**University of Antwerp**
Faculty of Arts

# Internship Presentation

Dimitris Boumparis

Supervisor: Prof. Dr. Walter Daelemans

Mentor: Jens Lemmens

08/12/2022

**DI** Suggest a couple of titles for a powerpoint presentation about an internship on natural language processing

1. "Exploring Natural Language Processing: My Internship Experience"
2. "Learning and Implementing NLP Techniques during My Internship"
3. "My Journey into Natural Language Processing: An Internship Experience"
4. "Applying Natural Language Processing in Real-World Settings: An Internship Perspective"
5. "Gaining Hands-On Experience with NLP through an Internship"
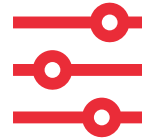
**DI** Be more creative; maybe add a pun

1. "NLP: It's a Natural Fit for My Internship"
2. "Intern-lingual Processing: My Experience with NLP"
3. "Interning in NLP Land: Adventures in Natural Language Processing"
4. "From Intern to NLP Pro: My Journey in Natural Language Processing"
5. "Interning on the Cutting Edge of NLP: My Experience with Natural Language Processing"

# General Information

## Objective

- Apply the knowledge and skills acquired in the MA DTA in concrete organizational contexts
- **Paths:** research, experiments, programming, application development, etc.

## Requirements

- **336 hours** minimum (240 work + 96 self-studying)
- Flexible **conditions** (on-site, WFH)
- combination with thesis

## Evaluation

- **Portfolio** (2,5k - 5k words)
  - More personal tone
  - Experiences and takeaways
- **Presentation:**
  - More formal tone
  - Results-oriented

# Motivation

- Already been an:
  - intern during my bachelor's (translation and localization) (summer of 2017)
  - associate at the Computational Stylistics Lab at NKUA (Prof. Dr. George Mikros) (2017-2019)
- Prior experience in programming (front-end web dev)
- New challenges:
  - Longest stay abroad to date
  - Multi-cultural environment
  - Research centre with international prestige
- Aspirations for a research-oriented career

# Tasks

- Main: **Author profiling**
- Predict:
  - Gender
  - Age group
  - Personality
  - Job (industry)
- **Classification-only**
- **Algorithms used:**
  - Experimented with traditional machine learning and deep learning (BERT)
  - (Linear) Support Vector Machines was preferred

- Other: **NRC Word-Emotion Association Lexicon** (Mohammad & Turney, 2010)
  - Used in (multilingual) sentiment analysis
  - Version 0.92
    - Translation into Greek
    - MTPE of the Greek MT

University of Antwerp
| Faculty of Arts

# NRC Lexicon

- **MTPE** (where column B already populated, ~5,500)
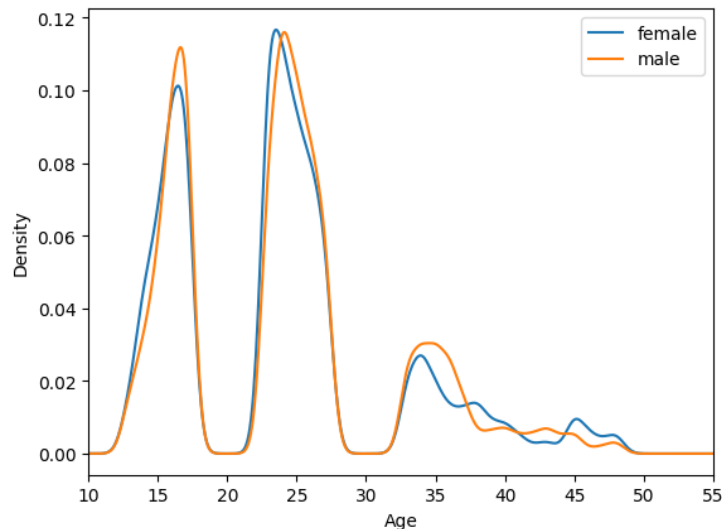
- Challenges:
  - English is analytical
  - The same form for different purposes (noun == verb) → translations in multiple PoS
  - Polysemous words → only emotional connotations
  - Greek is a cased language → prone to mistranslations

