003

004

006

007

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

029

030

DUAL VARIANCE REDUCTION WITH MOMENTUM FOR IMBALANCED BLACK-BOX DISCRETE PROMPT LEARNING

Anonymous authors

Paper under double-blind review

Abstract

Black-box prompt learning has proven to be an effective approach for customizing large language models (LLMs) offered as services to address various downstream tasks. Within this domain, policy gradient-based methods have garnered substantial attention as a prominent approach for learning discrete prompts. However, the highly imbalanced data distribution in the real world limits the applicability of such approaches by influencing LLMs' tendency to favor certain categories. To tackle the challenge posed by imbalanced data, this paper pioneers the integration of pairwise AUC loss into the policy gradient optimization of discrete text prompts and proposes learning discrete prompts with doubly policy gradient. Unfortunately, the doubly policy gradient estimation suffers from two variance components, resulting in unstable optimization. As a further improvement, we propose (1) a novel unbiased variance-reduced doubly policy gradient estimator and (2) incorporating the STORM variance reduction technique. Ultimately, we introduce a novel *momentum-based discrete prompt learning method* with doubly policy gradient (mDP-DPG). Crucially, we provide theoretical convergence guarantees for mDP-DPG within standard frameworks. The experimental results show that mDP-DPG surpasses baseline approaches across diverse imbalanced text classification datasets, emphasizing the advantages of our proposed approach for tackling data imbalance. Our code is available at the following URL: https://anonymous.4open.science/r/DPDPG-1ECB.

035

1 INTRODUCTION

Large language models (LLMs) have achieved milestone accomplishments on a wide range of nat-037 ural language processing (NLP) tasks (Brown et al., 2020; Devlin et al., 2018; Raffel et al., 2020). However, the increasing parameters pose challenges for tuning. Prompting has emerged as the parameter-efficient paradigm for adapting LLMs to specific NLP tasks (Li & Liang, 2021; Lester 040 et al., 2021; Liu et al., 2023). Well-crafted prompts can effectively enhance the performance of 041 LLMs on various downstream tasks instead of retraining. Recently, considering that most exist-042 ing LLMs, such as GPT-4, only provide cloud-based API services, researchers have introduced the 043 Language-Model-as-a-Service scenario, where users are limited to interacting with LLMs solely 044 through APIs, creating a black-box setting (Sun et al., 2022b;a). Black-box prompt learning has 045 been an effective strategy for adapting LLMs to downstream tasks due to the opacity of black-box LLMs (Deng et al., 2022; Prasad et al., 2022; Diao et al., 2022). 046

Currently, much attention has been focused on black-box discrete prompt learning, with policy gradient-based methods becoming highly influential (Diao et al., 2022; Lin et al., 2023). Within these investigations, discrete prompt learning is treated as a distribution optimization problem, using policy gradient to update the prompt distribution and overcome the issue of inaccessible gradients in black-box LLMs. However, these efforts focus solely on vanilla text classification without any additional handling of imbalanced data. These imply that adapting them to address the class imbalance problem to bridge the gap between LLM services and downstream tasks while providing theoretical convergence guarantees remains a significant challenge.

- Departing from idealized situations, real-world data usually features severe class distribution imbalances, where certain minority classes are markedly less prevalent than majority classes in classification problems (Henning et al., 2022). For example, non-hateful tweets are more prevalent than hateful ones on Twitter (Waseem & Hovy, 2016), and the positive and negative reviews on Amazon are also imbalanced. Specifically, the ratio of negative to positive movie reviews is approximately 6.24, while it is as high as 9.75 for book reviews (Gao et al., 2021).
- 060 Simultaneously, imbalanced class distributions hinder prompt learning by making LLMs more likely 061 to prioritize well-represented classes, which in turn lowers prompt performance. Dong et al. (2022) 062 initially identified this phenomenon within the computer vision (CV) domain, noting that prompts 063 could benefit models on long-tail data, but their performance still lagged far behind the state-of-064 the-art. Similarly, class imbalance issues also impact the NLP black-box prompt learning. When the downstream task is a text classification problem, cross-entropy is typically chosen to build the 065 objective function. However, this becomes unreasonable in the presence of imbalanced data since 066 the minority class data have little influence, leading to the aforementioned issues (Liu et al., 2020). 067
- 068 A common approach in previous research for addressing imbalanced class distributions is to use 069 AUC (Area Under the ROC Curve) maximization as the optimization objective. The AUC is defined as the probability that the prediction score for a positive example exceeds that for a negative example in statistical terms (Hanley & McNeil, 1982; 1983), and AUC maximization is proposed 071 to address the imbalanced data problem (Zhao et al., 2011; Yuan et al., 2020). However, high 072 variance emerges as a major challenge in the learning process of prompts in two respects. Firstly, 073 the sampling of prompts from the prompt distribution introduces inherent randomness, making the 074 optimization process unstable. Secondly, sampling both positive and negative examples for AUC 075 maximization will also introduce high variance. In scenarios with imbalanced classes, the sampling 076 of data for training not only amplifies the disparity between majority and minority classes but also 077 increases the variance of gradient estimation during optimization. This duality of variance poses sig-078 nificant difficulties in effectively learning prompts that generalize well across all classes, particularly 079 in settings with highly imbalanced data.
- In order to tackle the above problems, we propose a novel *momentum-based discrete prompt learn-*081 ing method with doubly policy gradient (mDP-DPG) that utilizes AUC maximization to adapt LLMs to downstream tasks with imbalanced data. By minimizing AUC's pairwise surrogate loss us-083 ing policy gradient, mDP-DPG prevents prompts from being degraded by the majority class, thereby 084 preserving performance. Moreover, due to the doubly sampling of examples during gradient estima-085 tion, we refer to the policy gradient in mDP-DPG as the doubly sampled policy gradient, abbreviated as doubly policy gradient. We further propose an unbiased, variance-reduced doubly policy gradi-087 ent estimator (VR-DPGE) to improve convergence in practice. While the estimator suffers from 088 variance, stochastic sampling of mini-batches from the dataset also introduces variance into gradient estimation. To reduce the dual variance, we introduce the momentum-based variance reduction 089 strategy STORM (Cutkosky & Orabona, 2019; Huang et al., 2020). STORM does not rely on con-090 structing variance-reduced gradients through giant batch sizes, as SVRG does (Johnson & Zhang, 091 2013; Xiao & Zhang, 2014). Instead, it employs a variant of momentum, making it can be seam-092 lessly integrated into the optimization of pairwise AUC loss for variance reduction. Additionally, unlike other policy gradient-based methods for black-box discrete prompt learning, mDP-DPG has 094 rigorous theoretical convergence guarantees. 095
- ⁰⁹⁶ The main contributions of this work are summarized as follows:

099

- 1. We propose a novel momentum-based discrete prompt learning method named mDP-DPG, which introduces VR-DPGE and STROM strategy to address challenges posed by dual variance. Using pairwise AUC loss as objective function, mDP-DPG preserves prompts' performance in downstream tasks with class imbalance. As far as we know, we are the first to discuss how to address imbalanced data in, NLP prompt learning.
- 2. We establish rigorous theoretical analysis of the mDP-DPG. Specifically, we provide proof that mDP-DPG achieves an oracle complexity of $O(1/\epsilon^3)$, validating its effectiveness in optimizing the pairwise AUC loss for black-box prompt learning in the context of imbalanced data.
- 107 3. Numerous experimental results on RoBERTa-large, GPT2-XL, and Llmma3 show that our method achieve state-of-the-art across various class imbalance datasets, demonstrating the

110 111 effectiveness of prompts learned through mDP-DPG on imbalanced data. Furthermore, our research findings confirm that imbalanced data negatively impacts prompt learning, emphasizing the importance of imbalanced prompt learning.

