# Human Values Behind Arguments
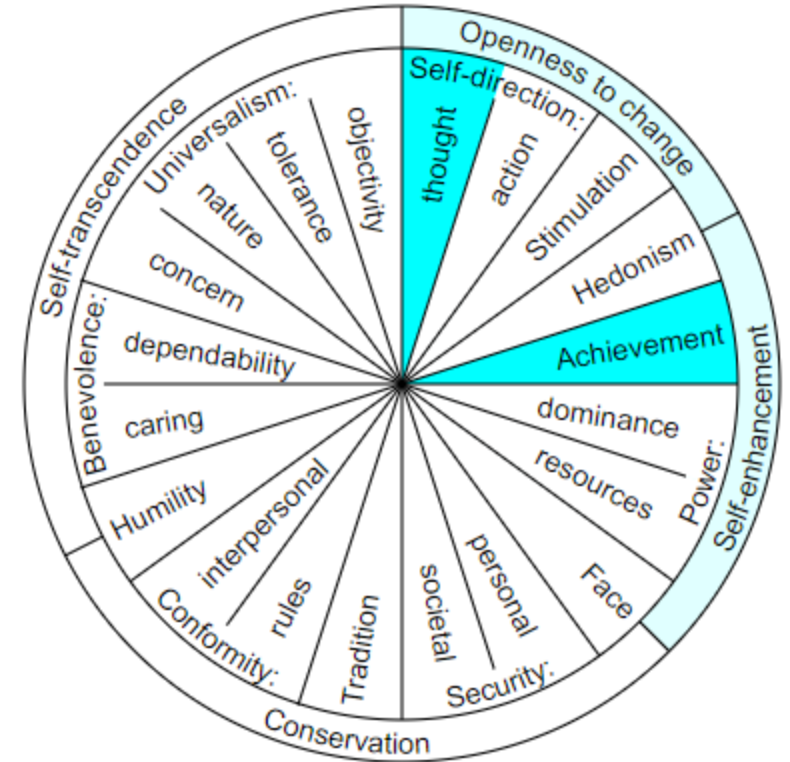
# The Big Picture

- Task: Any prompt should be automatically labeled with 1+ labels

  **Multi-label Classification Problem**

- Prompts/Inputs consist of triplets:

  **<premise>, <stance>, <justification>**

- Hierarchy of Label Levels:
  - 4 levels (original dataset, 2021)
  - 2 levels (revised dataset, 2023)
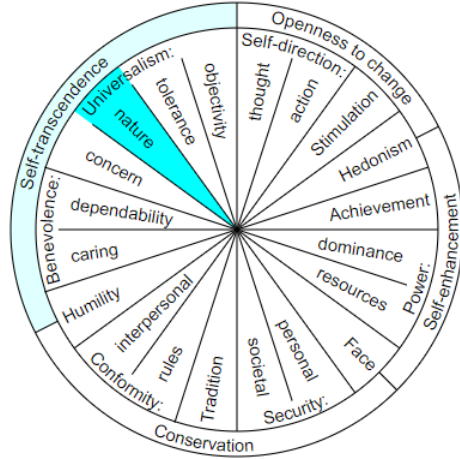  - Sub-levels do not necessarily fall under one fixed upper level (e.g., Humility, Face, Hedonism)

*Input: ("We should subsidize student loans", "against", "it isn't the obligation of any tax payer to give money to help someone else get an education. that is a personal choice, and if they want to go to college they need to pay for it on their own.")*
*Your response: (Achievement, resources, personal, dependability)*

