

Reprezentacje wektorowe tekstu i ich wizualizacja

Machine Learning Gdansk - Hackerspace

20.06.2017

Adam Wróbel

github.com/Adamage

adam.wrobel1@gmail.com

[linkedin.com/in/adamwrobel1](https://www.linkedin.com/in/adamwrobel1)

Dlaczego chcemy przedstawić słowa jako wektory?

- Większość rozwiązań w ML opiera się na operacjach mnożenia, dodawania, na macierzach liczb

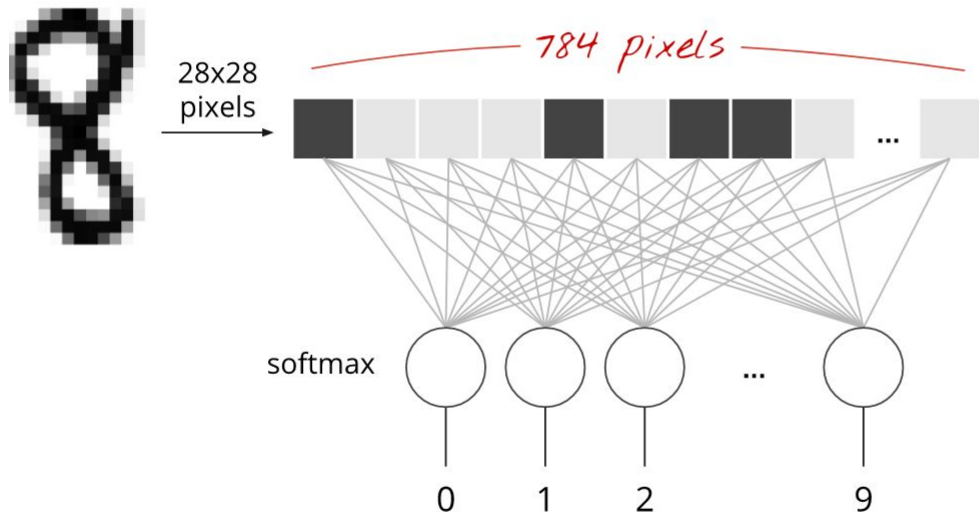
Np. wykrywanie pisanej cyfry to problem klasyfikacji:

