



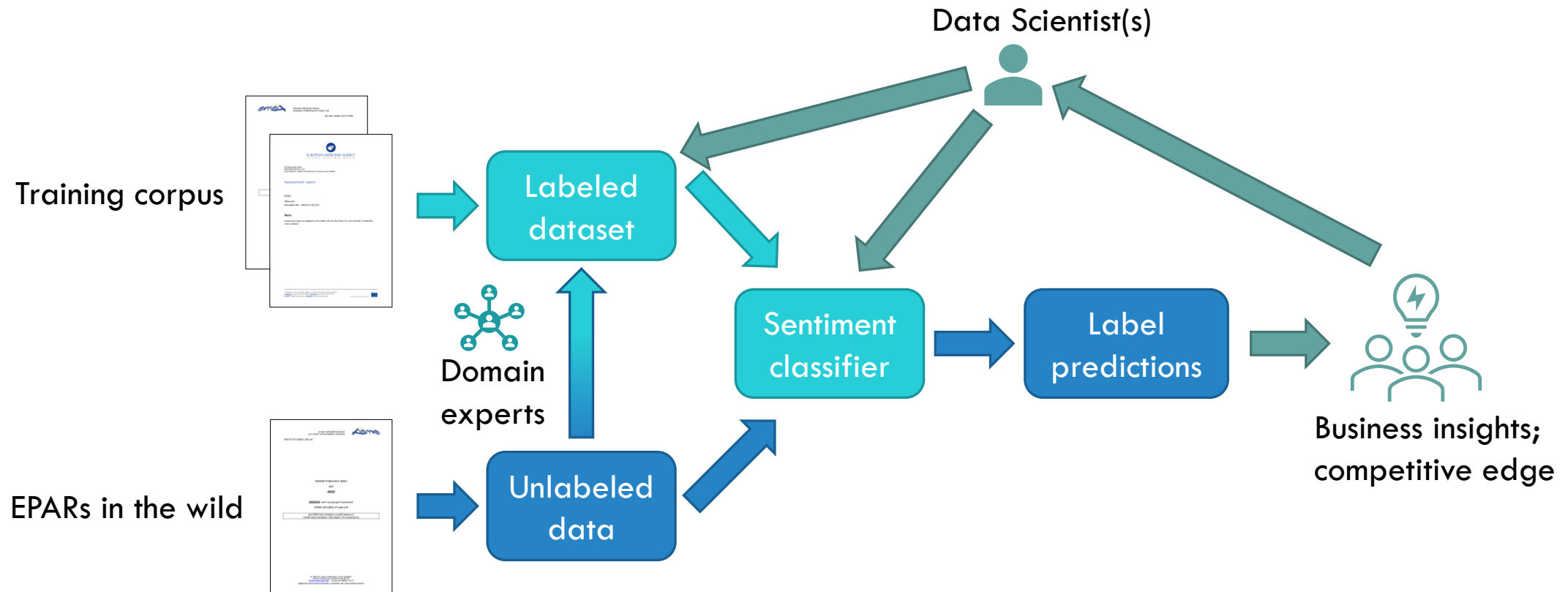
SENTIMENT ANALYSIS FROM EPAR DOCUMENTS

Panu Aho, Data Scientist
17.1.2021

AGENDA

- Business objective of research
- Training data overview
- Sentiment classification model pipeline
 - Data cleaning
 - Data preprocessing
 - Classifier training
- Results
- Discussion
 - Achievements
 - Limitations
 - Future research

BUSINESS OBJECTIVES



TRAINING DATASET

- An annotated dataset was provided in a 266×5 matrix containing an integer ID, Sentence and one-hot-encoded sentiment labels
- Out of the 266 rows, 236 were unique while 30 were at least once duplicated (not taking in to account the ID column)

ID			Sentence	Positive	Negative	Neutral
0	1	The results in 2nd line treatment show an ORR ...		1	0	0
1	2	The long duration of response and high durable...		1	0	0
2	3	The median OS time in the updated results exce...		0	0	1
3	4	Therefore, the clinical benefit in 2nd line tr...		1	0	0
4	5	The data provided in 1st line, although prelim...		1	0	0

DUPLICATED SENTENCES

- 30 rows contained sentences that were identical replicas of other rows (excluding the first occurrence)
 - In these instances, also the labels were identical
- In addition, some near-identically replicated sentences (~spelling mistakes) were found by examining the data visually
- Without further knowledge of the data gathering/annotation process, the decision was made to retain the replicated in data.