https://values.args.me/
(last acessed: 23.01.2024)
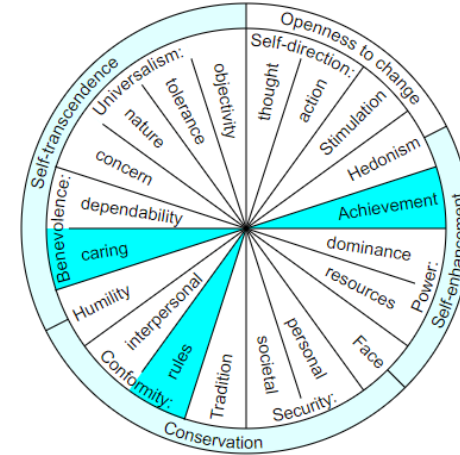
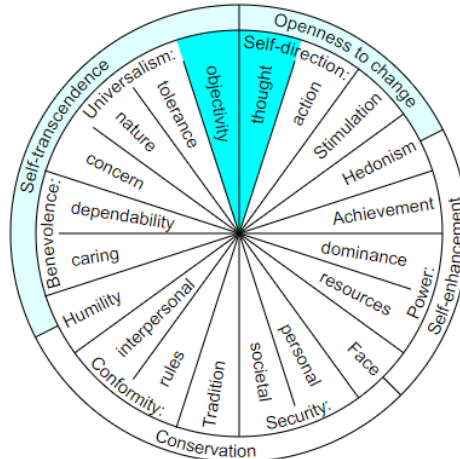# The Big Picture

We need to reduce our CO2 emissions to save the environment.

Submit

We should ban the use of child actors
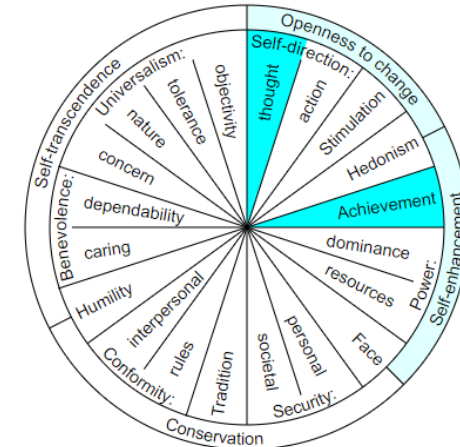
Submit

We should adopt atheism

Submit

Everybody should learn to code.

Submit

# An Overview – What Has Been Done?

- Research Questions
- Exploratory Data Analysis (EDA)
- Experiments on ChatGPT3.5
- Experiments on Llama-2-7b

# The Research Questions

- Are there any inconsistencies in the dataset?

- Does the presence of any label 'imply' another label?

- How well the general LLMs compare to the fine-tuned LMs (such as BERT, etc.) participating in SemEval-2023 competition?

- Will the result be any different if role-based prompts are used?
  (e.g., prompts asked from the perspective of a scientist or a psychologist)

# Exploratory Data Analysis (EDA) – What is that?

Formal Definition: *Exploratory data analysis [...] is used by data scientists to analyze [...] data sets and summarize their main characteristics [...]. It helps determine how best to manipulate data sources to get the answers you need, making it easier [...] to discover patterns, spot anomalies, test a hypothesis, or check assumptions.* *([https://www.ibm.com/topics/exploratory-data-analysis](https://www.ibm.com/topics/exploratory-data-analysis), Last access: 21.01.2024)*

# Initial Data Insights

| # of | Data Split | | | |
|---|---|---|---|---|
| | **Train** | **Validation** | **Test** | **Total** |
| **samples** | 5393 (60.83%) | 1896 (21.39%) | 1576 (17.78%) | 8865 (100%) |
| **missing values (initial assessment)** | 0 | 0 | 0 | 0 |

Label assignments for samples are done via a one-hot encoding like manner.
Therefore, it should also be checked **if all the samples are assigned to at least one label.**

# Feature Analysis – General Facts



| Argument | Value categories |
|---|---|
| ○ Con "We should end the use of economic sanctions": Economic sanctions provide security and ensure that citizens are treated fairly. | Security: societal, Universalism: concern |
| ○ Pro "We need a better migration policy.": Discussing what happened in the past between Africa and Europe is useless. All slaves and their owners died a long time ago. You cannot blame the grandchildren. | Universalism: concern |
| ○ Con "Rapists should be tortured": Throughout India, many false rape cases are being registered these days. Torturing all of the accused persons causes torture to innocent persons too. | Security: societal, Universalism: concern |

Labels on figure: Stance, Conclusion, Premise

| Statistics Val. | "Premise" Length | |
|---|---|---|
| | # of words | # of characters |
| Min. | 3 | 19 |
| Max. | 133 | 780 |
| Avg. | 21 | 126 |

