

UNIVERSITY OF BERGEN

Analysis of Word Embeddings: A Clustering and Topological Approach

Master's thesis defence, 30 June 2021

By: Jonas Folkvord Triki
Supervisor: Nello Blaser

UNIVERSITY OF BERGEN



Agenda

1. Introduction
2. Methods
3. Results
4. Conclusion
5. Future Work



Introduction



Introduction

- Natural language processing (NLP)
- Words to word embeddings
- Single vector representations
- Polysemy of words, e.g.
 - run
 - open
 - make



Word embeddings

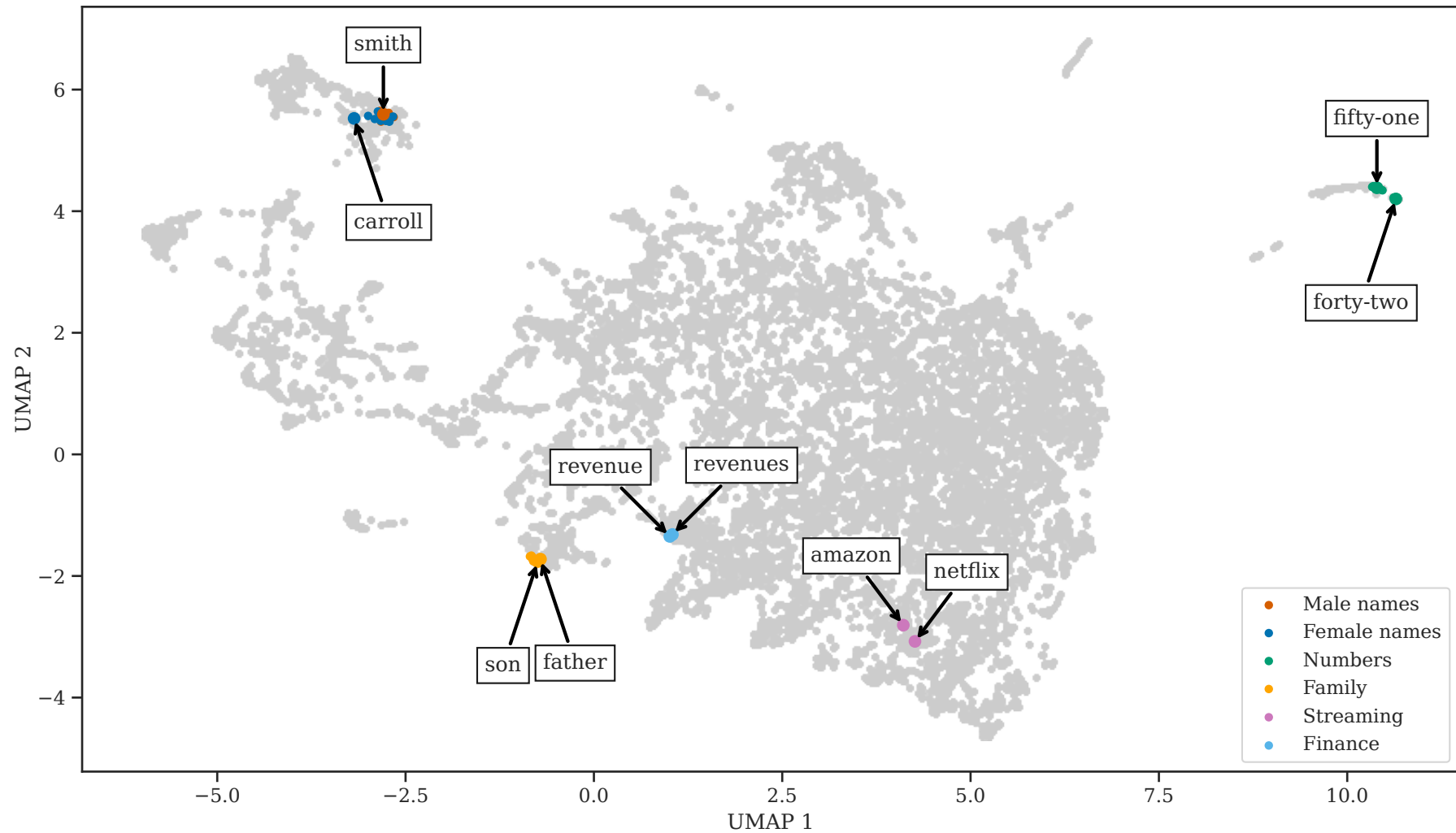
- Vector representations of words
 - Similar words have similar vector representations
 - 300 dimensions
- Word embedding models, e.g.
 - word2vec
 - GloVe
 - fastText