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | English (en) | Greek (el) - automatic | Greek (el) - manual | Positive | Negative |
| 39 | abrasion | τριβή | τριβή | 0 | 1 |
| 41 | abrogate | ακυρώνω | ακυρώνω | 0 | 1 |
| 44 | abscess | απόστημα | απόστημα | 0 | 1 |
| 45 | absence | απουσία | απουσία | 0 | 1 |
| 46 | absent | απών | απών | 0 | 1 |
| 47 | absentee | απών | απών | 0 | 1 |
| 48 | absenteeism | κατά συνήθεια απουσία | κατά συνήθεια απουσία | 0 | 1 |
| 50 | absolute | απόλυτος | απόλυτος | 1 | 0 |
| 51 | absolution | απαλλαγή | απαλλαγή | 1 | 0 |
| 52 | absorbed | απορροφηθεί | απορροφημένος | 1 | 0 |
| 61 | absurd | παράλογος | παράλογος | 0 | 1 |
| 62 | absurdity | παραλογισμός | παραλογισμός | 0 | 1 |
| 63 | abundance | αφθονία | αφθονία | 1 | 1 |
| 64 | abundant | άφθονος | άφθονος | 1 | 0 |
| 65 | abuse | κατάχρηση | καταχρώμαι, κατάχρηση | 0 | 1 |
| 68 | abysmal | αβυσσαλέες | αβυσσαλέος | 0 | 1 |
| 69 | abyss | άβυσσος | άβυσσος | 0 | 1 |

# Author Profiling Tasks

- Datasets:
  - **PAN15 Author Profiling Shared Task**
    - Tweets in **English**, Spanish, Italian and Dutch
    - Balanced in genders and ages (see graph)
  - **Blog Authorship Corpus**
    - Short and mid-length tweets/texts (English)



- Compatibility/Sanitization:
  - Age groups
  - Cleaning
  - Tokenization
- **Classifications:**
  - **5 or 10-fold** Grid Search Cross-Validation on F1-macro score
  - Run on CLiPS server or locally
  - Models saved as **joblib** files
  - Used sklearn's `GroupShuffleSplit()` for the PAN dataset

# Preprocessing

- Lowercasing (stdlib)
- Tokenization (spacy)
- Replace content words with POS tags while maintaining function words

- Different levels of cleaning (RegEx):
  - Remove https://urls, @mentions, #hastags
  - Remove repeated punctuation
  - Replace whitespace with single space

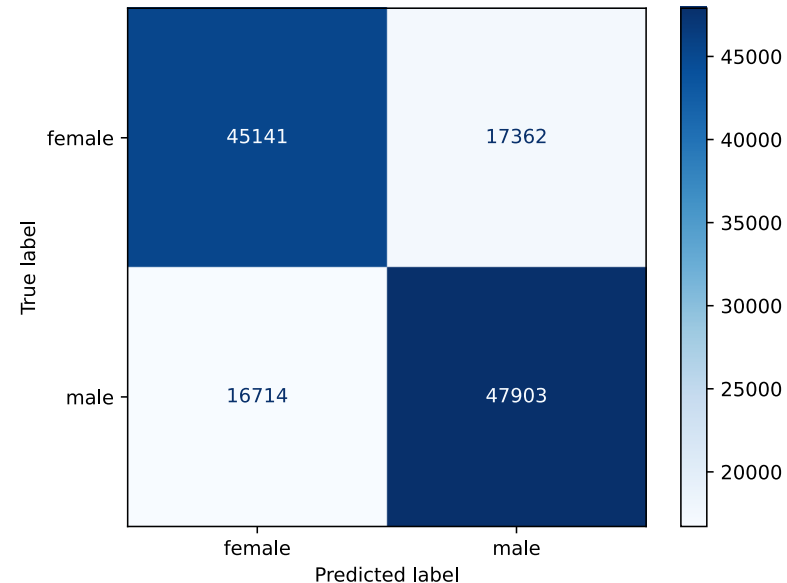| documents | semicleaned | tokenized | pos_fw |
| --- | --- | --- | --- |
| How to Test Your Startup Idea for $50 http://t... | how to test your startup idea for $50 | how to test your startup idea for $ 50 | how to VERB your NOUN NOUN for SYM NUM |
| @username @username @username @username @usern... | you've been quoted in my story "new story" | you 've be quote in my story " new story " | you ' AUX be INTJ in my NOUN " ADJ NOUN " |
| New Story http://t.co/Uu5AggZP #storify #cacer... | new story | new story | ADJ NOUN |
| @username @username @username @username You've... | you've been quoted in my story | you 've be quote in my story | you ' AUX be INTJ in my NOUN |
| @username @username @username @username @usern... | you've been quoted in my story | you 've be quote in my story | you ' AUX be INTJ in my NOUN |