# Feature Analysis – Digging Deeper

```python
# print rows containing premise value with 3 words or less in train data
df_train_simplified.loc[df_train_simplified['premise_preprocessed'].str.split().str.len() < 4]
```

✓ 0.0s                                                                                  Python

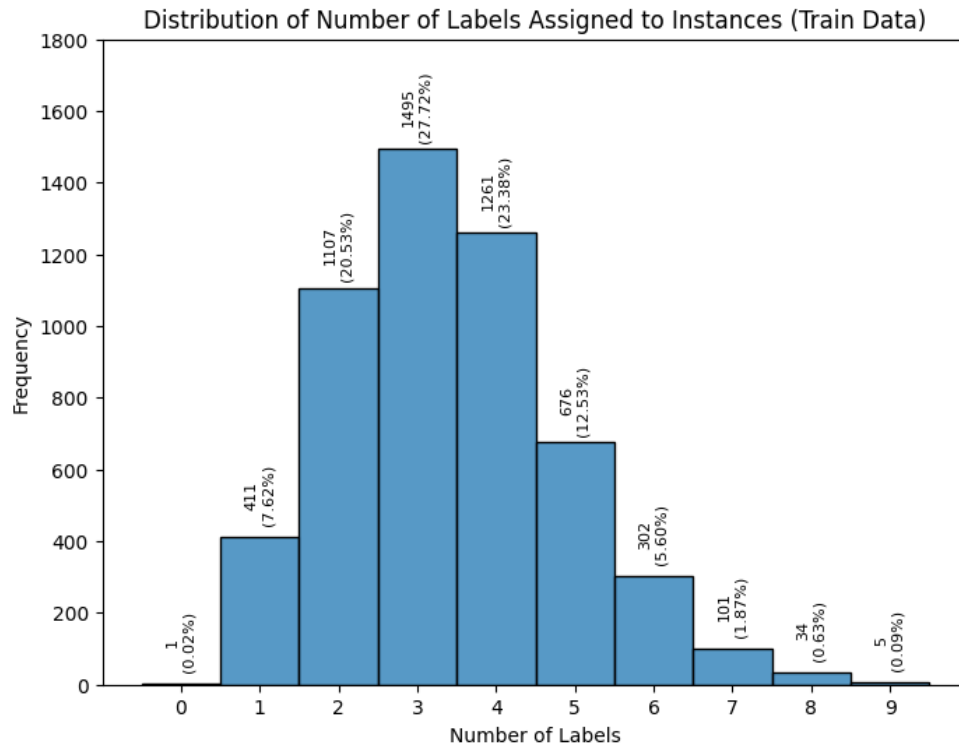| | argument_id | stance | conclusion | premise | labels | premise_preprocessed |
|---|---|---|---|---|---|---|
| 4792 | D27081 | against | Aadhaar should be made mandatory in India | Failed biometric authentication. | [Self-direction: action] | failed biometric authentication |
| 4951 | E04046 | against | We need legal ways of migration. | Migrants sell drugs. | [Security: societal, Conformity: rules] | migrants sell drugs |
| 5366 | E07262 | against | We need legal ways of migration. | Migrants sell drugs. | [Security: societal, Conformity: rules, Benevo... | migrants sell drugs |