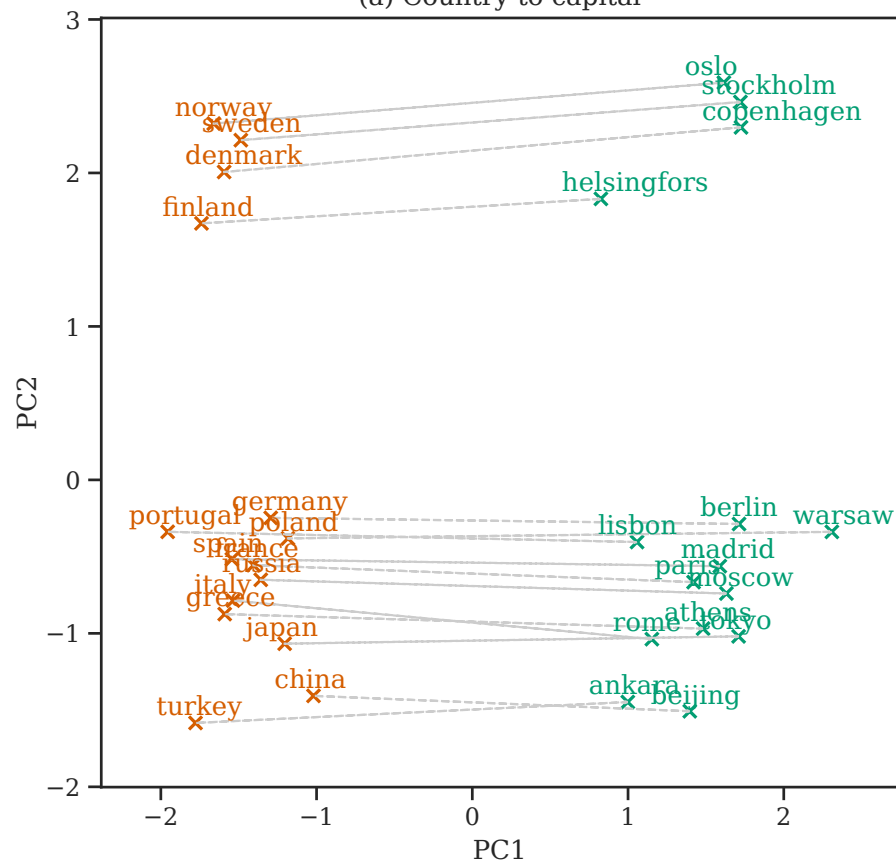
Source Text

Training Samples

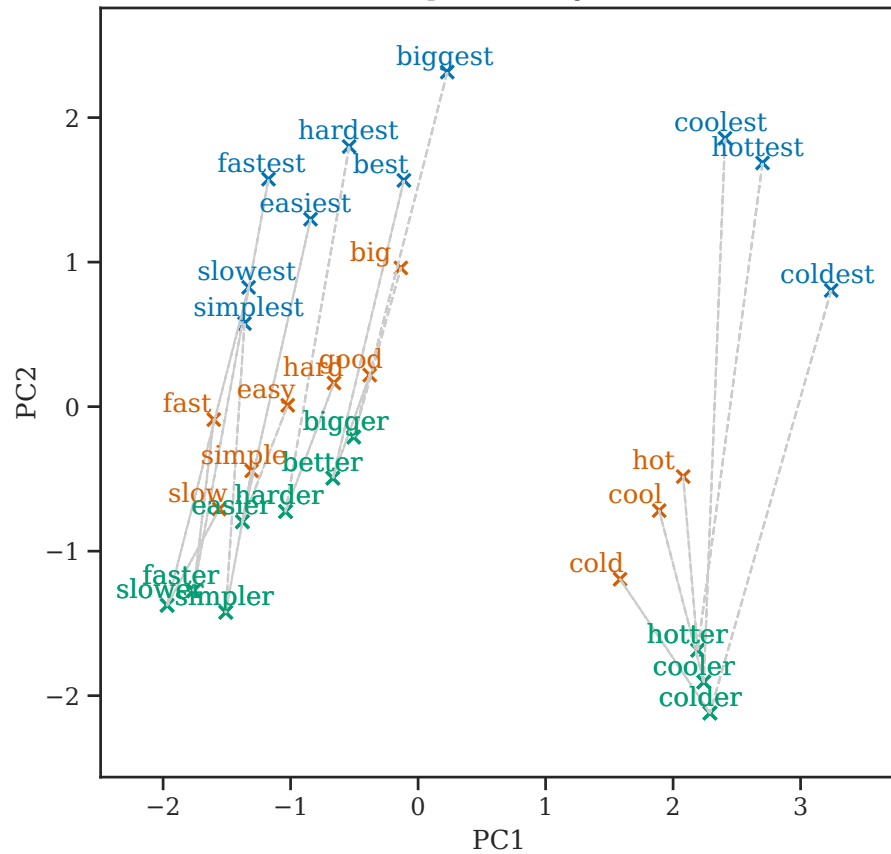
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)



(a) Country to capital



(b) Comparative adjectives



Motivation

- Finding hidden language relations
- Learning more from single vector representations
- Recent methods from topological data analysis (TDA)



Main goals

Deepen our understanding of word embeddings through

1. Cluster analysis
2. Prediction of word polysemy



Methods



Methods

- Analogy and cluster analysis
 - Classical word embedding model evaluation
 - Clustering and internal clustering validation algorithms
- Polysemy prediction
 - Topological polysemy
 - Geometric Anomaly Detection (GAD)
 - Intrinsic dimension estimation
 - Proposed supervised models



Analogy analysis

- Analogy test data sets
 - Syntactic, e.g. capital cities, currencies
 - Semantic, e.g. comparative, past tense
- Examples:
 - king – man \approx queen – woman
 - Paris – France \approx Norway – Oslo
 - good – better \approx cold – colder
 - work – works \approx speak – speaks



Analogy analysis

king – man \approx queen – woman

queen \approx king – man + woman

Converting to word embeddings:

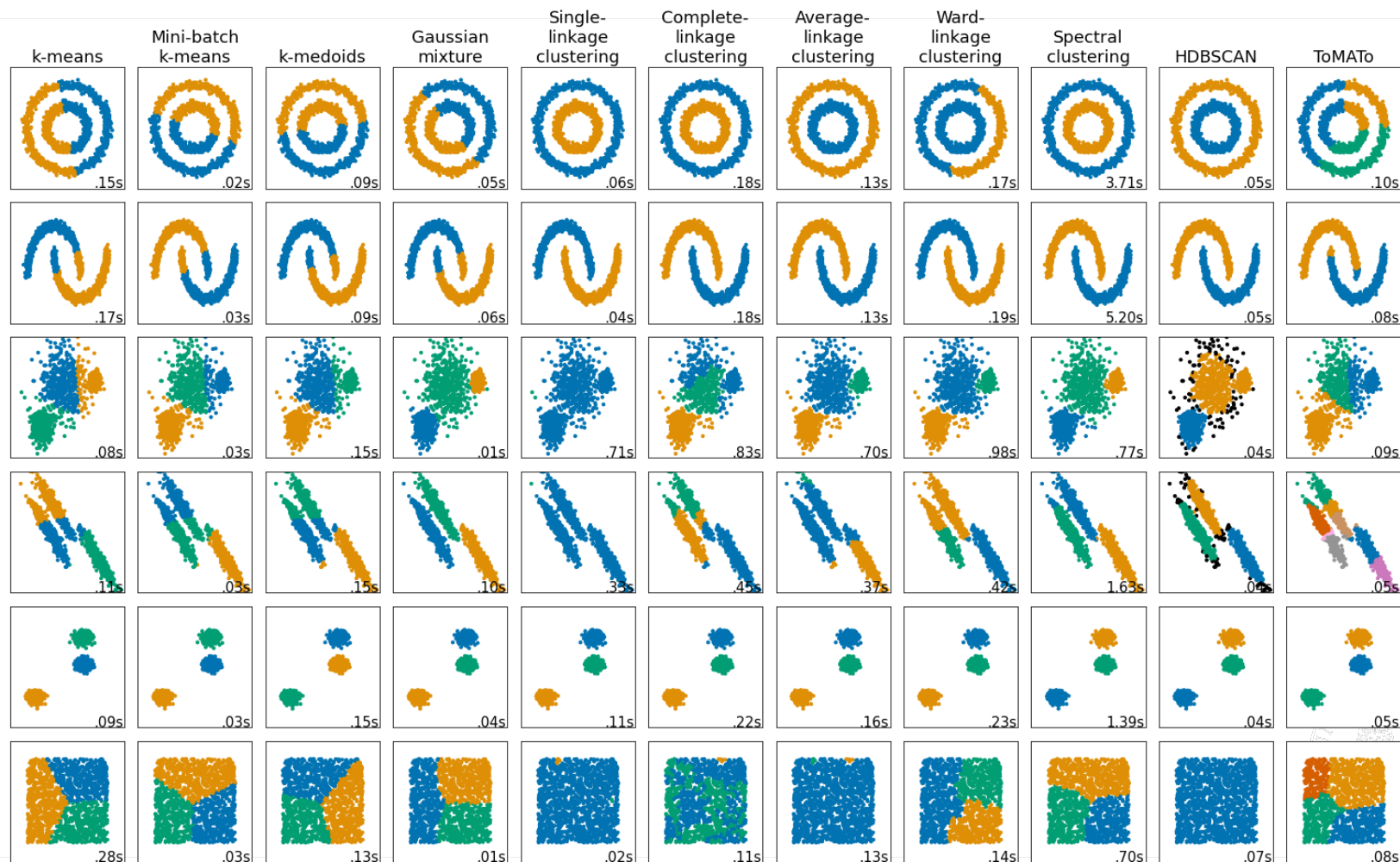
$$V_{\text{queen}} \approx V_{\text{king}} - V_{\text{man}} + V_{\text{woman}}$$