# Gender (LSVM / LogReg)

- Binary classification
- LogReg performed significantly worse

- Unigrams: .73
- **1-3grams: .74**
- 1-2grams + POS + Punct: .68
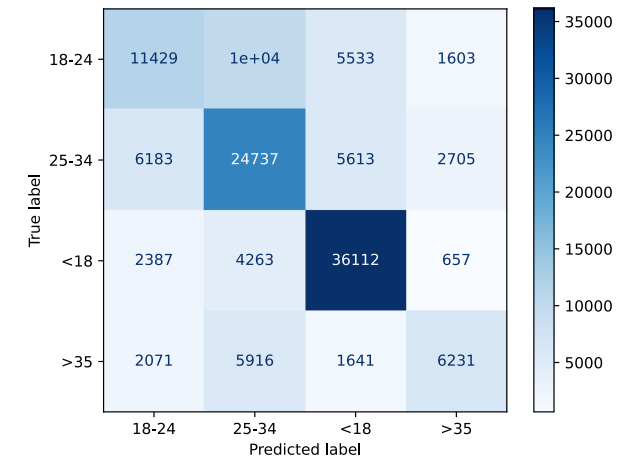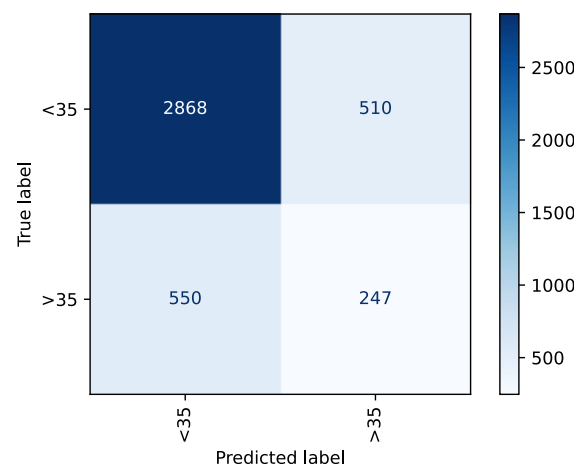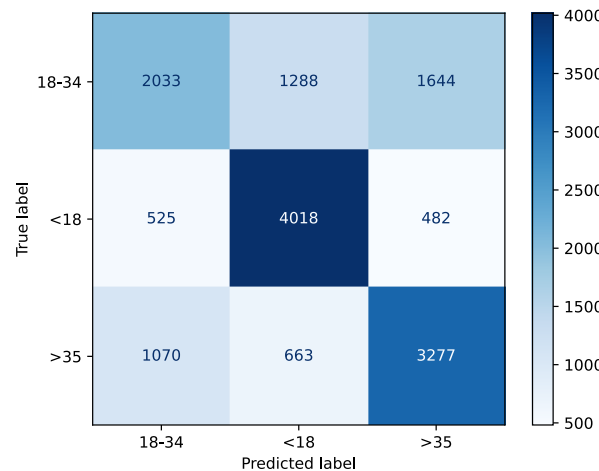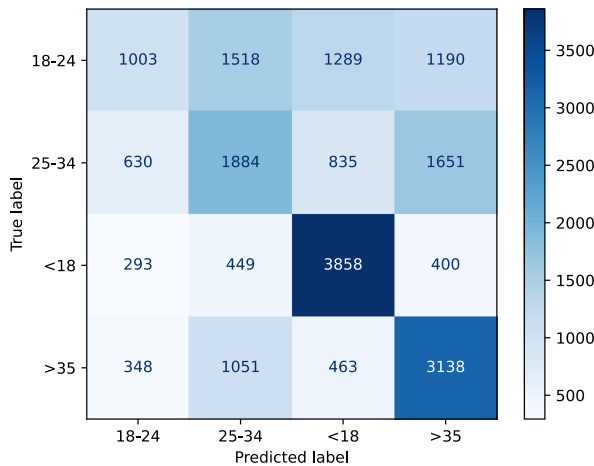- 1-2grams + POS + Punct + LIWC: .67

# Age group (LSVM)

- Groups:
  - Initially: 18-24, 25-34, 35-50, 50+
  - After sanitization: <18, 18-24, 25-34, 35+
  - Finally: <35, 35+ (binary)

- 4 classes: .35, .44, **.52 (1-2grams)**
- 3 classes: .51-**.57 (unigrams)**
- 2 classes: .54 (word/char 1-3grams + FU)
- **4 classes (all samples): .62 (1-3grams)**

University of Antwerp
Faculty of Arts

# Industry (SVM)

- 36 industries with various # of samples
- 8 with 1,000 samples each

- 1-2grams: .25 (left)
- **All features: .26 (right)**

# Personality (LSVM)

- Extrovert: -0.5...0.4
- Binary: 1 if score > 0, else 0
- Very imbalanced dataset → adjust weights

- 1-4grams (tokenized text): .56
- **1-3grams (tokenized text) + other personality traits + POS FW: .70**

```python
1  # Calculate the weights for each class
2  class_weights = dict(df['extroverted'].value_counts(normalize=True))
3  class_weights = {k: round(1-v, 3) for k, v in class_weights.items()}
```

University of Antwerp
Faculty of Arts

# Scikit-Learn **Classes** Spotlight

- **FeatureUnion()**
  - Combine several feature extraction mechanisms into a single transformer
    1. TF-IDF word Vectorizer
    2. TF-IDF POS Vectorizer
    3. Punctuation Vectorizer
    4. Numeric Transformers (punctuation and LIWC)