# Not Just Train Data...

```python
with pd.option_context('display.max_colwidth', None):
    display(df_test_simplified.loc[[1544, 1479]])
```
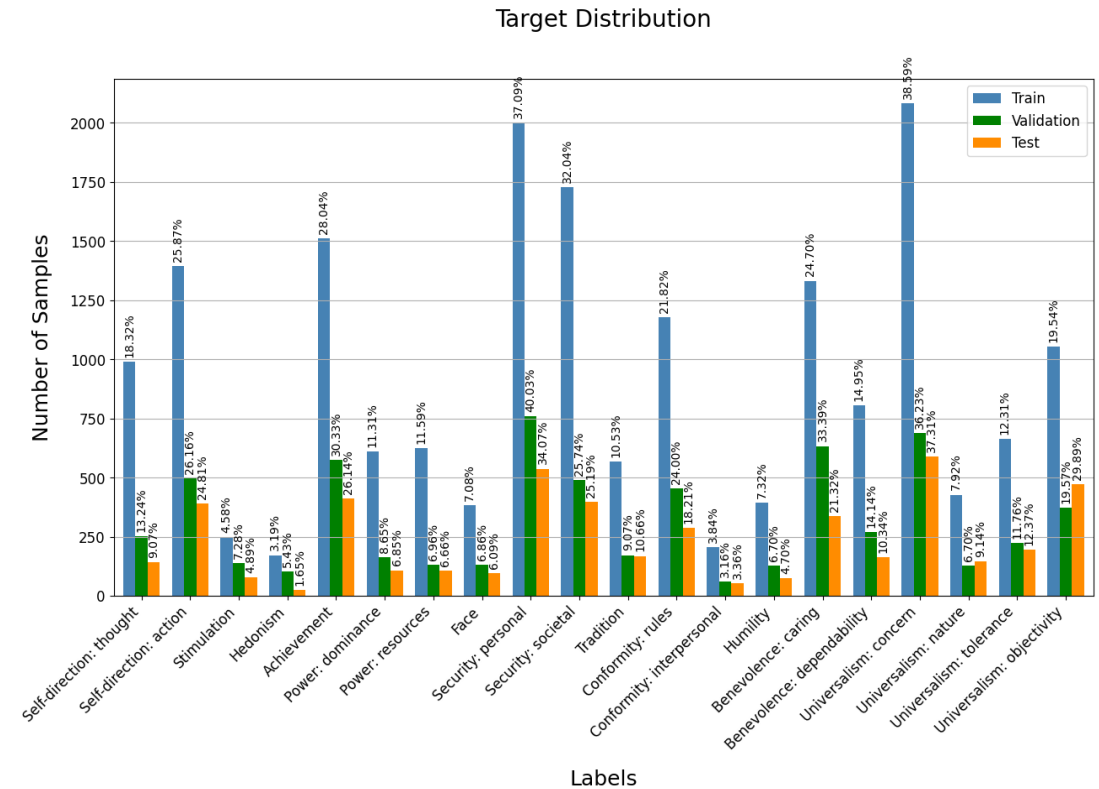Python

| | argument_id | stance | conclusion | premise | labels | premise_preprocessed |
|---|---|---|---|---|---|---|
| 1544 | E07097 | in favor of | We should be aware of migration due to climate change | There are no laws protecting climate refugees nowadays. They are facing all kinds of abuse, violence, and aggression, even from officials at the borders! | [Security: societal, Conformity: rules, Universalism: concern] | there are no laws protecting climate refugees nowadays they are facing all kinds of abuse violence and aggression even from officials at the borders |
| 1479 | E04116 | in favor of | We should be aware of migration due to climate change | There are no laws protecting climate refugees nowadays. They are facing all kinds of abuse, violence, and aggression, even from officials at the borders! | [Security: personal, Conformity: rules, Universalism: concern] | there are no laws protecting climate refugees nowadays they are facing all kinds of abuse violence and aggression even from officials at the borders |

# Target Analysis:
## Distribution of Labels in Data



Distribution of Number of Labels Assigned to Instances (Train Data)
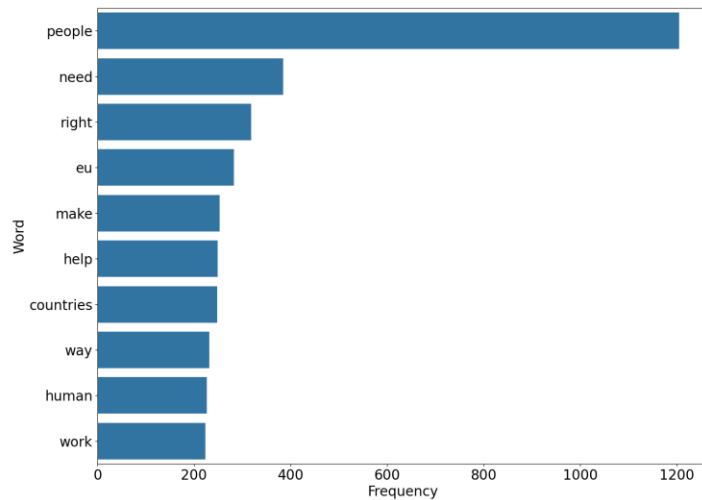


