# NLP Course Report

Hans Peter Lyngsøe, pvr448

September 5, 2024

## 1 Week 1

(a) Basic statistics:

- the data is quite evenly distributed across the 3 languages.
- We note that there are more answerable than unanswerable by a factor of 10-1.
- train_set: 15326
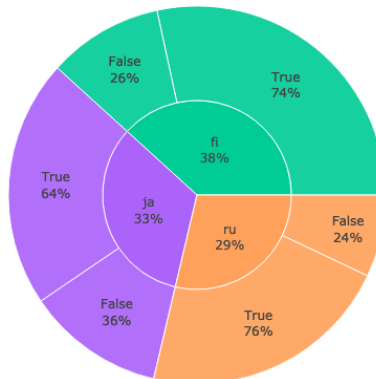- val_set: 3028

Training Data Language Distribution



Figure 1: Distribution of labels in the dataset

(b)
```python
import MeCab
def get_top_words(df: pd.DataFrame, lang: str,
    n=5):
    df_lang = df[df['lang'] == lang].copy()
```
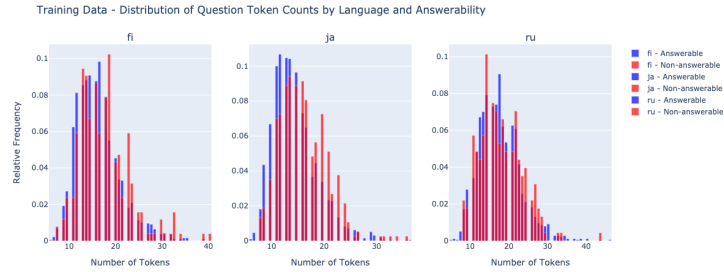
Figure 2: normalized Histogram for token count(using llama3 tokenizer) of answerable/unanswerable questions in the dataset

```python
if lang == 'ja':
    mecab = MeCab.Tagger("-Owakati")  #
        Initialize MeCab tokenizer
    df_lang.loc[:, 'words_question_tokens'
        ] = df_lang['question'].apply(
        lambda x: mecab.parse(x).split())
else:
    df_lang.loc[:, 'words_question_tokens'
        ] = df_lang['question'].apply(
        lambda x: x.split(' '))

all_tokens = np.concatenate(df_lang['
    words_question_tokens'].values)
unique, counts = np.unique(all_tokens,
    return_counts=True)
sorted_indices = np.argsort(counts)[::-1]
top_unique_tokens = unique[sorted_indices
    ][:n]
top_tokens_dict = {token: int(count) for
    token, count in zip(top_unique_tokens,
    counts[sorted_indices][:n])}
return top_tokens_dict
```

we get the following results:

qualitative analysis:

- Finnish:
    - mitä (1021) - what
    - on (774) - is
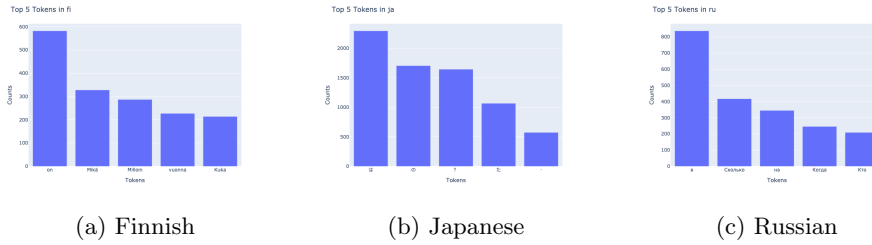    - mikä (441) - what/which

2

<div style="text-align:center">(a) Finnish      (b) Japanese      (c) Russian</div>

Figure 3: Top 5 tokens in Finnish, Japanese, and Russian

- – miten (293) - how
- – kuinka (214) - how

- Japanese:
  - – (3402) - of/in
  - – (2896) - topic marker
  - – (1827) - to/at
  - – (1604) - object marker
  - – (1328) - subject marker

- Russian:
  - – (1889) - in
  - – (812) - which
  - – (744) - what
  - – (628) - how
  - – (454) - on

From this analysis, we can observe:

- Finnish top words include question words (mitä, mikä, miten, kuinka) and a common verb (on).

- Japanese top words are primarily particles, which are essential for sentence structure but don't carry much meaning alone.

- Russian top words include both question words (, , ) and prepositions (, ).

This reflects differences in language structure and how questions are typically formed in each language.

(c)