# COMP348 — Document Processing and the Semantic Web

Week 05 Lecture 1: Processing Text Sequences

Diego Mollá

Department of Computer Science
Macquarie University

COMP348 2019H1

## Programme

1. Word Embeddings
   - Challenges of Text for Machine Learning
   - Word Embeddings

2. Text Sequences
   - Modelling Text Sequences
   - Sequence Labelling

### Reading

- Deep Learning book, chapter 6.
- Understanding LSTM Networks,
  https://colah.github.io/posts/2015-08-Understanding-LSTMs/.

## Programme

## Programme

# Words as Arbitrary Symbols

- Words are encoded as arbitrary symbols.
- Within one language there is no clear correspondence between a word symbol and its meaning.
  - "dig" vs. "dog"
  - "car" vs. "automobile"
- Different languages may use different representations of the same word.



https://en.wikipedia.org/wiki/File:Hello_in_different_languages_v

## Ambiguities Everywhere
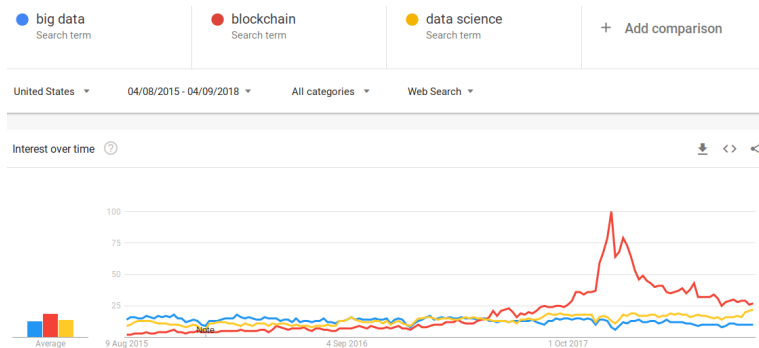
Language features ambiguity at multiple levels.

### Lexical Ambiguity

Example from Google's dictionary:

- bank (n): the land alongside or sloping down a river or lake.
- bank (n): financial establishment that uses money deposited by customers for investment, . . .
- bank (v): form in to a mass or mound.
- bank (v): build (a road, railway, or sports track) higher at the outer edge of a bend to facilitate fast cornering.
- . . .

# So many words!

- Any language features a large number of distinct words.
- New words are coined.
- Words change their use in time.
- There are also names, numbers, dates... an infinite number.



https://trends.google.com

## Long-distance Dependencies

- Sentences are sequences of words.
- Words close in the sentence are often related.
- But sometimes there are relations between words far apart.

grammatical: "The man living upstairs is very cheerful"
"The people living upstairs are very cheerful"

reference: "I bought a book from the bookshp and I liked it"

knowledge: "I was born in France and I speak fluent French"
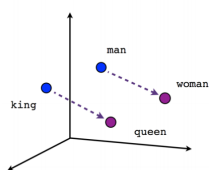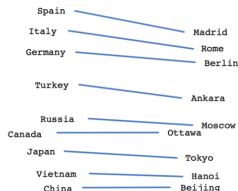
# Programme

# Word Embeddings

- First introduced in 2013, nowadays is one of the most common ingredients in text processing systems.
- Word embeddings squarely aim at addressing the issue of representing words as continuous vectors of integers.
- Words with similar context are mapped to similar vectors.
- Embeddings are learnt using large, unlabelled training data.



Male-Female          Verb tense          Country-Capital
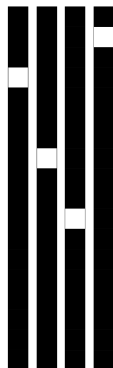
https://www.tensorflow.org/tutorials/representation/word2vec
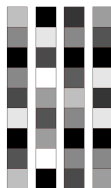
# One-hot vs. word embeddings

### One-hot

- Sparse
- Binary values (typically)
- High-dimensional
- Hard-coded

### Word embeddings

- Dense
- Continous values
- Lower-dimensional
- Learned from data

# Two Ways to Obtain Word Embeddings

1. Learn the word embeddings jointly with the task you care about (e.g. document classification).
2. Use pre-trained word embeddings.

## Learning Word Embeddings

- You can add a dense layer as the first layer of your network and let the system learn the optimal weights.
- This approach is so useful and common that many deep learning frameworks define an "embedding" layer that facilitates this.
- The input to the "embedding" layer is the word index.
- The output is the word embedding.

# Using pre-trained word embeddings

- Sometimes we have so little training data that many words are poorly represented.
- Often, words in the training data do not occur in the test data.
- For these words we would not be able to learn the embeddings.
- Several people have computed word embeddings for large vocabularies using large data sets.
- We can then use these pre-trained embeddings to map from the word index to the word embedding.