Target Distribution

- Average number of labels: 3

- Minimum number of labels: **0 (?)**

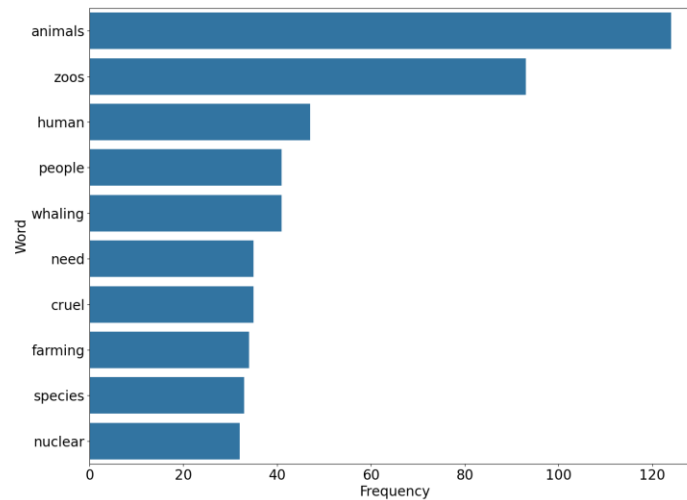- Maximum number of labels: 9

   **(out of 20 labels)**

Ratio of labelled samples don't follow the same distribution for training, validation, and test splits.
→ Any training or fine-tuning is prone to resulting in biased models!

# Do certain words in Premise imply a certain label?

**Frequency of word appearance in "Premise" column of Train Data**
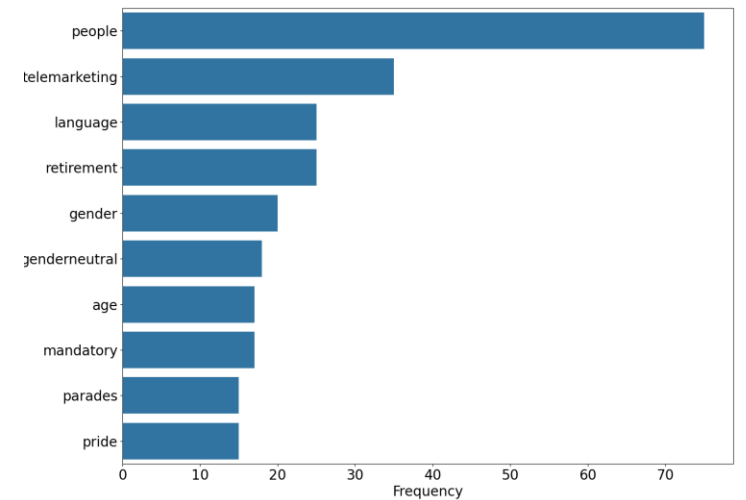**(stop-words are ignored: i.e., is, the, a, …)**



