

---

# 2R Assignment 2 Documentation

*Release 2025*

Constantina Maltezou

Sep 01, 2025

## 1 Introduction

This is the documentation for my second class assignment **2R Analysing Multimodal Language Data for Quantitative Social Science**.

1. `custom_packages.misc_modules.config`
2. `custom_packages.misc_modules.styling_and_animations`
3. `custom_packages.misc_modules.visualisations`

The code provided here is designed to:

- **Preprocess** the CSV data relating to extracted acoustic features and speaker demographics.
- **Classify** speech based on its intent (trustworthy vs. neutral).
- **Validate** the dataset's effectiveness in categorising speech according to speaker intent and demographics.
- **Visualise** the relationship between speaker demographics, acoustic features and trustworthy intent classifications.

For detailed instructions on how to set up the environment, run the analysis scripts, and interpret the results, refer to the **Getting Started** section.

## 2 Section 2

### 2.1 Module: `custom_packages.misc_modules.config`

```
get_openai_client
```

```
get_gemini_client
```

```
get_youtube_client
```

`custom_packages.misc_modules.config.get_openai_client`

`get_openai_client()`

## custom\_packages.misc\_modules.config.get\_gemini\_client

`get_gemini_client()`

## custom\_packages.misc\_modules.config.get\_youtube\_client

`get_youtube_client()`

## Module: custom\_packages.misc\_modules.styling\_and\_animations

<code>print_error_message</code>	Displaying a red-coloured error message in the console.
<code>print_success_message</code>	Displaying a green-coloured success message in the console.
<code>print_warning_message</code>	Displaying an orange-coloured warning message in the console.
<code>print_highlighted_message</code>	Displaying a yellow-highlighted message in the console.
<code>print_animated_ellipsis_message</code>	Displaying a message in the console, followed by a simple, animated ellipsis to indicate progress.

## custom\_packages.misc\_modules.styling\_and\_animations.print\_error\_message

`print_error_message(message, e)`

Displaying a red-coloured error message in the console.

### Parameters

- **message** (*str*)
- **e** (*Exception*)

### Return type

*None*

## custom\_packages.misc\_modules.styling\_and\_animations.print\_success\_message

`print_success_message(message)`

Displaying a green-coloured success message in the console.

### Parameters

**message** (*str*)

### Return type

*None*

## custom\_packages.misc\_modules.styling\_and\_animations.print\_warning\_message

`print_warning_message(message)`

Displaying an orange-coloured warning message in the console.

### Parameters

**message** (*str*)

### Return type

*None*

## custom\_packages.misc\_modules.styling\_and\_animations.print\_highlighted\_message

**print\_highlighted\_message** (*message*)

Displaying a yellow-highlighted message in the console.

### Parameters

**message** (*str*)

### Return type

*None*

## custom\_packages.misc\_modules.styling\_and\_animations.print\_animated\_ellipsis\_message

**print\_animated\_ellipsis\_message** (*message='Currently processing', duration=4.0, sleep\_interval=0.2*)

Displaying a message in the console, followed by a simple, animated ellipsis to indicate progress.

### Parameters

- **message** (*str, optional*) – The message to be displayed in the console. The default is “Currently processing”.
- **duration** (*float, optional*) – Indicates the total duration (seconds) of the animated ellipsis. The default is 4.0.
- **sleep\_interval** (*float, optional*) – Indicating how many seconds needed before each ellipsis update. The default is 0.2.

### Return type

*None*

## Module: custom\_packages.misc\_modules.visualisations

<i>plot_word_frequencies</i>	Shows a bar plot with the most frequent words across all transcripts.
<i>word_cloud_visualisation</i>	Visualises a word cloud of the top keywords across all transcripts.
<i>plot_sentiment_distribution</i>	Shows a bar plot of the count distribution of sentiment labels across all transcripts.
<i>plot_trust_distribution</i>	Shows a bar plot of the count distribution of trust labels across all transcripts.
<i>plot_topic_clusters</i>	Shows a scatter plot that groups transcripts with similar topics.
<i>plot_sentiment_by_topic</i>	Shows a bar plot that shows the distribution of trust scores from a model for each topic category.