Cluster analysis

- Clustering algorithms
- Internal cluster validation methods







Cluster analysis

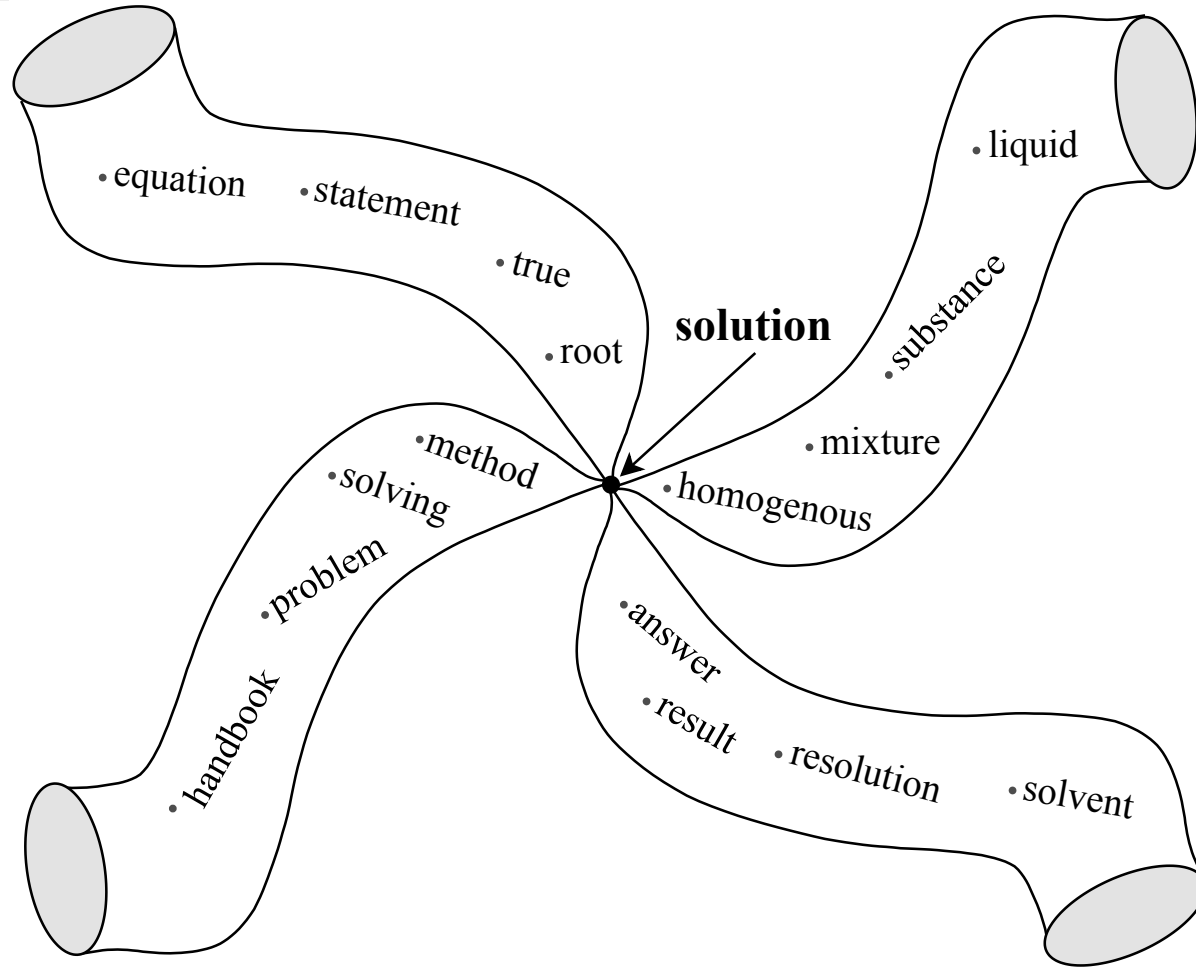
- Internal cluster validation methods
 - Silhouette Coefficient (SC)
 - Davies-Bouldin Index (DBI)
 - Caliński-Harabasz Index (CHI)

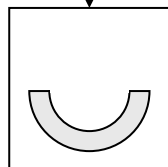
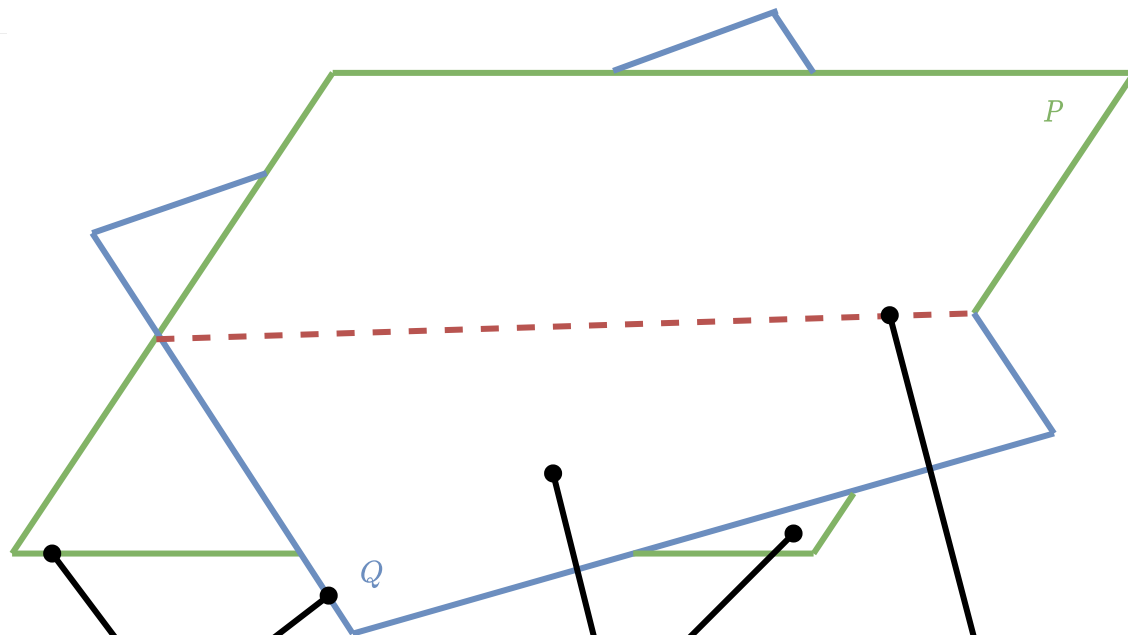


Polysemy prediction

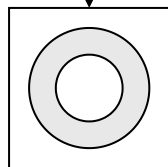
- Topological data analysis
 - Topological polysemy
 - Geometric Anomaly Detection
- Intrinsic dimension estimation
- Regression analysis
 - Lasso regression
 - Logistic regression