Most used words in dataset
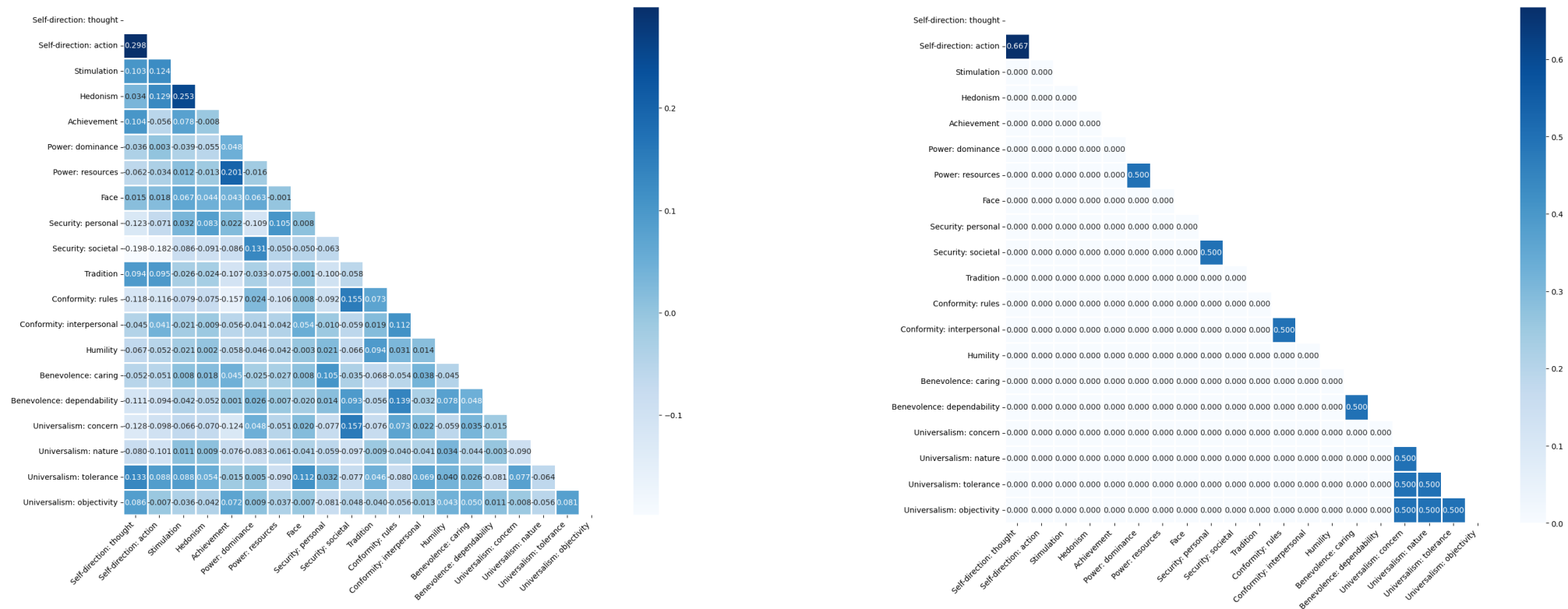for all instances

Most used words in dataset
with label:

Most used words in dataset
with label:

**Universalism – nature**

**Conformity: inter-personal**

# Correlation between labels?

...hardly



- **Left: Bi-variant Analysis; Right: Cosine Similarity**

- **Highest correlation between labels is "Self-direction: thought" and "Self-direction action" with correlation of 0.3 (=fairly low)**

- **However: Correlation ≠ Causation !**

# Summary of Inconsistencies in Dataset

## Missing Labels & Duplicate Inputs

|  | Train | Validation | Test | Total |
|---|---|---|---|---|
| # of missing labels | 1 | 0 | 0 | 1 |
| # of duplicate values | 89 | 35 | 32 | 156 |

## Target Label Distribution

Ratio of labeled samples per class in preset splits of Train, Validation & Test data do not necessarily follow the same distribution, sometimes off by a noticeable margin.

# Prompt - Techniques

- **Role setting**

- **Zero-shot**

- **Few-shot**

- **Chain-of-thought**

- **YAML for structured input, output, context**

# Prompt - Insight

*Imagine you are a scientist. It is your job to analyze the human values behind a list of arguments ...*

*An argument is thereby given as a triple (premise, stance, conclusion). The premise is the statement given to the person. The stance is ...*