## custom\_packages.misc\_modules.visualisations.plot\_word\_frequencies

**plot\_word\_frequencies** (*df, words\_num*)

Shows a bar plot with the most frequent words across all transcripts.

### Parameters

- **df** (*pd.DataFrame*) – A document-feature matrix (DFM) with keyword frequencies.
- **words\_num** (*int*) – How many highest-frequency keywords to show in the plot.

### Return type

*None*

## **custom\_packages.misc\_modules.visualisations.word\_cloud\_visualisation**

**word\_cloud\_visualisation** (*df*, *words\_num*)

Visualises a word cloud of the top keywords across all transcripts.

### **Parameters**

- **df** (*pd.DataFrame*) – A document-feature matrix (DFM) with keyword frequencies.
- **words\_num** (*int*) – How many highest-frequency keywords to show in the word cloud.

### **Return type**

*None*

## **custom\_packages.misc\_modules.visualisations.plot\_sentiment\_distribution**

**plot\_sentiment\_distribution** (*df*, *sentiment\_column*)

Shows a bar plot of the count distribution of sentiment labels across all transcripts.

### **Parameters**

- **df** (*pandas.DataFrame*)
- **sentiment\_column** (*str*)

### **Return type**

*None*

## **custom\_packages.misc\_modules.visualisations.plot\_trust\_distribution**

**plot\_trust\_distribution** (*df*, *sentiment\_column*)

Shows a bar plot of the count distribution of trust labels across all transcripts.

### **Parameters**

- **df** (*pandas.DataFrame*)
- **sentiment\_column** (*str*)

### **Return type**

*None*

## **custom\_packages.misc\_modules.visualisations.plot\_topic\_clusters**

**plot\_topic\_clusters** (*transcript\_topic\_distribution*, *dominant\_topic\_labels*)

Shows a scatter plot that groups transcripts with similar topics.

### **Parameters**

- **transcript\_topic\_distribution** (*numpy.ndarray*)
- **dominant\_topic\_labels** (*numpy.ndarray*)

### **Return type**

*None*

## **custom\_packages.misc\_modules.visualisations.plot\_sentiment\_by\_topic**

**plot\_sentiment\_by\_topic** (*df*, *trust\_column\_name*, *topic\_column\_name*)

Shows a bar plot that shows the distribution of trust scores from a model for each topic category.

### **Parameters**

- **df** (*pandas.DataFrame*)
- **trust\_column\_name** (*str*)
- **topic\_column\_name** (*str*)

### **Return type**

*None*

## 2.2 Module: custom\_packages.preprocessing\_modules.data\_collector

<code>store_all_video_data</code>	Uses the given dictionary of search query strings to retrieve relevant Youtube videos, and store them locally as json files.
<code>get_youtube_video_data</code>	
<code>get_json_data</code>	Get all json files from given directory and return a subset of their data in a dataframe format.
<code>get_transcripts</code>	Get all transcript .txt files from given directory and combine their text in a single string output.
<code>get_cleaned_csv_record</code>	Retrieves a CSV file as a dataframe from the given path.

### custom\_packages.preprocessing\_modules.data\_collector.store\_all\_video\_data

**store\_all\_video\_data** (*search\_query\_groups*, *max\_results=10*)

Uses the given dictionary of search query strings to retrieve relevant Youtube videos, and store them locally as json files.

#### Parameters

- **search\_query\_groups** (*dict*) – A dictionary of strings, for all search queries per search group (e.g., trust/distrust groups).
- **max\_results** (*int*, *optional*) – The maximum number of total video results to retrieve and store. The default is 10.

#### Raises

**ValueError** – Raised if no Youtube data retrieved / empty dataframe.

#### Returns

Contains all the extracted video data.

#### Return type

pd.DataFrame

### custom\_packages.preprocessing\_modules.data\_collector.get\_youtube\_video\_data

**get\_youtube\_video\_data** (*query*, *query\_group\_name*, *max\_results=10*)

1. Uses Youtube API v3 to retrieve the list of relevant videos based on the function parameters
2. Extract a list of video data to return.
3. Store all of the retrieved search data locally as a json file on the computer.

