

Natural Language Processing(NLP)

معالجة اللغات الطبيعية

تطبيق عملي-الدرس الاول

(تصفية الرسائل غير المرغوب فيها) SMS Spam Filtering

Natural Language Processing

NLP

Ahmad Shhadeh

معالجة اللغات الطبيعية

From the previous session

1. What is NLP?
2. Structured Data vs. Unstructured Data
3. NLP Components
4. NLP, Artificial intelligence & Machine Learning
5. Main approaches in NLP(TimeLine)
6. Why NLP is very important?
7. Natural Language Processing Applications
8. Areas That Leverage NLP Technology
9. Why is NLP so difficult?
10. What are the techniques used in NLP?
11. Libraries and tools
12. Arabic Natural Language Processing
13. The Future of NLP

Before we start !!!

What you should know before you start



- 1- Understanding of some of the key concepts in **natural language processing** and **machine learning algorithms**
- 2- Basic Knowledge in **Python**
- 3- Some experience using the **NumPy, pandas and scikit-learn** libraries.

Environment setup

Installing Python(IDE like anaconda) and **Jupyter Notebooks**
OR
Google Colab



Colab is a free cloud service based on Jupyter Notebooks for machine learning education and research. It provides a runtime fully configured for deep learning and **free-of-charge** access to a **robust GPU**.



01 Free-Of-Charge



02 Cloud Service



03 Jupyter Notebook Environment



04 Zero Configuration Required



05 Access to GPU /TPU

Python Natural Language Processing (NLP) libraries

Natural
Language
Toolkit (NLTK)

TextBlob

CoreNLP

Gensim

spaCy

polyglot

scikit-learn

Pattern

Libraries and tools

▶ NLTK

- Small but useful datasets with markup
- Preprocessing tools: tokenization, normalization...
- Pre-trained models for POS-tagging, parsing...

▶ Stanford parser

▶ spaCy:

- python and cython library for NLP

▶ Gensim

- python library for text analysis, e.g. for word embeddings and topic modeling

▶ MALLET

- Java-based library, e.g. for classification, sequence tagging, and topic modeling

• Ahmad Shhadeh

Comparison of Python NLP libraries Pros and Cons

	⊕ PROS	⊖ CONS
	<ul style="list-style-type: none"> + The most well-known and full NLP library + Many third-party extensions + Plenty of approaches to each NLP task + Fast sentence tokenization + Supports the largest number of languages compared to other libraries 	<ul style="list-style-type: none"> - Complicated to learn and use - Quite slow - In sentence tokenization, NLTK only splits text by sentences, without analyzing the semantic structure - Processes strings which is not very typical for object-oriented language Python - Doesn't provide neural network models - No integrated word vectors
	<ul style="list-style-type: none"> + The fastest NLP framework + Easy to learn and use because it has one single highly optimized tool for each task + Processes objects; more object-oriented, comparing to other libs + Uses neural networks for training some models + Provides built-in word vectors + Active support and development 	<ul style="list-style-type: none"> - Lacks flexibility, comparing to NLTK - Sentence tokenization is slower than in NLTK - Doesn't support many languages. There are models only for 7 languages and "multi-language" models
	<ul style="list-style-type: none"> + Has functions which help to use the bag-of-words method of creating features for the text classification problems + Provides a wide variety of algorithms to build machine learning models + Has good documentation and intuitive classes' methods 	<ul style="list-style-type: none"> - For more sophisticated preprocessing things (for example, pos-tagging), you should use some other NLP library and only after it you can use models from scikit-learn - Doesn't use neural networks for text preprocessing
	<ul style="list-style-type: none"> + Works with large datasets and processes data streams + Provides tf-idf vectorization, word2vec, document2vec, latent semantic analysis, latent Dirichlet allocation + Supports deep learning 	<ul style="list-style-type: none"> - Designed primarily for unsupervised text modeling - Doesn't have enough tools to provide full NLP pipeline, so should be used with some other library (Spacy or NLTK)
	<ul style="list-style-type: none"> + Allows part-of-speech tagging, n-gram search, sentiment analysis, WordNet, vector space model, clustering and SVM + There are web crawler, DOM parser, some APIs (like Twitter, Facebook etc.) 	<ul style="list-style-type: none"> - Is a web miner; can be not enough optimized for some specific NLP tasks
	<ul style="list-style-type: none"> + Supports a large number of languages (16-196 languages for different tasks) 	<ul style="list-style-type: none"> - Not as popular as, for example, NLTK or Spacy; can be slow issues solutions or weak community support

