

Detecting Machine-Generated Texts by Multi-Population Aware Optimization for Maximum Mean Discrepancy

Shuhai Zhang^{15*}, Yiliao Song^{2*}, Jiahao Yang¹, Yuanqing Li^{51†}, Bo Han^{4†}, Mingkui Tan^{13†}

South China University of Technology¹ The University of Adelaide²

Key Laboratory of Big Data and Intelligent Robot, Ministry of Education³

Department of Computer Science, Hong Kong Baptist University⁴ Pazhou Laboratory⁵

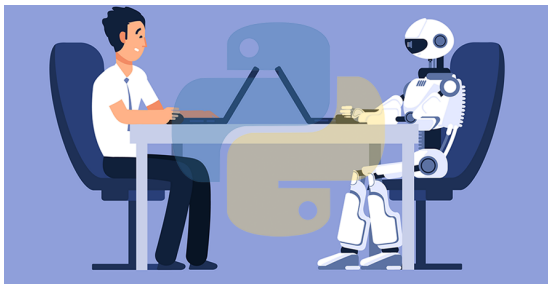
April 18, 2024

- 1 Background
- 2 Preliminaries and Motivations
- 3 MMD-MP for Text Detection
- 4 Experiments
- 5 Conclusions

- 1 Background
- 2 Preliminaries and Motivations
- 3 MMD-MP for Text Detection
- 4 Experiments
- 5 Conclusions

Background

Large language models (LLMs) have exhibited remarkable performance in generating human-like texts but may carry **critical risks**, *e.g.*, plagiarism issues, and hallucination issues.



Human or AI?

Background

It is challenging to distinguish machine-generated texts from human-written texts since **the distribution discrepancy between them is often very subtle** due to the advancement of LLMs.



Detector	HumanRec	MachineRec	AvgRec
GPT-3.5	96.98%	12.03%	54.51%
Human	61.02%	47.98%	54.50%

Li Y, Li Q, Cui L, et al. Deepfake text detection in the wild. arXiv, 2023.

Outline

- 1 Background
- 2 Preliminaries and Motivations**
- 3 MMD-MP for Text Detection
- 4 Experiments
- 5 Conclusions

Preliminaries of Maximum Mean Discrepancy

- *Maximum mean discrepancy* (MMD) aims to measure the distance between two distributions:

$$\begin{aligned}\text{MMD}(\mathbb{P}, \mathbb{Q}; \mathcal{H}_k) &= \sup_{f \in \mathcal{F}, \|f\|_{\mathcal{H}_k} \leq 1} |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]| \\ &= \sqrt{\mathbb{E}[k(X, X') + k(Y, Y') - 2k(X, Y)]}.\end{aligned}$$

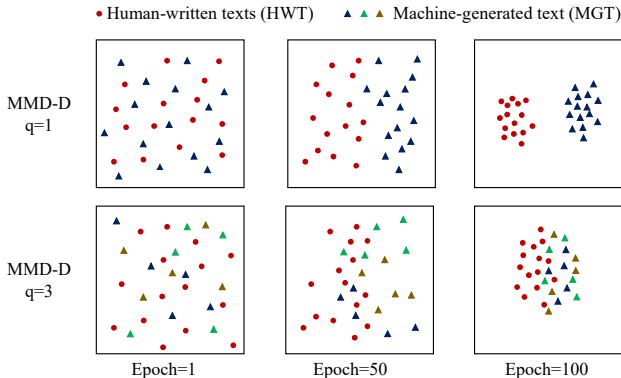
- Intuitively, we could view $k(X, X')$ or $k(Y, Y')$ as an *intra-class* distance and $k(X, Y)$ as an *inter-class* distance.
- Objective function for kernel-based MMD:

$$J(\mathbb{P}, \mathbb{Q}; k_\omega) = \text{MMD}^2(\mathbb{P}, \mathbb{Q}; k_\omega) / \sigma_{\mathcal{H}_1}(\mathbb{P}, \mathbb{Q}; k_\omega),$$

where $\sigma_{\mathcal{H}_1}^2 := 4 \left(\mathbb{E}[H_{ij}H_{i\ell}] - \mathbb{E}[H_{ij}]^2 \right)$ and
 $H_{ij} := k(\mathbf{x}_i, \mathbf{x}_j) - k(\mathbf{x}_i, \mathbf{y}_j) - k(\mathbf{y}_i, \mathbf{x}_j) + k(\mathbf{y}_i, \mathbf{y}_j)$.