#### Parameters

- **query** (*str*) – The search term to lookup on Youtube.
- **query\_group\_name** (*str*) – A unique identifier for this query - used in the exported json filename.
- **max\_results** (*int* (*optional*; *default value of 10*)) – Specifies specifies the maximum number of items(video) that should be returned.

#### Returns

Contains all the extracted video data.

#### Return type

pd.DataFrame

## custom\_packages.preprocessing\_modules.data\_collector.get\_json\_data

**get\_json\_data** (*files\_dir*)

Get all json files from given directory and return a subset of their data in a dataframe format.

### Parameters

**files\_dir** (*Path.PosixPath*) – The directory where all the json files are stored.

### Returns

Contains json data such as video\_id, url, publish\_time, etc.

### Return type

pd.DataFrame

## custom\_packages.preprocessing\_modules.data\_collector.get\_transcripts

**get\_transcripts** (*files\_dir, platform\_name*)

Get all transcript .txt files from given directory and combine their text in a single string output.

### Parameters

- **files\_dir** (*Path.PosixPath*) – The directory where all the transcription .txt files are stored.
- **platform\_name** (*str*)

### Returns

Contains the combined text from all the .txt files.

### Return type

*str*

## custom\_packages.preprocessing\_modules.data\_collector.get\_cleaned\_csv\_record

**get\_cleaned\_csv\_record** (*csv\_filepath*)

Retrieves a CSV file as a dataframe from the given path.

### Parameters

**csv\_filepath** (*Path*)

### Return type

pandas.DataFrame

## 2.3 Module: custom\_packages.preprocessing\_modules.text\_cleaning

<i>clean_transcripts</i>	Cleans transcripts per platform and returns a dataframe where each row contains a unique transcript identifier, the actual filename retrieved from Path, the platform name, the raw text, and the cleaned text.
<i>clean_text</i>	Preprocesses the string argument, by stripping unnecessary white spaces, converts to lowercase, decouples contractions and replaces dashes with an empty space.
<i>store_cleaned_data_as_csv</i>	Extracts the raw and cleaned data from the dataframe argument and stores them in the corresponding directory offered as an argument.
<i>store_as_csv</i>	Stores the dataframe argument as a CSV file in the given directory.

## custom\_packages.preprocessing\_modules.text\_cleaning.clean\_transcripts

**clean\_transcripts** (*filepaths\_per\_platform*)

Cleans transcripts per platform and returns a dataframe where each row contains a unique transcript identifier, the actual filename retrieved from Path, the platform name, the raw text, and the cleaned text.

### Parameters

**filepaths\_per\_platform** (*dict[str, list[Path]]*) – Contains the platform name (e.g., youtube) as key, and the full paths to the relevant transcripts as the value.

### Return type

pd.DataFrame

## custom\_packages.preprocessing\_modules.text\_cleaning.clean\_text

**clean\_text** (*transcript\_text*)

Preprocesses the string argument, by stripping unnecessary white spaces, converts to lowercase, decouples contractions and replaces dashes with an empty space.

### Parameters

**transcript\_text** (*str*) – The text to be preprocessed.

### Returns

The cleaned text after preprocessing.

### Return type

*str*

## custom\_packages.preprocessing\_modules.text\_cleaning.store\_cleaned\_data\_as\_csv

**store\_cleaned\_data\_as\_csv** (*transcript\_records, cleaned\_data\_dir, raw\_data\_dir*)

Extracts the raw and cleaned data from the dataframe argument and stores them in the corresponding directory offered as an argument.

### Parameters

- **transcript\_records** (*pandas.DataFrame*)
- **cleaned\_data\_dir** (*Path*)
- **raw\_data\_dir** (*Path*)

### Return type

*None*

## custom\_packages.preprocessing\_modules.text\_cleaning.store\_as\_csv

**store\_as\_csv** (*df, full\_filepath*)

Stores the dataframe argument as a CSV file in the given directory.