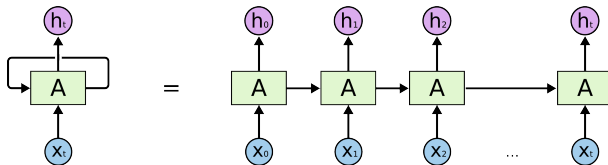# Programme

# Programme

# Handling Text Sequences

- A document is a sequence of words.
- Many document representations are based on a bag-of-words approach.
  - Word order is ignored.
- Even word embeddings ignore word order.
- A Recurrent Neural Network (RNN), however, is designed to process sequences.



https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# A Recurrent Neural Network

- A RNN is a neural network that is composed of RNN cells.
- Each RNN cell takes as input two pieces of information:
  1. A vector representing an item in the sequence.
  2. The state resulting from processing the previous items.
- The output of the RNN cell is a state that can be fed to the next cell in the sequence.
- All RNN cells are identical copies. In a sense, we can say that an RNN cell is the same for all words in the sequence.



https://colah.github.io/posts/2015-08-Understanding-LSTMs/

## Recurrent Neural Networks

- RNNs are designed to model long-distance dependencies.
- Vanilla RNN cells were used since 1990s.



https://colah.github.io/posts/2015-08-Understanding-LSTMs/

# LSTMs and GRUs

- Vanilla RNN cells are still too simple and they do not memorise long-distance dependencies easily.
- More complex RNN cells have been designed specifically to address this issue.
- These RNN cells have components that are programmed to memorise or forget past information.
- Current most popular RNN cells are:

    LSTM Long Short Term Memory (picture).
    GRU Gated Recurrent Unit; a more recent, simpler cell.

## RNNs in Practice

- Most deep learning frameworks include special layers for RNNs.
- When you use an RNN layer, you have the option to specify the type of RNN cell.
- You often have the option to use the state of the last cell, or the state of all cells.

# Programme

# What is Sequence Labelling?

- A sequence labelling problem is one where:
  - the input consists of a sequence $X = (X_1, \ldots, X_n)$, and
  - the output consists of a sequence $Y = (Y_1, \ldots, Y_n)$ of labels, where:
  - $Y_i$ is the label for element $X_i$
- Example: Part-of-speech tagging

$$
\left( \begin{array}{c} Y \\ X \end{array} \right) = \left( \begin{array}{ccc} \text{Verb,} & \text{Determiner,} & \text{Noun} \\ \text{spread,} & \text{the,} & \text{butter} \end{array} \right)
$$

- Example: Spelling correction

$$
\left( \begin{array}{c} Y \\ X \end{array} \right) = \left( \begin{array}{ccc} \text{write,} & \text{a,} & \text{book} \\ \text{rite,} & \text{a,} & \text{buk} \end{array} \right)
$$

# Other applications of sequence labelling

- Named entity recognition and classification (NER) involves finding the named entities in a text and identifying what type of entity they are (e.g., person, location, corporation, dates, etc.).
- Speech transcription can be seen as a sequence labelling task:
  - The input $X = (X_1, \ldots, X_n)$ is a sequence of acoustic frames $X_i$, where $X_i$ is a set of features extracted from a 50msec window of the speech signal.
  - The output $Y$ is a sequence of words (the transcript of the speech signal).
- Financial applications of sequence labelling:
  - Identifying trends in price movements.
- Biological applications of sequence labelling:
  - Gene-finding in DNA or RNA sequences.

# Sequence Labelling as Classification I

Can we just use a standard classifier?

- Standard classifiers (K-Nearest Neighbours, Naïve Bayes, Support Vector Machine, . . . ) assume independence between samples:
    - The probability of the label assigned to sample $i$ is independent to the probability of the label assigned to sample $j$.
- But in sequence labelling there is interdependence between the labels of different samples.

# Modelling Context

### Classifier with context features

- A (crude) approach to model interdependence between samples is to add context features.
- For example, we can use features based on previous words and following words.
- We can even incorporate the label of the previous word as a feature.
- But it is not so easy to incorporate the label of both the previous word and the following word.

# Using Recurrent Neural Networks for Sequence Labelling I

- We have seen how RNNs can be used to classify documents.
- Similarly, we can use RNNs to classify sequences of words.
- In Keras, we can define a dense layer that is repeated at the output of each recurrent cell.

```
model = Sequential()
model.add(Embedding(max_words, 32))
model.add(LSTM(32, return_sequences=True))
model.add(TimeDistributed(Dense(num_tags,
                        activation='softmax')))

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

# Using Recurrent Neural Networks for Sequence Labelling II

```
Layer (type)                    Output Shape            Param #
=================================================================
embedding_1 (Embedding)         (None, None, 32)        320000

_____
lstm_1 (LSTM)                   (None, None, 32)        8320

_____
time_distributed_1 (TimeDist    (None, None, 20)        660
=================================================================
```

## Take-home Messages

1. Explain some of the fundamental challenges that plain text represents to machine learning.
2. Apply word embeddings in deep learning.
3. Use recurrent neural networks for text classification.
4. Comment on the issues of sequence labelling.

## What's Next

Week 6

- Generating text.
- Reading: Deep Learning book, chapter 8.1