Imperial College London

Novelty Detection in Scientific Research Papers

Yousef Nami

Supervisor: Dr Loïc Salles

Contents

03 Problem Definition

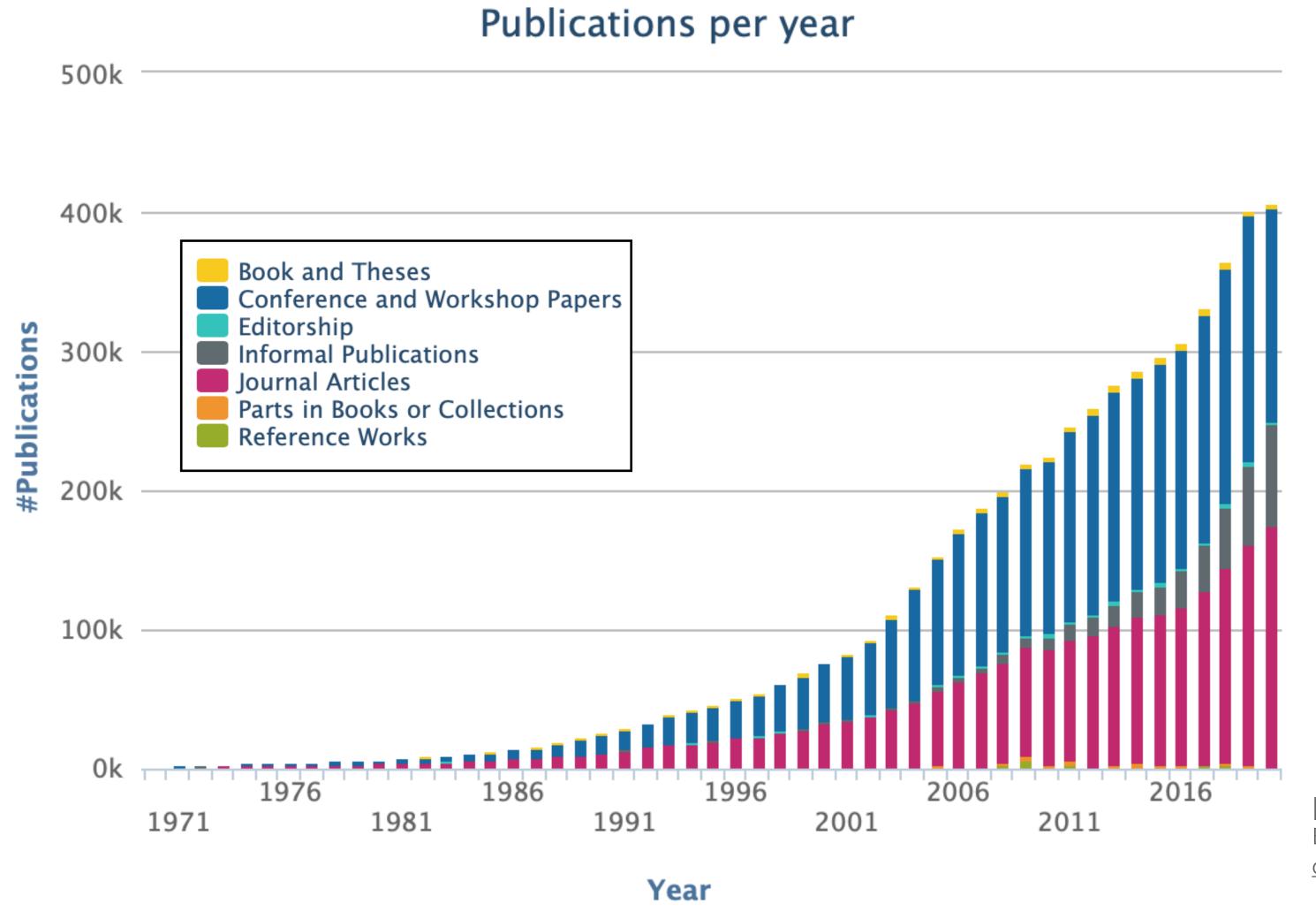
07 Data Processing

10 Methods

21 Outcomes

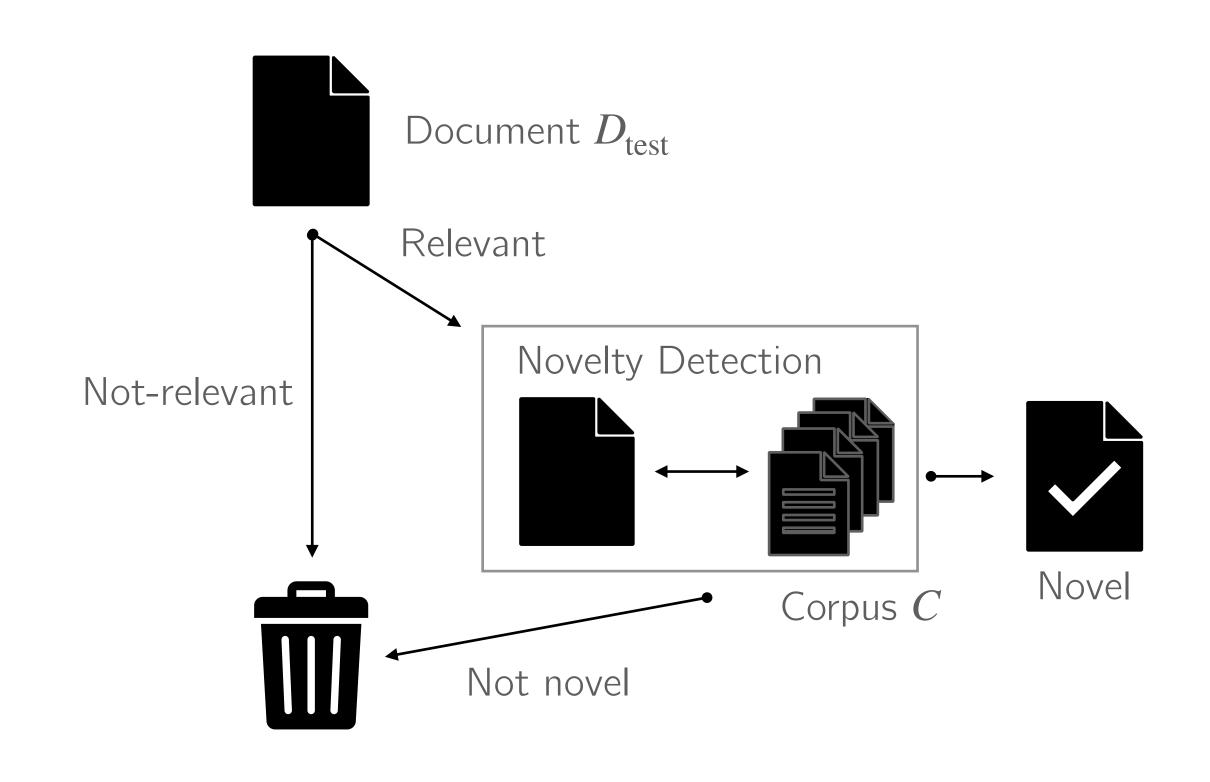
Concluding Remarks





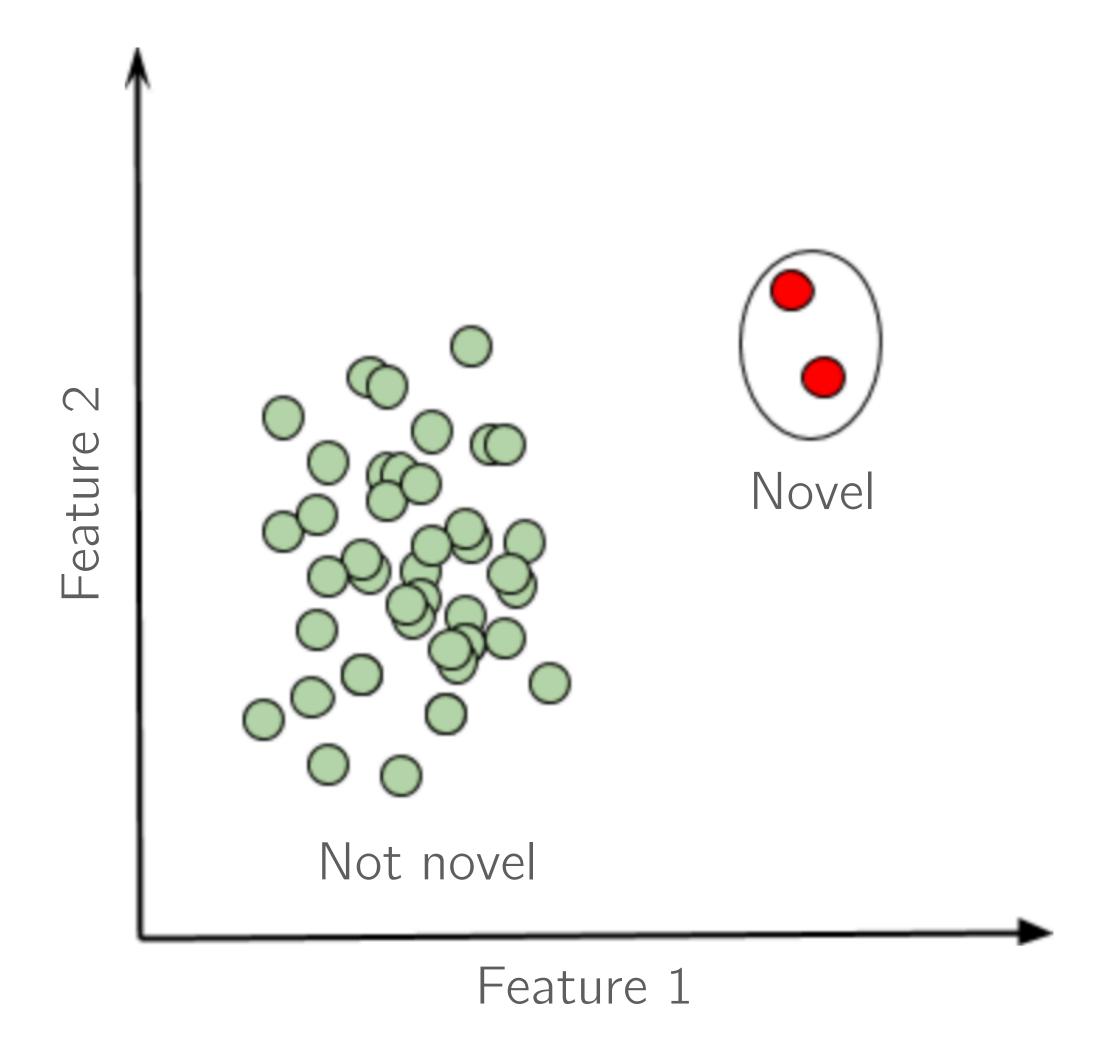
[1] *Publications per year.* Digital Bibliography and Library project. Link: dblp.org/statistics/publicationsperyear.html

