# Software project:Developing a decision-making assistant

Morozova Milena

July 26, 2022

# Table of Contents

# Basic Tasks

- Creating a training sample
- Preprocessing training sample
- Building and training the model
- Testing the Model

## Additional Tasks

- Rating calculation
- Creating GUI
- Testing the App on different platforms
- Setting the application icon

# Preprocessing the sample

Any text data in its raw material form cannot be analyzed by NLP libraries. This data must be cleaned using various data processing techniques.



Figure: Preprocessing the training sample

Our neural network will have 2 neurons at the output - the upper neuron will be responsible for the positive text, the lower - for the negative. For positive statements at the output, we will require the vector [1.0], respectively for negative - [0.1].

```python
X = data_pad
Y = np.array([[1, 0]]*count_true + [[0, 1]]*count_false)
print(X.shape, Y.shape)
```

Figure: Training

```python
t = "Не доверяй никому".lower()
data = tokenizer.texts_to_sequences([t])
data_pad = pad_sequences(data, maxlen=max_text_len)

print( sequence_to_text(data[0]) )

res = model.predict(inp)
print(res, np.argmax(res), sep='\n')
```

Figure: Testing

Figure: Working with a mathematical model

# Bibliography

- Machine Learning, Neural and Statistical Classification
  `https://www1.maths.leeds.ac.uk/~charles/statlog/`
- Cleaning Preprocessing Text Data for Sentiment Analysis
  `https://towardsdatascience.com/`
- Generating WordClouds in Python Tutorial
  `https://www.datacamp.com/tutorial/wordcloud-python`
- TripAdvisor
  `https://www.tripadvisor.ru/`
- Sentiment Analysis of Review Datasets
  `https://www.researchgate.net/publication`
- Keras Documentation
  `https://ru-keras.com/recurrent-layers/`