

---

# Aiming beyond the Obvious: Identifying Non-Obvious Cases in Semantic Similarity Datasets

Nicole Peinelt

ACL, Italy, 30 July 2019



# Authors

---



**Nicole Peinelt**

---

The Alan Turing Institute, UK  
University of Warwick, UK



**Dr. Maria Liakata**

---



**Dr. Dong Nguyen**

---

The Alan Turing Institute, UK  
Utrecht University, The Neverlands

---

# Semantic Similarity Detection for English

- Given a sentence pair, predict binary relatedness

Dataset	Task
MSRP	Paraphrase detection
Quora	Question paraphrase detection
STS	Similarity scoring
SemEval 2017 Task C	(A) CQA answer ranking (B) CQA paraphrase ranking (C) CQA answer ranking

# Examples from Quora

Positive label (paraphrase)	Negative label (not a paraphrase)
<p>Which is the best way to learn coding? ---<p>How do you learn to program?</p></p>	<p>What is meant by 'e' in mathematics? ---<p>What is meant by mathematics?</p></p>

# Examples from Quora

<b>Positive label (paraphrase)</b>	<b>Negative label (not a paraphrase)</b>
Which is the best way to learn coding? ---	What is meant by 'e' in mathematics? ---
How do you learn to program?	What is meant by mathematics?
What's the origin of the word o'clock? ---	What are the range of careers in Biotechnology in Indonesia? ---
What is the origin of the word o'clock?	How do you tenderize beef stew meat?

# Examples from Quora

Hard  
↓  
Easy

Positive label (paraphrase)	Negative label (not a paraphrase)
Which is the best way to learn coding? --- How do you learn to program?	What is meant by 'e' in mathematics? --- What is meant by mathematics?
What's the origin of the word o'clock? --- What is the origin of the word o'clock?	What are the range of careers in Biotechnology in Indonesia? --- How do you tenderize beef stew meat?

But equal weight for evaluation

---

## Background

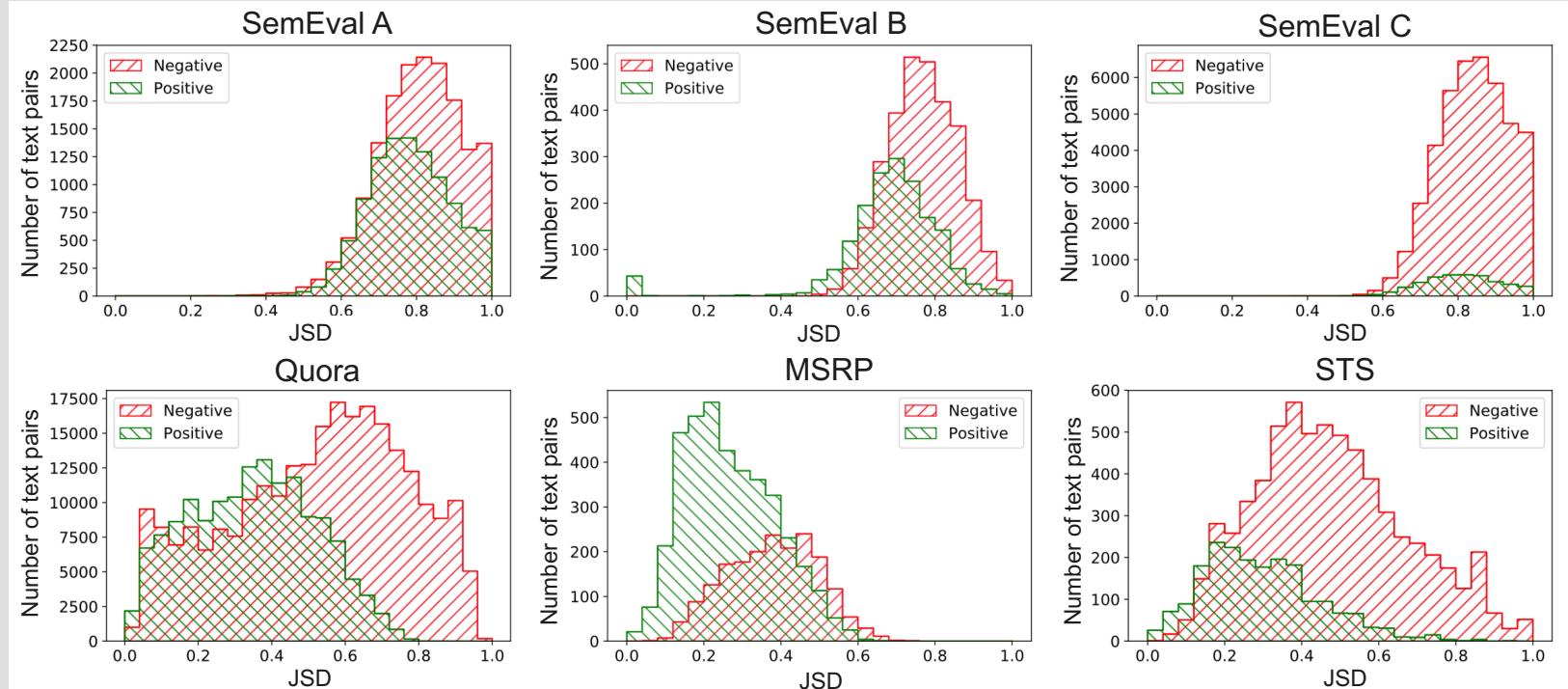
- Scrutiny of model and dataset biases in NLU  
(Wadhwa et al. 2018; Rajpurkar et al., 2018)
- High model performance even without crucial task information  
(Kaushik and Lipton 2018)
- Model performance inflated by annotation artefacts  
(Gururangan et al. 2018)

