# Applied Text Analytics &
# Natural Language Processing

with Dr. Mahdi Roozbahani
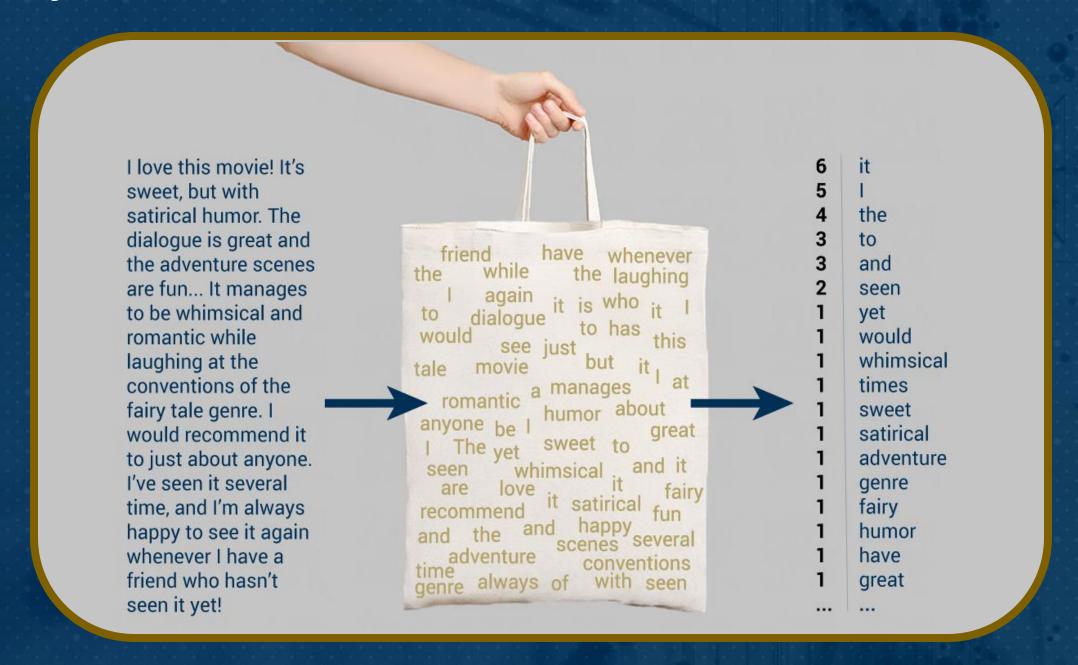& Wafa Louhichi

GT

# Learning Objectives

In this lesson, you will learn another discrete text representation

- Explain **text data** into numerical format using Bag of Words

- Understand the advantages and disadvantages of Bag of Words

GT

# Bag of Words Model

- Represent each **document** as a **bag of words**, ignoring words' ordering for **simplicity**.

# Bag of Words Model

- Represent a document as a column vector $X$ of word counts

- The bag of words is a fixed-length representation, which consists of a vector of word counts:

    Document: $\mathrm{D} = it\ was\ the\ best\ of\ times, it\ was\ the\ worst\ of\ times$

    Vocabulary: $V = [aardvak, \dots, it, \dots, best, \dots, times, \dots, zyther]$

    Bag of Words: $X = [0, \dots, 2, \dots, 1, \dots, 2, \dots, 0]$

- The size of $X$ is $1xd$ where $d$ is the size of the vocabulary.

- A collection of $n$ documents is represented with a matrix of size $nxd$

# Example of Bag of Words Using Sklearn

```python
1 from sklearn.feature_extraction.text import CountVectorizer
2 my_text=["it was the best of times, it was the worst of times"]
3 vectorizer = CountVectorizer()
4 X = vectorizer.fit_transform(my_text)
5 print(X.toarray())
6 vectorizer.get_feature_names()
```

```
[[1 2 2 2 2 2 1]]
['best', 'it', 'of', 'the', 'times', 'was', 'worst']
```

# Advantages and Disadvantages of Bag of Words

- Advantages of Bag of Words:

    - Simple and easy to implement

- Disadvantages of Bag of Words:

    - Every document is represented as a vector of the size of the vocabulary: not scalable for a large vocabulary (100,000 words)

    - High dimensional sparse matrix which can be memory & computationally expensive

    - The order of words is disregarded and thus the meaning coming from the context is not captured

**GT**

# Summary

- Use Bag of Words as a discrete text representation method to represent documents

- Understand the advantages and disadvantages of this method