### Parameters

- **df** (*pandas.DataFrame*)
- **full\_filepath** (*Path*)

### Return type

*None*

## 2.4 Module: custom\_packages.preprocessing\_modules.transcriber

<code>transcribe_all_videos</code>	Initiates the transcription process by retrieving the id and url of each video and passing it on to the audio transcription service.
<code>process_audio_stream</code>	
<code>transcribe_youtube_audio_from_videos</code>	
<code>check_file_size_exceeded_limit</code>	Check if the file size in the given path exceeds the expected maximum size limit.
<code>subset_audio_with_overlap</code>	Segment audio files (with pydub) larger than 25MB into smaller segments, to support OpenAI Whisper's transcription.
<code>transcribe_audio</code>	Retrieves the audio file from the given filepath, transcribes it using the whisper-1 model, and stores it based on the video id as a .txt file in the given filepath.
<code>transcribe_segmented_audio</code>	Retrieves the audio segment files from the given filepath based on the video id, transcribes them using the whisper-1 model, and stores them as a concatenated string as a .txt file in the given filepath with the video id argument as part of the transcript name.
<code>store_transcript</code>	Stores the transcribed text as a .txt file in the given filepath, using video id as part of the .txt filename.

### custom\_packages.preprocessing\_modules.transcriber.transcribe\_all\_videos

**transcribe\_all\_videos** (*video\_data*)

Initiates the transcription process by retrieving the id and url of each video and passing it on to the audio transcription service.

#### Parameters

**video\_data** (*pd.DataFrame*) – Contains json youtube data such as the video url, id, published\_date in a dataframe format for easier data manipulation.

#### Returns

The combined text across all transcription files.

#### Return type

*str*

### custom\_packages.preprocessing\_modules.transcriber.process\_audio\_stream

**process\_audio\_stream** (*video\_url*)

1. Downloading the best audio stream, starting from m4a if available, as its small size is convenient for my local resources. 2. Converting that audio stream to a mono 16kHz mp4 audio file. 3. Storing that audio file as a temporary resource.

#### Parameters

**video\_url** (*str*) – Used to identify and retrieve the online video from Youtube.

#### Returns

The path directory of the temporarily stored mp4 audio file. If failed, returns empty string.

#### Return type

*str*



## **custom\_packages.preprocessing\_modules.transcriber.transcribe\_youtube\_audio\_from\_videos**

**transcribe\_youtube\_audio\_from\_videos** (*video\_id*, *video\_url*, *trim\_duration\_secs*=2400)

1. Retrieving the Youtube video based on the URL in an m4a format.
2. Using OpenAI's whisper-1 model to transcribe all temporary-stored wav audio files.

### **Parameters**

- **video\_id** (*str*) – Used to mark the transcription file for identification purposes.
- **video\_url** (*str*) – Used to identify and retrieve the online video from Youtube.

### **Returns**

The full text of the transcription.

### **Return type**

String

## **custom\_packages.preprocessing\_modules.transcriber.check\_file\_size\_exceeded\_limit**

**check\_file\_size\_exceeded\_limit** (*path*, *max\_size\_limit*)

Check if the file size in the given path exceeds the expected maximum size limit.

### **Parameters**

- **path** (*str*) – The path where the file is located.
- **max\_size** (*float*) – The maximum size limit permitted for this file.
- **max\_size\_limit** (*float*)

### **Returns**

True if the size of the file exceeds the indicated maximum value; False if within the maximum size limit.

### **Return type**

*bool*

## **custom\_packages.preprocessing\_modules.transcriber.subset\_audio\_with\_overlap**

**subset\_audio\_with\_overlap** (*audio\_id*, *audio\_file\_path*, *output\_dir*, *chunk\_length\_ms*=2400000, *overlap\_ms*=1000)

Segment audio files (with pydub) larger than 25MB into smaller segments, to support OpenAI Whisper's transcription.