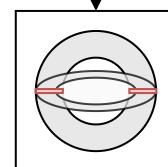




(a) Boundary point

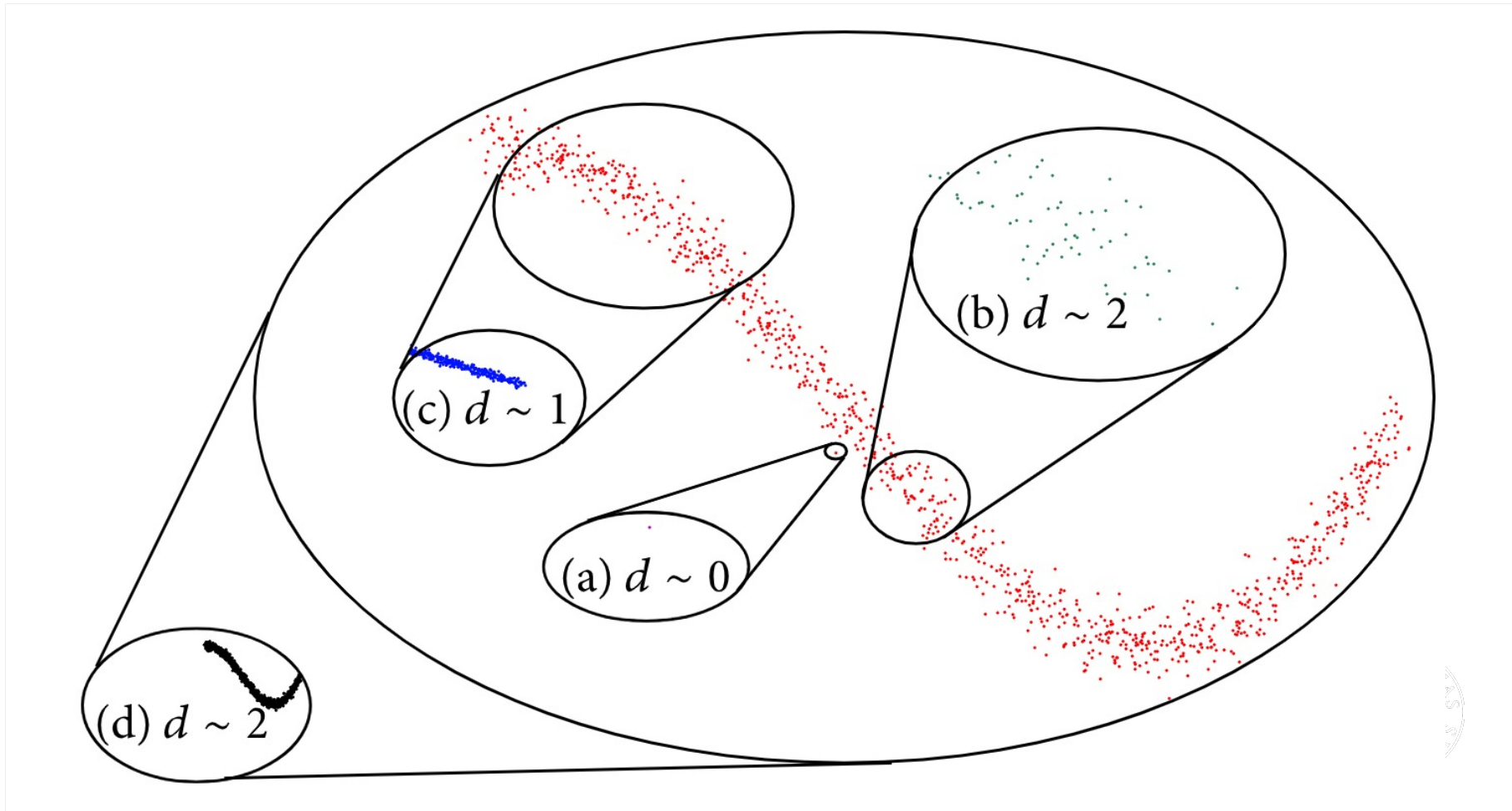


(b) Manifold point



(c) Singular point





Results



Training our word2vec models

- SGNS-enwiki (English Wikipedia)
 - ~4.4M words in vocabulary
- SGNS-semeval (SemEval-2010 task 14 data)
 - ~122k words in vocabulary
- 300 word embedding dimensionality



Analogy analysis

- Semantic-Syntactic Word Relationship test set (SSWR)
 - 19.5k analogy tests (semantic and syntactic)
- Microsoft Research Syntactic Analogies Dataset (MSR)
 - 8k analogy tests
- Phrase Analogy Dataset (PAD)
 - 3k phrase analogy tests



Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks



Model	SSWR		
	Semantic	Syntactic	Average
SG 300	55	59	57
SG 1000	66.1	65.1	65.6
NEG-15	61	61	61
RNN-1600	—	—	—
GloVe 300 42B	81.9	69.3	75.0
fastText	77.8	74.9	76
SGNS-enwiki	65.8	67.3	66.6