---

## Questions

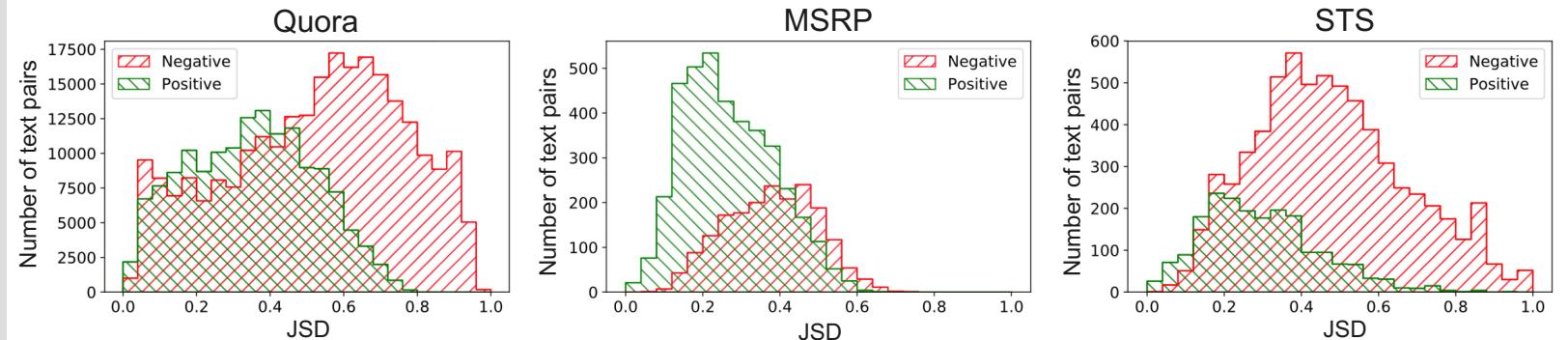
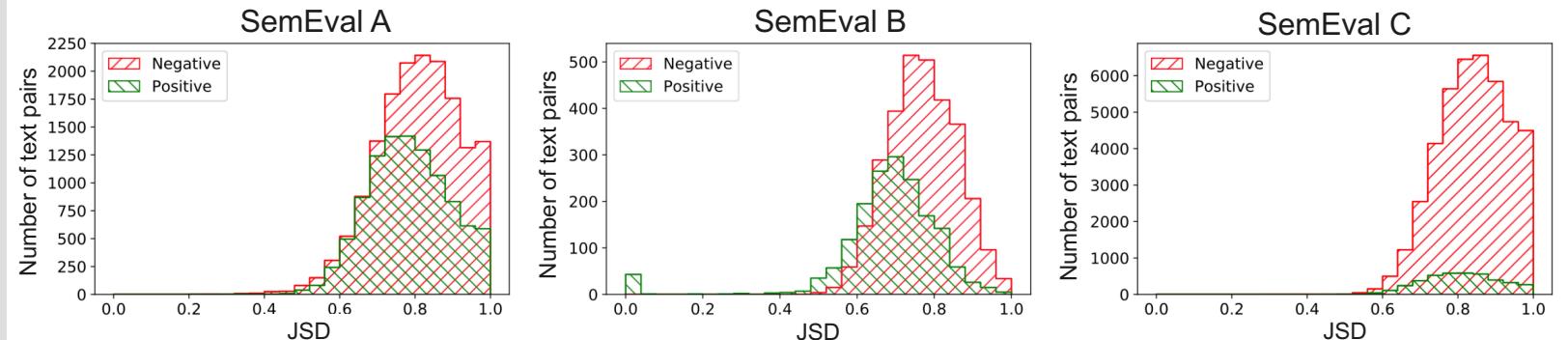
1. How to automatically distinguish obvious vs. non-obvious examples?
2. How common are obvious examples?
3. How to take this into account for evaluation?