### **Parameters**

- **video\_id** (*str*) – Used to mark the smaller audio file segments for identification purposes.
- **audio\_file\_path** (*str*) – To locate and retrieve the large audio file to be segmented.
- **output\_dir** (*str*) – Location path to save all smaller audio file segments.
- **chunk\_length\_ms** (*int* (*optional*)) – How long, in milliseconds, each segment duration should be. Default is 2.4m milliseconds (i.e., 40 minutes). 1 minute = 60k milliseconds.
- **overlap\_ms** (*int*) – How many milliseconds should the overlap be between segments.
- **audio\_id** (*str*)

### **Returns**

True if successfully segmented and stored in the *output\_dir*. Default is False.

### **Return type**

Boolean

## **custom\_packages.preprocessing\_modules.transcriber.transcribe\_audio**

**transcribe\_audio** (*video\_id*, *audio\_filepath*, *transcript\_filepath*)

Retrieves the audio file from the given filepath, transcribes it using the whisper-1 model, and stores it based on the video id as a .txt file in the given filepath.

**Parameters**

- **video\_id** (*str*)
- **audio\_filepath** (*str*)
- **transcript\_filepath** (*str*)

**Return type**

*str*

**custom\_packages.preprocessing\_modules.transcriber.transcribe\_segmented\_audio**

**transcribe\_segmented\_audio** (*video\_id, audio\_dir, transcript\_dir*)

Retrieves the audio segment files from the given filepath based on the video id, transcribes them using the whisper-1 model, and stores them as a concatenated string as a .txt file in the given filepath with the video id argument as part of the transcript name.

**Parameters**

- **video\_id** (*str*)
- **audio\_dir** (*Path*)
- **transcript\_dir** (*str*)

**Return type**

*str*

**custom\_packages.preprocessing\_modules.transcriber.store\_transcript**

**store\_transcript** (*video\_id, transcript\_text, transcript\_filepath*)

Stores the transcribed text as a .txt file in the given filepath, using video id as part of the .txt filename.

**Parameters**

- **video\_id** (*str*) – The id of the video's transcript to be used as part of the .txt filename.
- **transcript\_text** (*str*) – The full transcribed text.
- **transcript\_filepath** (*str*) – The full path to store the .txt transcript.

**Returns**

True if successfully stored locally, otherwise False.

**Return type**

*bool*

## 2.5 Module: custom\_packages.modelling\_modules.nlp\_modelling

<code>get_document_feature_matrix</code>	Creates a basic bag-of-words document-feature matrix (DFM) from the given cleaned text in the dataframe.
<code>sentiment_analysis_using_nltk</code>	
<code>sentiment_analysis_using_distilbert</code>	
<code>sentiment_analysis_using_gpt</code>	
<code>sentiment_analysis_using_google</code>	
<code>sentiment_model_comparison</code>	The dataframe argument should contain the trust sentiment values from all models.
<code>print_models_agreement_percentage</code>	Calculates and displays in the console the agreement rate between two models.
<code>map_sentiment_to_trust</code>	Converts classic sentiment labels to trust/distrust/neutral.
<code>get_dominant_topic</code>	Does topic modelling using FASTopic on the cleaned text from each transcript and retrieves the top keywords identified per topic (5 most popular topics in this instance).

### custom\_packages.modelling\_modules.nlp\_modelling.get\_document\_feature\_matrix

#### `get_document_feature_matrix(df)`

Creates a basic bag-of-words document-feature matrix (DFM) from the given cleaned text in the dataframe.

##### Parameters

`df` (`pd.DataFrame`) – Needs to contain a column named “cleaned\_text”.

##### Returns

`dfm`

##### Return type

`pd.DataFrame`

### custom\_packages.modelling\_modules.nlp\_modelling.sentiment\_analysis\_using\_nltk

#### `sentiment_analysis_using_nltk(df)`

1. **Retrieves the partly cleaned/preprocessed text (decoupling contractions, converting to lowercase and removing**  
dashes and unnecessary white spaces already done earlier) from the dataframe argument.
2. Finishes preprocessing by removing punctuation, stopwords, and applies tokenisation through NLTK.
3. **Does sentiment analysis using NLTK Vader SentimentIntensityAnalyzer on the tokenised text, and calculates the**  
polarity scores (positive, negative, neutral).
4. Visualises the sentiment and translated trust labels in a plot each.
5. Returns a dataframe with the scores and labels.

