

AuthentiGPT: Detecting Machine-Generated Text via Black-Box Language Models Denoising

Zhen Guo, Shangdi Yu

MIT Department of Electrical Engineering and Computer Science

1. Results

Methods	AUROC Scores
GPT-3.5 (zero-shot)	0.721
GPT-4 (zero-shot)	0.577
GPTZero	0.797
Originality AI	0.906
AuthentiGPT	0.918

GPT-3.5	0.53	0.93	0.80	0.96	0.97
GPT-4	0.19	0.99	0.93	0.95	1.00
GPTZero	0.99	0.51	0.14	0.71	0.74
Originality	0.89	0.85	0.99	0.75	0.97
AuthenticGPT	0.86	0.93	0.53	0.93	0.90
	PubMedQA	GPT3.5-re	GPT4-re	GPT3.5-new	GPT4-new

2. Algorithm

Algorithm 1 Detecting Machine-Generated Text

```

1: procedure GETSIMILARITY( $S$ , LLM)
2:   for  $i$  in  $[1, 2, \dots, \beta]$  do
3:      $M \leftarrow \text{maskSentences}(S, \alpha)$ 
4:      $D \leftarrow \text{denoiseSentences}(M, \text{LLM})$ 
5:      $E_S \leftarrow \text{computeEmbeddings}(S)$ 
6:      $E_D \leftarrow \text{computeEmbeddings}(D)$ 
7:      $D_{\text{sim},i} \leftarrow \text{getSimilarity}(E_S, E_D)$ 
8:   end for
9:   return mean( $[D_{\text{sim},1}, \dots, D_{\text{sim},\beta}]$ )
10: end procedure
11: procedure AUTHENTICGPT(LLM,  $S_{\text{test}}$ ,  $S_{\text{train}}$ , labels)
12:    $\mathcal{D}_{\text{train}} \leftarrow \text{GetSimilarity}(S, \text{LLM})$ 
13:    $gm \leftarrow \text{FindThreshold}(\mathcal{D}_{\text{train}}, \text{labels})$ 
14:    $\mathcal{D}_{\text{sim}} \leftarrow \text{GetSimilarity}(S_{\text{test}}, \text{LLM})$ 
15:   return  $gm.\text{classify}(\mathcal{D}_{\text{sim}})$ 
16: end procedure