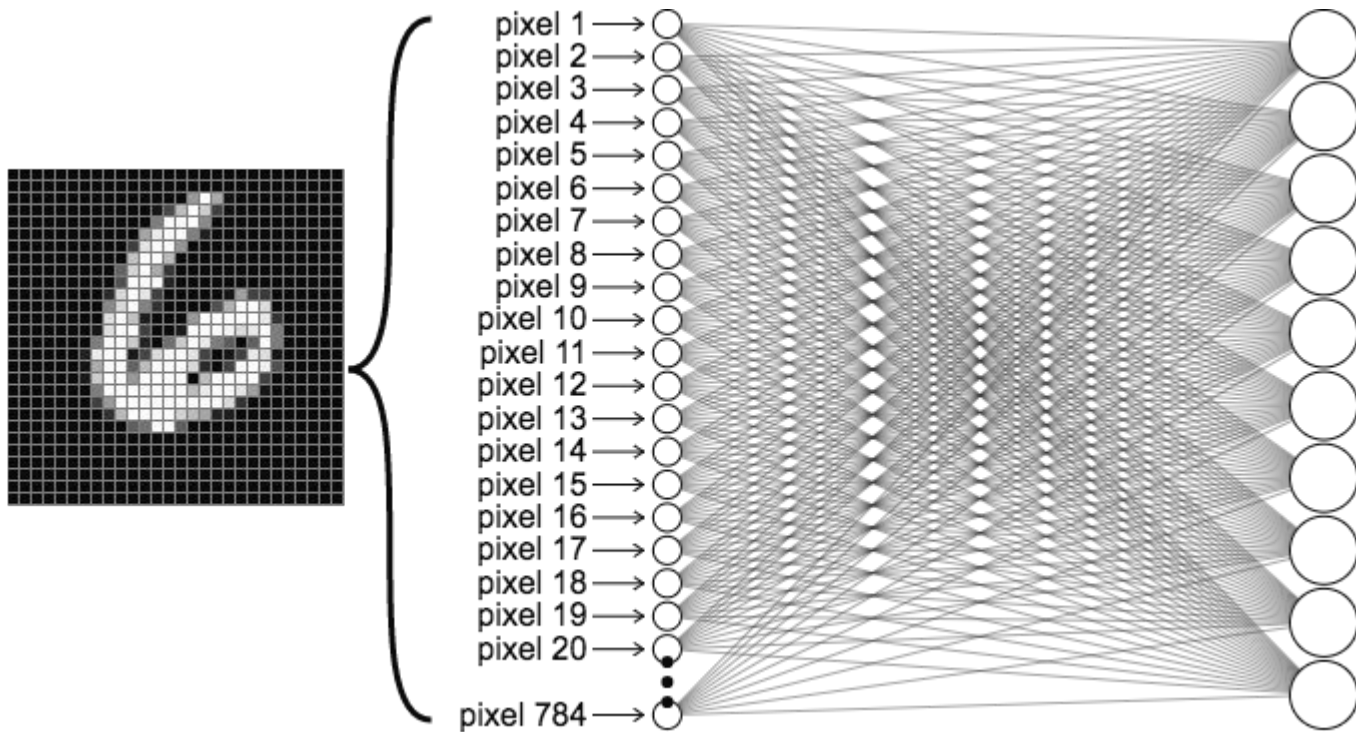
Obrazek 28x28, każdy piksel

Może mieć wartość 0 albo 1

Taki wektor musi dać się zmapować do jednej z 10 kategorii

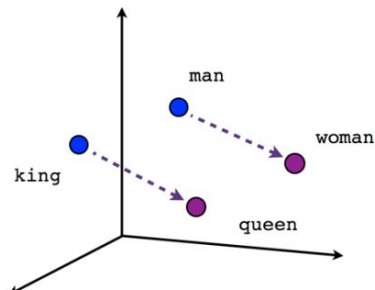
Poprzez np. wyuczenie sieci neur.



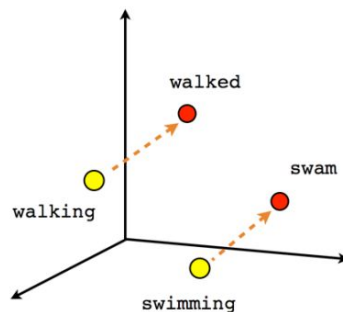


Dlaczego chcemy przedstawić słowa jako wektory?

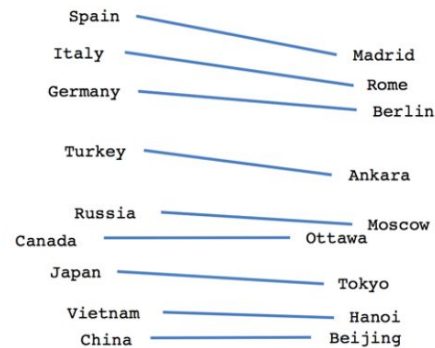
- Potrzebujemy reprezentacji, która jest w stanie uchwycić:
 - Wieloznaczność słów i synonimy
 - Slang, wulgaryzmy
 - Kontekst
 - Relacje między słowami, analogie, np:



Male-Female



Verb tense



Country-Capital

Jak przedstawić słowa jako wektory?

Intuicja:

Jeżeli słowa często występują w podobnym kontekście, to muszą mieć podobne znaczenie.

Mozemy “nauczyć” model ciągłej dystrybucji prawdopodobieństwa, że dla danego słowa, inne słowa będą w jego kontekście (zazwyczaj w odległości 2-8 słów).

Wyprowadzenie:

<http://web.stanford.edu/class/cs224n/lectures/cs224n-2017-lecture2.pdf>

Jak przedstawić słowa jako wektory?

Intuicja:

Jeżeli słowa często występują w podobnym kontekście, to muszą mieć podobne znaczenie.

Na rysunku:

“ruchome okno”

Wypisujemy wszystkie pary słów w środku i jego sąsiadów w oknie. Będą to pary input-output do trenowania modelu.