*Use following values:*

- *thought: it is good to have own ideas and interests*
- *action: it is good to determine one's own actions ...*
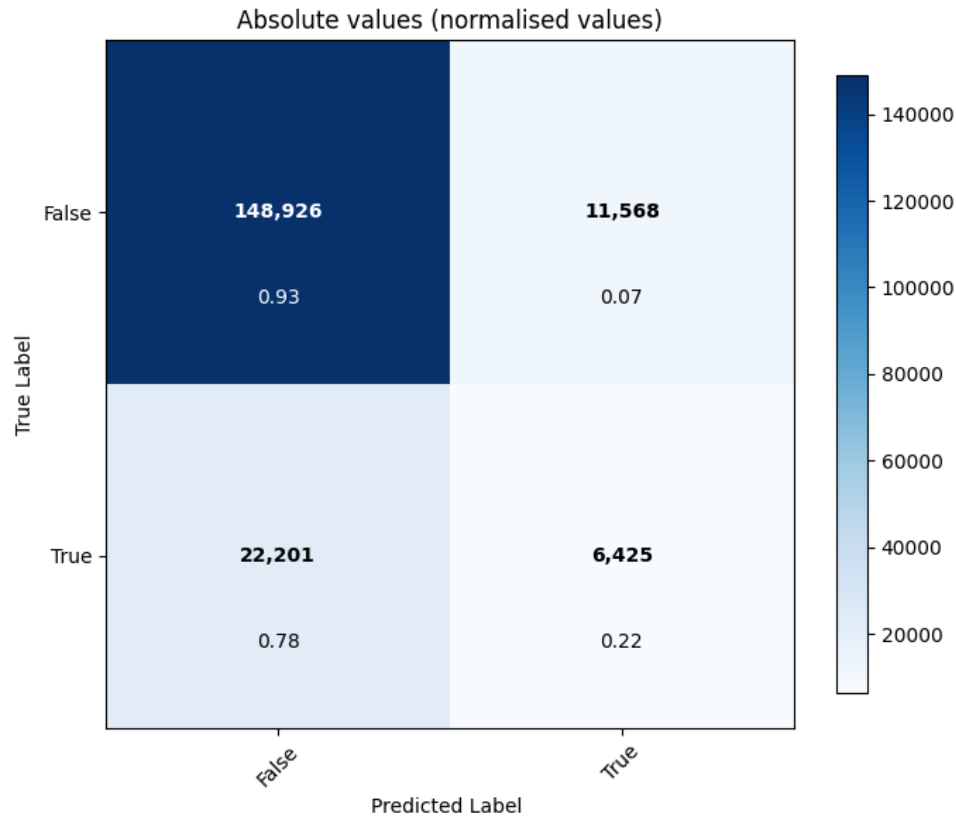
*Examples:*

*Input: ("We should cancel pride parades","in favor of","pride parades create a huge disturbance"),*