Identical training examples

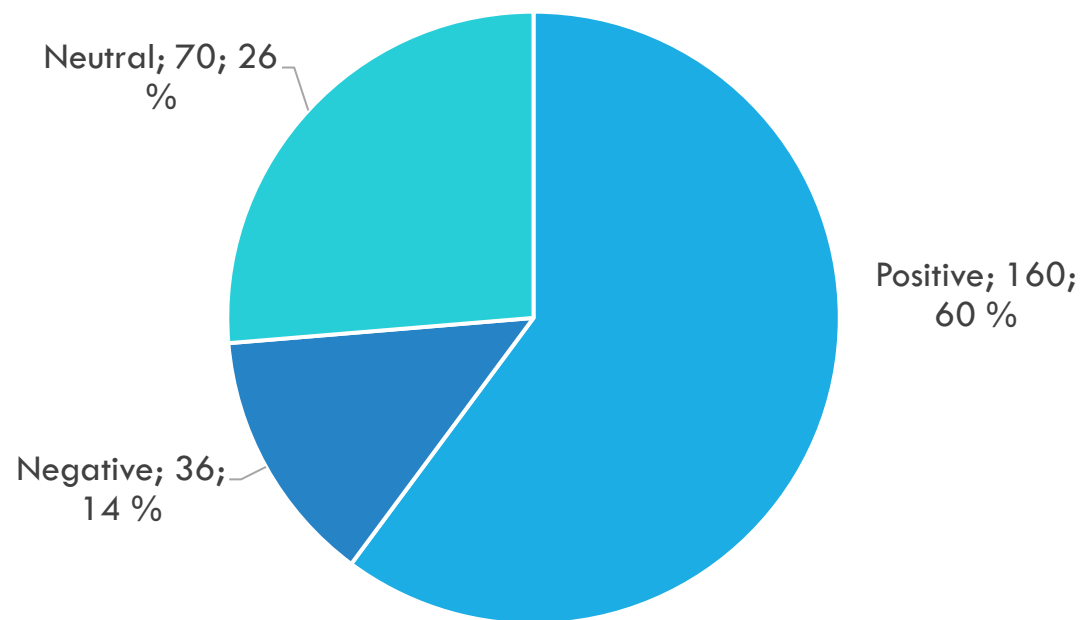
21	Biosimilarity of CT-P10 and MabThera is considered demonstrated based on the efficacy data.
137	Biosimilarity of CT-P10 and MabThera is considered demonstrated based on the efficacy data.
148	Biosimilarity of CT-P10 and MabThera is considered demonstrated based on the efficacy data.

Near-identical training examples

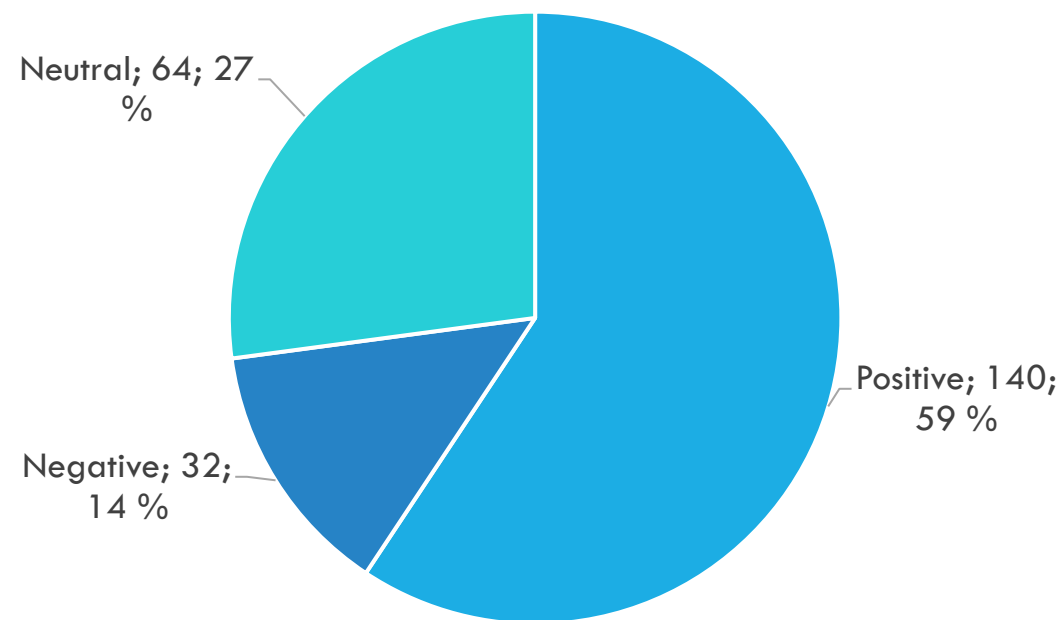
30	Additional safety data from maintenance study period CT-P10 3.3 and follow-up period should be provided (see RMP).
157	Additional safety data from Maintenance Study Period CT-P10 3.3. and Follow-up Period should be provided (see RMP).
146	Additional safety data from Maintenance Study Period CT-P10 3.3.and Follow-up Period should be provided (see RMP).