- 2 RELATED WORKS
- 113 114

112

Black-Box Prompt Learning. There is a significant amount of research focusing on black-box 115 prompt learning, which has achieved promising results in NLP tasks (Brown et al., 2020; Prasad 116 et al., 2022; Han et al., 2023; Hou et al., 2023). Prompts come in two formats: continuous and 117 discrete. The continuous prompt is a series of vectors, which concatenates with token embeddings. 118 Sun et al. (2022b) propose the black-box tuning (BBT) framework, which optimizes prompts in 119 low-dimensional subspace and obtains continuous prompts in the original space through a random 120 matrix. Sun et al. (2022a) present an improved version of BBT that adds continuous prompt prefixes 121 to each hidden layer of LLM. Zheng et al. (2023) point out the inappropriateness of random matric 122 in Sun et al. (2022b) and leverage meta-learning to identify the optimal subspace. On the other hand, 123 the discrete prompt is a sequence of tokens, which is more appropriate for real applications. Deng 124 et al. (2022) formulate discrete prompt learning as a reinforcement learning problem and generate 125 discrete prompts using policy network. Diao et al. (2022) utilize policy gradients to optimize the distribution of the discrete prompt. Zhao et al. (2023) introduce a genetic algorithm to search for 126 discrete prompts guided by LLMs predictive probabilities. 127

128 AUC Maximization. AUC maximization is a machine learning paradigm aimed at maximizing the 129 AUC score of models. A substantial amount of research has been dedicated to this topic, integrating 130 it with various contexts such as supervised learning (Joachims, 2005), semi-supervised learning 131 (Iwata et al., 2020), online learning (Gao et al., 2016), and federated learning (Yuan et al., 2021). Many algorithms frequently minimize the pairwise surrogate loss for AUC maximization. Zhao 132 et al. (2011) propose two online AUC maximization algorithms, which utilize the reservoir sampling 133 technique to avoid memorizing all training samples. Gao et al. (2016) focus on the challenge of 134 optimizing with only a single pass of training samples and propose a regression-based algorithm 135 using square surrogate loss. The issue of the pairwise surrogate loss lies in the necessity to construct 136 sample pairs from opposite classes. To overcome this challenge, Ying et al. (2016) demonstrate the 137 equivalence between AUC optimization and the saddle point problem. Liu et al. (2020) extend the 138 aforementioned equivalence to the case of deep neural network models. 139

Variance Reduction. In optimization problems, variance reduction is a frequently utilized improvement method (Cutkosky & Orabona, 2019). Numerous variance reduction techniques necessitate setting gradient checkpoints and calculating gradients at these checkpoints with giant batch sizes Johnson & Zhang (2013); Fang et al. (2018); Zhou et al. (2018). Cutkosky & Orabona (2019) propose a variance reduction technique based on the momentum variant, termed STORM, which facilitates variance reduction in non-convex optimization without giant batch sizes. Huang et al. (2020) propose a similar variance reduction technique for zero-order optimization to accelerate black-box minimization and minimax optimization problems.

147 148

3 METHODOLOGY

149 150

In this section, we first present the problem formulation in Section 3.1, where we define discrete prompt learning with the pairwise AUC loss as the objective function. Subsequently, in Sections 3.2 and 3.3, we discuss how to address the challenges posed by the dual variance in optimization. Specifically, in Section 3.2, we present VR-DPGE to reduce the variance introduced by prompt sampling in the doubly policy gradient estimation. In Section 3.3, we introduce a momentum-based variance reduction technique to reduce the dual variance.

Notations. Let $\mathcal{D} \triangleq \{(\mathbf{X}_1, y_1), (\mathbf{X}_2, y_2), \dots, (\mathbf{X}_M, y_M)\}$ denote a set of training data with cardinality M. For any $m \in \{1, 2, \dots, M\}$, \mathbf{X}_m represents an input training example (e.g., a piece of text), and $y_m \in \{-1, 1\}$ denotes its corresponding label. We use $T(\cdot)$ to represent a tokenizer that converts an input text to a token vector, and $\mathbf{S}_m \triangleq T(\mathbf{X}_m)$ as the *m*-th token vector. Let $\tilde{\mathcal{D}} \triangleq \{(\mathbf{S}_1, y_1), (\mathbf{S}_2, y_2), \dots, (\mathbf{S}_M, y_M)\}$ be the set of tuples composed of token vectors and labels. Discrete prompt is defined as a token vector with *n* discrete tokens $\mathbf{T} = [t_1, t_2, \dots, t_n]$. For any

162 $i \in \{1, 2, ..., n\}$, the word $t_i \in W$ and W is the set consisting of tokens from a given vocabulary. 163 Let $\mathbf{S} \in \{\mathbf{S}_1, \mathbf{S}_2, ..., \mathbf{S}_M\}$ denotes any token vectors, then the user query to a black-box LLM $h(\cdot)$ 164 is denoted as $h([\mathbf{T}, \mathbf{S}])$, i.e. $h([\mathbf{T}, \mathbf{S}])$ denote the prediction of the LLM $h(\cdot)$ on an input $[\mathbf{T}, \mathbf{S}]$. For 165 $p \in \mathbb{N}^*$, $\mathbf{1}_p$ denotes the vector of size p composed only of ones.

167 3.1 PROBLEM STATEMENT

166

168

176

187

192 193

194

195 196 197

199

200

201 202

169 Discrete Prompt Learning. Discrete prompt learning aims to learn a discrete textual prompt con- **170** sisting of *n* tokens, which is a more pragmatic scenario. Diao et al. (2022) formulate discrete **171** prompt learning as a distribution optimization problem. Specifically, they assume each token of the **172** prompt is sampled from the independent categorical distribution $t_i \sim \operatorname{Cat}(p_i)$, where $p_i \in C$ and **173** $C = \{x : ||x||_1 = 1, 0 \le x \le 1\}$. The component $p_{i,j}$ denotes the probability of sampling the *j*-th **174** token from the vocabulary W. By denoting C the subset of $\mathbb{R}^{|W| \times n}$ such that for any $p \in C$ (where p_i denotes a column of p), $p_i \in C$, the objective function during optimization is as follows:

$$\min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \mathbb{E}_{(\mathbf{S},y)} \mathbb{E}_{\mathbf{T}}[\ell(h([\mathbf{T},\mathbf{S}]),y)] = \min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \mathbb{E}_{(\mathbf{S},y)} \Sigma_{\mathbf{T}}\ell(h([\mathbf{T},\mathbf{S}]),y)P(\mathbf{T})$$
(1)

where $\ell(\cdot)$ is the loss function that evaluates the prediction of the black-box LLM $h(\cdot)$ based on the ground truth label y. $P(\mathbf{T}) = \prod_{i=1}^{n} P(t_i)$ is the joint probability of the discrete prompt **T**. To avoid confusion, we refer to p_i as the token distribution and the joint distribution p of n token distributions as the prompt distribution.

