# Text classification using Neural Networks

The goal of this notebook is to learn to use Neural Networks for text classification.

In this notebook, we will:

- Train a shallow model with learning embeddings
- Download pre-trained embeddings from Glove
- Use these pre-trained embeddings

## The BBC topic classification dataset

The BBC provides some benchmark topic classification datasets in English at: http://mlg.ucd.ie/datasets/bbc.html.

The raw text (encoded with the latin-1 character encoding) of the news can be downloaded as a ZIP archive:

```python
import os
import os.path as op
import zipfile
from urllib.request import urlretrieve

BBC_DATASET_URL = "http://mlg.ucd.ie/files/datasets/bbc-fulltext.zip"
zip_filename = BBC_DATASET_URL.rsplit('/', 1)[1]
BBC_DATASET_FOLDER = 'bbc'
if not op.exists(zip_filename):
    print("Downloading %s to %s..." % (BBC_DATASET_URL, zip_filename))
    urlretrieve(BBC_DATASET_URL, zip_filename)

if not op.exists(BBC_DATASET_FOLDER):
    with zipfile.ZipFile(zip_filename, 'r') as f:
        print("Extracting contents of %s..." % zip_filename)
        f.extractall('.')
```

Each of the five folders contains text files from one of the five topics:

```python
target_names = sorted(folder for folder in os.listdir(BBC_DATASET_FOLDER)
                      if op.isdir(op.join(BBC_DATASET_FOLDER, folder)))
target_names
```

```python
# Example of a file in the "business" category
with open(op.join(BBC_DATASET_FOLDER, 'business', '001.txt'), 'rb') as f:
    print(f.read().decode('latin-1')[:500] + '...')
```

Let's randomly partition the text files in a training and test set while recording the target category of each file as an integer:

```python
import numpy as np
from sklearn.model_selection import train_test_split

target = []
filenames = []
for target_id, target_name in enumerate(target_names):
    class_path = op.join(BBC_DATASET_FOLDER, target_name) # e.g. 'bbc/busine
    for filename in sorted(os.listdir(class_path)):
        filenames.append(op.join(class_path, filename))
        target.append(target_id)

target = np.asarray(target, dtype=np.int32)
target_train, target_test, filenames_train, filenames_test = train_test_spli
    target, filenames, test_size=200, random_state=0)
```

What we now have is pairs of target labels (which category the document belongs to) and filenames (where the document is stored on disk):

```python
target_train[:5], filenames_train[:5]
```

```python
size_in_bytes = sum([len(open(fn, 'rb').read()) for fn in filenames_train])
print("Training set size: %0.3f MB" % (size_in_bytes / 1e6))
```

This dataset is small so we can load everything into memory right now (which simplifies our code later). If we had substantially more data, we would need to use a `tf.data.Dataset` to stream it from disk in batches during training.

```python
texts_train = [open(fn, 'rb').read().decode('latin-1') for fn in filenames_t
texts_test = [open(fn, 'rb').read().decode('latin-1') for fn in filenames_te
```

# A first baseline model

For simple topic classification problems, one should always try a simple method first. Let's try using a `CountVectorizer` followed by `LogisticRegression` as a baseline. What this will do is:

- Convert the text documents to a matrix of token counts (each row is a document, each column is a word, each cell is the count of the word in the document)
- Train a logistic regression model on this matrix

It's a very efficient method and should give us a strong baseline to compare our deep learning method against.

```python
from sklearn.feature_extraction.text import CountVectorizer

# Understanding what the CountVectorizer does
vectorizer = CountVectorizer(max_features=2000) # only keep the 2000 most fr
X_train = vectorizer.fit_transform(texts_train)

# Compare the content of the first document with the vocabulary
print("Start of the first document:")
print(texts_train[0][0:100] + '...')
print('----------')
print("Sampling of vocabulary counts in the document:")
for word, count in zip(vectorizer.get_feature_names_out()[200:210], X_train.
    print(word, count)
```

```python
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline


text_classifier = make_pipeline(
    CountVectorizer(max_features=2000),
    LogisticRegression(),
)
```

```python
%time _ = text_classifier.fit(texts_train, target_train)
```

You may get a warning above that "lbfgs failed to converge". This means that the optimization algorithm did not reach the desired precision. This is not a big deal here, as we are not looking for the best possible accuracy, but just a baseline. We can check the accuracy of this model on the test set:

```python
text_classifier.score(texts_test, target_test)
```

Approximately 95 percent testing accuracy on a very simple baseline. It's quite unlikely that we can significantly beat that baseline with a more complex deep learning based model. This is simply not a complex task - we wouldn't expect to see this level of performance from a simple model on a real-world text classification problem. Let's move on, and see how well we can do with a simple neural network.