Python Natural Language Processing (NLP) libraries

- 1) Natural Language Toolkit (NLTK)
- 2) TextBlob
- 3) CoreNLP
- 4) Gensim
- 5) spaCy
- 6) polyglot
- 7) scikit-learn
- 8) Pattern

Natural Language Toolkit(NLTK)

- ▶ **The natural language toolkit** is the most utilized package for handling natural language processing tasks in Python. Usually called NLTK for short.
- ▶ NLTK is a **leading platform for building Python programs** to work with human language data. It provides easy-to-use interfaces to over 100 corpora and lexical resources such as WordNet.
- ▶ Providing a suite of text processing libraries for **classification, tokenization, stemming, tagging, parsing, and semantic reasoning**, wrappers for industrial-strength NLP libraries
- ▶ NLTK has been called “**a wonderful tool for teaching, and working in, computational linguistics using Python,**” and “**an amazing library to play with natural language.**”
- ▶ It is a suite of open-source tools originally created in **2001 at the University of Pennsylvania** for the purpose of making building NLP processes in Python easier.
- ▶ This package has been **expanded** through the extensive **contributions of open-source users** in the years since its original development.

NLP Tool Kit installation

- ▶ How to install NLTK on your local machine:
 - ▶ We assumed that python and anaconda both are installed
 - ▶ `pip install nltk` → installing nltk
 - ▶ `dir(nltk)` → check all installed packages under NLTK

