

## 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

However keep in mind:

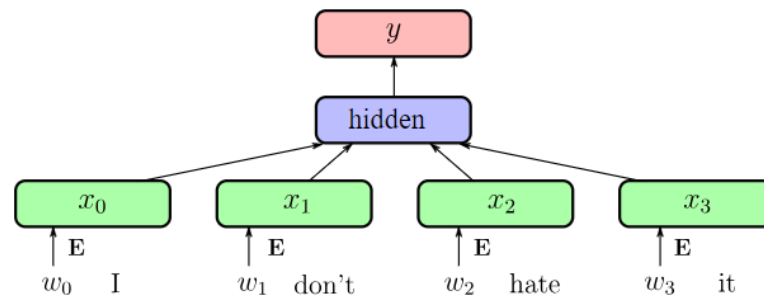
- Deep Learning can be better on text classification than simpler ML techniques, but only on very large datasets and well designed/tuned models.
- We won't be using the most efficient (in terms of computing) techniques, as Keras is good for prototyping but rather inefficient for training small embedding models on text.
- The following projects can replicate similar word embedding models much more efficiently: [word2vec](#) and [gensim's word2vec](#) (self-supervised learning only), [fastText](#) (both supervised and self-supervised learning), [Yowpal Wabbit](#) (supervised learning).
- Plain shallow sparse TF-IDF bigrams features without any embedding and Logistic Regression or Multinomial Naive Bayes is often competitive in small to medium datasets.

### 20 Newsgroups Dataset

The 20 Newsgroups data set is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups <http://qwone.com/~jason/20Newsgroups/>

### 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 20 classes (+ softmax)

### Building more complex models

#### Exercise

- From the previous template, build more complex models using:
  - 1d convolution and 1d maxpooling. Note that you will still need a GlobalAveragePooling or Flatten after the convolutions
  - Recurrent neural networks through LSTM (you will need to reduce sequence length before)

