TEXT MINING

L12. CONCLUSIONS

SUZAN VERBERNE 2021



TODAY'S LECTURE

- Course summary
 - Tasks
 - and evaluation
 - Models
 - and resources
 - > Applications
 - and domains
- Exercises to practice for the exam
- Schedule of tests (exam/assignment)



COURSE SUMMARY



MODELS, TASKS, APPLICATIONS









TASKS IN THE TM PIPELINE





PRE-PROCESSING

Cleaning

- PDF/docx/HTML to text
- Language filtering
- Encoding issues
- Regex patterns
- Spelling correction

Linguistic pipeline

- Tokenization
- Stop word removal
- Lemmatization/stemming
- POS-tagging

Minimal edit distance



CLASSIFICATION

- Multi-class vs multi-label
- Feature selection
- Term weighting: tf-idf



INFORMATION EXTRACTION

Named entity recognition

- Segmentaton & classification
- Sequence labelling task with IOB-labels
- Rule-based, feature-based, neural-network based
- Help of dictionaries

Relation extraction

- Co-occurrences, patterns, classification
- Distant supervision



SUMMARIZATION

- Extractive: sentence classification
- Abstractive: sequence-to-sequence (compare with translation)

> Challenges:

- Training data (ground truth, humans disagree)
- Evaluation



EVALUATION

Intrinsic: comparison against ground truth (=human) for task

- Metrics:
 - Accuracy
 - Precision, recall, F1
 - ROUGE for summarization
- Train-test split to prevent overfitting, or cross validation
 - Hyperparameter tuning on train-tune-set, or cross validation (GridSearchCV)

Extrinsic: effectiveness in context



MODELS, TASKS, APPLICATIONS







Neural language models (embeddings)

- Traditional: word2vec (a NN with 1 hidden layer) and others
- Transformer-based: BERT
- Goal: from high-dimensional sparse vectors (10,000s) to lowerdimensional (~100-800) dense vectors
- Pre-trained as a word/sentence prediction task
 - = Language modelling
 - The hidden layer has the dimensionality of the embeddings

Distributional hypothesis



Classification

- Vector space model, bag of words
- Dimensionality reduction
- Machine learning methods
 - Naïve Bayes (probabilistic)
 - Support Vector Machines (vector space)
 - Feedforward neural networks (compare with logistic regression, multi-node, multi-layer)
 - Transformer models (BERT)



Sequence labelling

- Conditional Random Fields
- Recurrent Neural Networks / Bi-LSTMs
- Transformer models: BERT

Sequence-to-sequence

- Encoder-decoder models
- Transformer models



TRANSFER LEARNING

Pre-trained neural language models allow for transfer learning

Inductive transfer learning: transfer the knowledge from pretrained language models to any text mining task

Fine-tuning

Current state-of-the-art for many text mining tasks



RESOURCES

Labeled data

- For training and evaluating task-specific models
- Typically small (1000s of examples, or even 100s)
- Supervised learning
- How to obtain labelled data:
 - Benchmark data
 - Existing expert labels
 - User-generated content
 - Data annotation (crowdsourcing)

Inter-rater agreement (Cohen's Kappa)



RESOURCES

Unlabeled data

- For pre-training language models (language modelling = word prediction/sentence prediction)
- General vs domain-specific
- Typically large (Millions of words)



RESOURCES

- Dictionaries (gazetteers)
- > Ontologies / taxonomies

Controlled vocabulary

General domain (names, places) or specific domain (e.g. bio)



APPLICATIONS



APPLICATIONS

Sentiment analysis

- Classification / ordinal regression / linear regression
- Extraction: aspect-based sentiment analysis (E,A,S,H,C)

Argument mining

- Sentence classification (classifying argument components into different types such as claims and premises)
- Structure identification focuses (linking arguments or argument components)



APPLICATIONS

CV-vacancy matching

- Extracting the structure of a CV
- > Extracting structured fields (e.g. name, education) from a CV
- Mapping words with similar meaning (e.g. function titles) between CV and vacancy



DOMAINS

- Social media analytics (classification, extraction, sentiment)
 - E-commerce (sentiment classification/extraction, ontologies)
- Biomedical text mining (classification/retrieval, extraction)
 - Clinical applications (EHRs)
- Humanities, historical documents (pre-processing, classification/retrieval, extraction)



EXERCISES



EXERCISES

- Minimal edit distance
- ➤ Tf-idf
- Naïve Bayes
- Inter-rater agreement
- Precision and recall



MINIMAL EDIT DISTANCE

Compute the Levenshtein distance between 'where' and 'hear'. Show your computation.