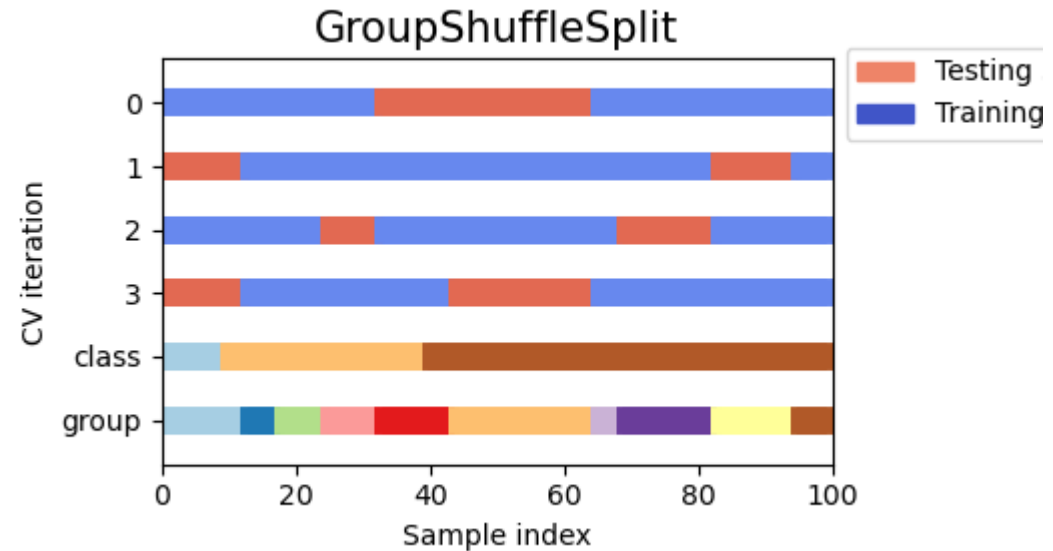  `from sklearn.pipeline import FeatureUnion`

```python
1   word_pipe = Pipeline([
2       ('selector', ItemSelector(key='tokenized')),
3       ('vect', TfidfVectorizer(analyzer='word'))
4   ])
5
6   pos_pipe = Pipeline([
7       ('selector', ItemSelector(key='pos')),
8       ('vect', TfidfVectorizer(analyzer='word'))
9   ])
10
11  num_cols = ['total_punct', 'punct_dist'] + [col for col in X.columns
12                                              if col.startswith('liwc_')]
13
14  num_pipe = Pipeline([
15      ('selector', ItemSelector(key=num_cols)),
16  ])
17
18  punct_pipe = Pipeline([
19      ('selector', ItemSelector(key='count_punct')),
20      ('iter', RowIterator()),
21      ('vect', DictVectorizer())
22  ])
23
24  pipe = Pipeline([(
25          'feats',
26          FeatureUnion([
27              ('word', word_pipe),
28              ('pos', pos_pipe),
29              ('num', num_pipe),
30              ('punct', punct_pipe)
31          ],
32      )),
33      ('clf', LinearSVC(random_state=97, class_weight='balanced'))
34  ])
```

# Scikit-Learn **Classes** Spotlight (cont'd)

- **GroupShuffleSplit()**
  - Shuffle-Group(s)-Out cross-validation iterator
  - Provides randomized train/test indices to split data according to a third-party provided group

`from sklearn.model_selection import GroupShuffleSplit`



GroupShuffleSplit

```
1  gss_cv = GroupShuffleSplit(n_splits=10, test_size=0.2, random_state=97)
2  train_idx, test_idx = next(gss_cv.split(X, y, groups))   # GSS is a generator, so we need to call next()
3
4  X_train, X_test, y_train, y_test = X[train_idx], X[test_idx], y[train_idx], y[test_idx]
5  X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

University of Antwerp
Faculty of Arts

# Conclusions & Takeaways

## Tasks

- There's always more to learn in the packages we use every day!
- Tokenization performs better than uncleaned and semicleaned text (kudos to NLP group project)
- Still room to improve models for secondary/tertiary author features

**Better data (+ better features) > better models**

## Overall experience

- Pros:
  - Met awesome people!
  - Participated in very intriguing meetings
  - Tested and expanded the knowledge acquired during the past year

- Cons:
  - Always (or mostly) on-site would be nice
  - Some technical difficulties
  - Put more time on other stuff

University of Antwerp
Faculty of Arts

Dank u wel!
Nog vragen?

When there's a task that can be done manually in 10 minutes but you find a way to automate it in 10 days

I'm gonna do what's called a programmer move.