LABEL DISTRIBUTION

All rows



Unique rows



DATA EXPLORATION

- Frequent (top-ten) 1-grams, 2-grams and 3-grams were extracted from the cleaned data
- Top-ten 1-grams have many terms replicated in both positive and negative classes, but the occurrence rates are different
- Negative class frequently has terms with obvious negative connotation

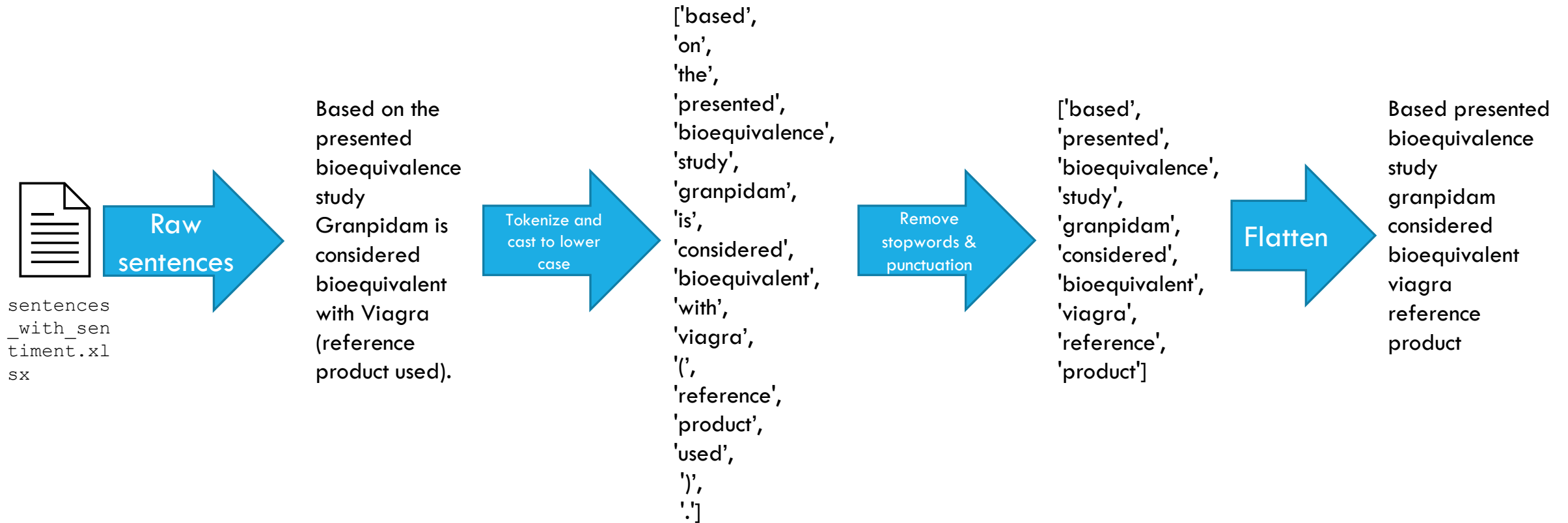
	Positive	Pos_rate	Negative	Neg_rate	Neutral	Neutr_rate
1	safety	0.29375	safety	0.472222	studies	0.300000
2	data	0.28125	data	0.388889	safety	0.242857
3	study	0.20625	patients	0.333333	study	0.214286
4	efficacy	0.18750	study	0.250000	ct-p10	0.171429
5	clinical	0.17500	should	0.222222	efficacy	0.157143
6	patients	0.16875	treatment	0.194444	data	0.157143
7	considered	0.16250	limited	0.194444	patients	0.142857
8	treatment	0.15000	further	0.166667	dose	0.142857
9	profile	0.14375	address	0.166667	insulin	0.142857
10	product	0.13125	efficacy	0.166667	product	0.128571

