

## Application of Deep Learning to Text and Images

### Module 2, Lab 3: GloVe Word Vectors

This notebook supports the topics presented on on the Word Embeddings lecture.

In this lab you will learn how to use word embeddings. Word embeddings, or word vectors, are a way of representing words as numeric vectors in a high-dimensional space. These embeddings capture the meaning of the words, the relationships between them, and can be used as inputs to machine learning models for a variety of natural language processing tasks.

The term **Word vectors** refers to a family of related techniques, first gaining popularity via Word2Vec which associates an *n*-dimensional vector to every word in the target language.

• Note: Normally n is in the range of 50 to 500. In this lab, you will set it to 50

You will learn:

- · What GloVe word vectors are
- · How to load GloVe word vectors
- How to use GloVe to produce word vectors
- What cosine Similarity is
- How to use cosine similarity to compare words

You will be presented with two kinds of exercises throughout the notebook: activities and challenges.



No coding is needed for an activity. You try to understand a concept, answer questions, or run a code cell.

Challenges are where you can practice your coding skills.

#### Index

- 1. GloVe Word Vectors
- 2. Cosine Similarity

First, install the latest versions of the libraries.

```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts.
autovizwidget 0.21.0 requires pandas<2.0.0,>=0.20.1, but you have pandas 2.0.3 which is incompatible. hdijupyterutils 0.21.0 requires pandas<2.0.0,>=0.17.1, but you have pandas 2.0.3 which is incompatible. sparkmagic 0.21.0 requires pandas<2.0.0,>=0.17.1, but you have pandas 2.0.3 which is incompatible.

In [2]: from torchtext.vocab import GloVe

#from torchtext.vocab import GloVe

GloVe.url['6B'] = 'https://huggingface.co/stanfordnlp/glove/resolve/main/glove.6B.zip'

from sklearn.decomposition import PCA
from sklearn.metrics.pairwise import cosine_similarity

%matplotlib inline
import matplotlib.pyplot as plt
```

#### **GloVe Word Vectors**

In [1]: |# installing libraries

!pip install -U -q -r requirements.txt

You learned about **Word2Vec** and **FastText** as word embedding techniques. Now you will use a set of pre-trained word embeddings. Pre-trained embeddings are created by someone else who took the time and computational power to train. This reduces your cost by not having to train the model yourself. One popular word embedding is **GloVe** embeddings. GloVe is a variation of a Word2Vec model. To learn more about GloVe, read the Project GloVe website.

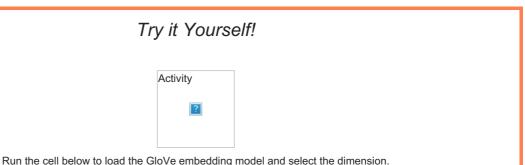
In this exercise, you will discover relationships between word vectors using the GloVe embeddings.

You can easily import GloVe embeddings from the Torchtext library. Here, you will get vectors with 50 dimensions.

The name parameter refers to the particular pre-trained model that should be loaded:

Matplotlib is building the font cache; this may take a moment.

- Wikipedia 2014 + Gigaword 5
  - 6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download: "6B"
  - This is the model that you will load.
- Common Crawl
  - 42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download: "42B"
- Common Crawl
  - 840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download: "840B"
- Etc
  - See documentation in Stanford link above



Now that the data is loaded, you can access it and print example word embeddings.

You might notice that the tensor has 50 values in it. This is related to the dimension flag ( dim=50 ) you set when you loaded the GloVe model. You can generate word embeddings for several words and use them to determine how closely related words are. This is a task that machine learning is really good at.