- Need for tools that distinguish novel subject matter from redundant content
- Novelty: dissimilarity provided that a document is relevant
- Needs deep semantic representation of documents



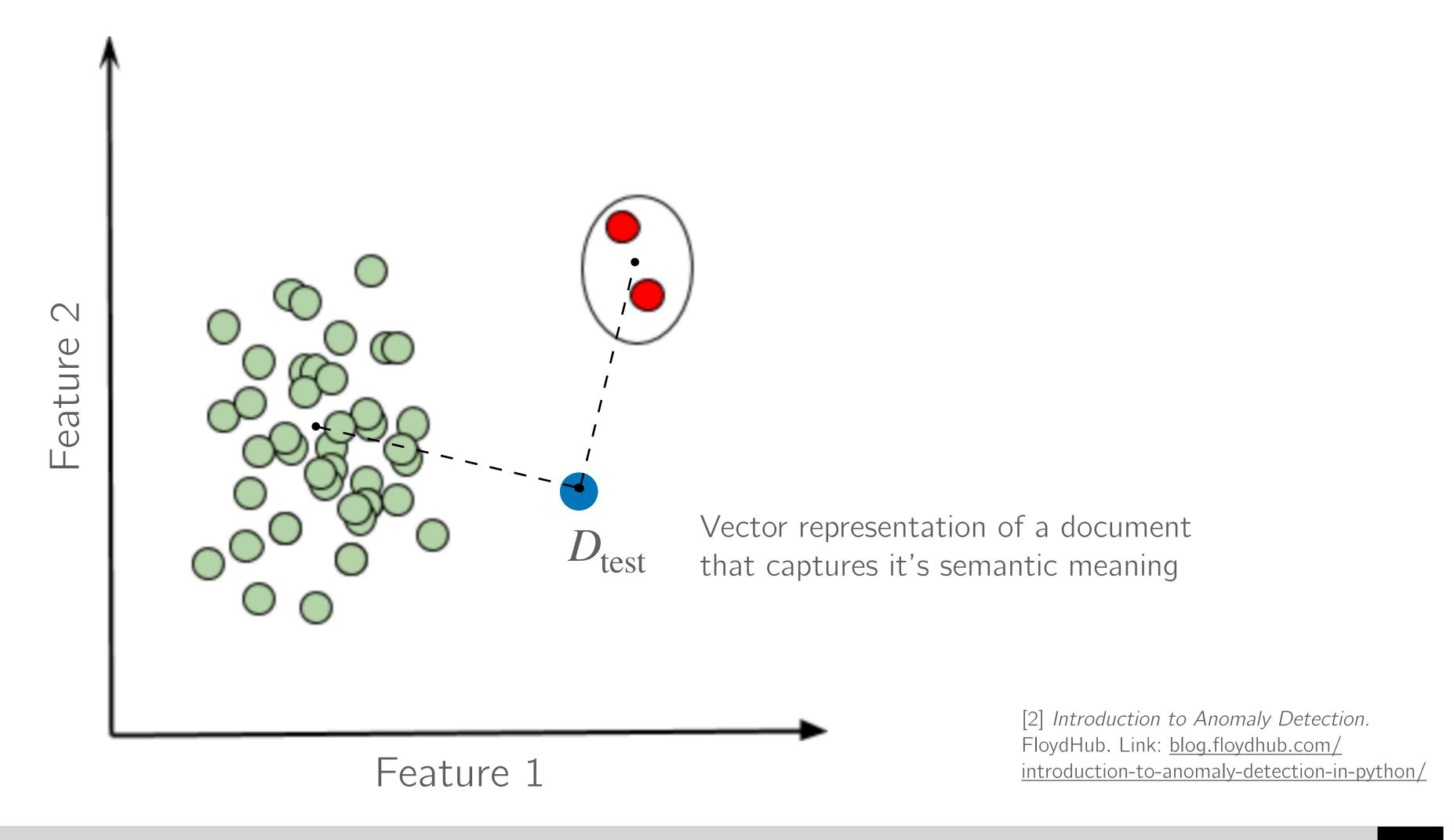
Novelty

Dissimilarity provided that a document is relevant



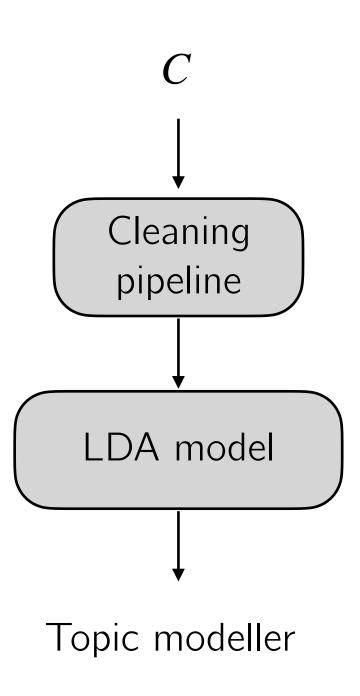
Novelty

Dissimilarity provided that a document is relevant

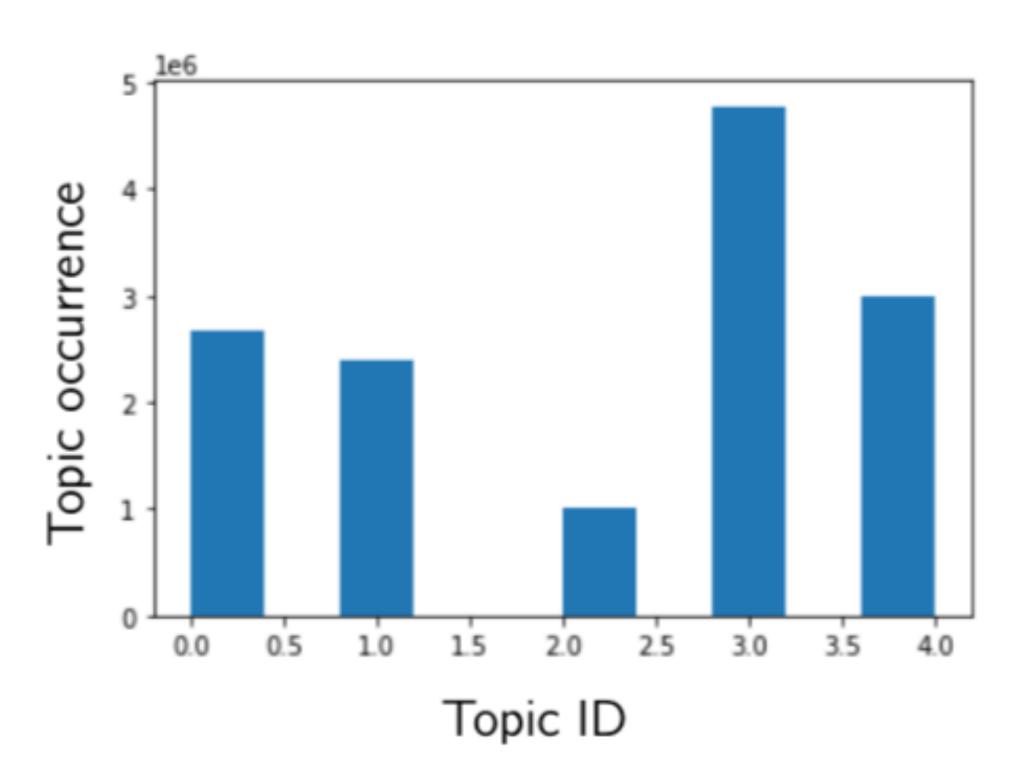


Exploratory Data Analysis

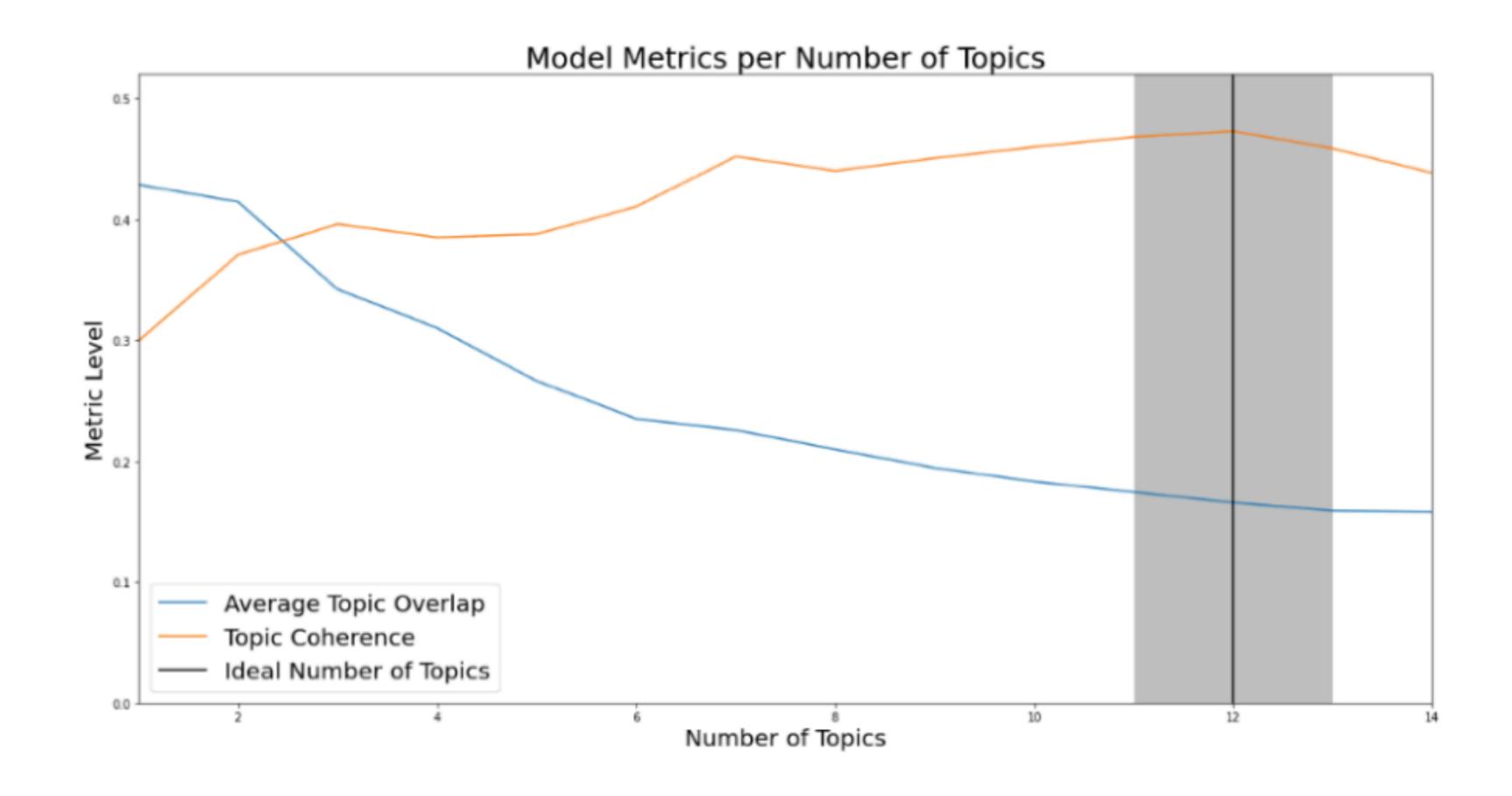
- Statistical model known as Latent Dirichlet Allocation (LDA)
- Learns the topic distribution of documents
- ullet Effectively finds the probability that a word w_i belongs to any of m topics selected by user



Exploratory Data Analysis



Exploratory Data Analysis



Jaccard score