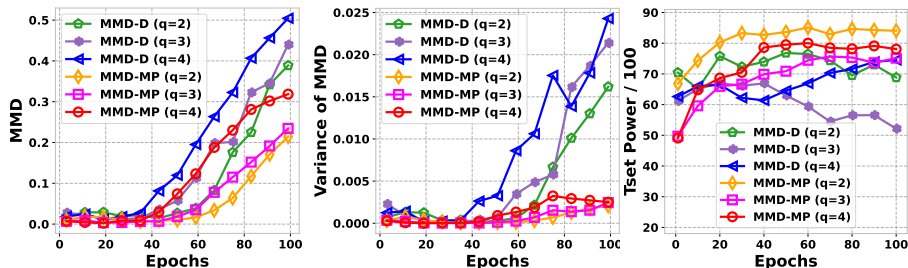
- For text detection, we define \mathbb{P} and \mathbb{Q} as the distribution of human-written texts (HWT) and machine-generated text (MGT), respectively.

High Variance Issue of MMD in Multiple Populations



- For MMD-D optimization, it tends to **separately aggregate** HWTs and **all possible** MGTs, such as decreasing the intra-class distance, and simultaneously push them away from each other like increasing inter-class distance.
- When the machine-generated texts population S_Q^{tr} **comprises different populations**, this optimization presents challenges due to **significant high variance issue of MMD**.

High Variance Issue of MMD in Multiple Populations



- During the optimization, as the number of S_Q^{tr} populations (*i.e.*, q) increases, kernel-based MMD (MMD-D) shows **an increase in MMD**, accompanied by a **sharp rise in variance**, resulting in **unstable test power** when testing.
- We propose a novel *multi-population* aware optimization method for MMD called MMD-MP, which exhibits **minimal variance of MMD values**, leading to **higher and more stable test power** when testing.

- 1 Background
- 2 Preliminaries and Motivations
- 3 MMD-MP for Text Detection**
- 4 Experiments
- 5 Conclusions

Intuition for MMD-MP: During the training, we **do not consider optimizing the intra-class distance of machine-generated text samples** in $S_{\mathbb{Q}}^{tr}$.

- MMD-MP maximizes the following objective with a proxy MPP:

$$J(\mathbb{P}, \mathbb{Q}; k_{\omega}) = \text{MPP}(\mathbb{P}, \mathbb{Q}; k_{\omega}) / \sigma_{\mathfrak{H}_1^*}(\mathbb{P}, \mathbb{Q}; k_{\omega}),$$

$$\text{MPP}(\mathbb{P}, \mathbb{Q}; \mathcal{H}_k) := \mathbb{E} [k_{\omega}(X, X') - 2k_{\omega}(X, Y)].$$

- Empirically, we can estimate it by:

$$\hat{J}(S_{\mathbb{P}}, S_{\mathbb{Q}}; k_{\omega}) = \frac{\widehat{\text{MPP}}_u(S_{\mathbb{P}}, S_{\mathbb{Q}}; k_{\omega})}{\sqrt{\hat{\sigma}_{\mathfrak{H}_1^*}^2(S_{\mathbb{P}}, S_{\mathbb{Q}}; k_{\omega}) + \lambda}},$$