AUC maximization. AUC is a common metric for evaluating model performance in imbalanced binary classification problems and AUC maximization is a form of pairwise learning that aims to maximize the AUC score during the model training. By employing the square loss as the surrogate, which is statistically consistent with AUC (Gao & Zhou, 2012), AUC maximization can be formulated as

$$\arg\min_{f} \mathbb{E}[(1 - f(x) + f(x'))^2 | y = +1, y' = -1]$$
(2)

Discrete Prompt Learning with AUC maximization. We propose employing pairwise AUC loss for black-box discrete prompt learning to address the challenge posed by class imbalance in prompt learning. Therefore, we can formulate the objective of black-box prompt learning that minimizes the expected risk $\mathcal{L}(p)$ as follows

$$\boldsymbol{p}^* = \arg\min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \mathcal{L}(\boldsymbol{p}) = \arg\min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \mathbb{E}_{(\mathbf{S},y)} \mathbb{E}_{(\mathbf{S}',y')} \mathbb{E}_{\mathbf{T}} L(h([\mathbf{T},\cdot]), (\mathbf{S},y), (\mathbf{S}',y'))$$
(3)

where $L(h([\mathbf{T}, \cdot]), (\mathbf{S}, y), (\mathbf{S}', y'))$ is the pairwise AUC square loss given a discrete prompt \mathbf{T} and a pair of samples $(\mathbf{S}, y), (\mathbf{S}', y')$.

$$L(h([\mathbf{T}, \cdot]), (\mathbf{S}, y), (\mathbf{S}', y')) = \begin{cases} (1 - h([\mathbf{T}, \mathbf{S}]) + h([\mathbf{T}, \mathbf{S}']))^2, \text{ if } y = +1 \text{ and } y' = -1, \\ 0, \text{ otherwise.} \end{cases}$$
(4)

In real-world applications, instead of minimizing the expected risk $\mathcal{L}(\boldsymbol{p})$, we consider an independent and identically distributed training set $\tilde{\mathcal{D}}$ and the empirical risk $\mathcal{L}_M(\boldsymbol{p})$ of the pairwise loss function on $\tilde{\mathcal{D}}$ is as follows

$$\boldsymbol{p}^{*} = \arg\min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \mathcal{L}_{M}(\boldsymbol{p}) = \arg\min_{\boldsymbol{p}\in\boldsymbol{\mathcal{C}}} \frac{1}{M(M-1)} \Sigma_{i,j\in\tilde{\mathcal{D}},i\neq j} \underbrace{\mathbb{E}_{\mathbf{T}} L(h([\mathbf{T},\cdot]), (\mathbf{S}_{i}, y_{i}), (\mathbf{S}_{j}, y_{j}))}_{F_{i,j}(\boldsymbol{p})}$$
(5)

203 204 205 206

213 214 215

3.2 REDUCE VARIANCE IN PROMPT SAMPLING WITH VR-DPGE

Now black-box discrete prompt learning with AUC maximization transforms into the task of solving the black-box pairwise learning problem. Given a pair of samples (\mathbf{S}_i, y_i) and (\mathbf{S}_j, y_j) , we can formulate doubly stochastic gradient $\nabla_{\mathbf{p}} F_{i,j}(\mathbf{p})$ as equation 6 with the aid of the policy gradient estimator for solving problem 5. Due to the double sampling, we refer to equation 6 doubly sampled policy gradient, abbreviated as doubly policy gradient.

$$\nabla_{\boldsymbol{p}} F_{i,j}(\boldsymbol{p}) = \nabla_{\boldsymbol{p}} \mathbb{E}_{\mathbf{T}} L(h([\mathbf{T}, \cdot]), (\mathbf{S}_i, y_i), (\mathbf{S}_j, y_j)) \\ = \sum_{\mathbf{T}} L(h([\mathbf{T}, \cdot]), (\mathbf{S}_i, y_i), (\mathbf{S}_j, y_j)) \nabla_{\boldsymbol{p}} P(\mathbf{T}) \\ = \sum_{\mathbf{T}} L(h([\mathbf{T}, \cdot]), (\mathbf{S}_i, y_i), (\mathbf{S}_j, y_j)) P(\mathbf{T}) \nabla_{\boldsymbol{p}} \log P(\mathbf{T}) \\ = \mathbb{E}_{\mathbf{T}} L(h([\mathbf{T}, \cdot]), (\mathbf{S}_i, y_i), (\mathbf{S}_j, y_j)) \nabla_{\boldsymbol{p}} \log P(\mathbf{T})$$
(6)



Figure 1: Overview of mDP-DPG. In each iteration, the positive and negative example batches are concatenated with the sampled prompts and input into the LLM to obtain predictions. S_{t+1}^{pos} and S_{t+1}^{neg} represent the sets of sampled positive and negative examples in the *t*-th iteration, respectively. \mathbf{T}_t are prompts sampled from distribution \mathbf{p}_t

where $P(\mathbf{T}) = \prod_{i=1}^{n} P(t_i)$ denotes the joint probability of the prompt \mathbf{T} . Considering $t_i = \mathcal{W}[j_i]$, i.e. the *i*-th token in \mathbf{T} is the j_i -th token in \mathcal{W}^{-1} , and $t_i \sim \operatorname{Cat}(\mathbf{p}_i)$, we can give explicitly $\nabla_{\mathbf{p}} \log P(\mathbf{T})$ as follow (proof in Appendix A),

$$\nabla_{\boldsymbol{p}_{i,j}} \log P(t_i) = \begin{cases} \frac{1}{\boldsymbol{p}_{i,j_i}} & j = j_i \\ 0 & j \neq j_i \end{cases}$$
(7)

However, due to the randomness introduced by prompt sampling, the doubly policy gradient suffers from high variance, which adversely affects convergence, just as in policy gradient estimation (Sutton et al., 1999; Rezende et al., 2014). We implement VR-DPGE, which incorporates a baseline subtraction term to reduce variance (Greensmith et al., 2004). Compared to the high variance of loss values, the difference between the loss value and the baseline term has a lower variance, which facilitates more stable gradient estimation. By defining mini-batch $S = \{(\mathbf{S}_q^{pos}, y_q^{pos}), (\mathbf{S}_q^{neg}, y_q^{neg})\}_{q=1}^B$ and $L_B(h([\mathbf{T}, \cdot]), S) = \frac{1}{B} \Sigma_q L(h([\mathbf{T}, \cdot]), (\mathbf{S}_q^{pos}, y_q^{pos}), (\mathbf{S}_q^{neg}, y_q^{neg}))$, we can replace $\nabla_p F_{i,j}(p)$ by VR-DPGE g_p as equation 8.

$$\boldsymbol{g}_{\boldsymbol{p}} := L_{avg} \boldsymbol{1}_{|\mathcal{W}|} \boldsymbol{1}_{n}^{\top} + \frac{1}{I-1} \Sigma_{\mathbf{k}} \nabla_{\boldsymbol{p}} \log P(\mathbf{T}^{(\mathbf{k})}) (L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) - L_{avg})$$
(8)

with $L_{avg} := \frac{1}{I} \Sigma_{\mathbf{k}} L_B(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S)$. $\mathbf{T}^{(\mathbf{k})}$ represents the k-th prompt obtained from I prompt samplings. $L_{avg} \mathbf{1}_{|W|} \mathbf{1}_n^{\top}$ is the unbiased correction term to ensure unbiasedness of the VR-DPGE and the theoretical guarantee is Lemma 1.

3.3 REDUCE DUAL VARIANCE WITH MOMENTUM TECHNIQUE

The sampling of positive and negative examples also introduces high variance, especially in the case of highly imbalanced data, and the dual variance from this and prompt sampling poses substantial difficulties for imbalanced discrete prompt learning. However, various variance reduction techniques typically require giant batch sizes to compute the gradient in the checkpoint, such as SVRG, which is difficult to apply to pairwise learning because the number of positive and negative sample

¹The vocabulary W is constructed following Diao et al. (2022)

combinations is substantial. Therefore, we further employ the momentum-based variance reduction strategy STORM. Specifically, mDP-DPG uses a variant of momentum m_{t+1} (line 16 in Algorithm 1) as the update direction. The content of Lemma 4 demonstrates that the momentum-based strategy effectively reduces the dual variance. The bound in Lemma 4 contains terms that decay with the momentum parameter θ_t , showing that variance of the moving average m_{t+1} is systematically reduced over iterations, leading to a more stable and accurate approximation of the true gradient.