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Coherence

Measures the semantic similarity

Bag-of-words and TF-IDF

- Consider $C = \{ "My \ name \ is \ Yousef", "NLP \ is \ awesome" \}$
- ullet Represented by matrix $V_{ij} \in \mathbb{R}^{N \times L}$ where N is # of documents, L vocab size
- Does not capture deep semantic information!

My name is Yousef
NLP is awesome

My	name	is	Yousef	NLP	awe some
1	1	1	1	0	0
0	0	1	0	1	1

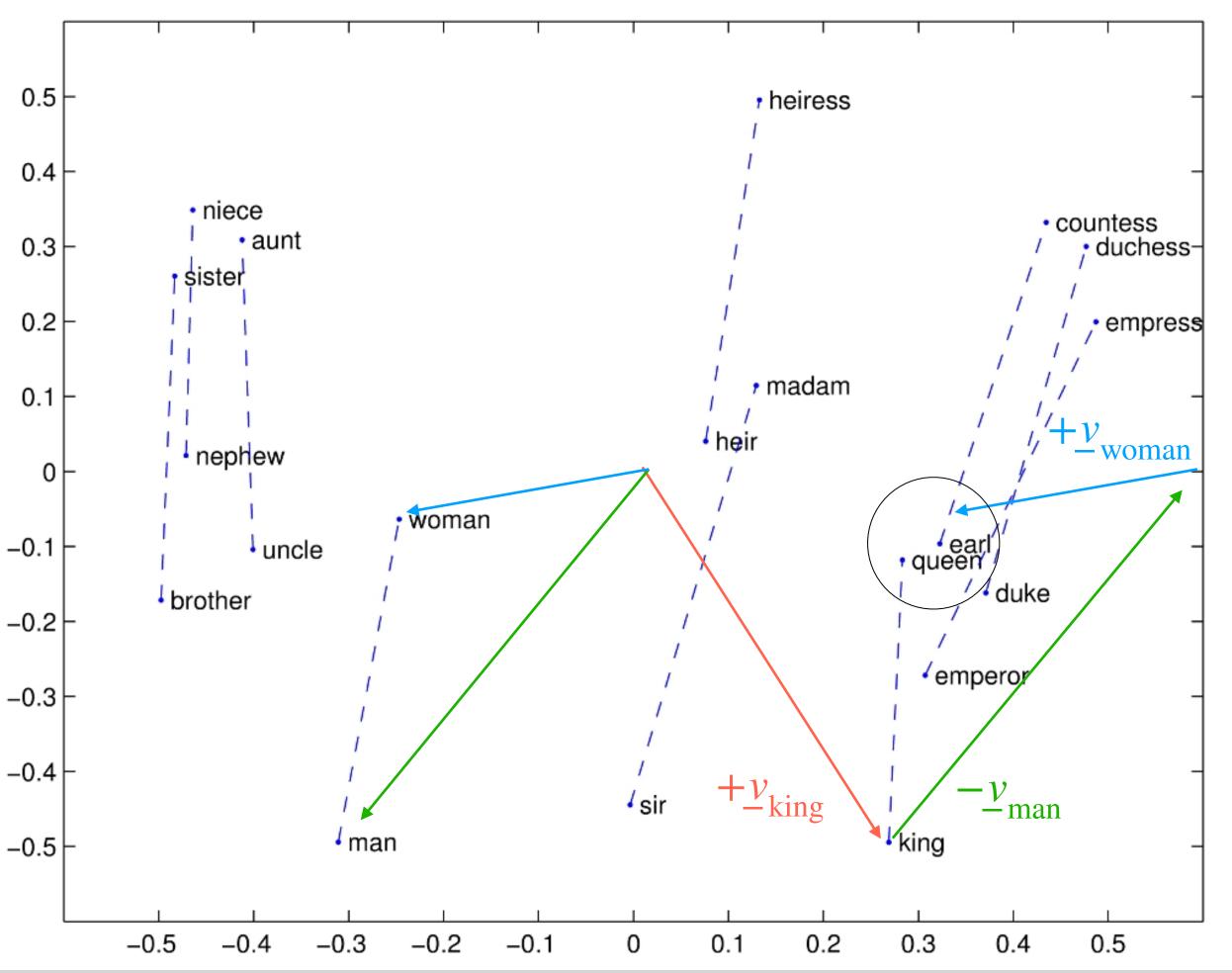
My name is Yousef
NLP is awesome

My	name	is	Yousef	NLP	awe some
0.5	0.5	0.25	0.5	0	0
0	0	0.33	0	0.67	0.67

GloVe

- Represent each word w_i as a vector $\underline{v} \in \mathbb{R}^l$, where l is an embedding length
- Typically, $l \in [50, 300]$
- Word vectors carry semantic meaning
- They are linear and additive, but also contain spatial information