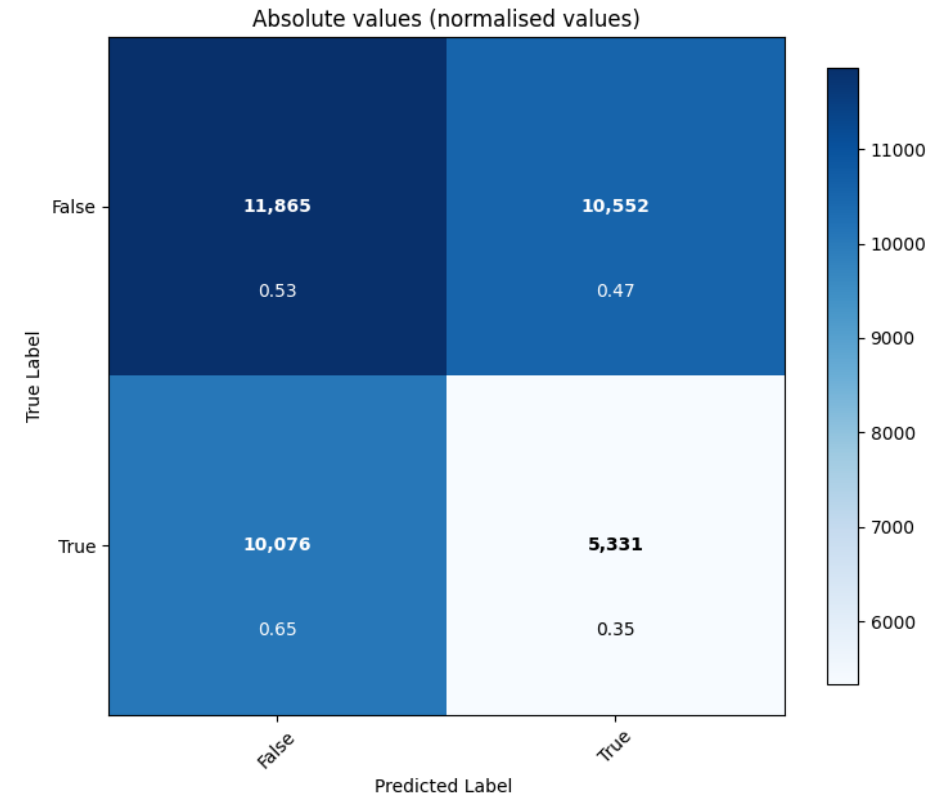*Output: (rules, interpersonal)...*

*Input: ("We should end mandatory retirement","in favor of","mandatory retirement is purely age discrimination")*
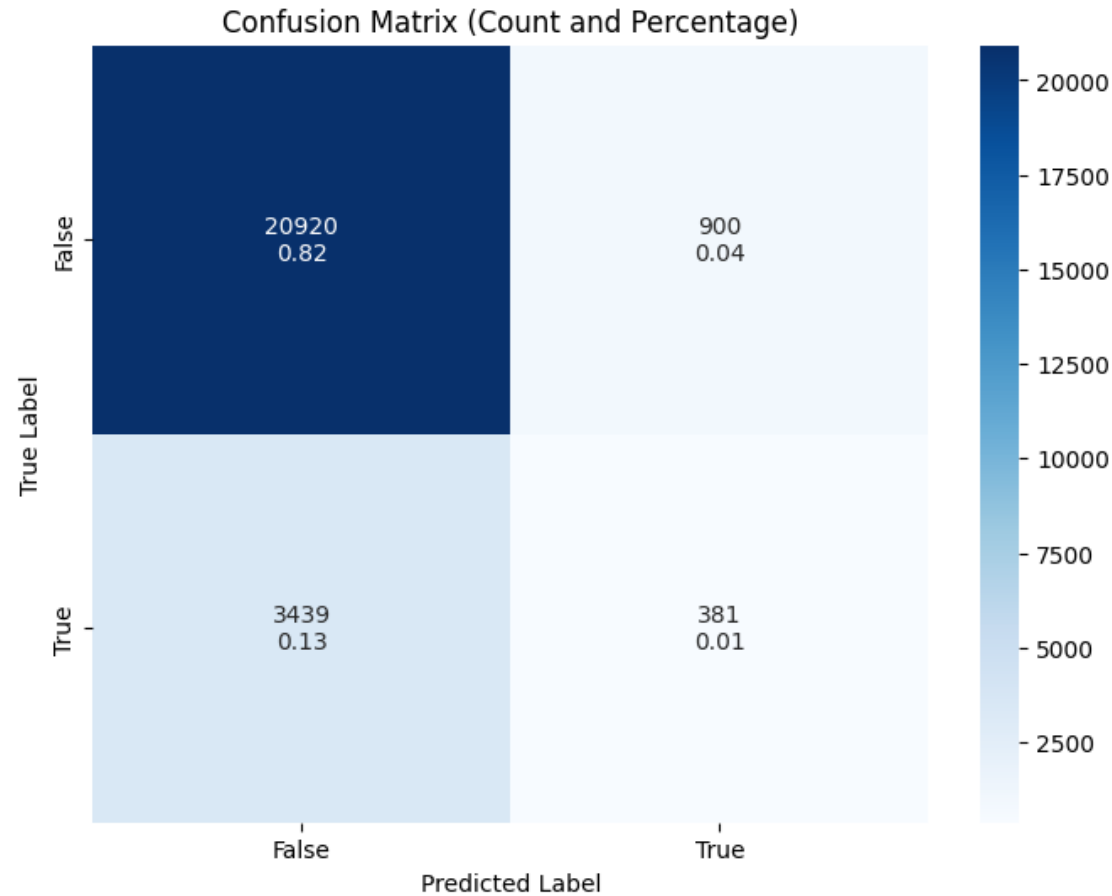
*Output:*

# ChatGPT - Results



Confusion matrix fine granularity (F1 ~ 0.275553)

Confusion coarse granularity (F1 ~ 0.79924)

# Llama2- Results



Confusion Matrix (Count and Percentage)

Confusion coarse granularity (F1 ~ 0.15)

# Conclusions

- Developers of ML-based solutions for problems should prioritize thorough analysis of the data
    - <span style="color:red">EDA is a MUST</span>

- Smaller language models prior to the LLMS are still capable if fine-tuned for a specific task
    - This makes them ideal for utilizing them into Production pipelines (faster inference time, less costly)

# QA



**Word Frequency Cloud Generated from "Premise" Column of Training Data**