## Preprocessing text for the (supervised) CBOW model

We will implement a simple classification model in Keras. Raw text requires (sometimes a lot of) preprocessing.

The following cells uses Keras to preprocess text:

- using a tokenizer. This converts the texts into sequences of indices representing the `20000` most frequent words

- sequences have different lengths, so we pad them (add 0s at the end until the sequence is of length `1000`). For example, if we were padding to three words, and we had a sequence of just "dog", we would pad it to "[<dog token>,0,0]".
- we convert the output classes as 1-hot encodings

```python
from tensorflow.keras.preprocessing.text import Tokenizer

MAX_NB_WORDS = 20000

# vectorize the text samples into a 2D integer tensor
tokenizer = Tokenizer(num_words=MAX_NB_WORDS, char_level=False)
tokenizer.fit_on_texts(texts_train)
sequences = tokenizer.texts_to_sequences(texts_train)
sequences_test = tokenizer.texts_to_sequences(texts_test)

word_index = tokenizer.word_index
print('Found %s unique tokens.' % len(word_index))
print(f'Example of word_index: {list(word_index.items())[:5]}')
```

Tokenized sequences are converted to list of token ids (with an integer code). We can convert them back to text to see what they now look like:

```python
index_to_word = dict((i, w) for w, i in tokenizer.word_index.items())
```

```python
print(f'Original text:          {texts_train[0].split(" ")[0:4]}')
print(f'Tokenized text:         {sequences[0][0:4]}')
print(f'Converted back to text: {[index_to_word.get(i, "UNK") for i in seque
```

Let's have a closer look at the tokenized sequences:

```python
seq_lens = [len(s) for s in sequences]
print("average length: %0.1f" % np.mean(seq_lens))
print("max length: %d" % max(seq_lens))
```

```python
import matplotlib.pyplot as plt

plt.hist(seq_lens, bins=50);
```

We can see that while we do have sequences up to 4355 words long, the vast majority of sequences are less than 1000 words long. We can use this information to truncate or pad all the sequences to 1000 symbols to build the training set. This will simplify our model and speed up training.

```python
from tensorflow.keras.preprocessing.sequence import pad_sequences

MAX_SEQUENCE_LENGTH = 1000

# Make all sequences exactly 1000 words long
```
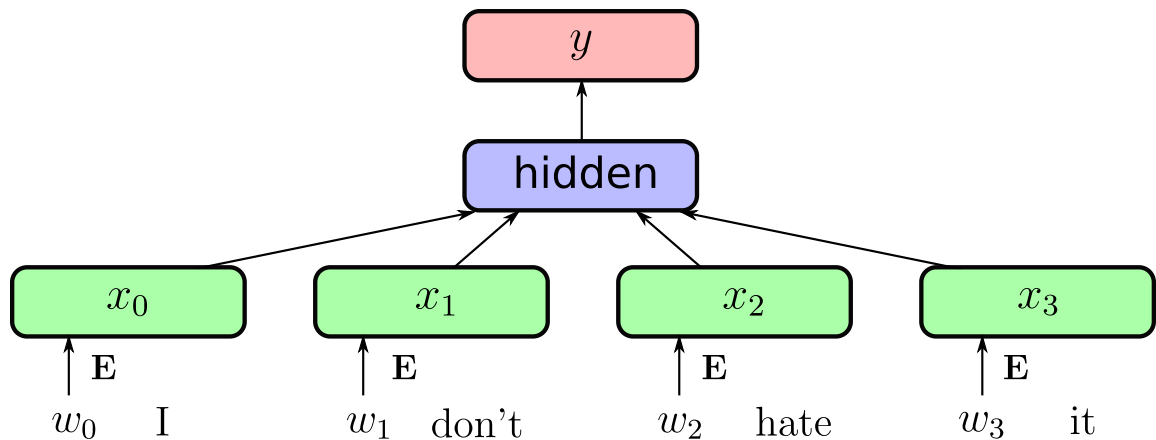
```python
x_train = pad_sequences(sequences, maxlen=MAX_SEQUENCE_LENGTH)
x_test = pad_sequences(sequences_test, maxlen=MAX_SEQUENCE_LENGTH)
print('Shape of data tensor:', x_train.shape)
print('Shape of data test tensor:', x_test.shape)
```

In [ ]:
```python
from tensorflow.keras.utils import to_categorical

y_train = to_categorical(target_train)
print('Shape of label tensor:', y_train.shape)
```

## A simple supervised CBOW model in Keras

The following computes a very simple model, as described in fastText:



- Build an embedding layer mapping each word to a vector representation
- Compute the vector representation of all words in each sequence and average them
- Add a dense layer to output 5 classes