GloVe



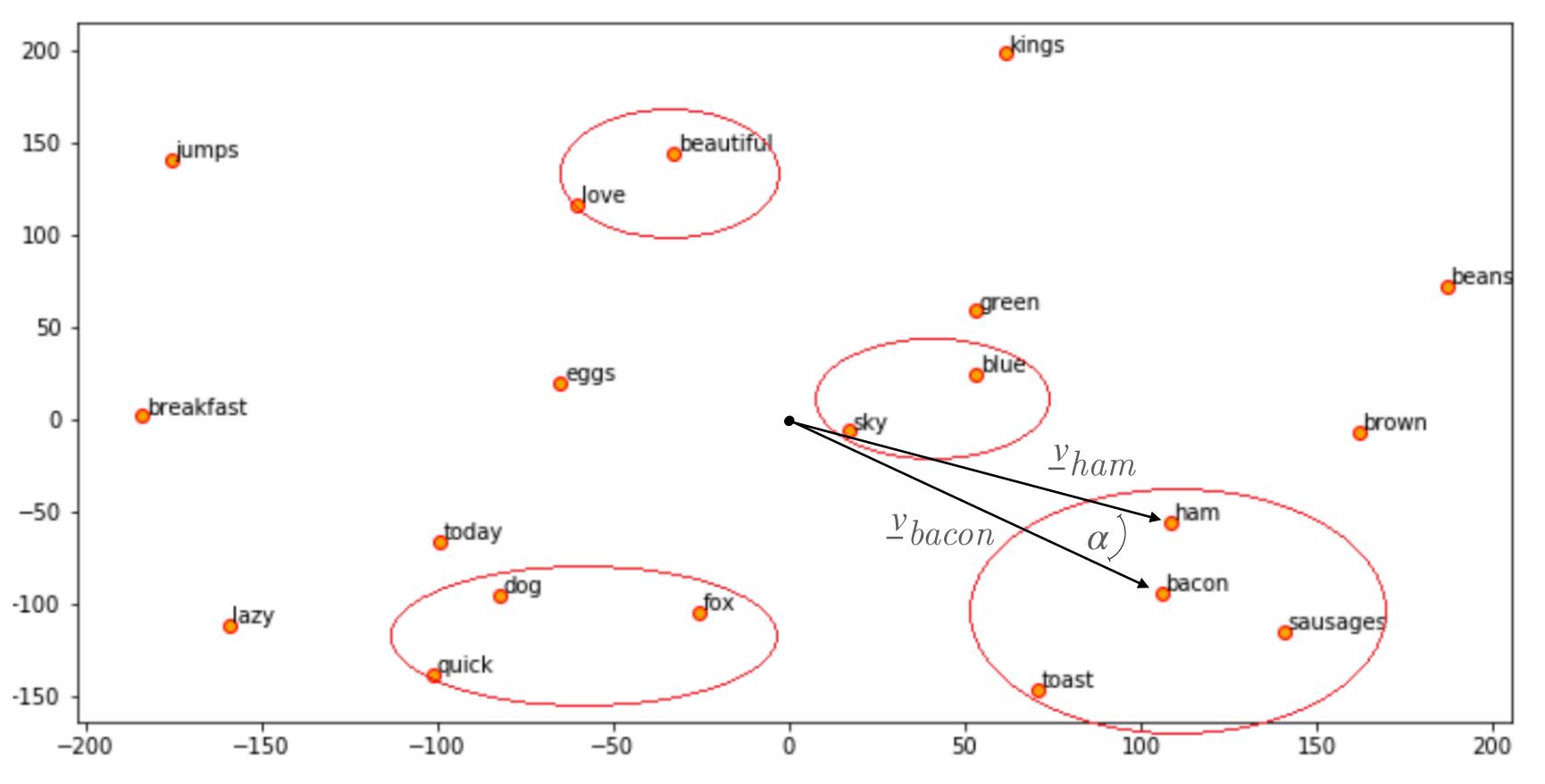
Additive meaning

Vectors can add to each other, mimicking real relationships

$$v_{king} - v_{man} + v_{woman} \approx v_{queen}$$

[3] GloVe: Global Vectors for Word Representation. Pennington et al. Link: nlp.stanford.edu/projects/glove/

GloVe



Spatial closeness

clusters indicate semantic family $C = \{ "quick", "dog", "fox" \} \approx "animal" \}$

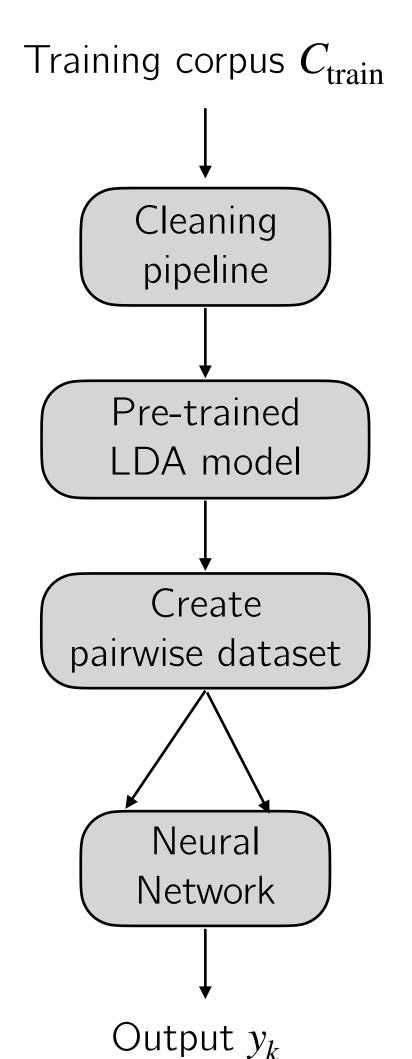
$$\cos \alpha = \underline{v}_{ham} \cdot \underline{v}_{bacon} \approx 1$$

[4] Hands on approach to Deep Learning methods for text data. Dipanjan Sarkar. Link: towardsdatascience.com/understanding-feature-engineering-part-4-deep-learning-methods-for-text-data-96c44370bbfa

Pseudo-Supervised Model

Training process

- Assigns a topic to each sentence
- Pairs all sentences together
- Trains a neural network to predict the probability that a pair belongs in the same class
- Model based on pairwise network by Qin et al. [6]



[6] Text Classification with novelty detection. Qin Q, Hu W, Liu B. Link: https://arxiv.org/abs/2009.11119

Pseudo-Supervised Model

Pairwise data

	Topic ID
My name is Yousef	0
NLP is awesome	1
Deep Learning	1

Topic ID	Meaning
0	People
1	Artificial Intelligence

Sen	tence pair Sent 1	Sent 2	Topic ID 1	Topic ID 2	Same topic?
	My name is Yousef	NLP is awesome	0	1	0
	My name is Yousef	Deep Learning	0	1	0
	NLP is awesome	Deep Learning	1	1	1

Pseudo-Supervised Model

Testing process

Topic 0: People

Topic 1: Artificial Intelligence

Test sentence	Sentences from corpus	Probability	Mean by topic
Machine learning is very broad	My name is Yousef	0.1	0.1
Machine learning is very broad	NLP is awesome	0.7	0.75
Machine learning is very broad	Deep Learning	0.8	0.75

Not-novel

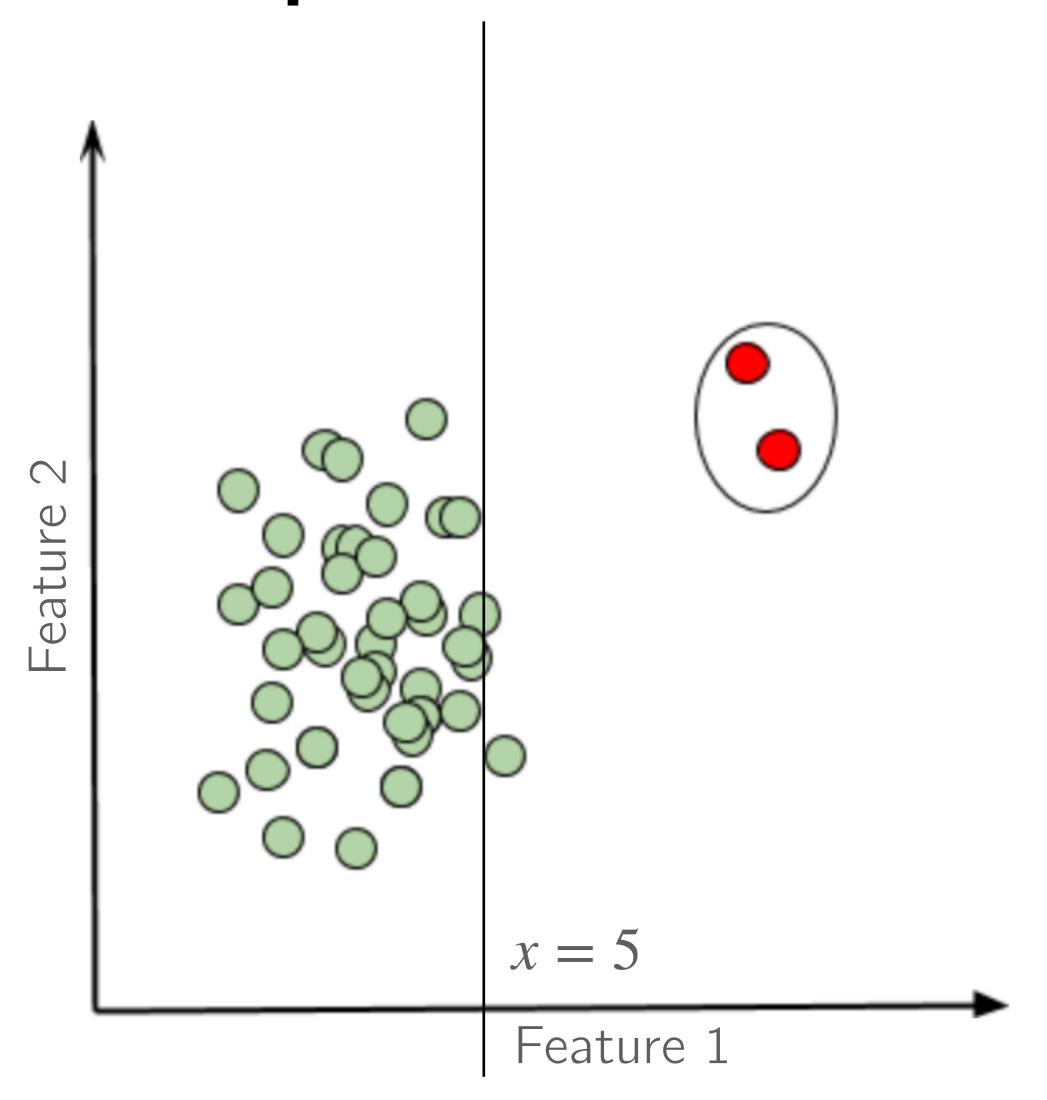
Test sentence	Sentences from corpus	Probability	Mean by topic
A turbine is a Turbomachine	My name is Yousef	0.1	0.1
A turbine is a Turbomachine	NLP is awesome	0.3	0.25
A turbine is a Turbomachine	Deep Learning	0.2	0.23