### Try it Yourself!



In the code block below, generate word embeddings for the words "computer" and "human" using pre-trained GloVe embedding.

```
########### CODE HERE ##############
print(f"computer -> {glove['computer']}\n")
print(f"human -> {glove['human']}\n")
computer -> tensor([ 0.0791, -0.8150, 1.7901, 0.9165, 0.1080, -0.5563, -0.8443, -1.4951,
                                                  0.1342, \quad 0.6363, \quad 0.3515, \quad 0.2581, \quad -0.5503, \quad 0.5106, \quad 0.3741, \quad 0.1209, \quad 0.1
                                              \hbox{-1.6166,} \quad 0.8365, \quad 0.1420, \ \hbox{-0.5235,} \quad 0.7345, \quad 0.1221, \ \hbox{-0.4908,}
                                                                                                                                                                                                                                                                                                                                                                                                                0.3253,
                                                  0.4531, -1.5850, -0.6385, -1.0053, 0.1045, -0.4298,
                                                                                                                                                                                                                                                                                                                                                               3.1810, -0.6219,
                                                  0.1682, -1.0139, 0.0641, 0.5784, -0.4556,
                                                                                                                                                                                                                                                                                                                                                                    0.3720, -0.5772,
                                                                                                                                                                                                                                                                                                                0.7378,
                                             0.6644, 0.0551,
-0.1143, 0.2071])
                                                                                                                                                      0.0379, 1.3275, 0.3099, 0.5070,
                                                                                                                                                                                                                                                                                                                                                               1.2357, 0.1274,
human -> tensor([ 0.6185,  0.1191, -0.4679,  0.3137,  1.0334,
                                                                                                                                                                                                                                                                                                                                                                 0.9596, 0.8780, -1.0346,
                                                  1.6322, 0.2935, 0.8084, -0.0589, 0.0213,
                                                                                                                                                                                                                                                                                                                0.4099.
                                                                                                                                                                                                                                                                                                                                                                0.5444, -0.3331,
                                                                                                                                                                                                                                                                                                                0.6939,
                                                  0.5371, -0.3582, 0.2937, 0.0902, -0.9205,
                                                                                                                                                                                                                                                                                                                                                                 0.3910, -0.6439,
                                                  0.7783, -1.7215, -0.4839, -0.5033, -0.2251,
                                                                                                                                                                                                                                                                                                                0.0992,
                                                                                                                                                                                                                                                                                                                                                                  3.2095, -0.3155,
                                              -0.7175, \ -1.6752, \ -1.3537, \quad 0.1520, \quad 0.0546, \ -0.1633, \ -0.0280, \quad 0.3917,
                                              -0.5501, \; -0.0792, \; \; 0.6339, \; \; 0.5145, \; \; 0.7012, \; \; 0.2764, \; -0.5344, \; \; 0.0648, \; -0.5344, \; \; 0.0648, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.5344, \; -0.
                                              -0.2197, -0.5205])
```

#### Cosine Similarity

'cat'

is closer to

'dog'

than

'sea'

You learned about cosine similarity in class, now let's look at an example. Use the cosine\_similarity() function from scikit-learn to easily calculate cosine similarity between word vectors.

## Try it Yourself!



Run the cell below to calculate cosine similarity between word vectors.

```
In [6]: # define the similarity between two words
def similarity(w1, w2):
    return cosine_similarity([glove[w1].tolist()], [glove[w2].tolist()])

# Say if w1 is closer to w2 than w3
def simCompare(w1, w2, w3):
    s1 = similarity(w1, w2)
    s2 = similarity(w1, w3)
    if s1 > s2:
        print(f"'{w1}'\tis closer to\t'{w2}'\tthan\t'{w3}'\n")
    else:
        print(f"'{w1}'\tis closer to\t'{w3}'\tthan\t'{w2}'\n")

In [7]: simCompare("actor", "pen", "film")
    simCompare("cat", "dog", "sea")
    'actor' is closer to 'film' than 'pen'
```

# Try it Yourself!



Write code to determine if "car" is closer to "truck" than "bike".

### Conclusion

You have now seen how to use word embeddings and determine relationships between word vectors using the GloVe embeddings.

# Next Lab: Word Embeddings

In the next lab of this module you will learn how to build a recurrent neural network (RNN) with PyTorch. It will also show you how to implement a simple RNN-based model for natural language processing.

Processing math: 100%