

# Title and Group:

Group members: Chow Jun Wei ([jun.chow.18@ucl.ac.uk](mailto:jun.chow.18@ucl.ac.uk), Brunei Darussalam, NIL, Data Science and NLP and Data Analyst), Esraa Sultan.

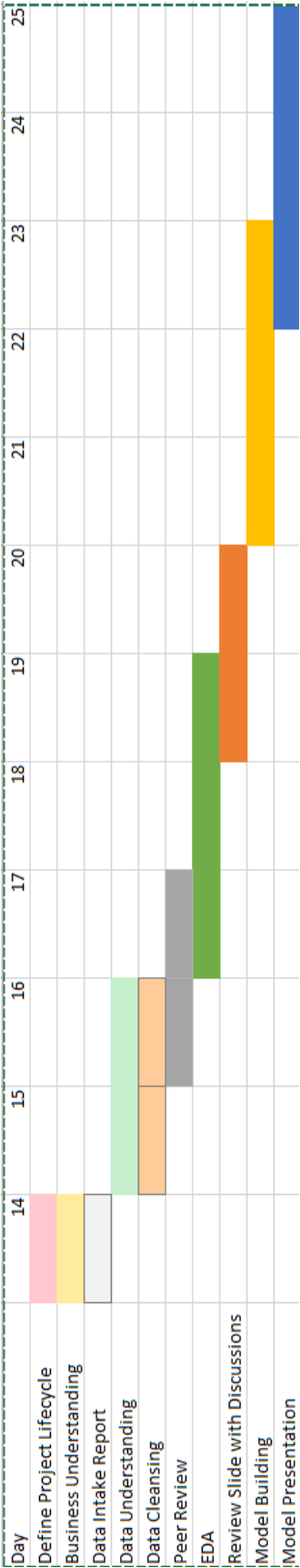
Internship Batch: NLP 02

# Problem Description

There are 20 different categories of articles, and our job is to classify, to which category does each article belongs to.

# Business Understanding

Working for a newspaper office, our boss task us with classification of the news article being written. Previously, it was being classified by human (particularly the article writer). Our boss thought that article writer should focus on writing articles than classification job, and a human classifier is too expensive to hire, so he would like to develop an ML model to do the classification job.



# Project Lifecycle

# Data Intake Report

Name: NLP Group Project

Report date: November 14, 2021

Internship Batch: NLP 02

Version:<1.0>

Data intake by: Chow Jun Wei, Esraa Sultan

Data intake reviewer:<intern who reviewed the report>

Data storage location: [https://github.com/Wabinab/NLP\\_GroupProject\\_DG](https://github.com/Wabinab/NLP_GroupProject_DG)

**Note: Since we have too many files, we will list their folders instead.**

**Note: At the end these files aren't used, rather we changed to this:**

```
from sklearn.datasets import fetch_20newsgroups  
  
all_xs, all_y = fetch_20newsgroups(subset="all",  
    remove=('headers', 'footers', 'quotes'),  
    shuffle=True, return_X_y=True)
```

The files are arranged such that we retain the original data method: each article are their own .txt files. There are more than 1000 files so if we list them here it would take too long. Rather, we would group by each category instead and mention how many files there are. Hence, “Total number of observations” and “Total number of features” would be the NIL for all. We changed “Total number of features” to “Base Folder”. And one file called “errors.txt” containing the files that cannot be processed due to reasons (mostly due to cannot decode with UTF-8 and we aren't sure about what encoding it uses so it's ignored).

## Tabular data details:

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	387
<b>Base Folder</b>	Alt.atheism
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.1MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	185
<b>Base Folder</b>	Comp.graphics
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.5MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	184
<b>Base Folder</b>	Comp.os.ms-windows.misc
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.0 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	195
<b>Base Folder</b>	Comp.sys.ibm.pc.hardware
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 924 kB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	132
<b>Base Folder</b>	Comp.sys.mac.hardware
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 624 kB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	249
<b>Base Folder</b>	Comp.windows.x
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.8 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	180
<b>Base Folder</b>	Misc.forsale
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 820 KB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	234
<b>Base Folder</b>	Rec.autos
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.1 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	168
<b>Base Folder</b>	Rec.motorcycles
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 744 KB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	271
<b>Base Folder</b>	Rec.sport.baseball

<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.3 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	310
<b>Base Folder</b>	Rec.sport.hockey
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.7M

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	321
<b>Base Folder</b>	Sci.crypt
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.0 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	193
<b>Base Folder</b>	Sci.electronics
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 900 KB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	277
<b>Base Folder</b>	Sci.med
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.7 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	272
<b>Base Folder</b>	Sci.space
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 1.6 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	442
<b>Base Folder</b>	Soc.religion.christian
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.5 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	400
<b>Base Folder</b>	Talk.politics.guns
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.2 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	530
<b>Base Folder</b>	Talk.politics.mideast
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 3.5 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	450
<b>Base Folder</b>	Talk.politics.misc
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.8 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	384
<b>Base Folder</b>	Talk.religion.misc
<b>Base format of the file</b>	.txt
<b>Size of the data</b>	Total: 2.1 MB

<b>Total number of observations</b>	NIL
<b>Total number of files</b>	1
<b>Total number of features</b>	NIL
<b>Base format of the file</b>	errors.txt
<b>Size of the data</b>	4.1 kB

**Note: Replicate same table with file name if you have more than one file.**

**Proposed Approach:**

- Mention approach of dedup validation (identification)
- Mention your assumptions (if you assume any other thing for data quality analysis)

**Note: Convert this doc in pdf and provide the link of pdf file in your dashboard.  
Please do not forget to remove this section while converting the file into pdf.**