NLP for Operationalizing (Viewpoint) Diversity

Myrthe Reuver CLTL, VU Amsterdam

VRIJE UNIVERSITEIT AMSTERDAM

Vrijenhoek et al. (2021)'s Metrics: NLP operationalization needed

- 1) **Fragmentation:** shared public sphere
- 2) **Representation:** diverse actors and **opinions**
- Alternative Voices: non mainstream opinions
- 4) Calibration: personalization
- 5) Affect: emotional content



Current methods: stance

- **Models:** Pre-trained Large Language Models such as BERT and RoBERTa
- **Data: Stance Benchmark** (Schiller et al., 2021) combines 10 different stance datasets:

Dataset	Domain	Topic	Comment	Stance
ibmcs	Encyclopedia	[] atheism is the only way	Atheism is a superior basis for ethics	PRO
semeval2019t7	Social media	(Charlie Hebdo)	"[] #CharlieHebdo gunmen have been killed" yayyy []	Support
semeval2016t6		Feminist Movement	[] every women should have their own rights !! #SemST	Favor
fnc1	News	Hugh Hefner Dead?	Hugh Hefner has denied reports that he is dead []	Disagree
snopes		Farmers feed their cattle candy []	[] padding out cow feed with waste candy is nothing new.	Agree
scd	Debating forums	(Obama)	I think Obama has been a great President. []	For
perspectrum		School Day Should Be Extended	So much easier for parents!	Support
iac1		existence of god	[] the Bible tells me that Jesus existed []	Pro
arc		Salt should have a place at the table	[] the iodine in salt is necessary to prevent goiter. []	Agree
argmin	Web search	school uniforms	We believe in freedom of choice.	CON

Topics in parentheses signal implicit information

Cross-topic, cross-domain stance

Main question of cross-topic stance detection:

can we detect stance (pro, con)

on topics or issues not seen in training?

(The news always has new topics coming up!)

Reuver, M. E., Verberne, S., Vallejo, R. M., & Fokkens, A. (2021, November). Is Stance Detection Topic-Independent and Cross-topic Generalizable?-A Reproduction Study. In Proceedings of the 8th Workshop on Argument Mining,

OpenClipart Vectors @ Pixabay

Reimers et. al. (2019): cross-topic stance classification

Train: 7 topics, test: 8th topic Fine-tuning BERT (base & large) Findings:



Marco Verch @ Flickr, Creative Commons 2.0. https://foto.wuestenigel.com/businessman-walking-from-a-to-b-point/

- avg. F1 (10 seeds) = 0.633
- +0.20 over reference model (LSTM)
- Results are *"very promising and stress the feasibility of the task"* (Reimers et al. 2019, p. 575)

Mean (stdv) over 10 seeds	F1
Reimers et. al. (2019)	
LSTM (baseline)	.424
BERT-base	.613 (-)
BERT-large	.633 (-)
Reuver et. al (2021)	
SVM+tf-idf (baseline)	.517
Reproduction BERT-base	.617 (.006)
Reproduction BERT-base Reproduction BERT-large (all)	.617 (.006)
Reproduction BERT-base Reproduction BERT-large (all) BERT-large - 5 good seeds	.617 (.006) .596 (.043)

Reproduction (Reuver et. al. 2021b)

- BERT-large under-performs in 50% of seeds
- SVM+tf-idf model outperforms the LSTM reference model from the original study (F1 of .517 > .424)



- /



Alessandra Polimeno's work: Clustering with SBERT (for the Fragmentation metric) MA thesis supervised by Sanne, myself, and prof. dr. Antske Fokkens

HLGD: News Story Chains dataset

#	Topic	Size	\overline{x} Characters	\overline{x} Tokens
1	Human Cloning	108	4354	805
2	International Space Station	215	3141	597
3	Ireland Abortion Vote	170	4134	787
4	US Bird Flu Outbreak	75	2266	422
5	Facebook Privacy Scandal	172	4098	763
6	Wikileaks Trials	153	7398	1390
7	Tunisia Protests	86	3201	593
8	Ivory Coast Army Mutiny	104	2231	417
9	Equifax Breach	156	4041	744
10	Brazil Dam Disaster	247	2970	564

Table 3.1: Topics of the news story chains in HLGD, the number of articles in each chain, and the mean number of characters and tokens of the articles per chain.

Operationalize 'fragmentation' with NLP, and **evaluate** different operationalizations



Clustering evaluation: SBERT clearly wins

Setup	${\bf H}\uparrow$	$\mathbf{C}\uparrow$	$\mathbf{V}\uparrow$	$\mathbf{S}\uparrow$	$\mathbf{DBI}\downarrow$
Baseline	0.166	0.156	0.161	-0.060	12.441
AHC*SBERT	0.921	0.844	0.881	0.290	1.933
AHC*GloVe	0.762	0.708	0.734	0.183	1.913
AHC*BoW	0.813	0.658	0.727	0.413	1.965
DB*SBERT	0.694	0.872	0.773	0.231	1.509
DB*GloVe	0.002	0.236	0.004	0.390	0.387
DB*BoW	0.993	0.283	0.441	0.213	0.218

Table 5.1: Evaluation of the different representation methods (Sentence-BERT, word embeddings, and Bag of Words) and clustering methods (agglomerative hierarchical clustering, and DB-Scan), and the baseline. The measures are abbreviated as follows: H (homogeneity), C (completeness), V (V-measure), S (Silhouette Score), and DBI (Davies-Bouldin Index). The arrow indicates whether a high or low score is more desirable.

How does this measure the concept fragmentation 7

Scenario	Chains per user	Fragmentation
Scenario 1	7	Low
Scenario 2	1	High
Scenario 3, profile 1 (70%)	5	Balanced
Scenario 3, profile 2 (15%)	2	Balanced
Scenario 3, profile 3 (15%)	7	Balanced

Table 4.3: Overview of the number of chains that are present in the recommendation sets per scenario. In each scenario, there are 1000 users who receive a recommendation set containing 7 articles. Scenario 3 is build with 3 distinct user profiles that differ in the amount of story chains users are exposed to.

Setup	Scen. 1 \downarrow	Scen. 2 \uparrow	Scen. 3	Variation
Gold	0.00	0.85	0.58	0.85
Baseline	0.67	0.73	0.70	0.06
AHC*SBERT	0.31	0.87	0.64	0.56
AHC*GloVe	0.38	0.84	0.63	0.46
AHC*BoW	0.62	0.85	0.63	0.23
DB*SBERT	0.16	0.74	0.48	0.58
DB*GloVe	0.01	0.01	0.00	0.01
DB*BoW	0.99	0.99	0.99	0.00

Table 5.4: Fragmentation Scores for each setup per scenario



Current work;

Fragmentation metric:

- Expanding this work to new datasets and labels;
- Implementing this on MIND dataset → needed: evaluation data!
 Labelled data!

Representation metric

- Expanding this to new datasets;
- new experimental set-up (same vs different stance or argument)
- Needed: evaluation data! Labelled data!