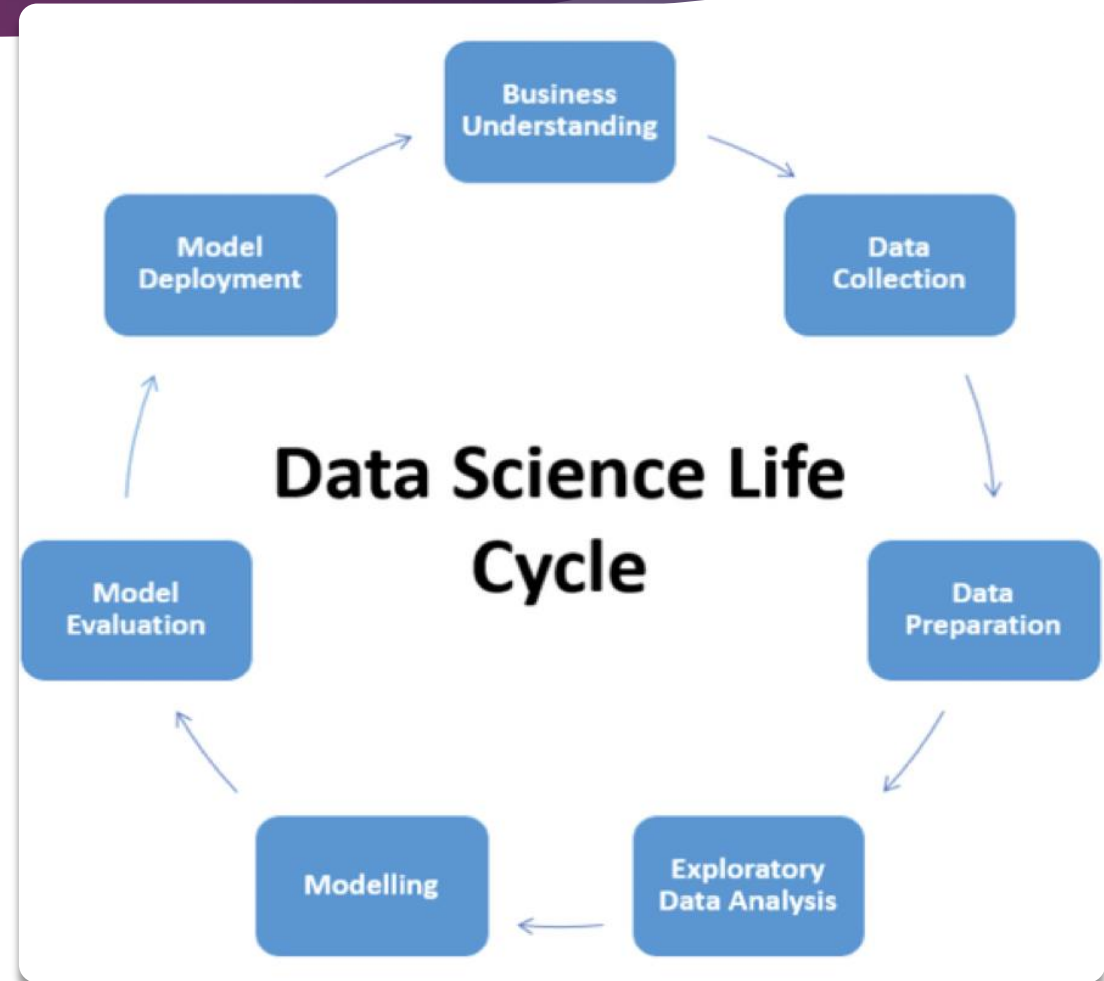


Email/SMS Spam Filtering Use case



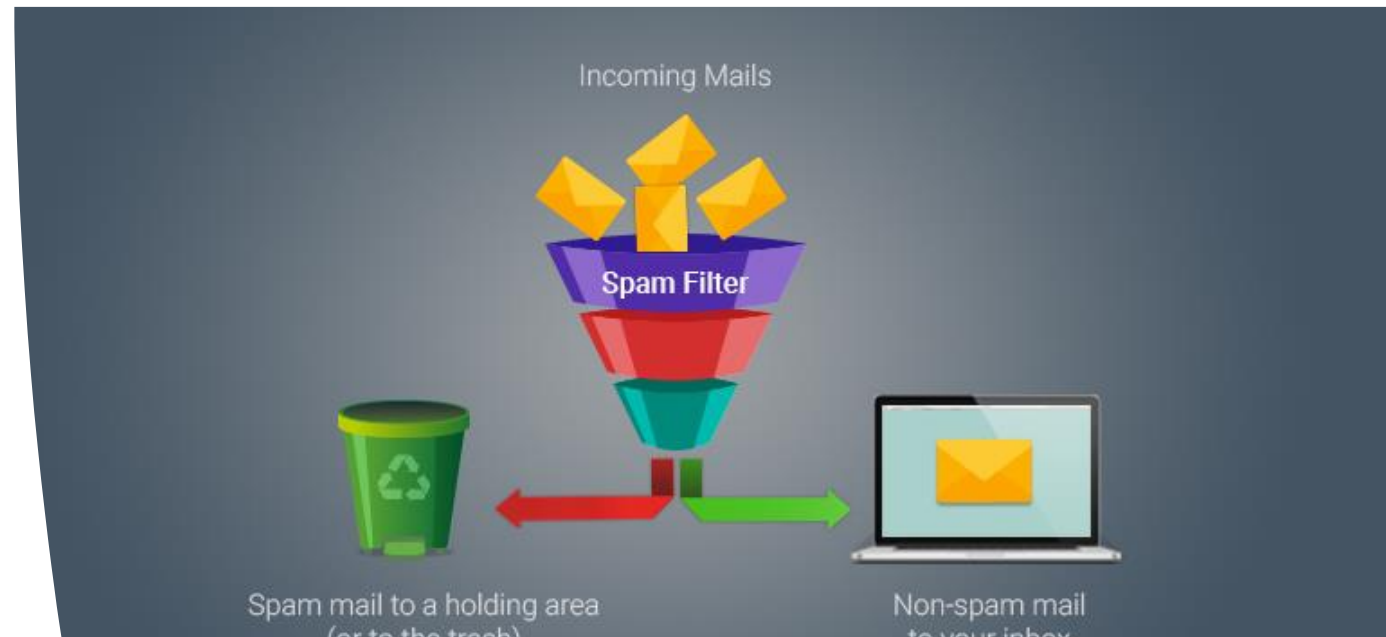
Data Science life cycle

1. **Business Understanding**
2. Data Collection
3. Data Preparation
4. Exploratory data analytics(EDA)
5. Data Modelling
6. Model Evaluation
7. Model Deployment



