

# Natural Language Processing (NLP)

The datasets in the early days of Machine Learning were mainly numeric.

This simplified an already difficult problem because numeric data is easily amenable to manipulation by a computer.

But the quantity of non-numeric data, such as Image and Text, is increasing rapidly thanks to the omnipresence of Internet connected devices.

*Natural Language Processing* is the set of techniques that facilitate using *unstructured* text as raw material.

What are some typical tasks

- Is a document *relevant*
- Does the document express a Positive or Negative sentiment
- Summarization of the document

We will provide a brief introduction, based on Neural Network techniques.

That is not to discount more "classical" methods for NLP

- Part of speech tagging
- Stemming
- Lemmatization

All of these are potentially useful as pre-processing steps for Deep Learning.

However, if our data sets are big enough, it may be counter-productive to preprocess.

- We may introduce our own biases
- "Let the data speak for itself"

# Notation

We define a finite set  $\mathbf{V}$  to be the universe of possible tokens that make up text

- We denote the items as *tokens* rather than *words*
- Tokens may involve non-words such as punctuation (e.g., exclamation) and special characters (emoji)
- Tokens may consist of sub-words "un-important"

We denote the  $j^{th}$  token in vocabulary  $\mathbf{V}$  as  $\mathbf{V}_j$ .

Text  $\mathbf{w}$  is a *sequence* of tokens

$$\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(n_{\mathbf{w}})}$$

using parenthesized subscripts to index into sequences, as usual.

$$\mathbf{w}_{(t)} \in \mathbf{V} \text{ for all } 1 \leq t \leq n_{\mathbf{w}}$$

Two non-word tokens are used to denote the start and end of the sequence of tokens in text.

- $\mathbf{w}_{(0)} = \langle \text{START} \rangle$
- $\mathbf{w}_{(n_{\mathbf{w}}+1)} = \langle \text{END} \rangle$

# What is special about NLP ?

We begin the study of NLP by discussing issues that are [particular to text](#)  
([NLP\\_Tokens.ipynb](#)).



# Sentiment classification notebook on Colab : simple model

Let's illustrate the traditional summarization of the sequence by some code for the following Classification task

- Input: Movie review, as sequence of characters
- Label: Positive/Negative

The first model will be extremely simple

- Use One Hot Encodings to represent tokens
- Use Global Pooling to reduce variable length sequences to fixed length
- Followed by Binary Classifier

NLP notebook: examine the data ([https://colab.research.google.com/github/kenperry-public/ML\\_Spring\\_2020/blob/master/Keras\\_examples\\_imdb\\_cnn.ipynb#scrollTo=shHO2I](https://colab.research.google.com/github/kenperry-public/ML_Spring_2020/blob/master/Keras_examples_imdb_cnn.ipynb#scrollTo=shHO2I))

NLP notebook: simple model ([https://colab.research.google.com/github/kenperry-public/ML\\_Spring\\_2020/blob/master/Keras\\_examples\\_imdb\\_cnn.ipynb#scrollTo=QtvUFJz](https://colab.research.google.com/github/kenperry-public/ML_Spring_2020/blob/master/Keras_examples_imdb_cnn.ipynb#scrollTo=QtvUFJz))

---

# Embeddings

One Hot Encoding of tokens is the obvious choice, which is why we used it in our simple model.

But we can do much better

- Learn shorter encoding vectors
- That also capture inter-token relationships

We will study these short vectors in the notebook on [Embeddings](#) ([NLP Embeddings.ipynb](#)).

# Learning embeddings notebook Colab: Sentiment Classification

Let's demonstrate

- How we can learn embeddings
- And use them to hopefully improve upon the simple model

[NLP notebook: learned embeddings](https://colab.research.google.com/github/kenperry-public/ML_Spring_2020/blob/master/Keras_examples_imdb_cnn.ipynb#scrollTo=f5XrUD)

[https://colab.research.google.com/github/kenperry-public/ML\\_Spring\\_2020/blob/master/Keras\\_examples\\_imdb\\_cnn.ipynb#scrollTo=f5XrUD](https://colab.research.google.com/github/kenperry-public/ML_Spring_2020/blob/master/Keras_examples_imdb_cnn.ipynb#scrollTo=f5XrUD)

---

# Using Convolutions for Text

Our models still fail to take advantage of features particular to text.

Because we used Global Pooling to summarize the sequence

- We lost all information about token (word) order

We don't have to resort to the "ultimate" Layer-type for sequences (RNN) to take at least some advantage of order.

A Convolutional Layer can capture some limited sense of the temporal ordering.

Each kernel in the Convolution can

- Create a synthetic feature
- That is a combination of consecutive tokens
- (What was a "spatial" dimension in, e.g., an image, is now a "temporal" dimension)

Let's learn about the concept of [neural n-grams \(NLP Convolution.ipynb\)](#).

# Sentiment classification notebook on Colab : Convolutional n-grams

Let's augment our Sentiment Classification notebook with neural n-grams.

[NLP notebook: neural n-grams \(https://colab.research.google.com/github/kenperry-public/ML\\_Spring\\_2020/blob/master/Keras\\_examples\\_imdb\\_cnn.ipynb#scrollTo=LPChvdI\)](https://colab.research.google.com/github/kenperry-public/ML_Spring_2020/blob/master/Keras_examples_imdb_cnn.ipynb#scrollTo=LPChvdI)

---

# Conclusion

Machine Learning on textual data is becoming an increasingly important source of insight in many domains.

NLP brought together many elements of the knowledge we gained in this course

- One Hot Encoding of tokens
- Reducing variable length sequence inputs to fixed length (RNN)
- Classical ML (Classification)

So, as useful as this lecture was in its own right, we hope it re-inforces everything you've learned throughout this course.



```
In [ ]: print("Done")
```