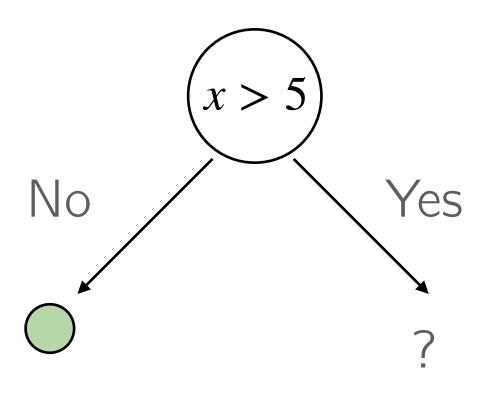
Novel

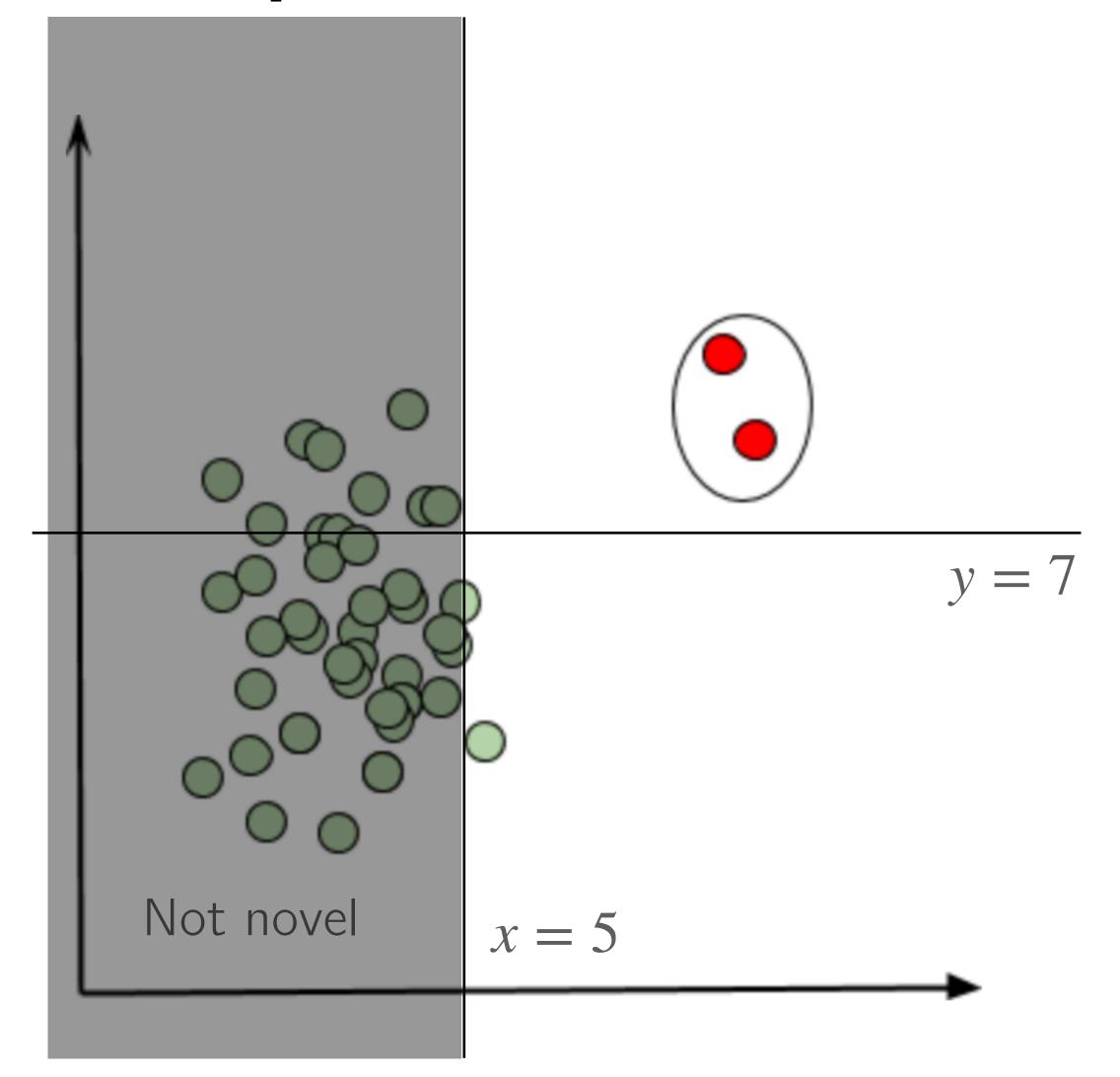
Document representation found using mean GloVe word vectors

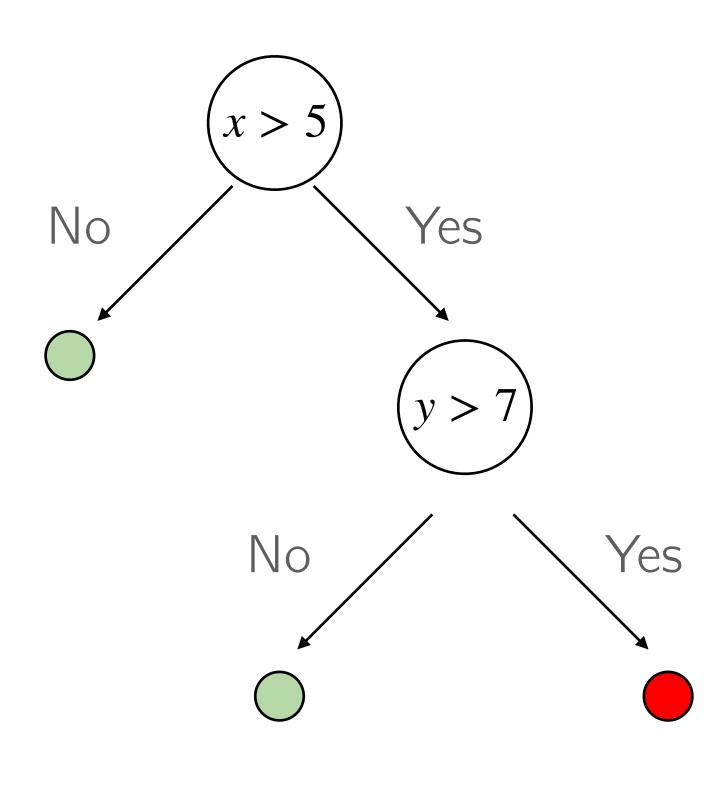
$$\bar{v} = \frac{\sum_{i}^{N} v_{i}}{N}$$
 N is the $\#$ of words per document v_{i} is the word vector for word i

