

On the Automatic Detection and Classification of Linguistic Bias

Richard Diehl Martinez, Sabri Eyuboglu



Dataset

Wikipedia Neutrality Corpus

50,000 biased sentences

Words responsible for bias are tagged

Schnabel himself did a **fantastic** reproduction of Basquiat's work.

No bias type labels!

Hand-labeled Data

500 biased sentences labeled as framing or epistemological

Schnabel himself did a **fantastic** reproduction of Basquiat's work. → **Framing**

Lower cost of living means a **possibly** higher standard of living → **Epistemoglical**

Training on Gold Labels

Bias Detection Model

Lower cost of living means a **possibly** higher standard of living

BERT

Bias Classification Model

Epistemological Framing

ATTENTION

Lower cost of living means a possibly higher standard of living

Lower cost of living means a **possibly** higher standard of living

BERT out of Box		BERT		Attention	
Accuracy	0.573	Accuracy	0.606	Accuracy	0.679
AUROC	0.618	AUROC	0.608	AUROC	0.756

Training on Weak Labels

✓ **Wikipedia Neutrality Corpus (WNC)**
 $n = 50,000$

✓ **Bias Detection Model**
BERT Trained on WNC

Schnabel himself did a **fantastic** reproduction of Basquiat's work.

Biased Sentence

Epistemological Bias

Framing Bias

Principal Challenge: **LACK OF LABELS**

Wikipedia Neutrality Corpus
50,000 unlabeled sentences

fantastic

Glove

Feed Forward

λ_1

fantastic

POS Tagger

Feed Forward

λ_2

fantastic

Linguistic Features *

Feed Forward

λ_3

$p(y|\lambda)$

y

Wikipedia Neutrality Corpus
50,000 **weakly labeled** sentences

	BERT		Attention
Accuracy	0.742	Accuracy	0.704
AUROC	0.847	AUROC	0.770