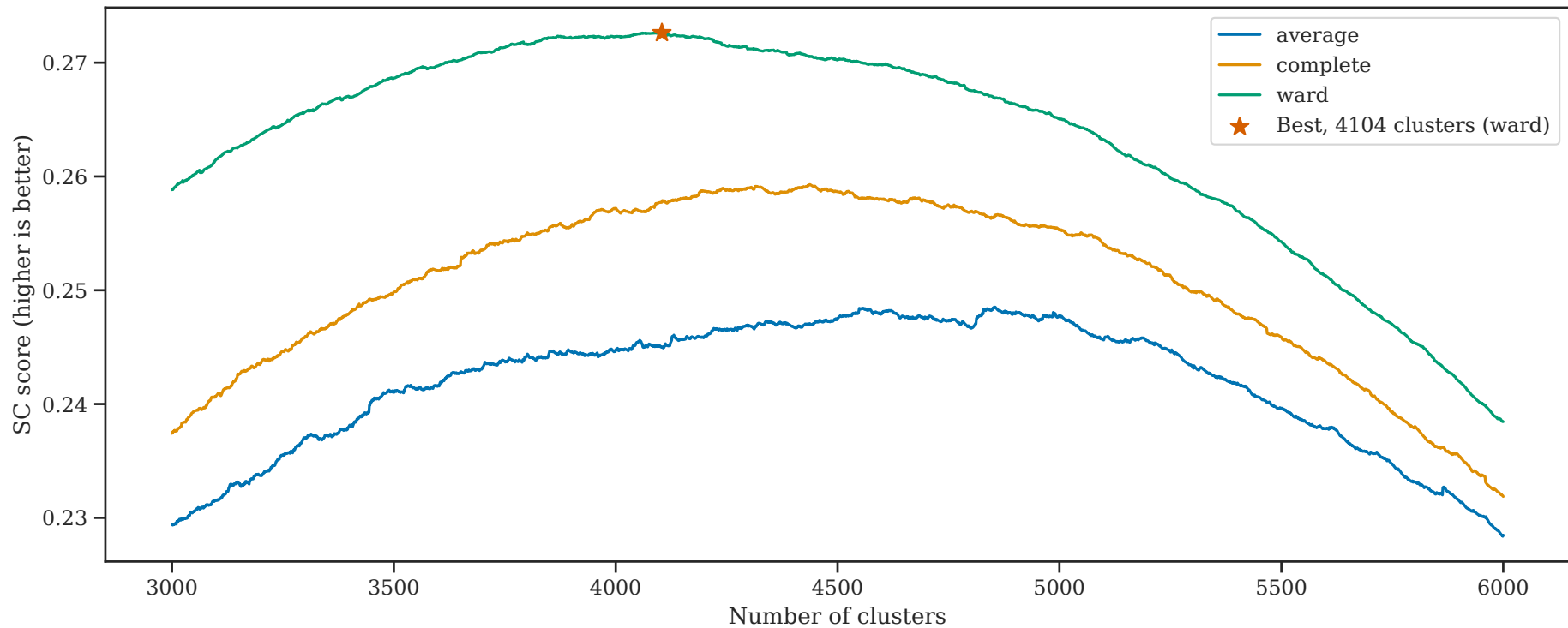




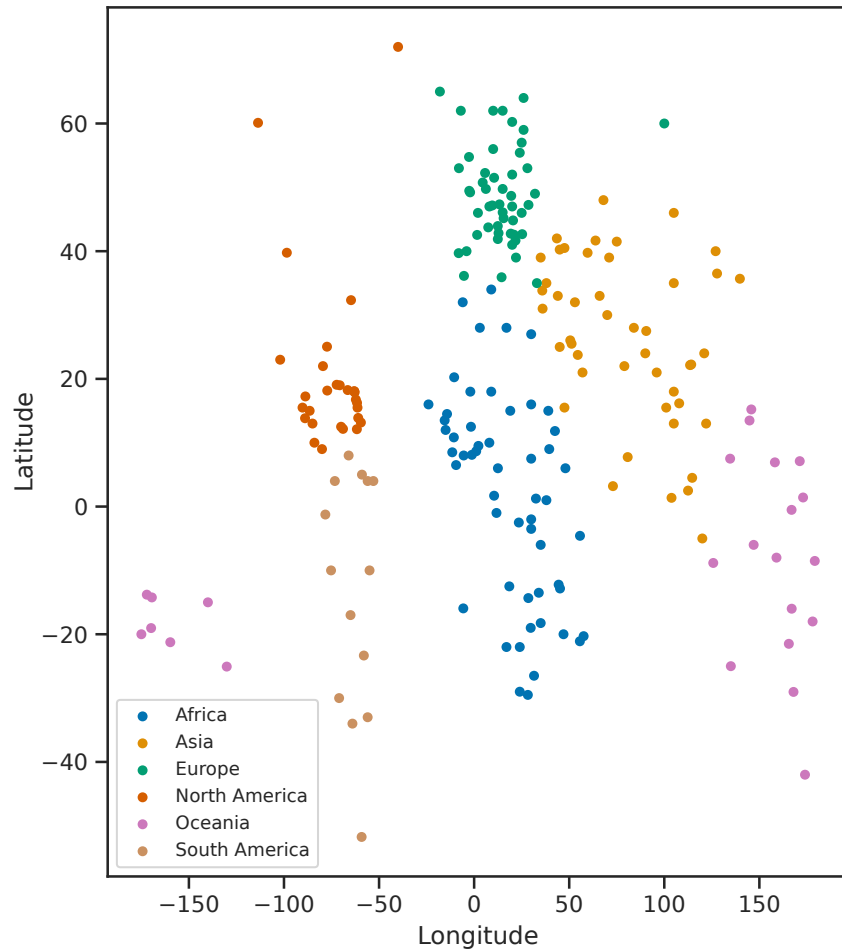
Cluster analysis

- Comparing cluster algorithms
- Internal cluster validation methods
- Analysis of
 - 4.4M \rightarrow 10k most common words
 - Distinct word groups
 - Countries/capitals
 - Number words

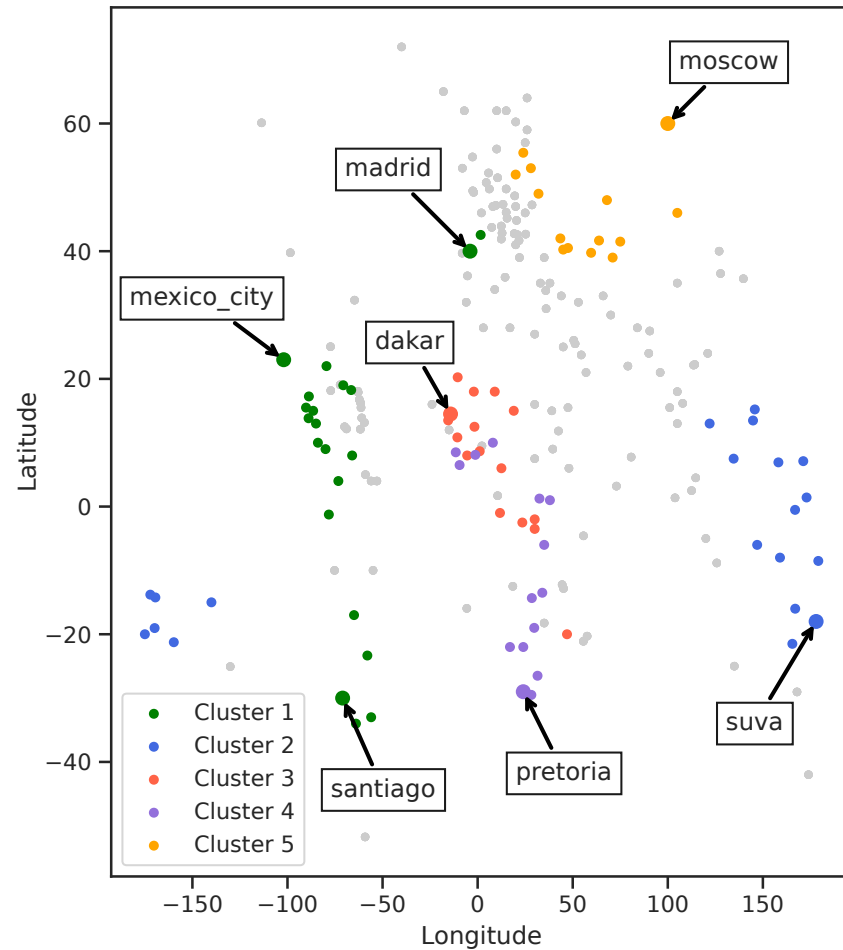




(a) Countries divided into six continents



(b) Top 5 clusters using capitals word group

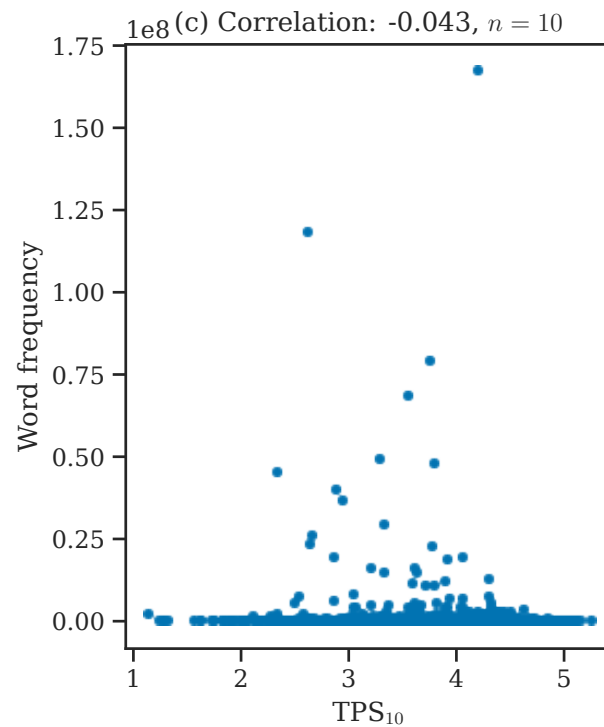
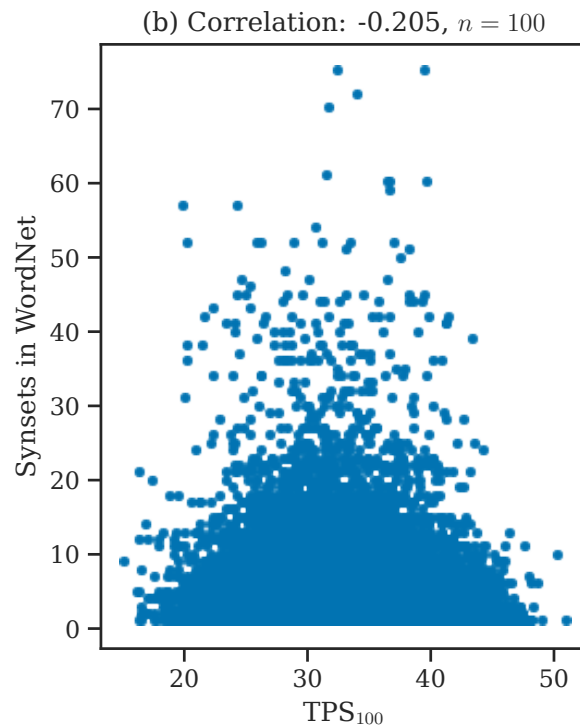
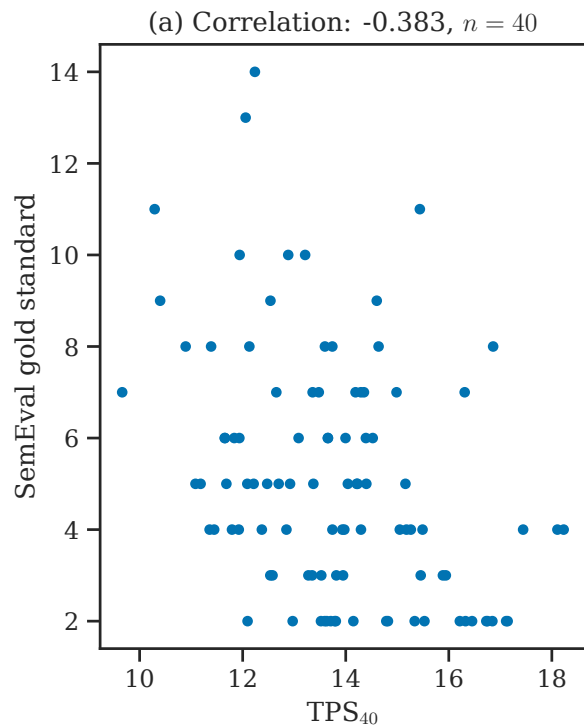