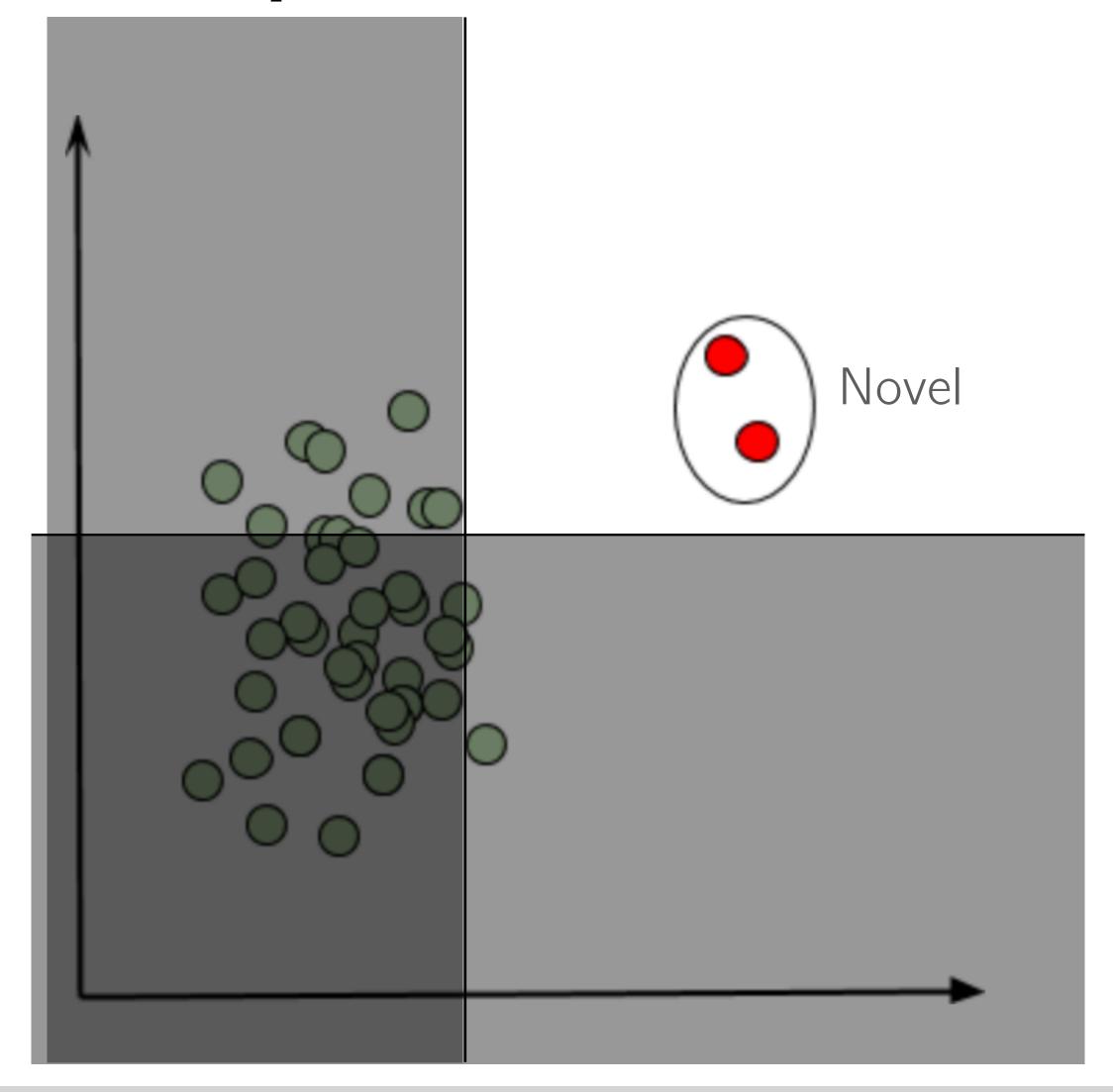
- IsolationForest algorithm was used to determine anomalies
- Results compared with TF-IDF embedded documents