##### Parameters

`df` (`pandas.DataFrame`)

##### Return type

`pandas.DataFrame`

## custom\_packages.modelling\_modules.nlp\_modelling.sentiment\_analysis\_using\_distilbert

`sentiment_analysis_using_distilbert(df)`

1. **Retrieves the partly cleaned/preprocessed text (decoupling contractions, converting to lowercase and removing**  
dashes and unnecessary white spaces already done earlier) from the dataframe argument. No need for any further preprocessing nor tokenisation as DistilBert handles that internally.
2. **Does sentiment analysis using the DistilBert sentiment pipeline and calculates the polarity scores**  
(positive, negative).
3. Visualises the sentiment and translated trust labels in a plot each.
4. Returns a dataframe with the scores and labels.

### Parameters

`df` (*pandas.DataFrame*)

### Return type

*None*

## custom\_packages.modelling\_modules.nlp\_modelling.sentiment\_analysis\_using\_gpt

`sentiment_analysis_using_gpt(df)`

1. **Retrieves the partly cleaned/preprocessed text (decoupling contractions, converting to lowercase and removing**  
dashes and unnecessary white spaces already done earlier) from the dataframe argument. No need for any further preprocessing nor tokenisation as gpt-4o-mini handles that internally.
2. **Does sentiment analysis using the gpt-4o-mini model (requires OpenAI API) and responds with json-formatted**  
results on polarity scores (positive, negative, neutral).
3. Visualises the sentiment and translated trust labels in a plot each.
4. Returns a dataframe with the scores, labels, and summary highlights.

### Parameters

`df` (*pandas.DataFrame*)

### Return type

*pandas.DataFrame*

## custom\_packages.modelling\_modules.nlp\_modelling.sentiment\_analysis\_using\_google

`sentiment_analysis_using_google(df)`

1. **Retrieves the partly cleaned/preprocessed text (decoupling contractions, converting to lowercase and removing**  
dashes and unnecessary white spaces already done earlier) from the dataframe argument. No need for any further preprocessing nor tokenisation as gemini-2.0-flash-lite handles that internally.
2. **Does sentiment analysis using the gemini-2.0-flash-lite model (requires Google API) and responds with**  
json-formatted results on polarity scores (positive, negative, neutral).
3. Visualises the sentiment and translated trust labels in a plot each.
4. Returns a dataframe with the scores, labels, and summary highlights.

### Parameters

`df` (*pandas.DataFrame*)

### Return type

*pandas.DataFrame*

## **custom\_packages.modelling\_modules.nlp\_modelling.sentiment\_model\_comparison**

**sentiment\_model\_comparison** (*df*)

The dataframe argument should contain the trust sentiment values from all models. Compares and displays in the console the agreement rate between all models.

### **Parameters**

**df** (*pandas.DataFrame*)

### **Return type**

*None*

## **custom\_packages.modelling\_modules.nlp\_modelling.print\_models\_agreement\_percentage**

**print\_models\_agreement\_percentage** (*first\_model, second\_model*)

Calculates and displays in the console the agreement rate between two models.

### **Parameters**

- **first\_model** (*pandas.Series*)
- **second\_model** (*pandas.Series*)

### **Return type**

*None*

## **custom\_packages.modelling\_modules.nlp\_modelling.map\_sentiment\_to\_trust**

**map\_sentiment\_to\_trust** (*sentiment\_label*)

Converts classic sentiment labels to trust/distrust/neutral.

### **Parameters**

**sentiment\_label** (*str*)

### **Return type**

*str*

## **custom\_packages.modelling\_modules.nlp\_modelling.get\_dominant\_topic**

**get\_dominant\_topic** (*df*)

Does topic modelling using FASTopic on the cleaned text from each transcript and retrieves the top keywords identified per topic (5 most popular topics in this instance). Displays a plot per sentiment model for a trust-sentiment score by topic assessment. It also displays another plot that clusters transcripts with similar topics.

### **Parameters**

**df** (*pandas.DataFrame*)

### **Return type**

*pandas.DataFrame*