**MS4S09**

**Topic Modelling**



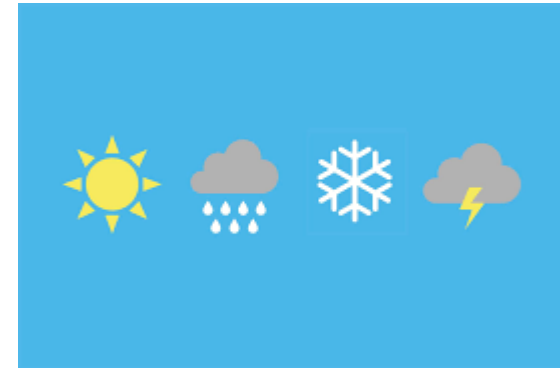https://learn.datacamp.com/courses/topic-modeling-in-r

# Topic Modelling

Topics give us an idea what text is about quickly.

A topic is a label for a collection of words that often occur together.

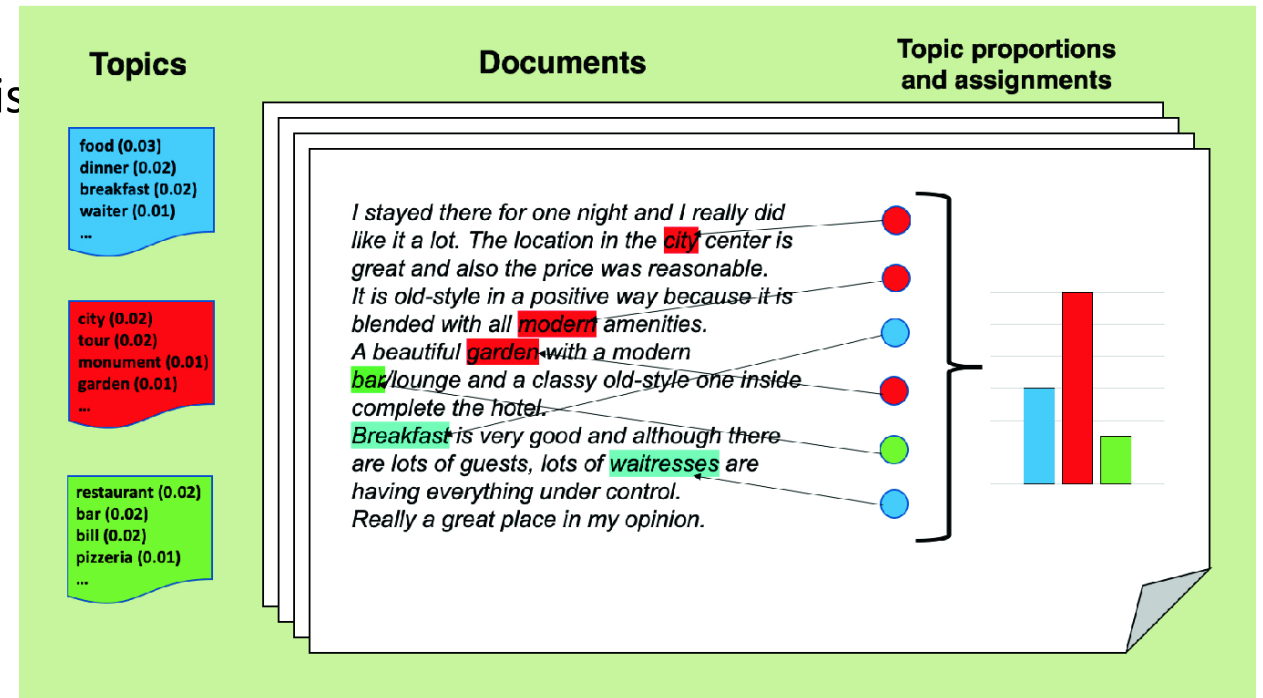What words would you associate with the **weather**?

# Algorithms for Topic Models

The most common algorithm for topic modelling is called the Latent Dirichlet Allocation (LDA).

LDA takes a document-term matrix as its input.

LDA returns two matrices:
- Prevalence of topics in documents
- Probability of words belonging to topics

# Topic Modelling

**Every document is a mixture of topics.**
For example, in a two-topic model we could say "Document 1 is 90% topic A and 10% topic B, while Document 2 is 30% topic A and 70% topic B."

**Every topic is a mixture of words.**
For example, we could imagine a two-topic model of the news, with one topic for "politics" and one for "entertainment."

# Topic Modelling

LDA is a mathematical method for estimating both at the same time: finding the mixture of words that is associated with each topic, while also determining the mixture of topics that describes each document.

Need to decide what words to keep.
        Remove common words, stop words etc.

We need to know beforehand how many topics we want. 'k' is pre-decided.
        This can be an iterative process.

# Topic Modelling Example

The AssociatedPress dataset provided by the topicmodels package, is an example of a DocumentTermMatrix.

This is a collection of 2246 news articles from an American news agency, mostly published around 1988.

We will use the LDA() function from the topicmodels package, setting k = 2, to create a two-topic LDA model.