276 As illustrated in Algorithm 1 and Figure 1, mDP-DPG addresses the class imbalance issue in down-277 stream tasks by minimizing the pairwise AUC surrogate loss. In each iteration, we randomly pick B278 positive and negative sample pairs from $\tilde{\mathcal{D}}$ to compute pairwise loss. To estimate the doubly policy 279 gradient, we sample I prompts according to the prompt distribution p_{t+1} in the current iteration. 280 Specifically, $p_{t+1,i}$ is the *i*-th token distribution of the prompt distribution p_{t+1} updated by momen-281 tum technique in the t-th iteration. For the sampled prompt $\mathbf{T}^{(k)}$, we concatenate it with all positive 282 and negative samples in the mini-batch to construct the input for the LLM. Then, we calculate 283 g_{p_{t+1},S_t} based on LLM predictions and similarly obtain g_{p_t,S_t} . Ultimately, we can determine the update direction m_{t+1} for the next iteration and $\operatorname{proj}_{\mathcal{C}}(\cdot)$ in the update step is a projection function 284 that projects updated prompt distribution to the constraint set C following Diao et al. (2022). 285

Algorithm 1 mDP-DPG

286

287

311

312 313

314

315

316

321

322

288 **Input:** Dataset \hat{D} , hyperparameters k, ξ, c, γ 289 **Initialization:** Construct $S_1 = \{(\mathbf{S}_q^{pos}, y_q^{pos}), (\mathbf{S}_q^{neg}, y_q^{neg})\}_{q=1}^B$ from $\tilde{\mathcal{D}}$ in the same way as Line 290 5-9, then compute $m_1 = g_{p_1, S_1}$. 291 1: for t = 1, ..., T do Learning rate $\eta_t = \frac{k}{(\xi+t)^{1/3}}$ Update $\tilde{p}_{t+1} = \text{proj}_{\mathcal{C}}(p_t - \gamma m_t), p_{t+1} = p_t + \eta_t(\tilde{p}_{t+1} - p_t)$ Compute $\theta_{t+1} = c\eta_t^2$ 292 2: 293 3: 294 4: 295 5: $S_{t+1} = \emptyset$ 296 6: for $q \leq B$ do Sample positive $(\mathbf{S}_q^{pos}, y_q^{pos})$ and negative $(\mathbf{S}_q^{neg}, y_q^{neg})$ examples from $\tilde{\mathcal{D}}$, respectively. $S_{t+1} = S_{t+1} \cup \{(\mathbf{S}_q^{pos}, y_q^{pos}), (\mathbf{S}_q^{neg}, y_q^{neg})\}$ 297 7: 298 8: 299 end for 9: for $\mathbf{k} < I$ do 300 10: Sample $j_1^{(\mathbf{k})} \sim \operatorname{Cat}(\boldsymbol{p}_{t+1,1}), \dots, j_n^{(\mathbf{k})} \sim \operatorname{Cat}(\boldsymbol{p}_{t+1,n})$ $\mathbf{T}^{(\mathbf{k})} = t_1^{(\mathbf{k})} \dots t_n^{(\mathbf{k})} = \mathcal{W}[j_1^{(\mathbf{k})}] \dots \mathcal{W}[j_n^{(\mathbf{k})}]$ 301 11: 302 12: end for 303 13: $L_{avg} = \frac{1}{T} \Sigma_{\mathbf{k}} L_B(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S_{t+1})$ 304 14: $g_{p_{t+1},S_{t+1}} = L_{avg} \mathbf{1}_{|W|} \mathbf{1}_n^\top + \frac{1}{I-1} \Sigma_{\mathbf{k}} \nabla_{p_{t+1}} \log P(\mathbf{T}^{(\mathbf{k})}) (L_B(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S_{t+1}) - L_{avg})$ Compute $m_{t+1} = g_{p_{t+1},S_{t+1}} + (1 - \theta_{t+1})[m_t - g_{p_t,S_{t+1}}]$ 305 15: 306 16: 307 17: end for 308 **Output:** p_R with R chosen uniformly at random in $\{1, \ldots, T\}$. (p_T in practice). 309 310

4 CONVERGENCE ANALYSIS

In this section, we provide theoretical convergence guarantees for the proposed mDP-DPG algorithm. We first introduce some necessary assumptions and definitions. Then, we analyze the convergence properties of mDP-DPG, showing that it achieves an oracle complexity of $O(1/\epsilon^3)$.

Lemma 1 (Unbiasedness of the VR-DPGE estimator, Proof in Appendix B.1). Consider g_p the VR-DPGE policy gradient estimator in equation 8. For any p in the constraint set (i.e. such that each p_i for i in [n] defines a categorical probability distribution), such an estimate is an unbiased estimate of the policy gradient, i.e.:

$$\mathbb{E}_{S}\mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{k=1}^{I}}\boldsymbol{g}_{\boldsymbol{p}} = \nabla_{\boldsymbol{p}}\mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)$$

Assumption 1 (Finiteness of the loss). We assume that there is a constant C > 0 such that for any prompt **T** and any batch of data S, we have $|L_B(h([\mathbf{T}, \cdot]), S)| \le C$.

Remark 1. Note that if one uses a loss $L_B(h([\mathbf{T}, \cdot]), S)$ which could be arbitrary large (for instance if L_B is the cross-entropy function), in practice one can always clip such value to ensure boundedness (indeed, since we consider black-box optimization through reinforcement learning in our paper, even if the clipping operation is non-differentiable, optimization of such loss function will still be possible).

Lemma 2 (Smoothness of the loss, Proof in Appendix B.2). Let us denote C the subset of $\mathbb{R}^{|\mathcal{W}| \times n}$ such that for any $p \in C$, $p_i \in C$ (where p_i denotes a column of p). Let us denote \mathcal{P}_p the probability distribution on prompts \mathbf{T} parameterized by $p \in C$. Then, $\mathbb{E}_S \mathbb{E}_{\mathbf{T} \sim \mathcal{P}_p} L_B(h([\mathbf{T}, \cdot]), S)$ is a smoothfunction of p on its domain, with constant $L = \sqrt{n|\mathcal{W}|}C$, that is, for any $(p, p') \in C^2$:

$$\|\nabla \mathbb{E}_{S} \mathbb{E}_{\mathbf{T} \sim \mathcal{P}_{p}} L_{B}(h([\mathbf{T}, \cdot]), S) - \nabla \mathbb{E}_{S} \mathbb{E}_{\mathbf{T} \sim \mathcal{P}_{p'}} L_{B}(h([\mathbf{T}, \cdot]), S)\| \leq L \|\boldsymbol{p} - \boldsymbol{p}'\|.$$
(9)

Assumption 2 (Boundedness of the variance of the gradient). We assume the following bound on the variance of the VR-DPGE gradient estimator. For any $p \in C$:

$$\mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}\left|\left|oldsymbol{g}_{oldsymbol{p}}-
abla_{oldsymbol{p}}\mathcal{L}(oldsymbol{p})
ight|
ight|^{2}\leq\sigma_{1}^{2}/I+\sigma_{2}^{2}/B,$$

where σ_1 denotes the variance introduced by the random selection of vocabularies and σ_2 denotes the variance introduced by the random selection of pairs of positive-negative examples.

