

Detecting the difficulty level of French texts



Joris Booij-Liewes
HEC Lausanne
Data Mining
Ahmed M'sabou



Introduction

01

Models

02

The Code

03

Conclusion

04



Presentation Outline

Introduction

01

Problematic

Given a text in French,
can you predict its
difficulty level (A1, A2, B1,
B2, C1, C2)?

02

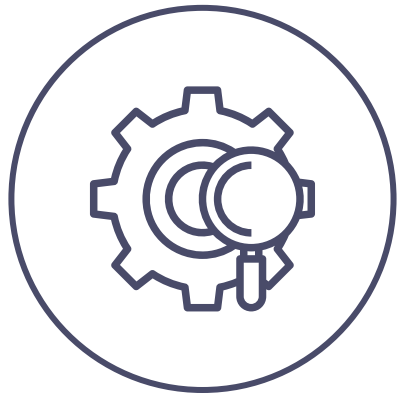
Why this project?

We have noticed that to
improve one's skills in a
new foreign language, it is
vital to read texts in a
suitable difficulty level.



Models

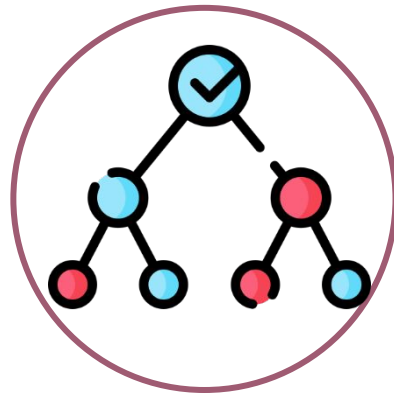
Classic models saw in course



Logistic regression



kNN



Decision Tree



Random Forest

Results

Using a validation set that is 20% of the data

| | Logistic regression | kNN | Decision Tree | Random Forests | CamemBERT |
|-----------------|---------------------|---------------|---------------|----------------|---------------|
| Precision | 0.4582 | 0.3997 | 0.2964 | 0.4187 | 0.5795 |
| Recall | 0.4604 | 0.3677 | 0.3021 | 0.4208 | 0.5740 |
| F1-score | 0.4561 | 0.3530 | 0.2975 | 0.4099 | 0.5717 |
| Accuracy | 0.4604 | 0.3677 | 0.3021 | 0.4208 | 0.5740 |

Techniques to improve our models?

What we did to improve the accuracy of our models



Ensemble

Ensemble models are machine learning models that combine the predictions of multiple individual models to make a final prediction.



Embedding

embeddings are a way of representing data such as words, phrases, or other discrete items as continuous, low-dimensional numerical vectors



Word2vec

Word2vec is a method for learning vector representations of words



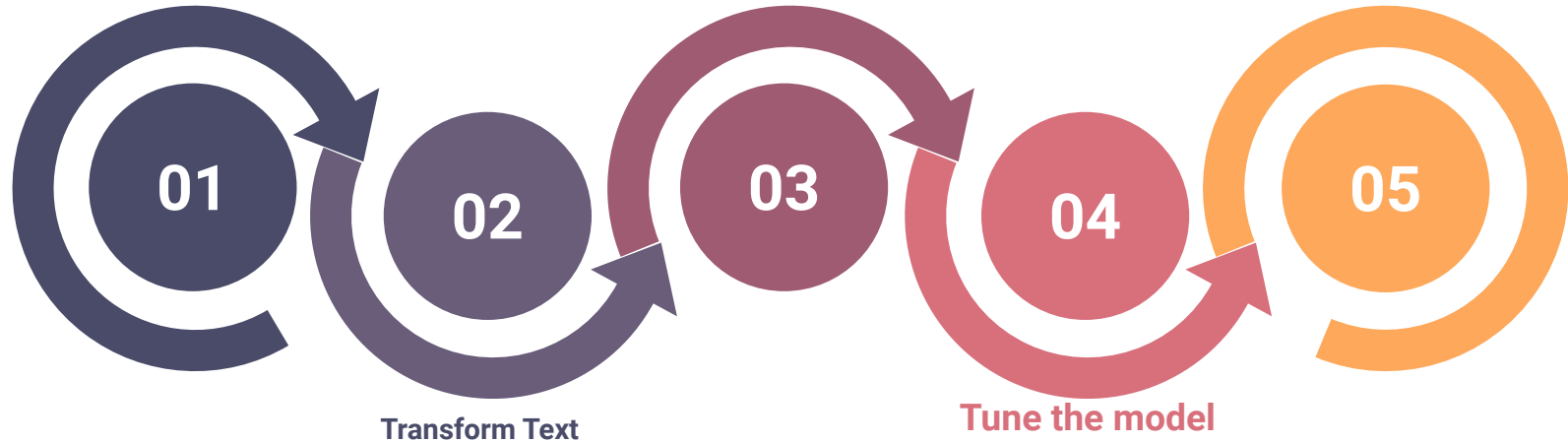
CamemBERT model

French version of the BERT model for NLP classification

Load training Sentences

Train model

Make predictions



Explaining CamemBERT & Code

The process of preparing the data

```
# train data
df_train = pd.read_csv('training_data.csv')

# Rename Labeling
df_train['difficulty'] = df_train['difficulty'].replace(['A1', 'A2', 'B1', 'B2', 'C1', 'C2'], [0, 1, 2, 3, 4, 5])

# test data
df_test = pd.read_csv('unlabelled_test_data.csv')
```

```
epochs = 6
MAX_LEN = 64
batch_size = 16
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')

# Initialize CamemBERT tokenizer
tokenizer = CamembertTokenizer.from_pretrained('camembert/camembert-large', do_lower_case=True)
```


Explaining CamemBERT & Code

Embedding

```
#user tokenizer to convert sentences into tokenizer'
input_ids = [tokenizer.encode(sent,add_special_tokens=True,max_length=MAX_LEN) for sent in text]

# Pad our input tokens
input_ids = pad_sequences(input_ids, maxlen=MAX_LEN, dtype="long", truncating="post", padding="post")

attention_masks = []

# Create a mask of 1s for each token followed by 0s for padding
for seq in input_ids:
    seq_mask = [float(i>0) for i in seq]
    attention_masks.append(seq_mask)
```

Words

| | cycle | car | road | tree | root | hotel | river |
|---------|-------|-----|------|------|------|-------|-------|
| Indices | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

↑
d
↓

| | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|
| 0.2 | 0.1 | 0.1 | 0.2 | 0.3 | 0.4 | 0.3 |
| 0.1 | 0.7 | 0.4 | 0.9 | 0.2 | 0.8 | 0.7 |
| 0.5 | 0.8 | 0.6 | 0.5 | 0.8 | 0.2 | 0.8 |
| 0.6 | 0.2 | 0.9 | 0.3 | 0.3 | 0.6 | 0.1 |

| | | |
|-----|------|------|
| 101 | 7592 | 2088 |
| 999 | 102 | 0 |

input_ids

attention_mask

| | | |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 0 |

- red tokens
- [PAD] tokens

WORD EMBEDDINGS

Explaining CamemBERT & Code

[illegible]

CONFUSION MATRIX:

Conclusion

What was our best model?

If we had more time?

What did we learn from
this project?

THANK YOU !