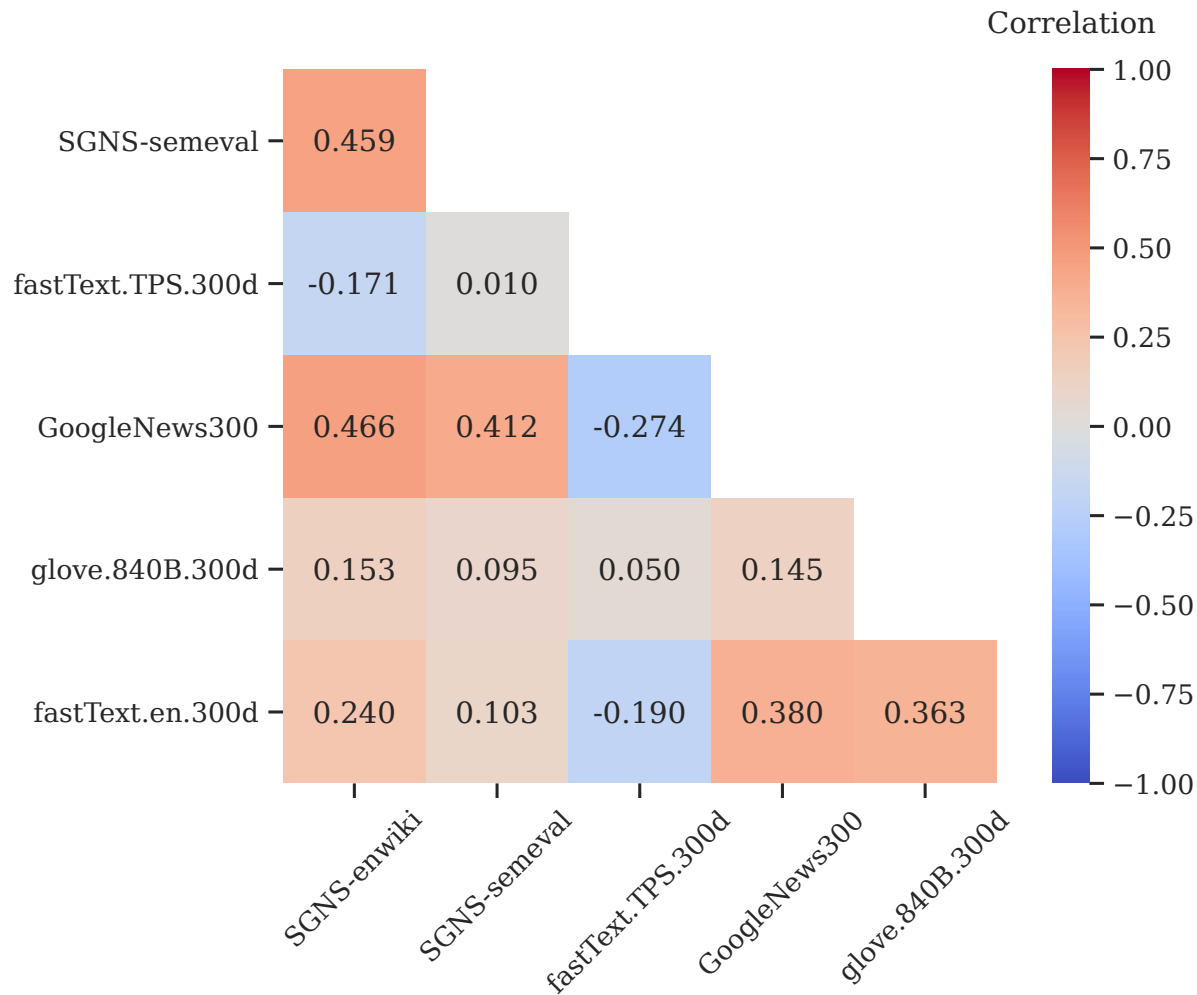


Topological polysemy

- Word embedding models
- SemEval-2010 task 14
 - 100 polysemous words
 - Gold standard (GS)
- Correlation with
 - Gold standard
 - Number of WordNet synsets
 - Word frequency







Geometric Anomaly Detection (GAD)

- WordNet SGNS-enwiki word embeddings
- Singular words reflected by singular GAD group



	GAD group			<i>Sum</i>
	Manifold	Boundary	Singular	
Number of monosemous words	4640	86731	4161	95532
Number of polysemous words	634	47902	344	48880
<i>Sum</i>	5274	134633	4505	144412

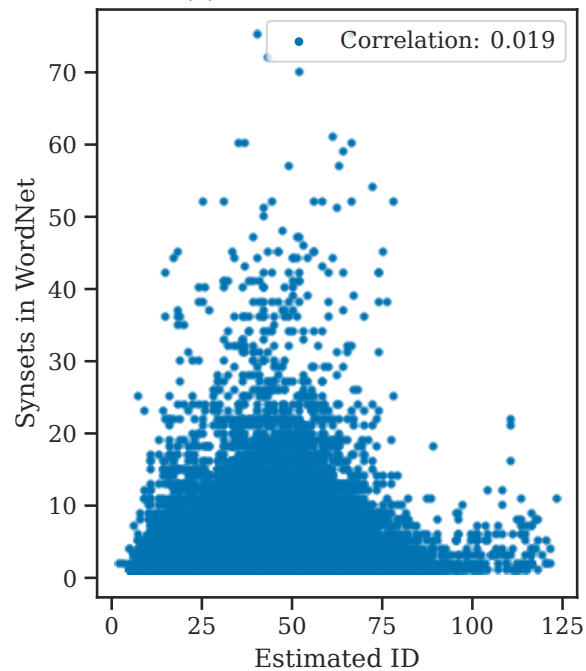


Intrinsic dimension (ID) estimation

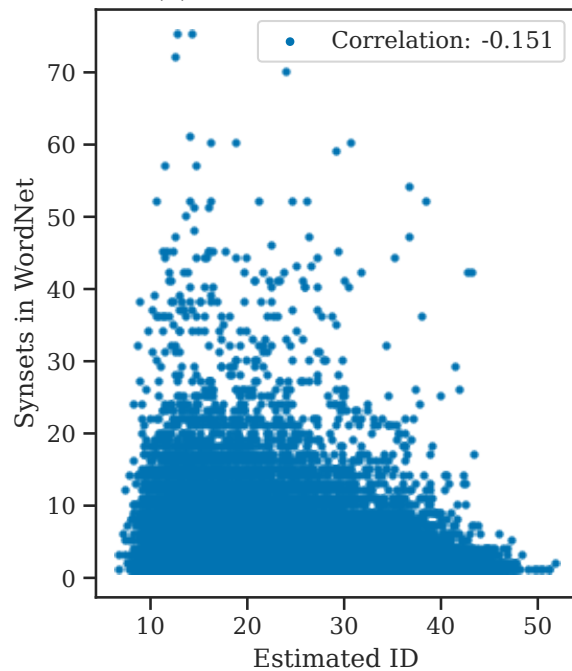
- Five ID estimation methods
- Relation to number of word meanings