Remark 2 (Proof in Appendix B.3). Even if Assumption 2 is not verified for $C = \{p \in \mathbb{R}^{|\mathcal{W}| \times n} : \forall i \in [n], \|p_i\|_1 = 1, \forall j \in [|\mathcal{W}|], 0 \leq p_{j,i} \leq 1\}$, it is actually verified if one ensures some lower bound on the values of p, i.e. it is verified on the set $C = \{p \in \mathbb{R}^{|\mathcal{W}| \times n} : \forall i \in [n], \|p_i\|_1 = 1, \forall j \in [|\mathcal{W}|], \nu \leq p_{j,i} \leq 1\}$, for some $\nu \in (0, 1]$, as we prove in Appendix B.3. In the experiments however, we could take $\nu = 0$ (i.e. we could keep the original constraint C), which still worked well in practice.

348 349

350

334 335

336

341

4.1 CONVERGENCE RESULTS

Convergence for projected stochastic gradient descent is usually measured in terms of the expected squared norm of the gradient mapping, which we will define below. Since we proved above that the function we consider is smooth and the stochastic gradient is bounded, and the set we project onto is bounded, we can establish the following convergence result:

Theorem 1 (Convergence rate of mDP-DPG, proof in Appendix B.4). Suppose that $\{p_t\}_{t=1}^T$ are generated from mDP-DPG. Let $\eta_t = \frac{k}{(\xi+t)^{1/3}}, 0 < \gamma \leq \min\{\frac{\xi^{1/3}}{2Lk}, \frac{1}{2\sqrt{2}L}\}, c \geq \frac{2}{3k^3} + \frac{5}{4}, \xi \geq \max\{2, k^3, c^3k^3\}$, then we have

359 360

361

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\|G_{\boldsymbol{\mathcal{C}}}(\boldsymbol{p}_t, \gamma)\| \le \frac{2\sqrt{2M}}{T^{1/2}}\xi^{1/6} + \frac{2\sqrt{2M}}{T^{1/3}},\tag{10}$$

362 363 where $G_{\mathcal{C}}(\boldsymbol{p}_t, \eta_t) := \frac{1}{\eta_t} (\boldsymbol{p}_t - proj_{\mathcal{C}}(\boldsymbol{p}_t - \eta_t \nabla \mathcal{L}(\boldsymbol{p}_t)))$ and $M = \frac{\mathcal{L}(\boldsymbol{p}_1) - \mathcal{L}^*}{\gamma k} + \frac{\xi^{1/3} \sigma^2}{k} + 2c^2 \sigma^2 k^3 \log(\xi + T)$, where σ^2 is an abbreviation of the upper bound in Assumption 2.

Remark 3. In order to obtain an ϵ -solution $(\frac{1}{T}\sum_{t=1}^{T} \mathbb{E} \| G_{\mathcal{C}}(\mathbf{p}_t, \gamma) \| \leq \epsilon)$, we need to choose $T = \frac{1}{\epsilon^3} I = O(1)$. Thus the total oracle complexity is $O(\frac{1}{\epsilon^3})$.

Remark 4. The theorems demonstrate that the proposed framework, with variance-reduction techniques and momentum-based updates, ensures convergence towards a prompt distribution that minimizes the empirical risk of the pairwise AUC loss. This implies that the learned prompts are expected to be optimal in terms of performance for the downstream task under the given black-box constraints.

371372373

369

370

5 EXPERIMENTS

374 375

In this section, we present the experiment setups and provide the results and analysis of both the main experiments and ablation studies. Due to space constraints, additional experimental results are provided in the Appendix D.

378 5.1 EXPERIMENT SETUPS379

380 **Datasets.** To evaluate the effectiveness of our methods, we conduct experiments on 4 datasets 381 including 2 widely used datasets from the GLUE benchmark (Wang et al., 2018), CoLA (Warstadt et al., 2018), MRPC (Dolan & Brockett, 2005), and 2 real-world class imbalanced datasets: Amazon 382 books and electronics reviews (McAuley et al., 2015). For GLUE benchmark, we perform down-383 sampling on datasets with a given imbalanced ratio (negative to positive samples ratio) $\tau = 20, 50$ 384 to construct the imbalanced scenarios. For 2 real-world datasets, we down-sample datasets with 385 $\tau = 10$, which facilitates experiments and closely approximates the original imbalance ratio. We 386 use AUC to measure the performance of handling imbalanced data. 387

Backbone Models. We select RoBERTa-large (Liu et al., 2019), GPT2-XL (Radford et al., 2019),
Llama3 (AI@Meta, 2024) as our backbone models and conduct experiments separately. These
models have approximately 355M, 1.5B, and 8B parameters, respectively.

391 Baselines. We compare our proposed methods with the following black-box prompt learning meth-392 ods under the same experimental settings: Manual Prompt performs the zero-shot evaluation on 393 the LLMs with human-written templates, and the results can serve as initial points. **BBT** optimizes the continuous prompts in a random low-dimensional subspace through covariance matrix adapta-394 tion evolution strategy (Sun et al., 2022b). GAP3 introduces a genetic algorithm that considers the 395 prompt as individual and employs auxiliary LLM to generate discrete prompts from the empty (Zhao 396 et al., 2023). BDPL utilizes policy gradients to optimize discrete prompt distribution as mentioned 397 in Section 3.1 (Diao et al., 2022). 398

Implementation Details. The proposed methods and all baselines are implemented using PyTorch
 and experimented on NVIDIA A40 GPUs with 48 GB memory. For all backbone models, we initial ize them with checkpoints provided by the HuggingFace. The details of the input template, output
 label words, and hyperparameters can be found in the Appendix C.

- 403
- 404

415

5.2 MAIN RESULTS AND ANALYSIS

405 Comparison on constructed imbalanced scenarios. We report the average AUC scores on CoLA 406 over 3 random seeds in Table 1. The results on MRPC can be found in Appendix D.1. and our 407 methods exhibit higher performance compared to all baselines. And in many cases, there are sig-408 nificant improvements. In particular, our methods achieve enhancements over BDPL, confirming 409 that minimizing the pairwise AUC loss can effectively address the class imbalance problem. On the 410 other hand, BBT, GAP3, and BDPL have almost no improvement compared to Manual Prompt, and 411 even there are substantial declines in some cases. We attribute this phenomenon to the fact that these 412 baselines do not have additional handling for imbalanced data. For instance, BDPL uses a cross-413 entropy loss function, which in imbalanced scenarios leads to the minority class having almost no 414 effect on the training process (Liu et al., 2020).

Table 1: Comparison of AUC scores (mean±std.) on constructed imbalanced scenarios of CoLA.
We conduct three groups of experiments on pre-trained RoBERTa-large, GPT2-XL, and Llama3 with a prompt length of 20. The best results are highlighted in **bold**.