Source Text	Training Samples						
<table border="1"><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. →	The	quick	brown	(the, quick) (the, brown)			
The	quick	brown					
<table border="1"><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. →	The	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)		
The	quick	brown	fox				
<table border="1"><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. →	The	quick	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)	
The	quick	brown	fox	jumps			
<table border="1"><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. →	The	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
The	quick	brown	fox	jumps	over		

Intuicja odnośnie działania algorytmu

Model **Skip-gram** - predykcja słów w oknie dookoła słowa

Model **Continuous Bag of Words** - predykcja słowa pośrodku, na podstawie słów w oknie

Użyc prostej sieci neuronowej, z jedną warstwą ukrytą:

- Wejście to wektor $[10000 \times 1]$, z jedną jedynką oznaczającą TO słowo
- Wyjście:
 - Przy trenowaniu, to wektor $[10000 \times 1]$, z jedną jedynką oznaczającą jednego z sąsiadów
 - Przy ewaluacji, wektor $[10000 \times 1]$ z rozkładem prawdopodobieństwa, które słowa mogłyby być sąsiadami

Aplikacja softmaxu na wyjściu z klasyfikacją na jedno z 10000 słów da nam najbardziej prawdopodobnego sąsiada.

Trenujemy konkretne słowo na inpuście tyle razy, ile ma sąsiadów w oknie, i w ten sposób wszystkie słowa.

W efekcie, jeżeli warstwa ukryta ma 300 jednostek, to uzyskamy za pomocą 300 różnych wag, liczby do reprezentacji wektorowej tego słowa!

Wymiar wektora, np. 300, to liczba “cech”, “features”, nauczonych automatycznie.

Intuicja odnosnie dzialania algorytmu

W trakcie trenowania:

- Input to wyraz zakodowany "1"
- Output to wyraz ktory jest na pewno w oknie dookola slowa

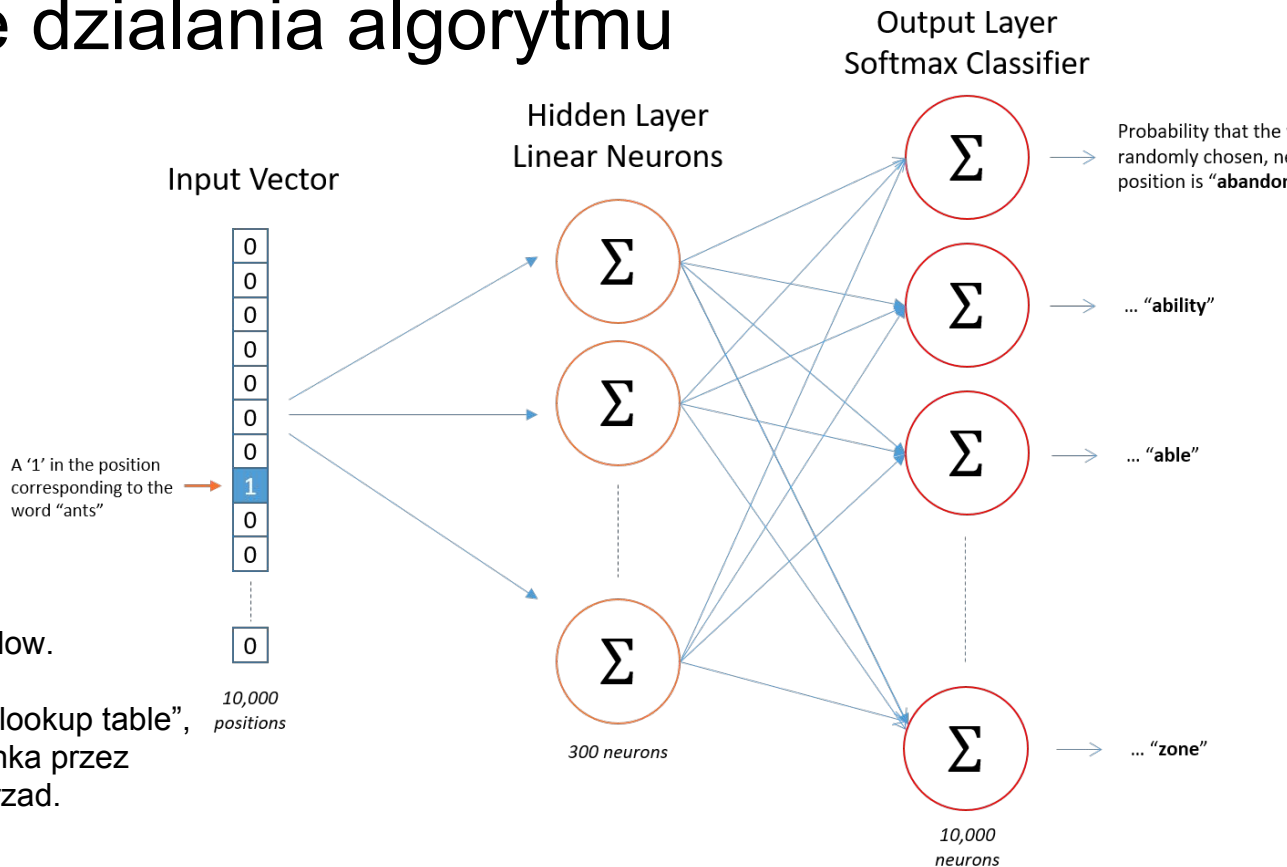
Po trenowaniu:

- Wagi 300 neuronow ukrytych sa jednocześnie wektorem slowa

Przy ewaluacji danego slowa, wynik to wlasciwie lista najbardziej prawdopodobnych sasiadow.