DATA EXPLORATION — FREQUENT 2- AND 3-GRAMS

	Positive	Pos_rate	Negative	Neg_rate	Neutral	Neutr_rate
1	(safety, profile)	0.12500	(chmp, considers)	0.111111	(insulin, glargine)	0.085714
2	(clinical, data)	0.06875	(considers, following)	0.111111	(safety, profile)	0.057143
3	(ct-p10, mabthera)	0.05625	(following, measures)	0.111111	(reference, products)	0.057143
4	(efficacy, data)	0.04375	(necessary, address)	0.111111	(safety, data)	0.057143
5	(reference, product)	0.04375	(measures, necessary)	0.083333	(pivotal, studies)	0.057143
6	(safety, data)	0.03750	(address, missing)	0.083333	(et, al)	0.057143
7	(comparable, between)	0.03750	(address, issues)	0.083333	(efficacy, safety)	0.057143
8	(between, ct-p10)	0.03750	(issues, related)	0.083333	(medicinal, product)	0.042857
9	(bioequivalence, study)	0.03750	(although, dataset)	0.083333	(overall, safety)	0.042857
10	(film-coated, tablets)	0.03750	(dataset, afl)	0.083333	(profile, ct-p10)	0.042857

	Positive	Pos_rate	Negative	Neg_rate	Neutral	Neutr_rate
1	(between, ct-p10, mabthera)	0.03750	(chmp, considers, following)	0.111111	(overall, safety, profile)	0.042857
2	(based, efficacy, data)	0.02500	(considers, following, measures)	0.111111	(safety, profile, ct-p10)	0.042857
3	(data, considered, supportive)	0.02500	(following, measures, necessary)	0.083333	(profile, ct-p10, appeared)	0.042857
4	(mg, film-coated, tablets)	0.02500	(measures, necessary, address)	0.083333	(ct-p10, appeared, roughly)	0.042857
5	(2nd, line, treatment)	0.01875	(necessary, address, issues)	0.083333	(appeared, roughly, similar)	0.042857
6	(biosimilarity, ct-p10, mabthera)	0.01875	(address, issues, related)	0.083333	(roughly, similar, reference)	0.042857
7	(ct-p10, mabthera, considered)	0.01875	(although, dataset, afl)	0.083333	(similar, reference, product)	0.042857
8	(mabthera, considered, demonstrated)	0.01875	(dataset, afl, patients)	0.083333	(reference, product, although)	0.042857
9	(considered, demonstrated, based)	0.01875	(afl, patients, updated)	0.083333	(product, although, pooled)	0.042857
10	(demonstrated, based, efficacy)	0.01875	(patients, updated, data)	0.083333	(although, pooled, incidences)	0.042857