Imbalanced Ratio	Method	RoBERTa-large	GPT2-XL	Llama3
	Manual Prompt	$.4586 \pm .0947$	$.5224 \pm .0180$	$.4917 \pm .0821$
	BBT	.4797±.1040	$.5000 \pm .0000$.4990±.0063
$\tau = 20$	GAP3	$.5042 \pm .0171$	$.5094 \pm .0162$	$.5089 \pm .0181$
	BDPL	$.4880 \pm .0316$	$.4963 \pm .0253$	$.5193 \pm .0171$
	mDP-DPG (ours)	.5615±.0486	.5271±.0064	.5453±.0906
	Manual Prompt	$.5288 {\pm} .0481$	$.5300 \pm .0017$	$.5289 {\pm} .0501$
	BBT	$.4094 \pm .0472$	$.4938 \pm .0016$	$.5111 \pm .0499$
$\tau = 50$	GAP3	$.4944 \pm .0035$	$.4989 \pm .0019$	$.4983 \pm .0017$
7 = 50	BDPL	$.4871 \pm .0105$	$.5394 \pm .1131$	$.5139 \pm .0580$
	mDP-DPG (ours)	.5700±.0351	.5589±.0139	.5466±.1314

429 430 431

Comparison on real-world imbalanced datasets. We conduct experiments with different prompt lengths and report the average AUC scores over 3 different seeds in Table 2. It should be noted that

Model	Mathad	Bo	ok	E	ec
WIGUEI	wiethou	len = 20	len = 50	len = 20	len = 50
	Manual Prompt	$.8491 {\pm} .0038$	$.8491 {\pm} .0038$	$.8225 \pm .0061$	$.8225 \pm .000$
POBERT ₂	BBT	$.8525 {\pm} .0032$	$.8514 {\pm} .0065$	$.8098 \pm .0172$.8480±.03
lorgo	GAP3	.8372±.0115	.8372±.0115	$.6581 \pm .0239$.6581±.02
-large	BDPL	$.8628 \pm .0066$.8611±.0174	.8431±.0147	.8559±.02
	mDP-DPG (ours)	.8678±.0084	.8623 ± .0047	.8569±.0365	.8588±.02
	Manual Prompt	$.7377 {\pm} .0068$	$.7377 {\pm} .0068$	$.6696 \pm .0544$.6696±.05
	BBT	$.7406 \pm .0133$	$.6078 \pm .0172$	$.6284 \pm .0450$.5196±.01
GPT2-XL	GAP3	$.7785 \pm .0661$	$.7785 \pm .0661$	$.5459 \pm .0270$.5459±.02
	BDPL	$.7884 \pm .0475$	$.7276 \pm .0040$	$.6941 \pm .0256$.7343±.02
	mDP-DPG (ours)	.8721 ±.0297	.7931±.0520	.7157 ±.0384	.7353±.03
	Manual Prompt	$.7502 {\pm} .0072$	$.7502 {\pm} .0072$	$.6549 \pm .0446$.6549±.04
	BBT	$.5164 \pm .0142$	$.5283 \pm .0127$	$.5137 \pm .0221$	$.5275 \pm .01$
Llama3	GAP3	OOM	OOM	OOM	OOM
	BDPL	$.7858 {\pm} .0363$	$.8009 \pm .0179$	$.5216 \pm .0162$.5422±.04
	mDP-DPG (ours)	.8098±.0129	.8151±.0376	.6804±.0814	.6657±.10

Table 2: Comparison of AUC values (mean±std.) on real-world imbalanced datasets Amazon books
 and electronics based on 3 backbone models. *len* represents the prompt length. OOM represents
 out-of-memory.

since GAP3 generates prompts from empty, under our query limit (in Appendix C.2), the lengths of generated prompts are always less than 20. Therefore, the results are the same for maximum prompt lengths of 20 and 50. We have observed that mDP-DPG surpasses all other baselines, which demonstrates our methods remain effective on real-world data. It is worth noting that in some results, mDP-DPG significantly outperforms all baselines, such as the results in the Book with the GPT2-XL. This demonstrates the potential of our method to significantly enhance performance on real-world imbalanced data distributions. Additionally, BBT, GAP3, and BDPL exhibit much poorer performance than Manual Prompt in many results, confirming the deficiencies of these methods in handling imbalanced data.

Table 3: Comparison of AUC values on 3 backbone models. "len" represents the prompt length. τ denotes the imbalanced ratio (negative to positive samples ratio).

Model	Method	CoLA (CoLA (len=20)		Book ($\tau = 10$)	
Mouel	Michiou	$\tau = 20$	$\tau = 50$	len=20	len=50	
POBERT ₂	BDPL-oversample	$.4474 \pm .0681$.4706±.0883	.8541±.0466	$.8479 \pm .0600$	
large	BDPL-reweight	$.5083 {\pm} .0033$	$.4706 \pm .0883$	$.8370 \pm .0096$.8221±.0034	
-large	mDP-DPG(ours)	.5615±.0486	.5700±.0351	.8678±.0084	.8623±.0047	
	BDPL-oversample	$.5172 \pm .0596$	$.5156 \pm .0706$.8171±.0386	.7413±.0731	
GPT2-XL	BDPL-reweight	$.5057 \pm .0213$	$.5428 \pm .0634$	$.8290 \pm .0367$	$.7628 \pm .0138$	
	mDP-DPG(ours)	.5271±.0064	.5589±.0139	.8721±.0297	.7931±.0520	
	BDPL-oversample	$.4943 \pm .0709$	$.5333 \pm .0076$	$.7826 \pm .0533$.7699±.0364	
Llama3	BDPL-reweight	$.5151 \pm .0765$	$.5082 \pm .0467$	$.7973 \pm .0126$.7169±.0471	
	mDP-DPG(ours)	.5453±.0906	.5466±.1314	.8098±.0129	.8151±.0376	

479 Comparison with Simple Techniques. In handling imbalanced data, simple techniques like oversampling minority class samples are among the solutions. To demonstrate the superiority of our proposed methods on imbalanced datasets, we have included such techniques as baselines for comparison. Consequently, we enhance the BDPL approach by incorporating over-sampling and reweighting. We augment the BDPL method as baseline because it formulates black-box prompt learning as a distribution optimization problem and updates the distribution using policy gradients similar to our methods. The results are presented in Table 3. The experimental results demonstrate that our methods outperform simple techniques for handling imbalanced data.

5.3 ABLATION STUDY

Ablation study about the gradient estimator. To lead to a more stable optimization process, we introduce the variance reduction technique and an unbiased correction term into the gradient esti-mator. As shown in Figure 2, we provide the performance comparison figure after removing both components on CoLA ($\tau = 20$) and Amazon books with a prompt length of 20. The VR-DPGE gradient estimator exhibits even stronger performance on 3 backbone models. These results indi-cate that the incorporation of both components in the gradient estimator allows for a more accurate estimation of the gradient.



Figure 2: Ablations of variance reduction and unbiased correction term in VR-DPGE.

Ablation study about the loss function. We incorporate a hinge loss and compare the results with those obtained using the square loss. The results in Figure 3 indicate that the experimental performance generally decreased with the hinge loss. We believe this is because the square loss is statistically consistent with AUC when used as the surrogate loss, whereas the hinge loss does not have this property (Gao & Zhou, 2012).



Figure 3: Ablations of loss function. Hinge loss vs Square loss.

CONCLUSION

In this paper, we propose a momentum-based imbalanced black-box discrete prompt learning frame-work mDP-DPG to handle imbalanced data in downstream tasks. Within this framework, we propose VR-DPGE and introduce the STORM technique for variance reduction to achieve more stable op-timization. We demonstrate the effectiveness mDP-DPG on constructed imbalanced scenarios and real-world imbalanced datasets, showing performance improvements in class imbalance problems. Although the AUC loss in our framework is specifically tailored for binary classification, we discuss in Appendix D.2 how to overcome the limitations of binary classification and provide additional experimental results. In addition, minimizing pairwise AUC loss in our framework suffers from the challenge of constructing sample pairs from opposite classes. Formulating AUC maximization as an equivalent saddle point problem has become dominant in addressing this challenge. However, this technique cannot be directly applied to our problem, as our objective function requires taking the expectation over prompt \mathbf{T} , which would invalidate the existing theoretical derivations. In future work, we will investigate how to introduce it to our framework.