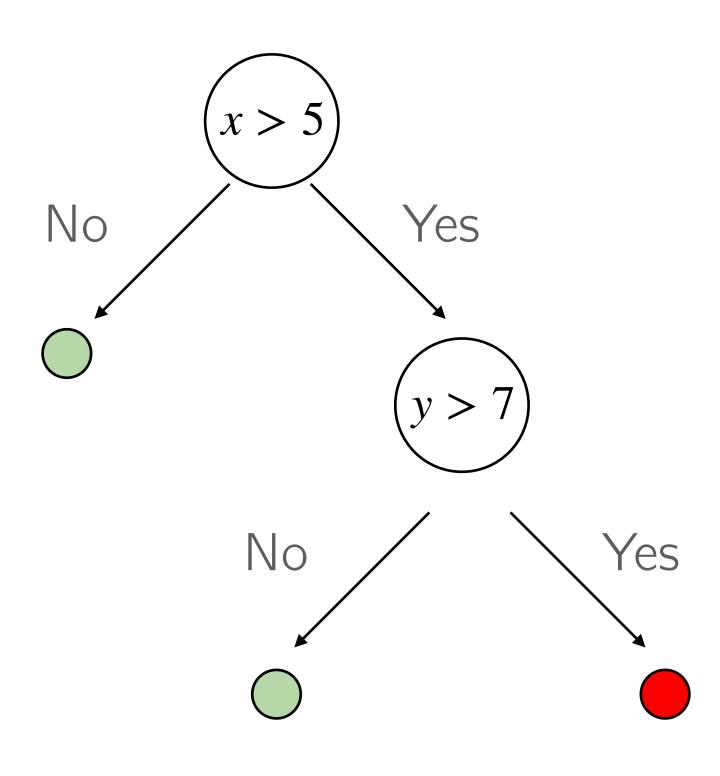




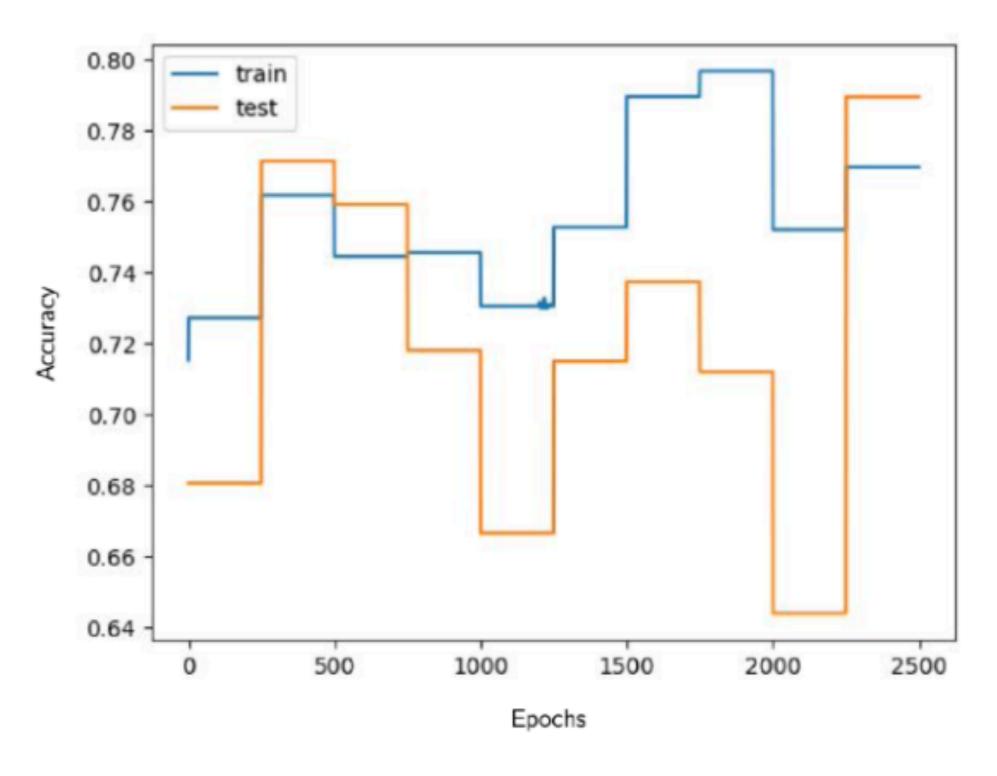




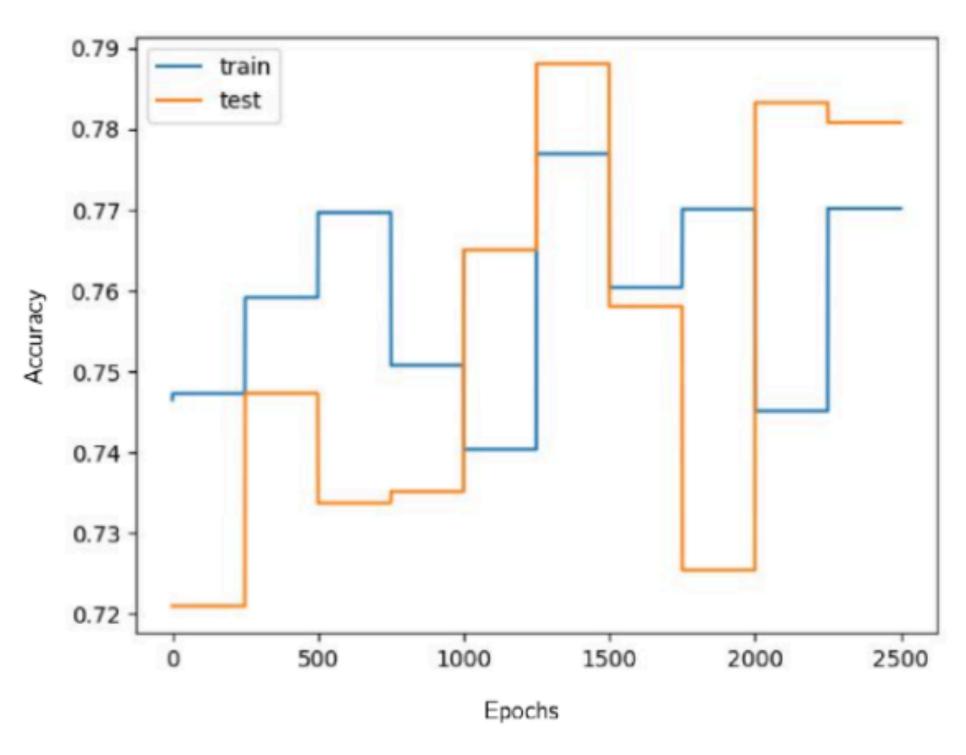




Pseudo-Supervised Model



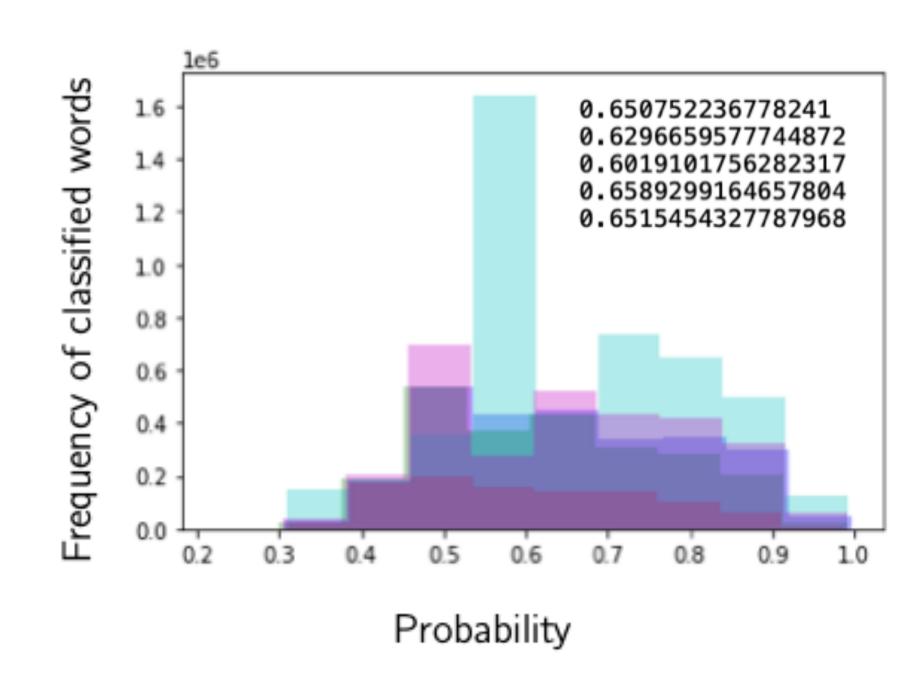
n = 200 sentences



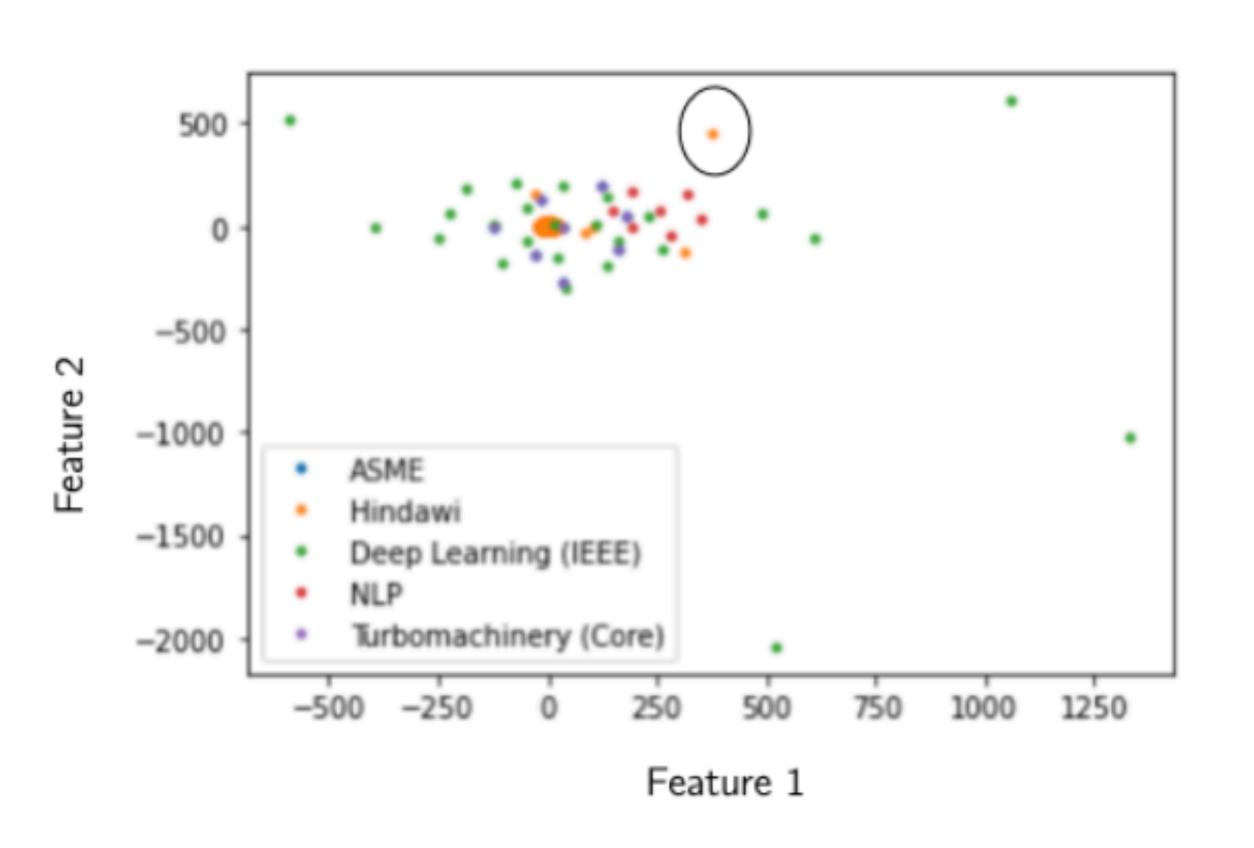
n = 1000 sentences

Pseudo-Supervised Model

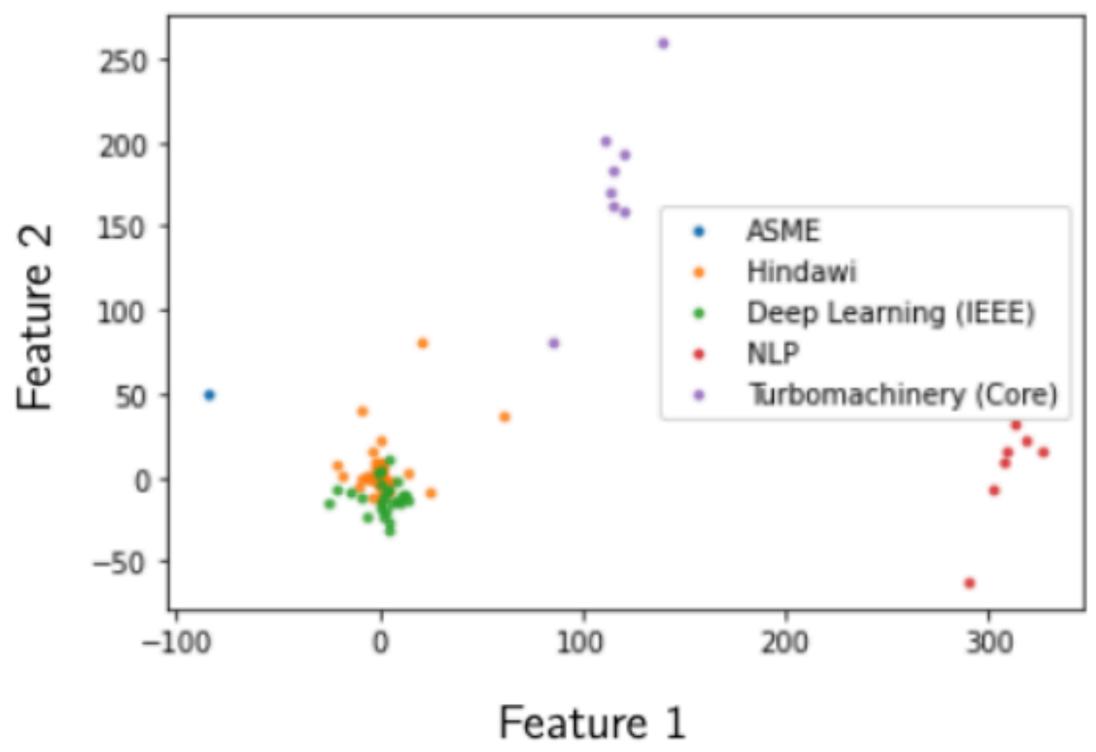
- Pseudo-supervised nature of the model not rigorously analysed
- Assumed that the LDA model can generate topics both at corpus (macro) and sentence (micro) level
- Degree of prediction confidence not considered
- Should set a threshold (e.g. 0.8)
- Would decrease computational cost