DATA CLEANING



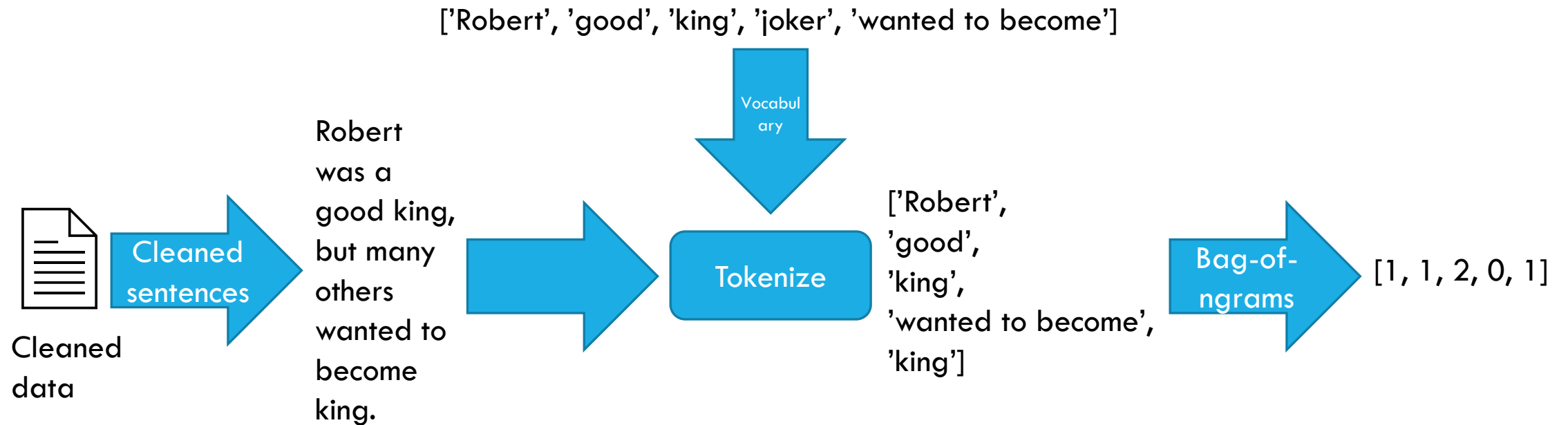
STOP WORDS

- A starting point for stop words list was obtained from the [nltk](#) Python library
- `nltk.corpus.stopwords.words('english')`
- The list was further hand-tweaked in an attempt to reduce the noise present by meaningless words such as 'a', 'the', 'it' etc. while keeping acceptable discriminative power between classes.

```
1 print(STOP_WORDS)
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves',  
'you', "you're", "you've", "you'll", "you'd", 'your', 'your  
s', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself',  
'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it  
s', 'itself', 'they', 'them', 'their', 'theirs', 'themselv  
s', 'what', 'which', 'who', 'whom', 'this', 'that', "that'l  
l", 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be',  
'been', 'being', 'have', 'has', 'had', 'having', 'do', 'doe  
s', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'o  
r', 'as', 'of', 'at', 'by', 'for', 'with', 'about', 'into',  
'through', 'during', 'to', 'from', 'in', 'out', 'on', 'off',  
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'ho  
w', 'both', 'each', 'other', 'such', 'own', 'so', 's', 't',  
'can', 'will', 'just', 'now', 'd', 'll', 'm', 'o', 're', 'v  
e', 'y']
```

FEATURE EXTRACTION



Illustrative example – not actual training data

VOCABULARY?

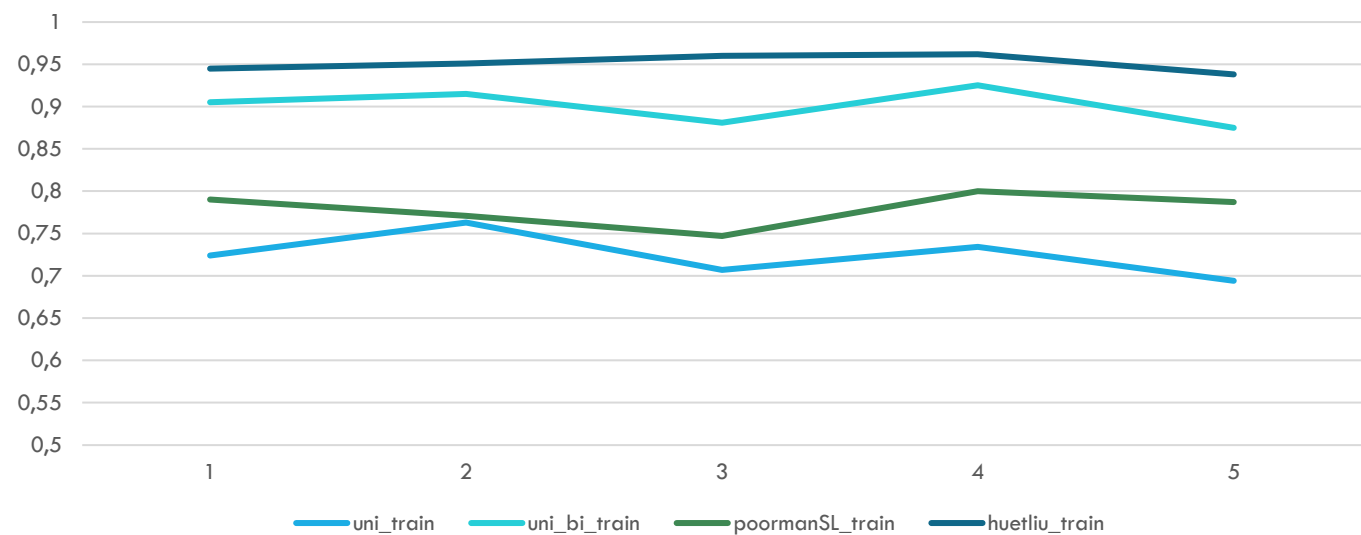
- Configurations tried:
 - Top ten 1-grams
 - Top ten 1-grams + 2-grams
 - Top ten 1-grams + hand-engineered "poor man's" *Sentiment Lexicon (SL)*
 - Contains 11 "low-hanging fruit" phrases obviously associated¹ with one of the classes and repeated multiple times in dataset, e.g.:
 - 'These objectives have been met'
 - 'Data are considered very limited'
 - The phrases were gathered by arranging sentences alphabetically and visually examining them without looking at the labels
 - Top ten 1-grams + poor man's SL + ~6800 word open source English SL (Hu & Liu 2004) [1]

