	KO	1254	1019	≈2.039
EN	NL	1260	1010	≈3.31
EN	RU	1260	1109	≈1.798
EN	ZHS	1248	689	≈3.36
EN	ZHT	1260	717	≈3.447
EN	ID	1248	1145	≈1.404
DE	EN	1239	958	≈1.727
DE	ES	1251	896	≈2.723
DE	FR	1239	912	≈2.611
DE	IT	1247	882	≈2.723
DE	JA	1254	836	≈3.232
DE	KO	1239	815	≈3.383
DE	NL	1253	802	≈3.757
DE	RU	1241	833	≈2.144
DE	ZHS	1233	600	≈3.923
DE	ZHT	1243	599	≈4.03
DE	ID	1239	869	≈2.169



Zgłoszenie patentowe

Zestaw Inter-language Vector Space był przedmiotem zgłoszenia patentowego w United States Patent and Trademark Office w 2020 roku:

“Inter-Language Vector Space: Effective assessment of cross-language semantic similarity of words using word-embeddings, transformation matrices and disk based indexes. “. Nr wniosku: 17/064,620.



Podsumowanie

- Algorytmy zrównoleglania i przechowywania wielojęzycznych zanurzeń słów pozwalają na efektywne i dokładne oszacowanie prawdopodobieństwa tłumaczenia słów.
 - Techniki te znajdują zastosowanie w algorytmach wspomaganego tłumaczenia.
 - Wdrożenie w systemie XTM Cloud potwierdza skuteczność rozwiązania.
-



ADAM MICKIEWICZ UNIVERSITY IN POZNAŃ

Dziękuję za uwagę!

Seminarium ZIL, IPI PAN, Warszawa

www.amu.edu.pl