

Locally Differentially Private In-context Learning

C. Zheng, K. Sun, W. Zhao, H. Zhou, L. Jiang, S. Song, Chunlai Zhou

Renmin Unversity of China

LREC-COLING' 2024

Contents



03.



Preliminar knowledgey

04. Main Contribution





an Era of ChatGPT

ChatGPT attracts attention!



Emergent Ability: In-context Learning





中國人民大學 RENMIN UNIVERSITY OF CHINA

Privacy Leakage



In February 2023, ChatGPtenhanced Bing was tricked by users into revealing project codes and other technical secrets about its programming. In March 2023, Samsung was exposed that employees used ChatGPT to leak chip secrets, including two "device information leaks" and one "conference content leak".

SAMSUNG



In March 2023, ChatGPT was revealed to have leaked user data and payment information. In November of that year, researchers discovered that prompt injections could enable ChatGPT to reveal a large amount of private information in its training data.





Legislation Laws about Privacy



protect large language models





in-context learning privacy protection methods

It is computationally inefficient ,computationally expensive, technically demanding, and difficult to apply non-open source models

DP-SGD

Multiple API access: high cost and slow, Assumption: Third party LLMS are trusted Requires open source datasets, multiple access to LLM

DP-ICL Prompt**PATE**

This paper adopts locally differential privacy technology as privacy protection for in-context learning.

LDP-ICL



Why choose Localized differential privacy technology?

- **1** Treat a third-party LLM as untrustworthy, with a higher level of protection
- 2 With only one access to the API, the calculation is fast and the cost is low
- 3 Easy to operate, LLM as a black box, no need to master the technology of LLM training or fine tuning
- 4 In practice, the set of examples for in-context learning is relatively small

5 No need to find open source datasets, no need to access multiple LLMS





Large language model privacy protection methods - Related papers

DP-SGD

□ Abadi, Martin, et al. "Deep learning with differential privacy." Proceedings of the 2016 ACM SIGSAC conference on computer and communications security. 2016.

□[85]Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. 2022. Large language models can be strong differentially private learners. In International Conference on Learning Representations.

□[86]Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. 2023. Differentially private bias-term only fine-tuning of foundation models.

□[87]Jiyan He, Xuechen Li, Da Yu, Huishuai Zhang, Janardhan Kulkarni, Yin Tat Lee, Arturs Backurs, Nenghai Yu, and Jiang Bian. 2023. Exploring the limits of differentially private deep learning with group-wise clipping. In The Eleventh International Conference on Learning Representations.

PromptPATE :

□ Duan, Haonan, et al. "Flocks of Stochastic Parrots: Differentially Private Prompt Learning for Large Language Models." arXiv preprint arXiv:2305.15594 (2023).

DP-ICL:

Ashwinee Panda, Tong Wu, et al. "Differentially Private In-Context Learning." arXiv preprint arXiv:2305.01639 (2023).





Large language model

Emergent ability
 In-context learning ability:
 Reasoning ability
 Cross-modal understanding

• Security issues

Model theftPrompt attackData theftMembers infererece attacksBackdoor attacksData poisoningadversary-example attack

Demonstrations

ALL.	FUSITIVE
\n	Neutral
\n	Negative
\n	<u></u>
	\n \n \n

03 Preliminary knowledge



中国人民大學

Avoid fine-tuning the weights of task-specific model altogether, and instead rely entirely on contextual information

In-context learning

• **Definition**: *In-context learning* is a paradigm that allows a language model to learn a task when given only a description of the task and a small number of samples (in demonstration form). It estimates the likelihood of a potential answer under demonstrative conditions by using a trained language model.

• Formula representation

$$P(y_j|x) \triangleq f_{\mathcal{M}}(y_j, C, x)$$

among $C = \{I, s(x_1, y_1), ..., s(x_k, y_k)\}$ or $C = \{s(x_1, y_1), ..., s(x_k, y_k)\}$





Differential private data processing framework



Locally differential privacy (left) and centralized Differential Privacy (right) data processing frameworks



K-RR randomized response mechanism

- <u>Definition</u>: *K-RR* mainly overcomes the problem of the random response technique (RR) for binary variables, which can directly perform a random response if the variable contains K selected values.
- Formula representation

 $Q_{k-\mathbf{R}\mathbf{R}}(y' \mid y) = \frac{1}{K-1+e^{\epsilon}} \begin{cases} e^{\epsilon} & \text{if } y' = y, \\ 1 & \text{if } y' \neq y. \end{cases}$



K-RR is a typical noising method in differential privacy technology. When k=2, k–RR is the famous Warner mechanism. In this paper, we focus on binary classification in context learning, using Warner mechanism to protect privacy in labels.



04. Main Contributions

04



Threat model:

Prompt-leakage Attack





Tong Wu, Ashwinee Panda, Jiachen T. Wang, Prateek Mittal: Privacy-preserving ICL, 2024

04



Design of privacy protection mechanism



Figure 1: Pipeline of LDP-ICL



Theoretical Foundations

Assumption

 By treating ICL as a dual form of optimization based on gradient descent,

This paper deduces a localized differential privacy in-context learning formula based on previous results in the literature

 $P(y_{\text{test}} \mid Q_{k-\text{RR}}(\mathcal{E}_n), \boldsymbol{x}_{\text{test}}) = \sigma(\boldsymbol{W}_0 \boldsymbol{x}_{\text{test}} - \eta \sum_{i=1}^n (\sigma(\boldsymbol{W}_0 \boldsymbol{x}_i) - Q_{k \cdot \text{RR}}(\boldsymbol{y}_i)) \boldsymbol{x}_i^T \boldsymbol{x}_{\text{test}})$





Experiment 1: LDP-ICL does classification task

Experiment 2: LDP-ICL is successfully applied to discrete distribution estimation problem





中國人民大學 RENMIN UNIVERSITY OF CHINA

Experimental design: Model dataset selection

Language Model: GPT-3.5-turbo

Data set:

- SST-2 and Subj for sentiment classification;
- Ethos is a hate speech detection dataset;
- SMS_Spam is used to identify spam





Experimental design: parameter setting

Privacy budget: $\varepsilon = \{0, 0.5, 1, 2, 3, 8, \infty\}$

Test examples: 150 test examples are selected each time from the evaluation verification set to evaluate the performance is the average of 6 runs under the same parameter configuration.