1) According to my subjective perception as a non-domain expert

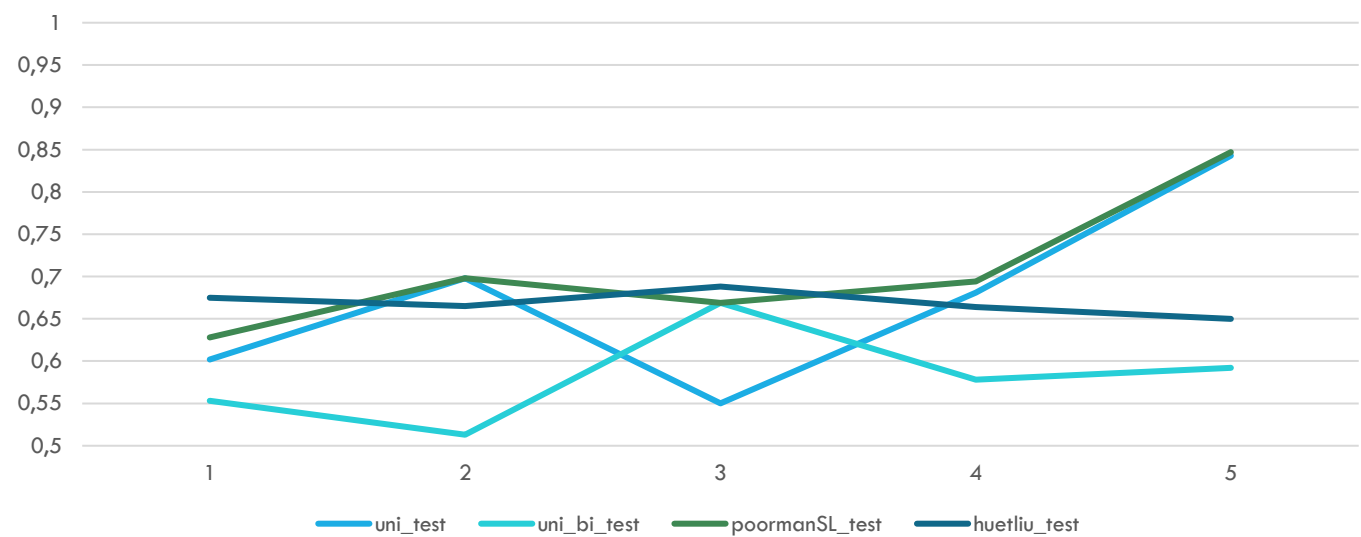
TRAINING

- A supervised classification model was trained using the [scikit-learn](#) library with the different vocabularies
- A support vector machine with linear kernel was chosen as initial approach due to its well-reported performance in similar tasks [2, 3, 4]
- Default parameters
- Random train/test split, 5-fold cross validation → 20 % used for testing per fold