```

Algorithm 2 Determine classification threshold

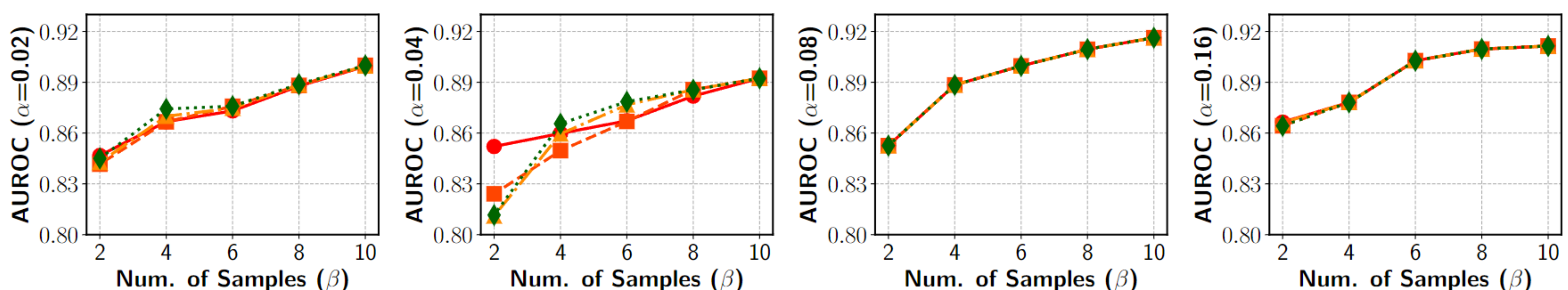
```

1: procedure FINDTHRESHOLD( $\mathcal{D}_{\text{sim}}$ , labels)
2:   for  $\lambda$  in  $[\lambda_1, \lambda_2, \dots, \lambda_n]$  do
3:      $\tilde{\mathcal{D}}_{\lambda} \leftarrow \text{Box-Cox}(\mathcal{D}_{\text{sim}}, \lambda)$ 
4:      $gm_{\lambda} = \text{GaussianMixture}(\tilde{\mathcal{D}}_{\lambda}, n_{\text{class}}=2)$ 
5:      $\text{score}_{\lambda} = \text{AUROC}(gm_{\lambda}, \text{labels})$ 
6:   end for
7:   return the  $gm_{\lambda}$  and corresponding  $\lambda$  that yield the maximum AUROC score
8: end procedure

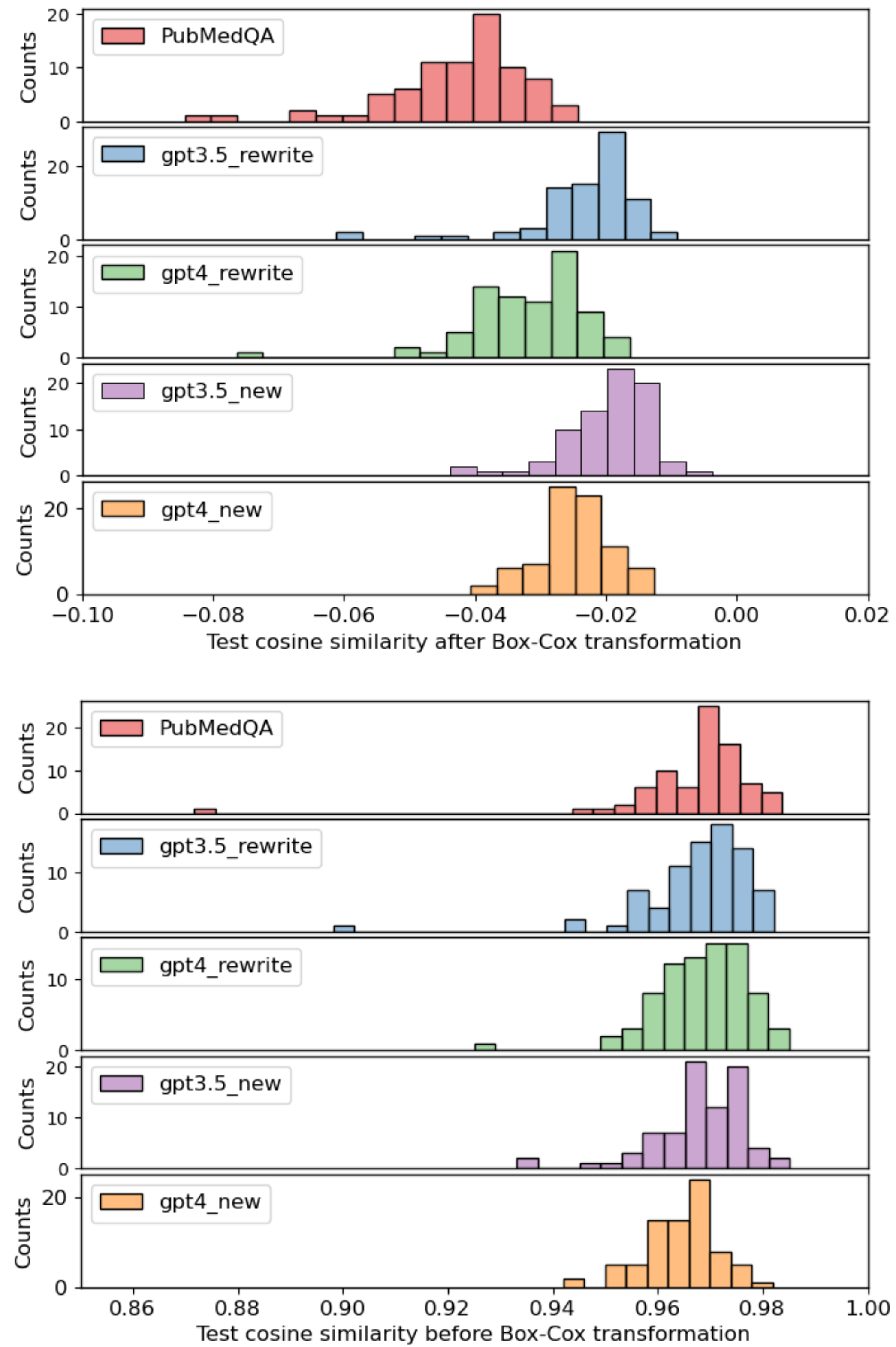
```

3. Training Sample Analysis

● training samples: 5 ■ training samples: 10 ▲ training samples: 15 ◆ training samples: 20



4. Distribution Analysis



5. References

- Zellers, Rowan, et al. "Defending against neural fake news." *Advances in neural information processing systems* 32 (2019).
- Gehrmann, Sebastian, Hendrik Strobelt, and Alexander M. Rush. "Gltr: Statistical detection and visualization of generated text." *arXiv preprint arXiv:1906.04043* (2019).
- Ippolito, Daphne, et al. "Automatic detection of generated text is easiest when humans are fooled." *arXiv preprint arXiv:1911.00650* (2019).
- Badaskar, Sameer, Sachin Agarwal, and Shilpa Arora. "Identifying real or fake articles: Towards better language modeling." *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-II*. 2008.
- Guo, Zhen, et al. "Dr. LLaMA: Improving Small Language Models in Domain-Specific QA via Generative Data Augmentation." *arXiv preprint arXiv:2305.07804* (2023).
- Kirchenbauer, John, et al. "On the Reliability of Watermarks for Large Language Models." *arXiv preprint arXiv:2306.04634* (2023).
- Kirchenbauer, John, et al. "On the Reliability of Watermarks for Large Language Models." *arXiv preprint arXiv:2306.04634* (2023).
- Yang, Xianjun, et al. "DNA-GPT: Divergent N-Gram Analysis for Training-Free Detection of GPT-Generated Text." *arXiv preprint arXiv:2305.17359* (2023).