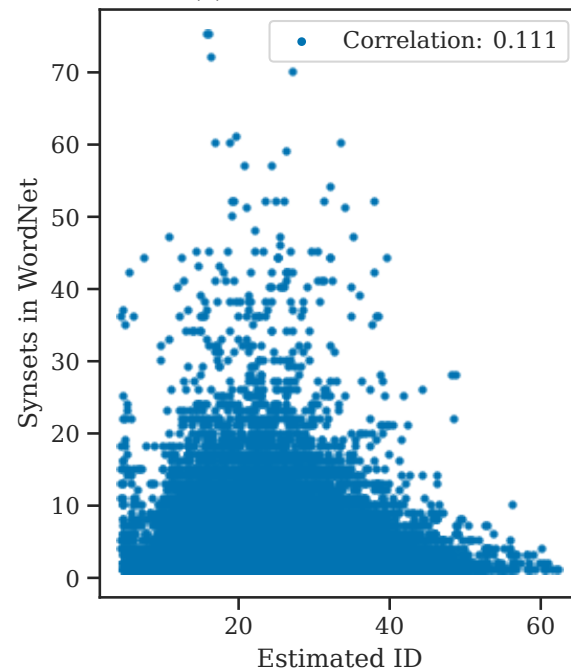
(a) Estimated ID w/IPCA



(b) Estimated ID w/TWO-NN



(c) Estimated ID w/TLE



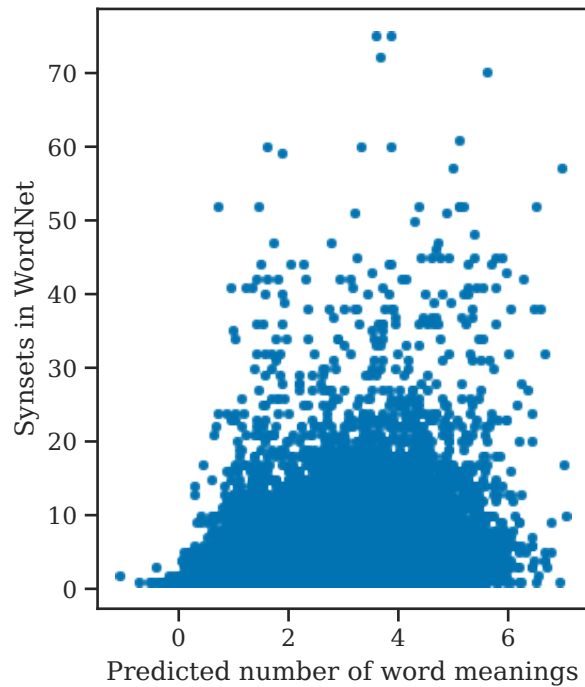


Supervised polysemy prediction

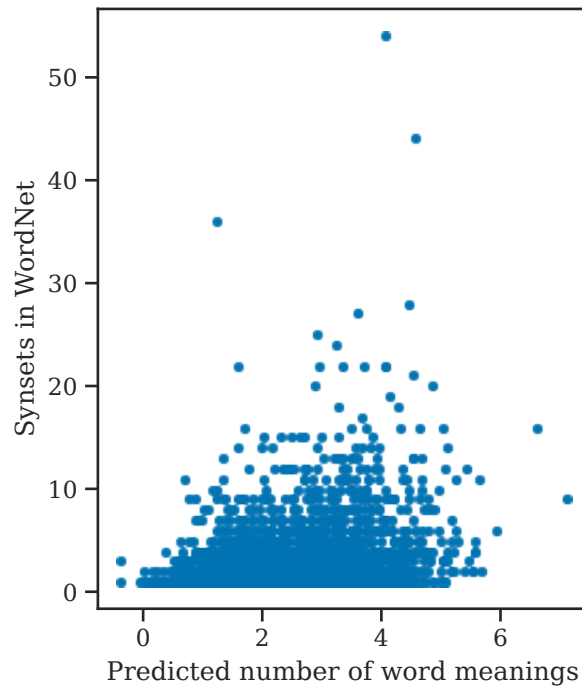
- WME- and BWME-enwiki models
 - Lasso and logistic regression
 - WordNet SGNS-enwiki word embeddings
- Features from TPS, GAD and ID estimation
- 95% train / 5% test split
- 20-fold cross validation on ℓ_1 -hyperparameter



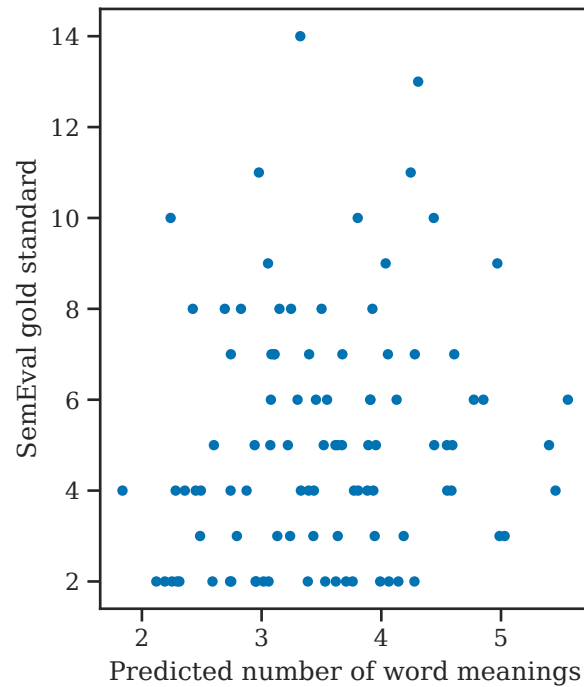
(a) Correlation: 0.364 (Train)