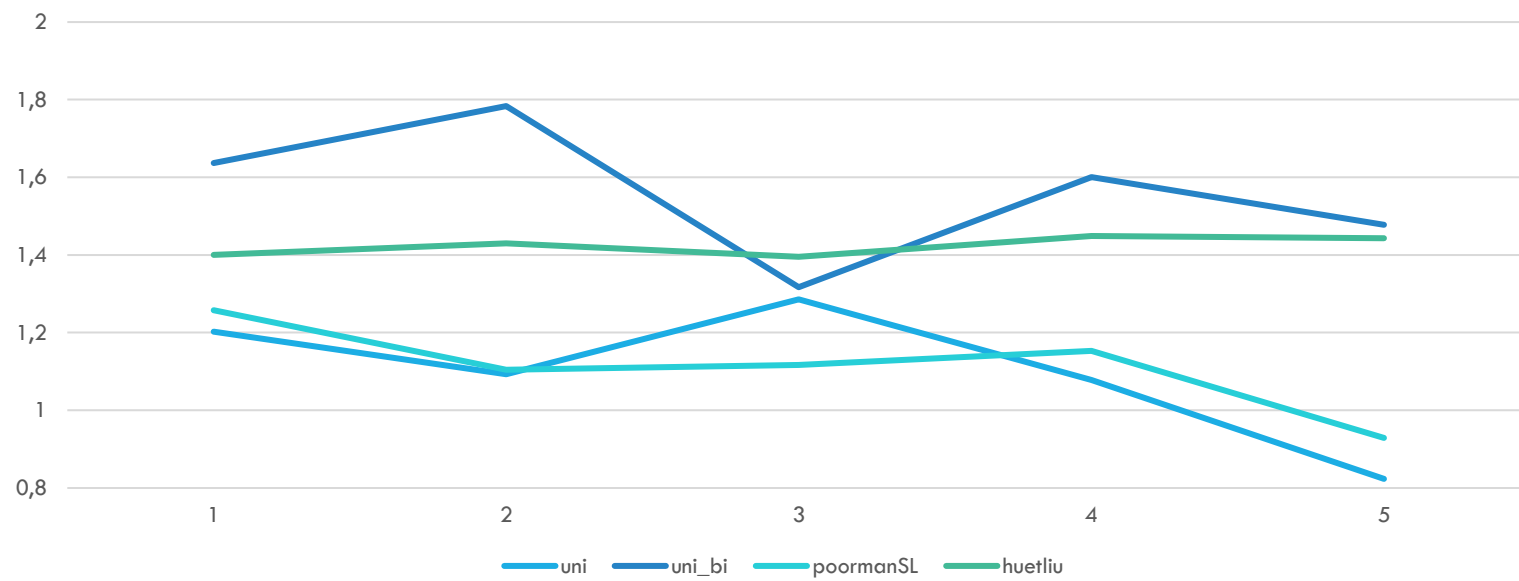
Training error per fold (Balanced Accuracy Score)



Test error per fold (Balanced Accuracy Score)

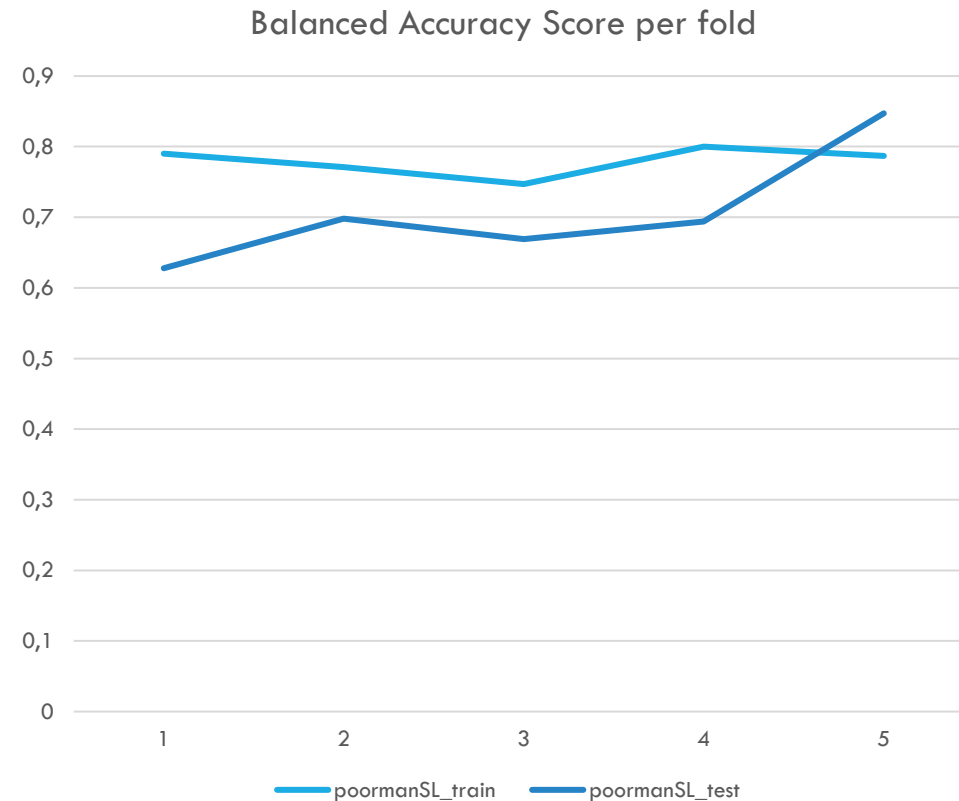


Overfitting (BAS-Train / BAS-Test)