Aplikujac softmax, dostaniemy jedno ze slow.

Tak naprawde warstwa ukryta dziala jak "lookup table", gdyz przemnozenie wektora z jedna jedynka przez macierz wag "wybiera" z macierzy jeden rzad.



Demo: użycie wytrenowanego modelu

Git clone <https://github.com/Adamage/neural-nets-tutorials>

<https://github.com/Adamage/neural-nets-tutorials/tree/master/tensorflow-tutorials/natural-language/notebook-check-similarities.ipynb>

Ten notebook ściągnie dane i wytrenuje model. Następnie możemy podglądać najbliższe słowa, pod względem częstości występowania w podobnych kontekstach.

\$ jupyter notebook

W przeglądarce -> localhost:8888

Wizualizacja

Jak pokazac 300 wymiarowe wektory na kartce, ekranie?

Jest wiele metod redukcji wymiarow, zachowujacych strukture:

- PCA - Principal Component Analysis
- SNE - Stochastic Neighbour Embedding
- Sammon mapping
- CCA - Curvilinear Component Analysis
-



Przykładowa wizualizacja w dwóch wymiarach wektorów 300 wymiarowych. Skupiska oznaczają zbiory wyrazów podobnych, lub występujących w podobnych miejscach.

Wizualizacja

- t-SNE = t-distributed stochastic neighbour embedding
- Algorytm nie jest nowy, ale konkretnie t-SNE opracowali G. Hinton i L. Maaten
- Intuicyjnie:
 - Skonstruuj rozkład prawdopodobieństwa między parami obiektów wysoko-wymiarowych, podobne mają dużą szansę być wybranymi na "sasiada", i na odwrot
 - Skonstruuj podobny rozkład prawdopodobieństwa, ale w 2 lub 3 wymiarach
 - Numerycznie minimalizuj dywergencję Kullbacka-Leiblera między dwoma dystrybucjami względem lokalizacji na mapie 2d/3d

https://en.wikipedia.org/wiki/T-distributed_stochastic_neighbor_embedding

Demo:

2D t-SNE wygenerowane z TED talks + Word2Vec
przy użyciu biblioteki bokeh

https://github.com/Adamage/neural-nets-tutorials/blob/master/tensorflow-tutorials/natural-language/oxford-course-practical1/part4_tsne_of_ted.py

Skrypt należy uruchomić, powinienn pociągnąć dane, wytrenować word2vec a następnie prostą mapę 2D t-sne, która jest interaktywna.

Outputs: webpage using **bokeh** library - interactive 2d map

Demo:

2D t-SNE wygenerowane z TED talks + Word2Vec,
Przy uzyciu tensorflow i tensorboard

https://github.com/Adamage/neural-nets-tutorials/blob/master/tensorflow-tutorials/natural-language/word2vec/tensorboard_word2vec_TED.py

Wymaga scipy, numpy, tensorflow 1.2+

Outputs: tensorflow checkpoints i plik z metadanymi - nalezy uzyc tensorboard

Zrodla, tutoriala, ciekawe linki

Deep learning and machine learning:

- <https://github.com/oxford-cs-deepnlp-2017/lectures>
- <http://cs224d.stanford.edu/syllabus.html>
- <https://www.quora.com/Deep-Learning-What-is-meant-by-a-distributed-representation>
- <http://www.cs.toronto.edu/~bonner/courses/2014s/csc321/lectures/lec5.pdf>
-

Word vectors:

- <https://www.tensorflow.org/tutorials/word2vec>

T-SNE

- Publikacja autorow metody: <http://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>
- <https://lvdmaaten.github.io/tsne/>
- <https://github.com/oreillymedia/t-SNE-tutorial>
- <http://alexanderfabisch.github.io/t-sne-in-scikit-learn.html>
- <http://distill.pub/2016/misread-tsne/>