$$\widehat{\text{MPP}}_u(S_{\mathbb{P}}, S_{\mathbb{Q}}; k_{\omega}) = \frac{\sum_{i \neq j} H_{ij}^*}{n(n-1)}, \quad \hat{\sigma}_{\mathfrak{H}_1^*}^2 := \frac{4}{n^3} \sum_{i=1}^n \left(\sum_{j=1}^n H_{ij}^* \right)^2 - \frac{4}{n^4} \left(\sum_{i=1}^n \sum_{j=1}^n H_{ij}^* \right)^2,$$

$$\text{where } H_{ij}^* := k_{\omega}(\mathbf{x}_i, \mathbf{x}_j) - k_{\omega}(\mathbf{x}_i, \mathbf{y}_j) - k_{\omega}(\mathbf{y}_i, \mathbf{x}_j).$$

Testing with MMD-MP for Text Detection

- Algorithm for paragraph-based detection, which can be considered as a two-sample test:

Algorithm 1 Testing with MMD-MP for 2ST

Input: Testing texts $S_{\mathbb{P}}^{te}, S_{\mathbb{Q}}^{te}, \hat{f}, k_{\omega}$;

$est \leftarrow \widehat{\text{MMD}}_u^2(S_{\mathbb{P}}^{te}, S_{\mathbb{Q}}^{te}; k_{\omega})$;

for $i = 1, 2, \dots, n_{perm}$ **do**

Shuffle $S_{\mathbb{P}}^{te} \cup S_{\mathbb{Q}}^{te}$ into S_X and S_Y ;

$perm_i \leftarrow \widehat{\text{MMD}}_u^2(S_X, S_Y; k_{\omega})$;

end for

Output: p -value $\frac{1}{n_{perm}} \sum_{i=1}^{n_{perm}} \mathbf{1}(perm_i \geq est)$

Testing with MMD-MP for Text Detection

- Algorithm for sentence-based detection, which can be considered as a single-instance detection task:

Algorithm 2 Testing with MMD-MP for 2ST

Input: Referenced HWT $S_{\mathbb{P}}^{re}$, testing texts

$S_{\mathbb{P}}^{te}, S_{\mathbb{Q}}^{te}, \hat{f}, k_{\omega};$

for $\mathbf{x}_i, \mathbf{y}_j$ in $S_{\mathbb{P}}^{te}, S_{\mathbb{Q}}^{te}$ **do**

$P_i \leftarrow \widehat{\text{MMD}}_b^2(S_{\mathbb{P}}^{re}, \{\mathbf{x}_i\}; k_{\omega});$

$Q_j \leftarrow \widehat{\text{MMD}}_b^2(S_{\mathbb{P}}^{re}, \{\mathbf{y}_j\}; k_{\omega});$

end for

Output: AUROC with two sets $\{P_i\}, \{Q_j\}$

Outline

- 1 Background
- 2 Preliminaries and Motivations
- 3 MMD-MP for Text Detection
- 4 Experiments**
- 5 Conclusions

Comparisons on Paragraph-based Detection

Our MMD-MP achieves **higher test power** than state-of-the-art text detection methods on HC3, given 3, 100 processed paragraphs in training data.

Method	ChatGPT	GPT3-S	Neo-S	ChatGPT Neo-S	ChatGPT GPT3-S
C2ST-S	62.83 \pm 0.90	43.64 \pm 5.92	30.68 \pm 2.37	34.62 \pm 2.73	46.66 \pm 2.95
C2ST-L	89.82 \pm 1.02	75.74 \pm 4.90	60.97 \pm 1.87	68.50 \pm 1.81	78.22 \pm 3.12
MMD-O	26.43 \pm 1.40	21.17 \pm 3.12	19.83 \pm 2.81	25.23 \pm 0.47	25.18 \pm 1.41
MMD-D	91.76 \pm 1.58	86.98 \pm 2.53	75.45 \pm 4.96	86.44 \pm 1.07	91.46 \pm 0.47
MMD-MP (Ours)	93.21\pm1.35	89.36\pm2.91	79.68\pm2.42	89.63\pm1.94	91.96\pm0.62