RESULTS

- Hand-engineered domain specific sentiment lexicon together with most popular 1-grams seems like a promising approach w.r.t.:
 - Good model robustness across folds
 - Not grossly overfitting
- Using a generic open-source SL also yielded surprisingly good results in terms of robustness, but with significant model overfitting
- In other approaches tried, robustness was questionable (~widely varying performance accross folds in test set)



DISCUSSION

- Achievements
 - A simple bag-of-ngrams approach with hand-selected features and linear SVN shows promising results
 - Naive approaches (e.g. guessing the largest class) are outperformed by a comfortable margin
- Limitations
 - Dataset is too small to effectively make use of complex features such as 2-grams and beyond (→ Risk overfitting)
 - Class proportions slightly skewed
 - BoW does not take word ordering and other context in to account, so model does not probably misses important information

DISCUSSION

- Suggestions for future work
 - Build a CI/CD pipeline, deploy the model to production and see how it performs in the wild!
 - Improve the Sentiment Lexicon vocabulary
 - Additional feature engineering – look into *Scientific Citation Sentiment Analysis (SCSA)* [5, 6]
 - If possible to get a LOT of data, could try more complex approaches e.g. bidirectional LSTM

Example		<i>Our data are generally consistent with that of other studies [TC], but not with studies where a single dose of paracetamol was administered [OC].</i>	<i>These values are lower than those reported by French et al. [20].</i>
Step 1		<i>[Our data]CITINGWORK are generally [consistent with]POSITIVE that of other studies [TC]TC , [but]CONTRAST [not]NEGATION with studies where a single dose of paracetamol was administered [OC]oc</i>	<i>[These values]CITINGWORK are [lower]COMPARATIVE [than]THAN those reported by French et al. [TC]TC</i>
Step 2	Unigram	'CITINGWORK', 'POSITIVE', 'TC', 'CONTRAST', 'NEGATION', 'OC'	'CITINGWORK', 'COMPARATIVE', 'THAN', 'TC'
	Bigram	'CITINGWORK_POSITIVE', 'POSITIVE_TC', 'TC_CONTRAST', 'CONTRAST_NEGATION', 'NEGATION_OC'	'CITINGWORK_COMPARATIVE', 'COMPARATIVE_THAN', 'THAN_TC'
	Trigram	'CITINGWORK_POSITIVE_TC', 'POSITIVE_TC_CONTRAST', 'TC_CONTRAST_NEGATION', 'CONTRAST_NEGATION_OC'	'CITINGWORK_COMPARATIVE_THAN', 'COMPARATIVE_THAN_TC'
	Direction	'CITINGWORK_CONTRAST_DIR', 'TC_CONTRAST_DIR', 'CONTRAST_OC_DIR'	

Examples of extracting *structure features* from citation data related to clinical trials, as proposed in [6]

REFERENCES

- [1] Mingqing Hu and Bing Liu. "Mining and Summarizing Customer Reviews." Proceedings of the ACM SIGKDD International Conference on Knowledge. Discovery and Data Mining (KDD-2004), Aug 22-25, 2004, Seattle, Washington, USA.
- [2] Walaa Medhat, Ahmed Hassan, Hoda Korashy. Sentiment analysis algorithms and applications: A survey. Ain Shams Engineering Journal. Volume 5, Issue 4. 2014
<https://doi.org/10.1016/j.asej.2014.04.011>
- [3] Mullen, Tony & Collier, Nigel. (2004). Sentiment Analysis using Support Vector Machines with Diverse Information Sources. 412-418.
- [4] Vinodhini, G., and R. M. Chandrasekaran. "Sentiment analysis and opinion mining: a survey." International Journal 2.6 (2012): 282-292.
- [5] Yousif, A., Niu, Z., Tarus, J.K. et al. A survey on sentiment analysis of scientific citations. Artif Intell Rev 52, 1805–1838 (2019). <https://doi.org/10.1007/s10462-017-9597-8>
- [6] Xu J, Zhang Y, Wu Y, Wang J, Dong X, Xu H. Citation Sentiment Analysis in Clinical Trial Papers. AMIA Annu Symp Proc. 2015;2015:1334-1341. Published 2015 Nov 5.

THANK YOU!

Panu Aho

Data Scientist
+358 50 4129 150
panu.aho@gmail.com
[linkedin.com/in/panuaho/](https://www.linkedin.com/in/panuaho/)