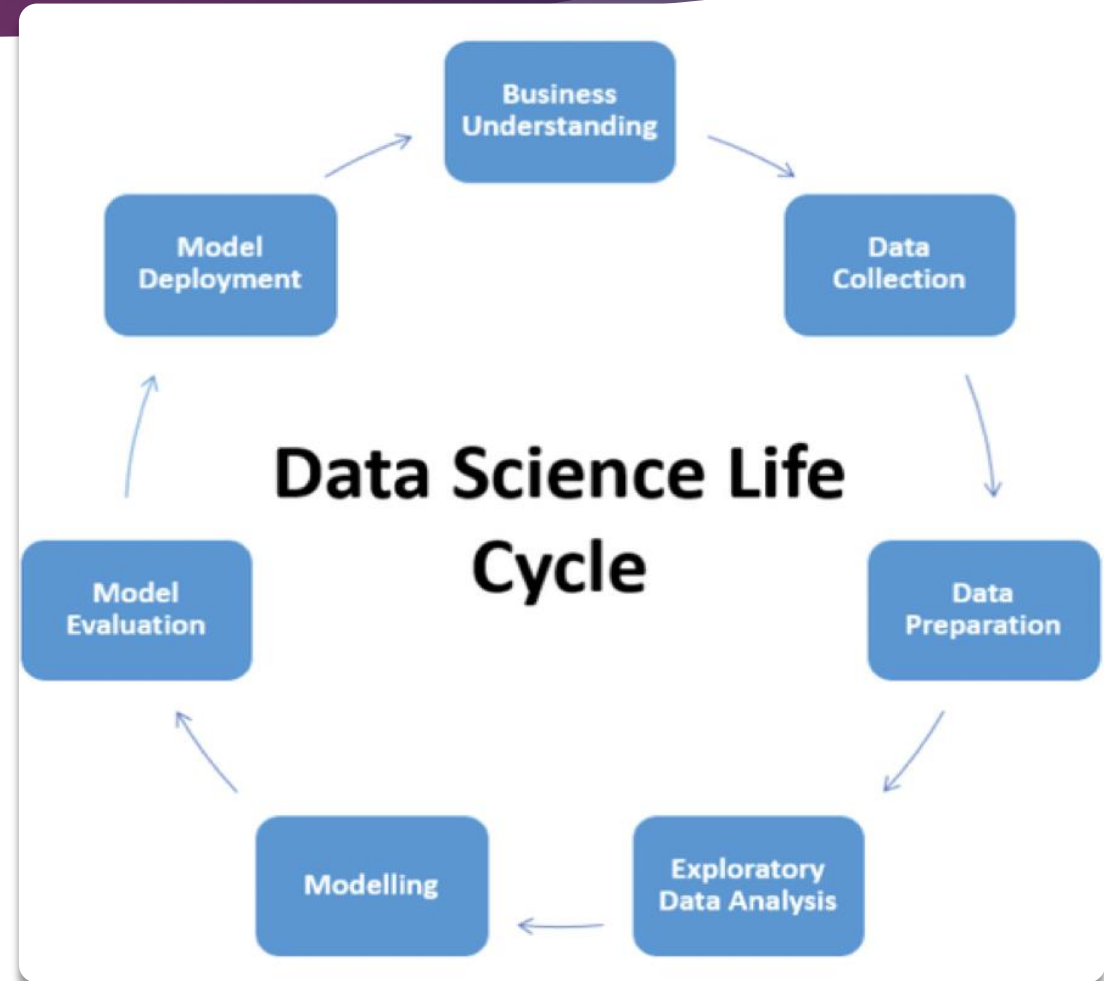
Email/SMS Spam Filtering

- ▶ Spam filtering is a beginner's example of document classification task which involves classifying an email as spam or non-spam (a.k.a. ham) mail.
- ▶ Spam box in your Gmail account is the best example of this.
- ▶ So let's get started in building a spam filter on a publicly available mail corpus.