Unsupervised Model

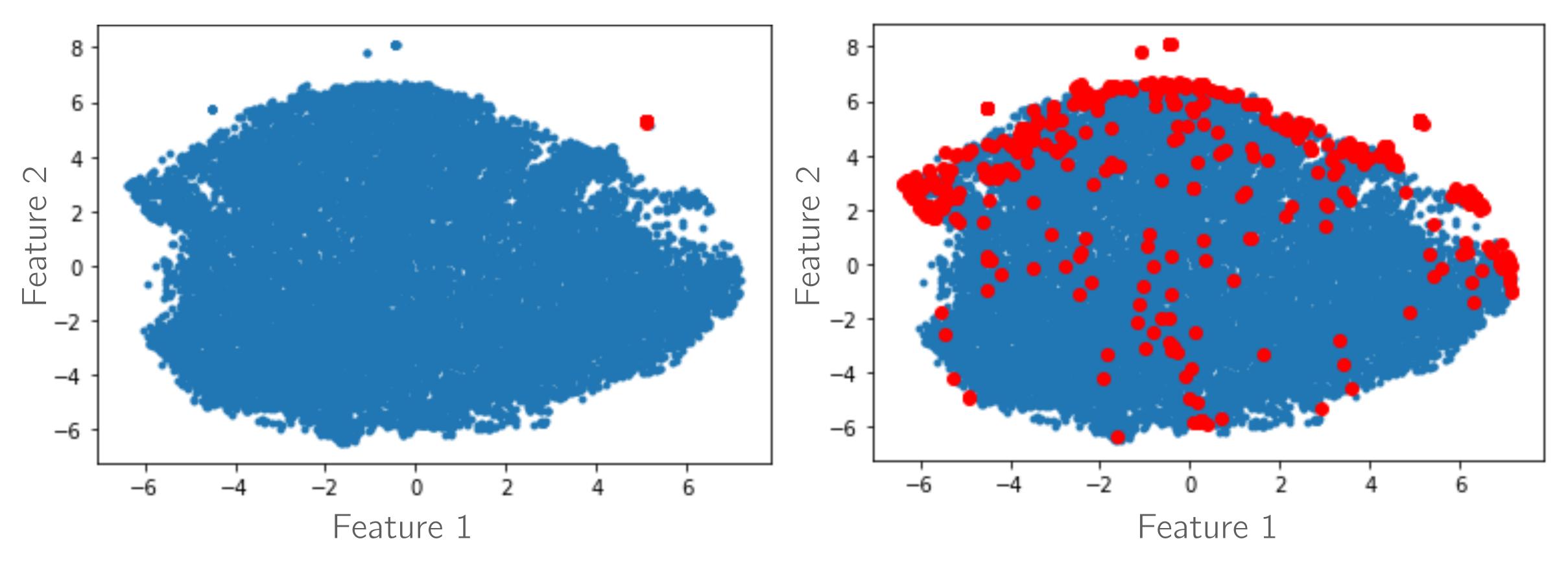


l=L length document vectors projected in 2D plane where L is the unique number of words in corpus



l=300 length document vectors projected in 2D plane

Unsupervised Model



l=300 length document vectors projected in 2D plane

Data Quality

ASME corpus data is not great

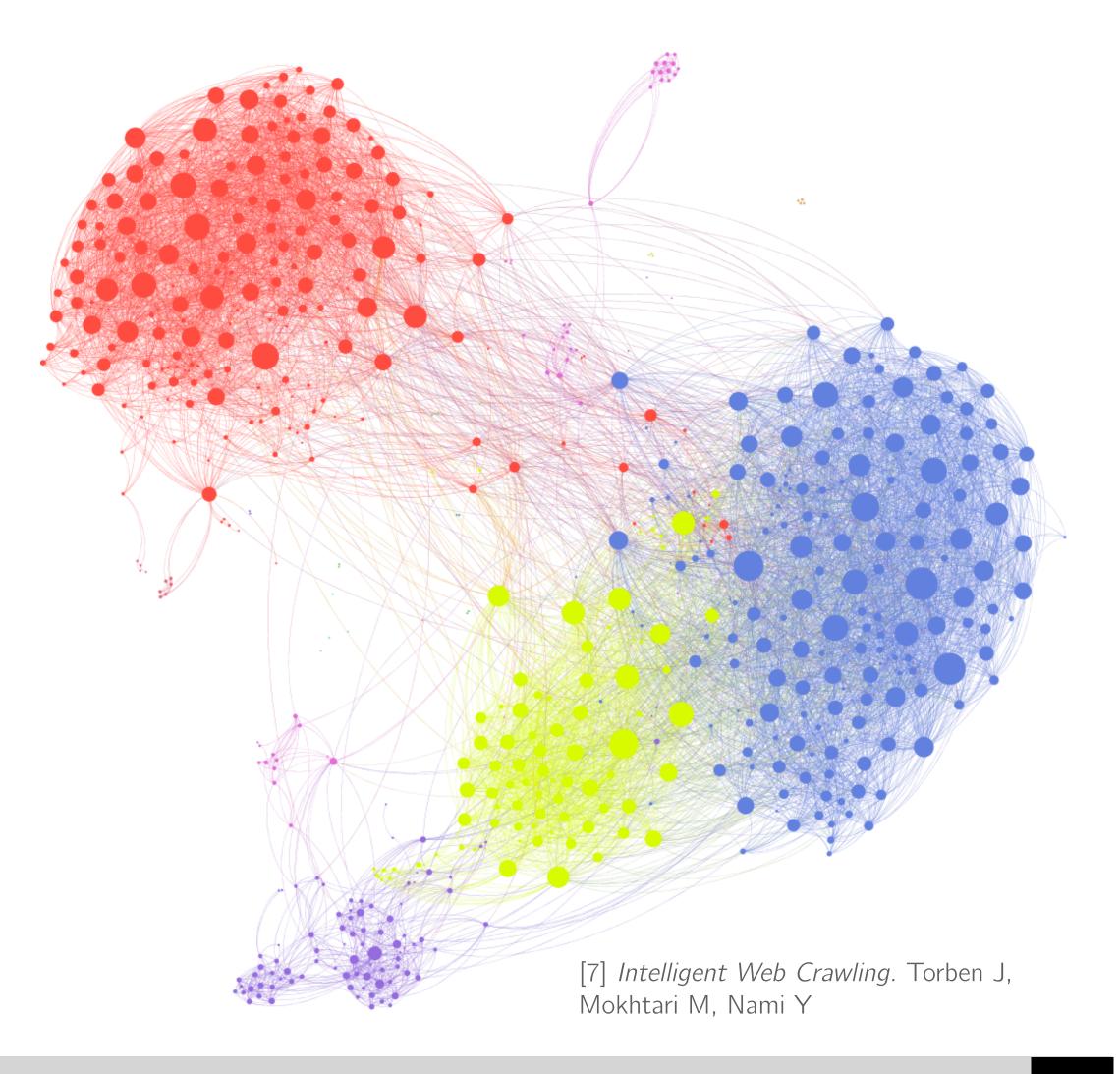
Type	Actual word	Unrecognized word
Spelling mistakes	capturing	cap <u>a</u> turing
Incorrectly read characters	copyright	c <u>a</u> pyright
Foreign words		brennkammersystem
Maths symbols	b sin y	bsiny
Phrases captured without whitespace	burner loudspeaker setup	burnerloudspeakersetup
Named entities that weren't removed		mohammadpour
Technical terms		eulerian

Conclusion

- Two models: 1 pseudo-supervised, 1 unsupervised
- Pseudo-supervised model: highly computationally expensive, best left for prelabelled and small sized documents
- Unsupervised model: Performance predictable, not consistent at detecting novelty
- Document representations not informative enough due to low quality data and choice of embeddings (GloVe, TF-IDF)
- TF-IDF is useful for detecting relevance, but not novelty

Future Work

- Redefining novelty: considering graphical representations of documents for extracting information
- A novel research trend could be a branch that leads to another branch from a different cluster
- Can also examine novelty in graphical structures



Future Work

- Data enhancement: use of TF-IDF as filtration step for noise removal, spell checking
- Improved document representations: use of state-of-the-art document representation techniques, including document encoder-decoder models
- Computational considerations: use of tensorflow's tf.Dataset for better GPU acceleration
- References: using citations as extra features for training

Imperial College London

Thank you for your attention

Any Questions

Reference List

- [1] Digital Bibliography and Library project. *Publications per year.* Available from: dblp.org/statistics/publicationsperyear.html [Accessed 5th June 2021]
- [2] FloydHub. Introduction to Anomaly Detection. Available from: blog.floydhub.com/introduction-to-anomaly-detection-in-python/ [Accessed 5th June 2021]
- [3] Pennington J, Socher R, Manning CD. Global Vectors for Word Representation. Available from: nlp.stanford.edu/projects/ glove/ [Accessed 5th June 2021]
- [4] Sarkar D. Hands on Approach to Deep learning methods for text data. Available from: towardsdatascience.com/understanding-feature-engineering-part-4-deep-learning-methods-for-text-data-96c44370bbfa [Accessed 5th June 2021]
- [5] Hui J. Word Embedding and GloVe. Available from: https://jonathan-hui.medium.com/nlp-word-embedding-glove-5e7f523999f6 [Accessed 5th June 2021]
- [6] Qin Q, Hu W, Liu B. Text Classification with Novelty Detection. Arxiv [Preprint] 2020. Available from: https://arxiv.org/abs/2009.11119 [Accessed 5th June 2021]
- [7] Torben J, Mokhtari R, Nami Y. Intelligent Web Crawling [Accessed 6th June 2021]