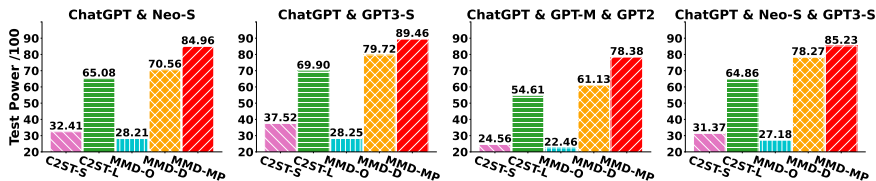
Comparisons on Sentence-based Detection

Our MMD-MP **outperforms state-of-the-art** methods for sentence-based detection on HC3 given 3, 100 processed paragraphs in training data.

Method	ChatGPT	GPT3-S	Neo-S	ChatGPT Neo-S	ChatGPT GPT3-S
Likelihood	89.82 \pm 0.03	60.56 \pm 1.32	61.18 \pm 1.25	75.81 \pm 0.51	75.05 \pm 0.25
Rank	73.20 \pm 1.49	71.96 \pm 1.01	72.09 \pm 0.51	72.74 \pm 0.74	72.34 \pm 1.38
Log-Rank	89.58 \pm 0.07	63.78 \pm 1.29	64.92 \pm 1.04	77.57 \pm 0.55	76.47 \pm 0.12
Entropy	31.53 \pm 0.90	54.34 \pm 1.33	56.19 \pm 0.33	44.08 \pm 0.24	42.08 \pm 2.01
DetectGPT-d	77.92 \pm 0.74	53.41 \pm 0.41	52.07 \pm 0.38	66.01 \pm 0.29	65.70 \pm 1.14
DetectGPT-z	81.07 \pm 0.77	53.45 \pm 0.53	52.28 \pm 0.31	67.54 \pm 0.19	67.32 \pm 1.02
OpenAI-D	78.57 \pm 1.55	84.05 \pm 0.71	84.86 \pm 0.87	81.20 \pm 0.95	80.68 \pm 1.64
ChatGPT-D	95.64 \pm 0.13	61.89 \pm 1.04	54.45 \pm 0.10	75.47 \pm 0.63	78.95 \pm 1.00
CE-Classifer	96.19 \pm 0.17	92.44 \pm 0.63	88.88 \pm 0.19	90.93 \pm 0.72	92.97 \pm 0.28
MMD-O	56.34 \pm 0.66	59.90 \pm 0.87	63.19 \pm 0.76	60.46 \pm 1.28	57.79 \pm 1.25
MMD-D	95.83 \pm 0.37	94.86 \pm 0.48	91.12 \pm 0.38	91.39 \pm 0.86	93.49 \pm 0.46
MMD-MP (Ours)	96.20\pm0.28	95.08\pm0.32	92.04\pm0.58	92.48\pm0.37	94.61\pm0.22

Comparisons on Unbalanced Training Data

Our MMD-MP exhibits **significantly superior performance** compared with other methods, *e.g.*, surpassing the test power of 6.96%~14.40% \uparrow than MMD-D, highlighting its stability under unbalanced training data scenarios.



Outline

- 1 Background
- 2 Preliminaries and Motivations
- 3 MMD-MP for Text Detection
- 4 Experiments
- 5 Conclusions**

Conclusions

- We delve into the optimization mechanism of MMD and reveal that **high variance of the MMD** when handling training data from **multiple different populations** can result in an **unstable discrepancy estimation** for MGT detection.
- We propose a novel multi-population aware optimization method for training kernel-based MMD (called MMD-MP), which can **alleviate the poor optimization of MMD-D** and **improve the stability of discrepancy measures**.
- Relying on MMD-MP, we develop two methods for **paragraph-based** and **sentence-based detection**, respectively. Extensive experiments across numerous LLMs, including ChatGPT, GPT2 series, GPT3 series, GPT-Neo series, demonstrate **superior detection performance**.

Closing Remarks

Thank you for your attention. Scan for more details.