MINIMAL EDIT DISTANCE





MINIMAL EDIT DISTANCE

cost	operation	input	output
1	delete	W	
0	(сору)	h	h
0	(сору)	e	e
1	substitute	r	а
1	substitute	e	r
Total: 3			



TF-IDF

> We have a collection of 100,000 movie scripts.

- a. The term *Elsa* occurs in 10 scripts. What is the inverse document frequency for *Elsa*? Show your computation.
- b. We have a film script *s* in which *Elsa* occurs 2 times. What is the tf-idf weight for *Elsa* in *s*?



TF-IDF

- a. $idf = log_{10}(100,000/10) = log_{10}(10,000) = 4$
- b. $tf = 1 + \log_{10}(2) \approx 1.3$ $tf^*idf = 1.3^*4 = 5.2$



NAÏVE BAYES

Consider this toy training set for a text classification task with Naïve Bayes:

Doc id	Content	Class
1	make our garden grow	relevant
2	we make the best of it	not relevant
3	together we can grow	not relevant
4	we make the best plans	not relevant



Doc id	Content	Class
1	make our garden grow	relevant
2	we make the best of it	not relevant
3	together we can grow	not relevant
4	we make the best plans	not relevant

- a. What is the prior probability of the 'relevant' class?
- b. What is the vocabulary size of the training set? Assume that we do not remove stop words.
- c. Estimate *P('make', not relevant)* using the maximum likelihood estimate on the train set.
- d. Why is add-one smoothing needed when we estimate the probability of an unseen document? Provide an example test document given the toy training set for which add-one smoothing is needed.



Doc id	Content	Class
1	make our garden grow	relevant
2	we make the best of it	not relevant
3	together we can grow	not relevant
4	we make the best plans	not relevant

a. 1/4

- b. 12
- c. (2+1)/(15+12)
- d. Because a word in the test document that does not occur in the training set will have a zero probability and the multiplication of zero probabilities will lead to a combined probability of zero; a correct example would be any text with a word that does not occur in training set.



INTER-RATER AGREEMENT

Compute Cohen's Kappa for this agreement table. Show your computation. (You can keep the last fraction of your computation as it is, without estimating the decimal numbers.)

Agreement table		Annotator 2			
		Positive	Negative	Neutral	
Annotator 1	Positive	25	10	5	
	Negative	0	25	15	
	Neutral	5	5	10	



INTER-RATER AGREEMENT

Agreement table		Annotator	2		
		Positive	Negative	Neutral	
	Positive	25	10	5	40
	Negative	0	25	15	40
	Neutral	5	5	10	20
		30	40	30	100
Pr(a	a) 60) / 100		0.6	
Pr(e	e,pos) 40	/ 100 * 30 /	100	0.12	
Pr(e	e,neg) 40	/ 100 * 40 /	100	0.16	
Pr(e	e,neut) 30	/ 100 * 20 /	100	0.06	
Pr(e	e) (si	um of the 3 r	ows above)	0.34	
kap	pa (0	5 - 0.34) / (1 - 0.34) = 0.3		26/0.66	



PRECISION AND RECALL

- Consider the following output table of an automatic classifier for 10 documents. Compute (please show the fractions):
- a. the recall for the A class
- b. the precision for the A class

doc	class assigned	ground
id	by classifier	truth class
1	А	С
2	А	А
3	В	В
4	В	А
5	С	С
6	А	А
7	D	А
8	А	D
9	В	В
10	С	А



PRECISION AND RECALL

- Consider the following output table of an automatic classifier for 10 documents. Compute (please show the fractions):
- a. Recall(A) = 2/5
- b. Precision(A) = 2/4

doc	class assigned	ground
id	by classifier	truth class
1	А	С
2	А	А
3	В	В
4	В	А
5	С	С
6	А	А
7	D	А
8	А	D
9	В	В
10	С	А





SUZAN VERBERNE 2021



HOMEWORK

Work on the final assignment

Prepare for the exam

- Exam Text mining, Thursday January 13, 10.15 – 13.15
 - (Students with a provision card get an addition 30 minutes)
- Location: GORL / 04/5
- The exam is closed book, individual
- Practice materials are on Brightspace



THANK YOU

Thank you for the participation in this course,

- > And a big thanks to the TAs:
 - Michiel van der Meer
 - Juan Bascur Cifuentes
 - Cheyenne Health
 - Hainan Yu



TIME FOR XKCD