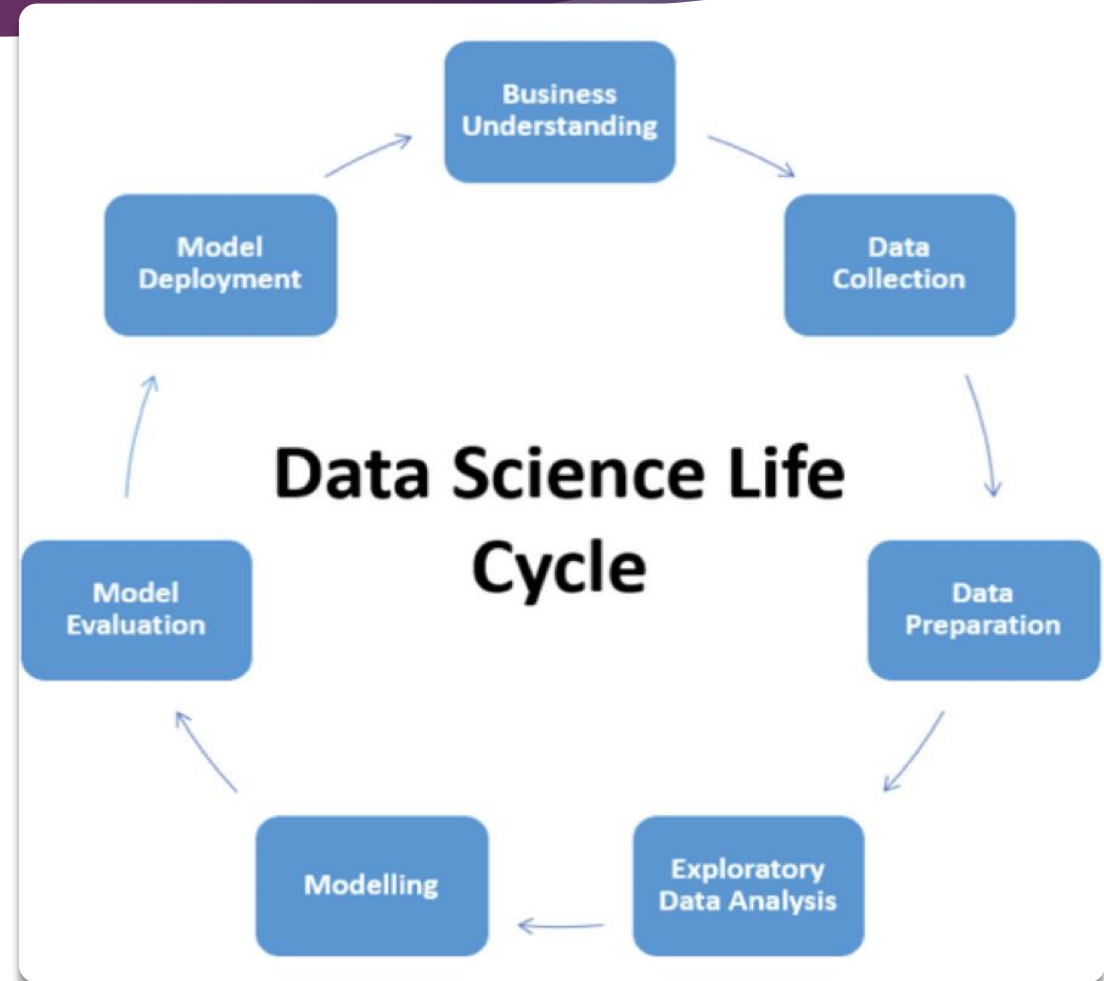
Data Science life cycle

- ▶ Business Understanding
- ▶ **Data Collection**
- ▶ Data Preparation
- ▶ Exploratory data analytics(EDA)
- ▶ Data Modelling
- ▶ Model Evaluation
- ▶ Model Deployment



Data Science life cycle

- ▶ Business Understanding
- ▶ **Data Collection**
- ▶ Data Preparation
- ▶ Exploratory data analytics(EDA)
- ▶ Data Modelling
- ▶ Model Evaluation
- ▶ Model Deployment



Public and Private Data

- ▶ **Private data**: it is **private and belongs to an organization**, and there are certain **security and privacy** concerns attached to it. It is used for the companies' **internal analysis** purposes in order to gain business and growth insights. Some examples of such organizational private data are telecom data, retail data, and banking and medical data.
- ▶ **Public data**: This is the data that is **available for public** use and is offered by many sites such as government websites and public agencies for the purpose of research. Accessing this data does not require any special permission or approval.
- ▶ **Open data** is the idea that some data should be freely available to everyone to use and republish as they wish, **without restrictions from copyright**, the goals of the open-source data movement are similar to those of other "open(-source)" movements such as **open-source software**

Data Collection

► SMS Spam Collection Data Set

Download: [Data Folder](#), [Data Set Description](#)

- Abstract: The SMS Spam Collection is a public set of SMS labeled messages that have been collected for mobile phone spam research.



Data Set Characteristics:	Multivariate, Text, Domain-Theory	Number of Instances:	5574	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	N/A	Date Donated	2012-06-22
Associated Tasks: Ahmad Shhadeh	Classification, Clustering	Missing Values?	N/A	Number of Web Hits:	331230

NEXT LESSON

Reading and exploring the data