In [ ]:
```python
from tensorflow.keras.layers import Dense, GlobalAveragePooling1D, Embedding
from tensorflow.keras.models import Sequential
from tensorflow.keras.optimizers import Adam

EMBEDDING_DIM = 50
N_CLASSES = len(target_names)

model = Sequential([
    Embedding(MAX_NB_WORDS, EMBEDDING_DIM, trainable=True), # Just like we'v
    GlobalAveragePooling1D(), # This layer averages the embeddings of all wo
    Dense(N_CLASSES, activation='softmax') # This layer outputs a probabilit
])

model.compile(
    loss='categorical_crossentropy',
    optimizer=Adam(learning_rate=0.01),
    metrics=['accuracy']
)
```

```
In [ ]:   model.fit(x_train, y_train, validation_split=0.1,
                    epochs=10, batch_size=32)
```

**Exercise**

- Evaluate the model on the test set
- Identify an example of a mis-classified document and display the text of the document

```
In [ ]:   # Evaluate the model
```

## Building more complex models

**Exercise**

- Copy the previous model, and add more complexity to it. You can try adding more layers, or using a different type of layer (e.g. LSTM, Conv1D, etc.)
- Some examples of what you could do:
  - Add a LSTM layer before the dense layer (LSTM documentation)
  - Add a Conv1D layer after the embedding layer (Conv1D documentation)
  - Add more dense layers

```
In [ ]:   from tensorflow.keras.layers import Conv1D, MaxPooling1D, Flatten, Dense, Em
          from tensorflow.keras.models import Sequential

          model = Sequential([
              # Your code here
          ])

          # Compile the model
          model.compile(loss='categorical_crossentropy',
                        optimizer='adam', metrics=['acc'])
```

```
In [ ]:   model.fit(x_train, y_train, validation_split=0.1,
                    epochs=15, batch_size=32)

          output_test = model(x_test)
          test_casses = np.argmax(output_test, axis=-1)
          print("Test accuracy:", np.mean(test_casses == target_test))
```

## Loading pre-trained embeddings

The file `glove100K.100d.txt` is an extract of Glove Vectors, that were trained on English Wikipedia and the Gigaword 5 corpus. They differ from word2vec in the way the vectors are trained, but the idea is the same: each word is represented as a vector of `100` numbers.

We extracted the `100 000` most frequent words for you, and the code below downloads them.

In [ ]:
```python
# Get pretrained Glove Word2Vec
URL_REPRESENTATIONS = "https://github.com/m2dsupsdlclass/lectures-labs/relea
ZIP_REPRESENTATIONS = "glove100k.100d.zip"
FILE_REPRESENTATIONS = "glove100K.100d.txt"

if not op.exists(ZIP_REPRESENTATIONS):
    print('Downloading from %s to %s...' % (URL_REPRESENTATIONS, ZIP_REPRESE
    urlretrieve(URL_REPRESENTATIONS, './' + ZIP_REPRESENTATIONS)

if not op.exists(FILE_REPRESENTATIONS):
    print("extracting %s..." % ZIP_REPRESENTATIONS)
    myzip = zipfile.ZipFile(ZIP_REPRESENTATIONS)
    myzip.extractall()
```

In [ ]:
```python
embeddings_index = {}
embeddings_vectors = []
with open('glove100K.100d.txt', 'rb') as f:
    word_idx = 0
    for line in f:
        values = line.decode('utf-8').split()
        word = values[0]
        vector = np.asarray(values[1:], dtype='float32')
        embeddings_index[word] = word_idx
        embeddings_vectors.append(vector)
        word_idx = word_idx + 1

inv_index = {v: k for k, v in embeddings_index.items()}
print("found %d different words in the file" % word_idx)
```

In [ ]:
```python
# Stack all embeddings in a large numpy array
glove_embeddings = np.vstack(embeddings_vectors)
glove_norms = np.linalg.norm(glove_embeddings, axis=-1, keepdims=True)
glove_embeddings_normed = glove_embeddings / glove_norms
print(glove_embeddings.shape)
```

In [ ]:
```python
def get_emb(word):
    idx = embeddings_index.get(word)
    if idx is None:
        return None
    else:
        return glove_embeddings[idx]


def get_normed_emb(word):
    idx = embeddings_index.get(word)
    if idx is None:
        return None
    else:
        return glove_embeddings_normed[idx]
```

```
In [ ]:  get_emb("computer")
```