Distribution estimation: For the distribution estimation scenario, we selected the number of queries R=1000 for the SST-2 dataset and the number of searches R=500 for the Ethos dataset





Experimental design: Sample sample selection Selection criteria:

 The number of samples in each group is the same, and 32 samples are selected in the text;
 Each group of sample samples has a complete label space, such as emotion classification label space {positive, negative};

3. The number of sample samples for each type of label is the same. For example, the emotion classification task, 16 positive examples and 16 negative examples;

4, sample examples in the same order. For different privacy parameter configurations, the same group of query predictions, use the same sample sample order;

5, uniform input and output formats. All sample examples use the same format, such as input: label, or input => label;

6, full correspondence. For example each sample example has a complete input corresponding to the label;

7. All sample examples are from the training data set.





Experimental design: LDP-ICL

Experiment 1: LDP does the classification task

Algorithm 1: LDP-ICL

Input: Private data D, query q, model **LLM**, privacy budget ϵ , number of demonstration examples n.

Output: Model prediction O(q)

- 1: Subsample of size n from \mathcal{D} and obtain \mathcal{E}_n
- 2: **Perturb** \mathcal{E}_n using k-RR and obtain $Q_{k-RR}(\mathcal{E}_n)$
- 3: Concatenate query and form $I(q) = Q_{k-\text{RR}}(\mathcal{E}_n) \cup q$
- 4: Obtain model output O(q) = LLM(I(q))

Experiment 2: Distribution estimation

Algorithm 2: LDP-ICL for distribution estimation

- **Input:** Private data D, model **LLM**, privacy budget ϵ , number of demonstration examples n, number of round(queries) R
- Output: Proportion estimation
- 1: Subsample of size R from \mathcal{D} , obtain $\{(x_{\text{test}}^i, y_{\text{test}}^i)\}_{i=1}^R$ and construct queries $\{q^i\}_{i=1}^R = \{(x_{\text{test}}^i, \underline{?})\}_{i=1}^R$
- 2: **Partition** \mathcal{D} into classes with size n: $\mathcal{D}_n^1, \dots, \mathcal{D}_n^l \leftarrow \mathcal{D}$
- 3: for $i \in \{1, \ldots, R\}$ do
- 4: **Perturb** \mathcal{D}_n^i using k-RR and obtain $Q_{k-\mathrm{RR}}(\mathcal{D}_n^i)$
- 5: Concatenate corresponding query and form $I(q^i) = Q_{k-\mathsf{RR}}(\mathcal{D}_n^i) \cup q^i$
- 6: Obtain *i*-th model output $O(q^i) =$ LLM $(I(q^i))$
- 7: end for
- 8: Calculate estimated rate (Eq.(8))





Experiment design: baseline setting

Baseline one: non-privacy-protected context learning (Non-ICL/Glod-ICL), that is, $\epsilon = \infty$ in LDP-ICL, is equivalent to n samples in the demonstration set using real labels without perturbation.

Baseline two: Zero sample learning (ZSL), which is the same as single sample learning, except that demo examples are not allowed and only natural language instructions describing the task are given to the model.

Baseline three: Flipped label Context Learning (FL-ICL), is a flipped example of all labels, indicating a divergence between semantic prior knowledge and input label mapping. Performance accuracy is inversely proportional to the ability to learn input-label mappings and override semantic priors.



RENMIN UNIVERS

Experimental results: LDP-ICL







Experiment result: Distribution estimation



LDP-ICL estimates are closer to the true distribution and maintain a higher level of stability, showing better utility even with smaller privacy budgets. 04

PART

中国人民大學

Results: Controlled experiment

模型	方法	ε =3	€=8	€≕∞
RoBERTa-large	DP-SGD[85]	93.04	93.81	96.2
	DP-SGD[86]	94.6	94.7	95.5
	DP-SGD[87]	94.23	94.87	96.2
RoBERTa-base	pormptPATE[25]	86.35	92.32	93.23
GPT-3 Babbage	DP-ICL(n=4)[24]	95.8	95.92	96.05
	DP-ICL(n=16)[24]	91.64	96.32	96.13
GPT-3.5 Turbo	LDP-ICL(n=16)	94.45	94.9	95.77
	LDP-ICL(n=32)	94.11	94.12	94.12

Compared with the previous experiments, the research method in this paper shows a comparable level of performance, and further proves the feasibility of applying localized differential privacy technology PART



Experimental results: Ablation experiment



Ablation studies were conducted to analyze how changes in the number of demo examples affected the performance of the task.

Through experiments, we found that the curves showed the same trend for different number of demo examples. The consistency of this trend is consistent with our previous conclusions.

D Summary

- ICL-LDP: a privacy protection mechanism for in-context learning of large language models by using locally differential privacy technology.
- a formula is proposed, and the correctness and feasibility of the theoretical analysis is proved by experiments.

Future Research

- Research on the protection of sensitive information in a wider range of fields (such as text summary generation, inference, question and answer, etc.). (Our IJCAI2024 papger)
- The methods of sensitive information protection under multi-modal input scenarios are studied.