# SMS Sentiment Classification based on Stylometric Features, Emoticons, Informal abbreviations and other Text Features

Branislava Šandrih

branislava.sandrih@fil.bg.ac.rs

University of Belgrade, Faculty of Philology, Serbia

JeRTeh – Society for Language Resources and Technology

# Motivation

- Sentiment analysis / Opinion mining / Sentiment classification:
  - contextual mining of text which identifies and extracts subjective information
- Analysis of social media is usually restricted to just basic sentiment analysis and count based metrics
  - what about SMS messages? They are even shorter!
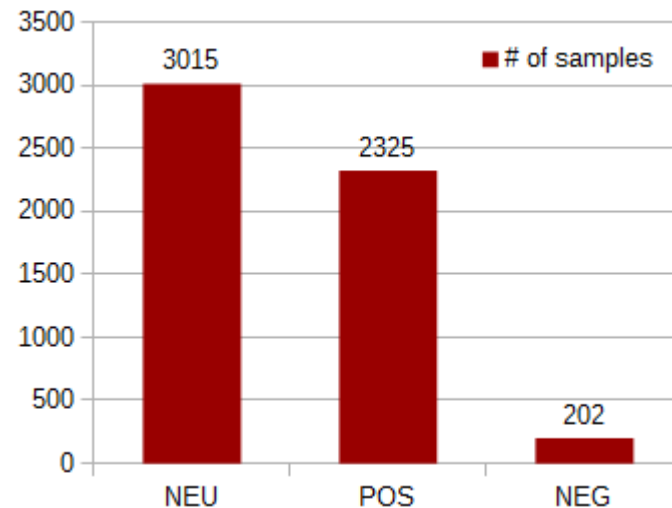
# Challenge

- Restrictions:
  - 160 characters
    - 70 if diacritics are used
  - Small keyboards, hard to type
    - messages contain only most important information
  - Need to express attitude, mood, voice tone, facial expression, gesture...
    - the only available tool: characters!?

# **Workaround**

- Authors
  - use sh-s for common used phrases
  - EMPHASIZE IMPORTANT INFORMATION WITH UPPERCASE
  - do not type whole ws
  - omit diacritics
  - excessively use emoticons :)  :(  :-P
- Consequence?
  - Hard to analyze using standard approaches

# A different approach

- But first dataset
  - modest ~ 5,500 SMS messages in Serbian/German/English, Cyrillic + Latin
  - hard to gather, because SMS are too personal!
  - Manual annotation

# Features

- Lexical
  - Character based
    - counts of lowercase and uppercase letters, total # of characters, ratios…
  - Word based
    - average sentence length, average length of tokens etc.
- Stylistic
  - sentence starts with uppercase, spaces after punctuation etc.
- Emoticons
  - kiss, confused, heart etc.
- Abbreviatons
  - ty, tnx, k, fb, cu, u, l8r etc.

# API and Web Interface

# Results

- Playing with features, incrementally adding:
  - Lexical
  - Lexical + Stylistic
  - Lexical + Stylistic + Emoticons
  - Lexical + Stylistic + Emoticons + Abbreviations
- Accuracy 94.4% in the last case

# Conclusion

- For short messages, sentiment classification can be performed by exploring stylometry and other important characteristics