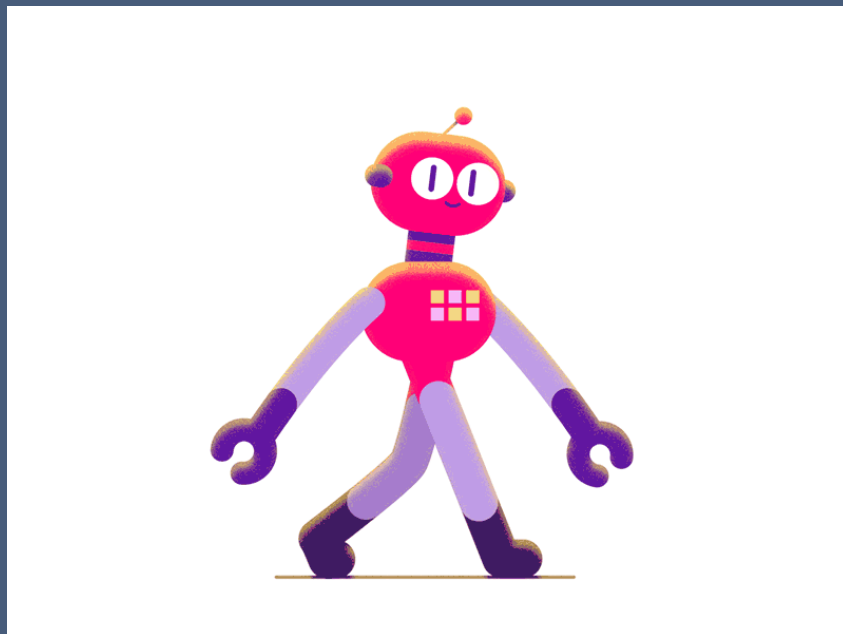


Bilbo the bot



A fancy NLP project by

[@leosanchezsoler](#)

The Bridge: Digital Talent Accelerator

March 2021

Contents

Prelude	3
Goals	4
Specifications	4
Hardware	4
Software	5
Requirements	5
Steps	5
Research the context	5
Get the data	5
Data Wrangling	5
Data Mining/Data Cleaning	6
Text preprocessing	6
Training the model	6
Deploying to production	6
Conclusions	7
Code conclusions	7
Personal conclusions	7
Sources	8
Data Source	8
Solve code doubts	8
Organization Tools	8
Git repository	9

Software

Python 3.8.8. has been the main software used.

Requirements

To execute the program, all tools and datasets are inside the program, in order to make things as easy as possible for the user.

Steps

1. Research the context

In this particular case, there was no need to do an exhaustive research, as recommendation engines documentation is easy to find on the Internet.

2. Get the data

The data was obtained from this Github repository: <https://github.com/NeelShah18/arxivData>. The dataset contains thousands of Data Science papers that were collected using web scraping.

Despite the fact that some variables were not properly formatted, after inspecting its structure it was considered valid for training a model.

3. Data Wrangling

For this project, Data Wrangling took some time to be done. As the info was collected using web scraping, most of them still had the html structure. In order to structure the data, it was necessary to iterate through dictionaries and extract the relevant tags.