## Finding most similar words

Here we define a function to find the most similar words to a given word. The similarity is computed using the cosine similarity between the word embeddings. It can also accept multiple words, and it will take the average of the embeddings of the words to find the most similar words.

```
In [ ]:  def most_similar(words, topn=10):
             query_emb = 0
             # If we have a list of words instead of one word
             # (bonus question)
             if type(words) == list:
                 for word in words:
                     query_emb += get_emb(word)
             else:
                 query_emb = get_emb(words)

             query_emb = query_emb / np.linalg.norm(query_emb)

             # Large numpy vector with all cosine similarities
             # between emb and all other words
             cosines = np.dot(glove_embeddings_normed, query_emb)

             # topn most similar indexes corresponding to cosines
             idxs = np.argsort(cosines)[::-1][:topn]

             # pretty return with word and similarity
             return [(inv_index[idx], cosines[idx]) for idx in idxs]
```

```
In [ ]:  most_similar("cpu")
```

```
In [ ]:  most_similar("nvidia")
```

```
In [ ]:  most_similar("1")
```

```
In [ ]:  # bonus: sum of two word embeddings
         most_similar(["toronto", "leaf"])
```

## Displaying vectors with t-SNE

```
In [ ]:  from sklearn.manifold import TSNE

         word_emb_tsne = TSNE(perplexity=30).fit_transform(glove_embeddings_normed[:1
```

```
In [ ]:  import matplotlib.pyplot as plt

         plt.figure(figsize=(40, 40))
         axis = plt.gca()
```

```python
np.set_printoptions(suppress=True)
plt.scatter(word_emb_tsne[:, 0], word_emb_tsne[:, 1], marker=".", s=1)

for idx in range(1000):
    plt.annotate(inv_index[idx],
                 xy=(word_emb_tsne[idx, 0], word_emb_tsne[idx, 1]),
                 xytext=(0, 0), textcoords='offset points')
plt.savefig("tsne.png")
plt.show()
```

## Using pre-trained embeddings in our model

We want to use these pre-trained embeddings for transfer learning. This process is rather similar than transfer learning in image recognition: the features learnt on words might help us bootstrap the learning process, and increase performance if we don't have enough training data.

- We initialize embedding matrix from the model with Glove embeddings:
- take all unique words from our BBC news dataset to build a vocabulary ( MAX_NB_WORDS = 20000 ), and look up their Glove embedding
- place the Glove embedding at the corresponding index in the matrix
- if the word is not in the Glove vocabulary, we only place zeros in the matrix

```python
EMBEDDING_DIM = 100

# prepare embedding matrix
nb_words_in_matrix = 0
nb_words = min(MAX_NB_WORDS, len(word_index))
embedding_matrix = np.zeros((nb_words, EMBEDDING_DIM))
for word, i in word_index.items():
    if i >= MAX_NB_WORDS:
        continue
    embedding_vector = get_emb(word)
    if embedding_vector is not None:
        # words not found in embedding index will be all-zeros.
        embedding_matrix[i] = embedding_vector # Place the Glove embedding a
        nb_words_in_matrix = nb_words_in_matrix + 1

print("added %d words in the embedding matrix" % nb_words_in_matrix)
```

Build a layer with pre-trained embeddings:

```python
pretrained_embedding_layer = Embedding(
    MAX_NB_WORDS, EMBEDDING_DIM,
    weights=[embedding_matrix],
)
```

## A model with pre-trained Embeddings

Now we can build a model with pre-trained embeddings. We will use the same architecture as before, but we will use the `pretrained_embedding_layer` as the first layer of the model.

```python
model = Sequential([
    # Add the pre-defined and pre-trained embedding layer
    pretrained_embedding_layer,
    GlobalAveragePooling1D(),
    Dense(N_CLASSES, activation='softmax')
])

# Set the embedding layer's trainable attribute to False to not fine-tune th
model.layers[0].trainable = False

# Compile the model
model.compile(loss='categorical_crossentropy',
              optimizer=Adam(learning_rate=0.01), metrics=['accuracy'])
```

```python
model.fit(x_train, y_train, validation_split=0.1,
          epochs=15, batch_size=32)
```

## Reality check

On small/medium datasets (few 10,000s) of reasonably large documents (e.g. more than a few paragraphs), simpler classification methods usually perform better, and are much more efficient to train and use. Here are two resources to go further, if you are curious:

- Naive Bayes approach, using scikit-learn http://scikit-learn.org/stable/datasets/twenty_newsgroups.html
- Alec Radford (OpenAI) gave a very interesting presentation, showing that you need a VERY large dataset to have real gains from GRU/LSTM in text classification https://www.slideshare.net/odsc/alec-radfordodsc-presentation

Training deep architectures from random init on text classification is usually a waste of time.

However, when looking at features, one can see that classification using simple methods isn't very robust, and won't generalize well to slightly different domains (e.g. forum posts => emails)

Nowadays, the strategy would be to use pre-trained deep network (BERT) to extract features and fit a linear classifer on top of this. This is especially useful when classifying short texts (e.g. one or a few sentences) as this kind of tasks can be very sensitive to understanding the meaning resulting from intra-sentence interactions between words. The next session on attentional mechanisms and pre-trained transformer-based word models will explain this in more details.