# Lexical Divergence by Labels



# Lexical Divergence by Labels

similar distributions



distinct distributions

# Obvious vs. Non-Obvious Examples

	<b>Positive label</b>	<b>Negative label</b>
<b>Low lexical divergence (JSD≤median)</b>	What's the origin of the word o'clock? --- What is the origin of the word o'clock?	What is meant by 'e' in mathematics? --- What is meant by mathematics?
<b>High lexical divergence (JSD&gt;median)</b>	Which is the best way to learn coding? --- How do you learn to program?	What are the range of careers in Biotechnology in Indonesia? --- How do you tenderize beef stew meat?

# Obvious vs. Non-Obvious Examples

	<b>Positive label</b>	<b>Negative label</b>
<b>Low lexical divergence (JSD≤median)</b>	<p>What's the origin of the word o'clock?</p> <p>P<sub>o</sub></p> <p>What is the origin of the word o'clock?</p>	<p>What is meant by 'e' in mathematics?</p> <p>---</p> <p>What is meant by mathematics?</p>
<b>High lexical divergence (JSD&gt;median)</b>	<p>Which is the best way to learn coding?</p> <p>---</p> <p>How do you learn to program?</p>	<p>What are the range of careers in Biotechnology in Indonesia?</p> <p>N<sub>o</sub></p> <p>How do you tenderize beef stew meat?</p>

# Obvious vs. Non-Obvious Examples

	<b>Positive label</b>	<b>Negative label</b>
<b>Low lexical divergence (JSD≤median)</b>	What's the origin of the word o'clock? <b>P<sub>o</sub></b> What is the origin of the word o'clock?	What is meant by 'e' in mathematics? <b>N<sub>n</sub></b> What is meant by mathematics?
<b>High lexical divergence (JSD&gt;median)</b>	Which is the best way to learn coding? <b>P<sub>n</sub></b> How do you learn to program?	What are the range of careers in Biotechnology in Indonesia? <b>N<sub>o</sub></b> How do you tenderize beef stew meat?

# Human Annotation Statistics on Quora

	Fleiss' Kappa	Avg. time per pair	Instances
P <sub>o</sub>	0.6429	11.58s	35
P <sub>n</sub>	0.0878	11.68s	15
N <sub>o</sub>	0.3886	12.50s	34
N <sub>n</sub>	0.0892	13.83s	16
Total			100

# Obvious Cases in Datasets

	SemEval			Quora	MSRP	STS
	A	B	C			
P <sub>o</sub>	5893	1162	2492	107612	2398	1597
P <sub>n</sub>	4428	531	1590	41691	1502	409
N <sub>o</sub>	8842	1843	22155	160410	1398	3900
N <sub>n</sub>	7377	1213	21253	94632	503	2719
Obvious	56%	63%	52%	66%	65%	64%
Median	0.80	0.79	0.82	0.53	0.52	0.52

# Evaluation

- Report non-obvious F1
- Performance decrease & different ranking across SemEval tasks

	IIT UHH	Bunji	KeLP	EICA
F1	0.144 (2 <sup>nd</sup> )	<b>0.197</b> (1 <sup>st</sup> )	0.121 (3 <sup>rd</sup> )	0.008
F1 <sub>n</sub>	0.047 (2 <sup>nd</sup> )	0.028 (3 <sup>rd</sup> )	<b>0.054</b> (1 <sup>st</sup> )	0.000

Evaluation for SemEval Task C

---

## Conclusion

1. Criterion to distinguish obvious vs. non-obvious items
2. Existing datasets contain >50% of obvious examples
3. Non-obvious F1 score to evaluate on difficult cases

---

# Thank you



@wuningxi



n.peinelt@warwick.ac.uk