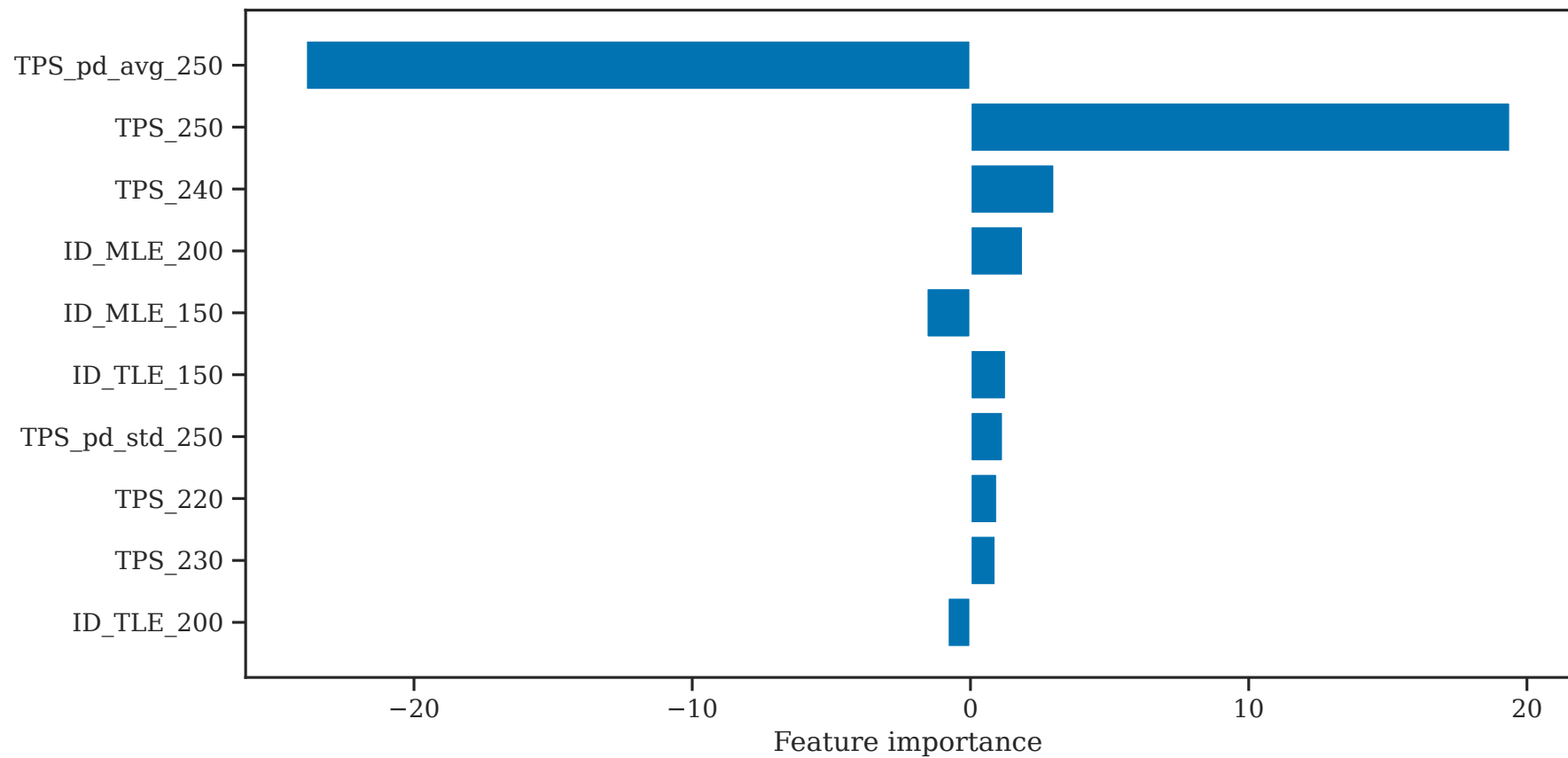


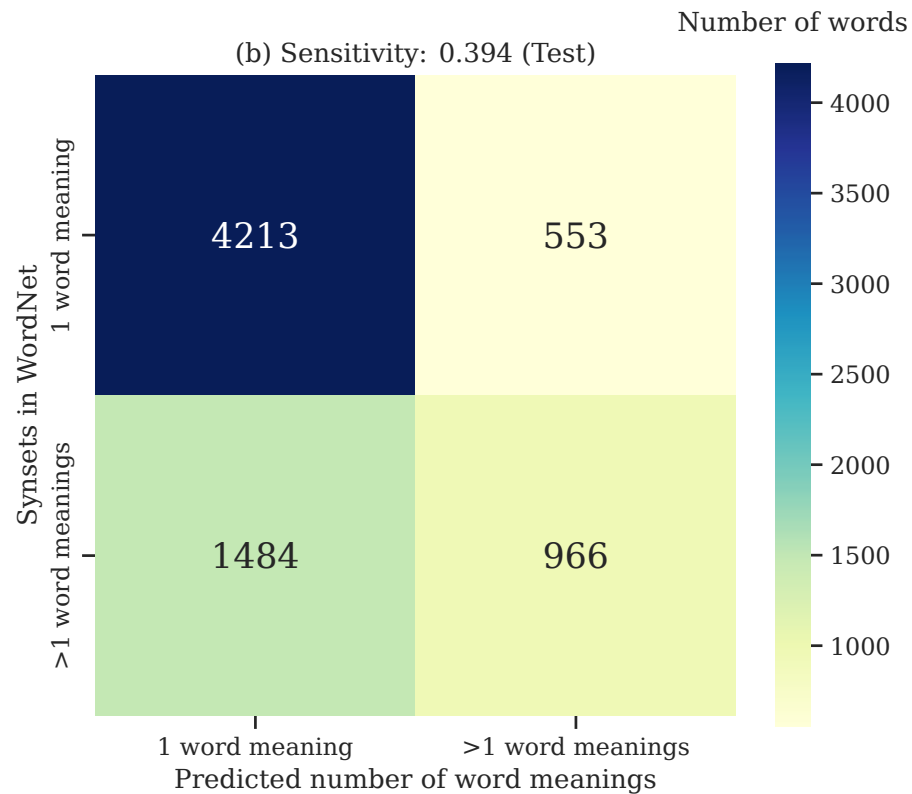
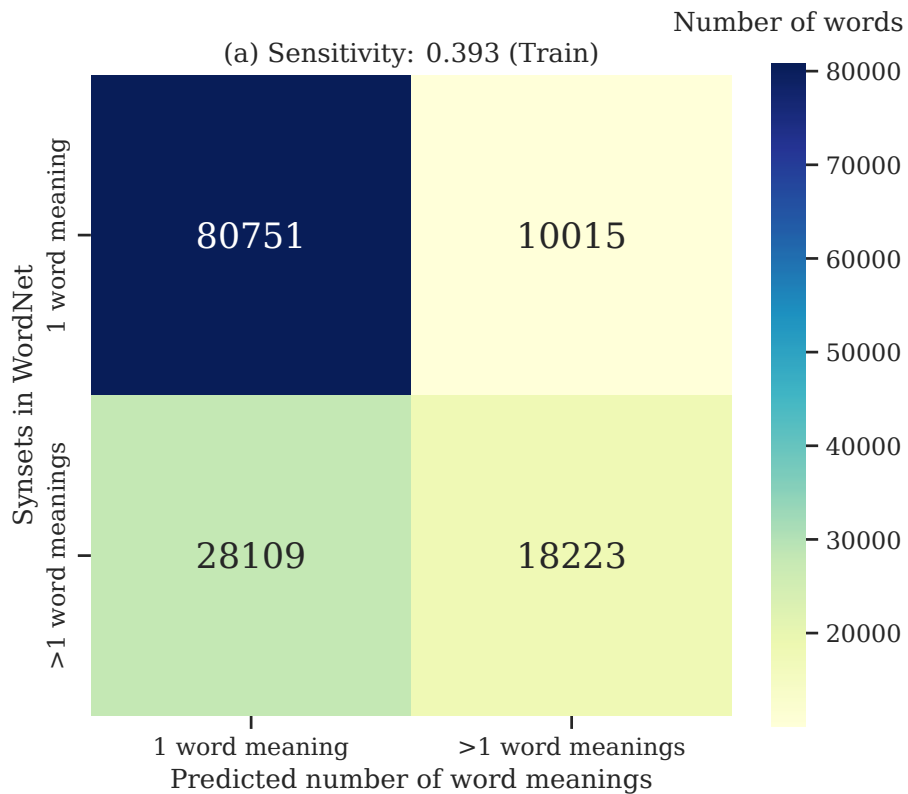
(b) Correlation: 0.370 (Test)



(c) Correlation: 0.154 (SemEval test)







Conclusion and Future Work





Conclusion

- Clustering of word embeddings deepen our understanding
- Topological polysemy yield inconsistent results

Future work

- Looking at other word embedding models
- Vectorization of GAD persistence diagrams
- GAD manifold dimension parametrization
- Word sense disambiguation problem

