# Homework 1 A

**Yusupha Juwara**
Sapienza University of Rome
juwara.1936515@studenti.uniroma1.it

## Abstract

This report walks through two tasks: task 0 (EmotivITA) and subtask A of task 1 (HODI). It provides a brief explanation of the two tasks, their input and output formats, how the data are formated, the prompts and the motivation for those specific prompt choices. It also explains how to run the scripts succesfully.

## 1 Task 0 – EmoITA

### 1.1 Overview

**EmotivITA** (Gafà et al., 2023) is a task of emotion regression that aims to advance research in emotion detection within the Italian language. It is based on EmoITA , the Italian version of the EmoBank dataset, annotated based on the dimensional model of emotions, where each of the three dimensions (Valence, Arousal, and Dominance) is a real-value associated with a sentence, allowing for nuanced comparisons. Note that in this task, however, we are asked to convert the real-values to categorical labels.

### 1.2 Input Format of the dataset

The dataset is in a **.csv** format with five fields: **id** (int), **text** (string), **V** (float), **A** (flaot), and **D** (float).

An example of the dataset with the headers and id 1 and 11 is:

```
1   id,text,V,A,D
2   1,Auguriamo a voi e alla vostra
    ↪  famiglia un nuovo anno pieno di
    ↪  gioia e amore.,5,3,4.25
3   11,"Ci sono ristoranti, piscine e
    ↪  spogliatoi.",4,3.33,3.67
```

### 1.3 Mapping Float to Categorical

The values of each dimension is a real-value. But since we are explicitly required to use **categorical** labels, we have to map them.

One simple option is to split the range (0-5) **equally** between categories. With this choice, I

found that the values are skewed toward the 'Alta' and 'Molto Alta' as shown below.

```
1   {'Bassa': 1, 'Media': 799, 'Alta': 14307,
    ↪   'Molto Alta': 8893}
```

This is problematic since the model can **always** predict the majority class and get a high score.

To mitigate this problem, another option is to split in this non-equal way:

```
1   [0.0,3.2) -> "Bassa"
2   [2.5, 3.8) -> "Media"
3   [4.0, 5.0] -> "Alta"
4
5   # frequencies
6   {'Bassa': 5315, 'Media': 11396, 'Alta':
    ↪   7289}
```

### 1.4 Reframing the Dataset

The input data is reformated in this way:

- Use pandas to read the **.csv** file, get the 4 fields (see the Input Format 1.2), convert the string of floats to floats, and map the floats to categorical labels (see the Categorical Mapping 1.3)

- Shuffle the choices/labels so that the model learns not just their ordering, but their actual semantic meaning. Note that there is a paper on this, that models sometimes memorize the ordering of the labels, and that if the ordering is changed, there may be some (drastic) drop in accuracy.

- Create a jsonl entry for dimension in their respective files, add the 'dimension' field to the jsonl entries so that the model can differentiate between the dimensions (said by one of the TAs on Google Classroom). E.g.,

```
1   # Output format for each jsonl entry
2   {'text': "...", 'choices': [...],
    ↪   'label': int, 'dimension': '...'}
```

## 1.5 Run the Code

To run the script for the **Development set.csv**, **cd** to the Home directory 'HM1_A-1936515' and run:

```
python ".\EmotivITA\scripts.py"
  --shuffle_labels --verbose
  --map_option=0
```

For the **Test set - Gold labels.csv**, do the same but add the 'test' arg. Type 'help' after the file name to see what these args are.

## 1.6 Prompts



Figure 1: Prompts of task 0



Figure 2: Motivation of the prompt choices for task 0

## 2 Task 1 – HODI

### 2.1 Overview

HODI (Homotransphobia Detection in Italian) is a shared task introduced at EVALITA 2023 that focuses on automatically identifying homotransphobic content in Italian text (Nozza et al., 2023). The HODI dataset includes messages directed at Trans, and the goal is to develop sytems like LLMs for detecting homotransphobic language in Italian.

HODI consists of two subtasks. **Subtask A** (binary classification) – classifies a given text as homotransphobic or not. This part is what concerns this section. **Subtask B** (explainability) – identifies which parts of the text are homotransphobic.

### 2.2 Input Format of the dataset

The dataset is in a **.tsv** format with 3 fields: **id** (int), **text** (string), **homotransphobic** (int). The homotransphobic field defines if a text is homotransphobia (1) or not (0).

An example of the dataset with the headers and id 10 is:

```
id,text,homotransphobic
10       @user_ab io ho una faccia
  proprio da frocia       0
```

### 2.3 Reframing the Dataset

The input data is reformated in a similar way as already mentioned in the **EvaITA** section. The difference is that this file is a **.tsv** (tab separated) and the choices are already in categorical format.

- Use pandas to read the **tsv** file, get the fields, shuffle the choices/labels just as explain in the **EvaITA** section, and create a jsonl entry that looks like this:

```
# Output format for each jsonl entry
{"text": "...", "choices": ["Vero",
  "Falso"], "label": int}
```

### 2.4 Run the Code

To run the script, **cd** to the Home directory 'HM1_A-1936515' and run:

```
python "./HODI_2023/scripts.py"
  --shuffle_labels --verbose
```

### 2.5 Prompts



Figure 3: Prompts of task 1, subtask A



Figure 4: Motivation of the prompt choices for task 1, subtask A

# References

Giovanni Gafà, Francesco Cutugno, and Marco Venuti. 2023. Emotivita at evalita2023: Overview of the dimensional and multidimensional emotion analysis task. In *Proceedings of the Eighth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2023)*. CEUR.org.

Debora Nozza, Alessandra Teresa Cignarella, Greta Damo, Tommaso Caselli, and Viviana Patti. 2023. HODI at EVALITA 2023: Overview of the Homotransphobia Detection in Italian Task. In *Proceedings of the Eighth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2023)*, Parma, Italy. CEUR.org.