540	REFERENCES
541	REI EREI(CES

AI@Meta. llama3/}	Llama 3 model card. 2024. URL https://github.com/meta-llama/ plob/main/MODEL_CARD.md.
Dainis A. Bor ture for der Calzolari, I (eds.), Pro Language I 2024. ELR	Imber, Fatima Zahra Qachfar, and Rakesh Verma. Domain-agnostic adapter architec- ception detection: Extensive evaluations with the DIFrauD benchmark. In Nicoletta Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue ceedings of the 2024 Joint International Conference on Computational Linguistics, Resources and Evaluation (LREC-COLING 2024), pp. 5260–5274, Torino, Italia, May A and ICCL. URL https://aclanthology.org/2024.lrec-main.468.
Tom Brown, Arvind New few-shot le	Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, elakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are arners. <i>Advances in neural information processing systems</i> , 33:1877–1901, 2020.
Clint Burfoot 2009 Confe 1667633. U	and Timothy Baldwin. Automatic satire detection. In <i>Proceedings of the ACL-IJCNLP erence Short Papers on - ACL-IJCNLP '09</i> , pp. 161, Jan 2009. doi: 10.3115/1667583. JRL http://dx.doi.org/10.3115/1667583.1667633.
Ashok Cutkos Advances i	sky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. <i>n neural information processing systems</i> , 32, 2019.
Mingkai Den Eric P Xin learning. a	g, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, g, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement <i>rXiv preprint arXiv:2205.12548</i> , 2022.
Jacob Devlin bidirection	Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep al transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> , 2018.
Shizhe Diao, Black-box 2022.	Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. prompt learning for pre-trained language models. <i>arXiv preprint arXiv:2201.08531</i> ,
Bill Dolan an Third inter	d Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In <i>national workshop on paraphrasing (IWP2005)</i> , 2005.
Bowen Dong image clas 2022.	Pan Zhou, Shuicheng Yan, and Wangmeng Zuo. Lpt: long-tailed prompt tuning for sification. In <i>The Eleventh International Conference on Learning Representations</i> ,
Cong Fang, Cong Fang, Cong Fang, Cong Fang, Cong Cong Cong Cong Cong Cong Cong Cong	Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex in via stochastic path-integrated differential estimator. <i>Advances in neural information systems</i> , 31, 2018.
Wei Gao and arXiv:1208	Zhi-Hua Zhou. On the consistency of auc pairwise optimization. <i>arXiv preprint</i> 2.0645, 2012.
Wei Gao, Lu Artificial In dx.doi.o	Wang, Rong Jin, Shenghuo Zhu, and Zhi-Hua Zhou. One-pass auc optimization. <i>itelligence</i> , pp. 1–29, Jul 2016. doi: 10.1016/j.artint.2016.03.003. URL http://org/10.1016/j.artint.2016.03.003.
Yang Gao, Ya learning fro	Fan Li, Yu Lin, Charu Aggarwal, and Latifur Khan. Setconv: A new approach for om imbalanced data. <i>arXiv preprint arXiv:2104.06313</i> , 2021.
Gene H Golu	b and Charles F Van Loan. Matrix computations. JHU press, 2013.
Evan Greensr estimates in	nith, Peter L Bartlett, and Jonathan Baxter. Variance reduction techniques for gradient a reinforcement learning. <i>Journal of Machine Learning Research</i> , 5(9), 2004.
Qingyan Guo and Yujiu Y prompt opt	, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, <i>Y</i> ang. Connecting large language models with evolutionary algorithms yields powerful imizers. <i>arXiv preprint arXiv:2309.08532</i> , 2023.

- 594 Chengcheng Han, Liqing Cui, Renyu Zhu, Jianing Wang, Nuo Chen, Qiushi Sun, Xiang Li, and 595 Ming Gao. When gradient descent meets derivative-free optimization: A match made in black-596 box scenario. arXiv preprint arXiv:2305.10013, 2023. 597 J A Hanley and B J McNeil. The meaning and use of the area under a receiver operating character-598 istic (roc) curve. Radiology, pp. 29-36, Apr 1982. doi: 10.1148/radiology.143.1.7063747. URL http://dx.doi.org/10.1148/radiology.143.1.7063747. 600 601 J A Hanley and B J McNeil. A method of comparing the areas under receiver operating char-602 acteristic curves derived from the same cases. Radiology, 148(3):839-843, Sep 1983. doi: 603 10.1148/radiology.148.3.6878708. URL http://dx.doi.org/10.1148/radiology. 604 148.3.6878708. 605 Sophie Henning, William Beluch, Alexander Fraser, and Annemarie Friedrich. A survey of methods 606 for addressing class imbalance in deep-learning based natural language processing. arXiv preprint 607 arXiv:2210.04675, 2022. 608 609 Bairu Hou, Joe O'connor, Jacob Andreas, Shiyu Chang, and Yang Zhang. Promptboosting: Blackbox text classification with ten forward passes. In International Conference on Machine Learning, 610 pp. 13309-13324. PMLR, 2023. 611 612 Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Accelerated zeroth-order and first-order 613 momentum methods from mini to minimax optimization. Cornell University - arXiv, Cornell 614 University - arXiv, Aug 2020. 615 Tomoharu Iwata, Akinori Fujino, and Naonori Ueda. Semi-supervised learning for maximizing the 616 partial auc. Proceedings of the AAAI Conference on Artificial Intelligence, 34(04):4239–4246, 617 Jun 2020. doi: 10.1609/aaai.v34i04.5846. URL http://dx.doi.org/10.1609/aaai. 618 v34i04.5846. 619 620 Thorsten Joachims. A support vector method for multivariate performance measures. In Proceedings 621 of the 22nd international conference on Machine learning - ICML '05, Jan 2005. doi: 10.1145/ 622 1102351.1102399. URL http://dx.doi.org/10.1145/1102351.1102399. 623 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance 624 reduction. Advances in neural information processing systems, 26, 2013. 625 626 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt 627 tuning. arXiv preprint arXiv:2104.08691, 2021. 628 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. arXiv 629 preprint arXiv:2101.00190, 2021. 630 631 Zihao Lin, Yan Sun, Yifan Shi, Xueqian Wang, Lifu Huang, Li Shen, and Dacheng Tao. 632 Efficient federated prompt tuning for black-box large pre-trained models. arXiv preprint 633 arXiv:2310.03123, 2023. 634 Mingrui Liu, Zhenyu Yuan, Yiming Ying, and Tianbao Yang. Stochastic auc maximization with 635 deep neural networks. International Conference on Learning Representations, International Con-636 ference on Learning Representations, Apr 2020. 637 638 Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt 639 understands, too. AI Open, 2023. 640 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike 641 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining 642 approach. arXiv preprint arXiv:1907.11692, 2019. 643 644 Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based rec-645 ommendations on styles and substitutes. In Proceedings of the 38th international ACM SIGIR 646 conference on research and development in information retrieval, pp. 43–52, 2015. 647
 - Yurii Nesterov et al. Lectures on convex optimization, volume 137. Springer, 2018.

Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. Grips: Gradient-free, edit-based in-649 struction search for prompting large language models. arXiv preprint arXiv:2203.07281, 2022. 650 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 651 models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019. 652 653 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yangi 654 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text 655 transformer. Journal of machine learning research, 21(140):1-67, 2020. 656 657 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In International conference on machine learning, 658 pp. 1278-1286. PMLR, 2014. 659 660 Tianxiang Sun, Zhengfu He, Hong Qian, Yunhua Zhou, Xuanjing Huang, and Xipeng Qiu. Bbtv2: 661 Towards a gradient-free future with large language models. arXiv preprint arXiv:2205.11200, 662 2022a. 663 664 Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning 665 for language-model-as-a-service. In International Conference on Machine Learning, pp. 20841– 20855. PMLR, 2022b. 666 667 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-668 ods for reinforcement learning with function approximation. Advances in neural information 669 processing systems, 12, 1999. 670 671 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings* 672 of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for 673 NLP, Jan 2018. doi: 10.18653/v1/w18-5446. URL http://dx.doi.org/10.18653/v1/ 674 w18-5446. 675 676 Alex Warstadt, Amanpreet Singh, and SamuelR. Bowman. Neural network acceptability judgments. 677 arXiv: Computation and Language, arXiv: Computation and Language, May 2018. 678 679 Zeerak Waseem and Dirk Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In Proceedings of the NAACL Student Research Workshop, Jan 2016. 680 doi: 10.18653/v1/n16-2013. URL http://dx.doi.org/10.18653/v1/n16-2013. 681 682 Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduc-683 tion. SIAM Journal on Optimization, 24(4):2057-2075, 2014. 684 685 Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online auc maximization. Neural Information Processing Systems, Neural Information Processing Systems, Dec 2016. 686 687 Zhenyu Yuan, Yan Yan, Milan Sonka, and Tianbao Yang. Robust deep auc maximization: A new 688 surrogate loss and empirical studies on medical image classification. arXiv: Learning, arXiv: 689 Learning, Dec 2020. 690 691 Zhenyu Yuan, Zhishuai Guo, Yi Xu, Yiming Ying, and Tianbao Yang. Federated deep auc maximiza-692 tion for hetergeneous data with a constant communication complexity. International Conference 693 on Machine Learning, International Conference on Machine Learning, Jul 2021. 694 Jiangjiang Zhao, Zhuoran Wang, and Fangchun Yang. Genetic prompt search via exploiting lan-695 guage model probabilities. In Proceedings of the Thirty-Second International Joint Conference 696 on Artificial Intelligence, pp. 5296–5305. IJCAI, 2023. 697 Peilin Zhao, Rong Jin, Tianbao Yang, and StevenC.H. Hoi. Online auc maximization. International 699 Conference on Machine Learning, International Conference on Machine Learning, Jun 2011. 700 Yuanhang Zheng, Zhixing Tan, Peng Li, and Yang Liu. Black-box prompt tuning with subspace 701 learning. arXiv preprint arXiv:2305.03518, 2023.

702 703 704	Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic nested variance reduced gradient descent for nonconvex optimization. <i>Neural Information Processing Systems, Neural Information Processing Systems</i> , Jan 2018.
705	
706	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan,
707	and Jimmy Ba. Large language models are human-level prompt engineers. arXiv preprint
708	arXiv:2211.01910, 2022.
709	
710	
711	
712	
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

APPENDIX

PROOFS FOR SECTION 3.2 А

We now prove the explicit derivation for the policy gradient estimator. The j-th component of $\nabla_{\boldsymbol{p}_i} \log P(t_i)$ is:

$$\nabla_{\boldsymbol{p}_{i,j}} \log P(t_i) = \nabla_{\boldsymbol{p}_{i,j}} \log \boldsymbol{p}_{i,j_i}.$$

When $j = j_i$, we have

$$\nabla_{\boldsymbol{p}_{i,j}} \log P(t_i) = \nabla_{\boldsymbol{p}_{i,j}} \log \boldsymbol{p}_{i,j_i} = \frac{\nabla_{\boldsymbol{p}_{i,j}} \boldsymbol{p}_{i,j_i}}{\boldsymbol{p}_{i,j_i}}|_{j=j_i} = \frac{1}{\boldsymbol{p}_{i,j_i}}$$

When $j \neq j_i$, we have

$$\nabla_{\boldsymbol{p}_{i,j}} \log P(t_i) = \nabla_{\boldsymbol{p}_{i,j}} \log \boldsymbol{p}_{i,j_i} = \frac{\nabla_{\boldsymbol{p}_{i,j}} \boldsymbol{p}_{i,j_i}}{\boldsymbol{p}_{i,j_i}}|_{j \neq j_i} = 0.$$

PROOFS FOR SECTION 4.1 В

B.1 PROOF OF LEMMA 1

Proof. First, it is easy to show that the expectation of L_{avg} over the sampling of the I prompts is equal to the expected loss with expectation taken over the distribution over prompts:

$$\mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} L_{avg} \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top} = \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \frac{1}{I} \Sigma_{\mathbf{k}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top}$$

$$= \frac{1}{I} \Sigma_{\mathbf{k}} \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top}$$

$$= \frac{1}{I} \Sigma_{\mathbf{k}} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S) \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top} = \frac{I}{I} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S) \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top}$$

$$= \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S) \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top}$$

$$(11)$$
Now, let us analyze the full expression for the expectation of \boldsymbol{q}_{n} :

Now, let us analyze the full expression for the expectation of g_p :

$$\mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \boldsymbol{g}_{\boldsymbol{p}} = \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \left(L_{avg} \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top} + \frac{1}{I-1} \Sigma_{k} \nabla \log P(\mathbf{T}^{(\mathbf{k})}) (L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) - L_{avg})) \right)$$

$$= \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} L_{avg} \mathbf{1}_{|\mathcal{W}|} \mathbf{1}_{n}^{\top} + \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \underbrace{\left(\frac{1}{I-1} \Sigma_{k} \nabla \log P(\mathbf{T}^{(\mathbf{k})}) (L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) - L_{avg})\right)}_{(\mathbf{A})}$$

$$(12)$$

We can rewrite (A) as follows:

812 Now, first, we have:

$$\mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \left[(1) \right] = \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})})$$
$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \frac{\nabla_{p} P(\mathbf{T}^{(\mathbf{k})})}{P(\mathbf{T}^{(\mathbf{k})})}$$

 $= \nabla_{\boldsymbol{p}} \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \sum_{\mathbf{T}^{(\mathbf{k})}} P(\mathbf{T}^{(\mathbf{k})}) L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S)$

 $= \nabla_{p} \mathbb{E}_{S} \frac{I}{I} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S) = \nabla_{p} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S)$

(13)

 $= \nabla_{\boldsymbol{p}} \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S)$

$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \frac{\mathbf{P}^{I}(\mathbf{x})}{P(\mathbf{T}^{(\mathbf{k})})}$$
$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{T}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \frac{\nabla_{\mathbf{p}} P(\mathbf{T}^{(\mathbf{k})})}{P(\mathbf{T}^{(\mathbf{k})})}$$

$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \sum_{\mathbf{T}^{(\mathbf{k})}}^{I} P(\mathbf{T}^{(\mathbf{k})}) L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \frac{\nabla_{\mathbf{p}} P(\mathbf{T}^{(\mathbf{k})})}{P(\mathbf{T}^{(\mathbf{k})})}$$
$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \sum_{\mathbf{T}^{(\mathbf{k})}}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{k})}, \cdot]), S) \nabla_{\mathbf{p}} P(\mathbf{T}^{(\mathbf{k})})$$

Then, we have:

$$\begin{split} \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \left[\left[\widehat{\mathbf{2}} \right] &= \mathbb{E}_{S} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j} \neq \mathbf{k}}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}} \frac{1}{I-1} \sum_{\mathbf{j} \neq \mathbf{k}}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{j})}\}_{\mathbf{j}=1,\mathbf{j}\neq \mathbf{k}}^{I}} \frac{1}{I-1} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\{\mathbf{T}^{(\mathbf{j})}\}_{\mathbf{j}=1,\mathbf{j}\neq \mathbf{k}}^{I} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \nabla_{p} \log P(\mathbf{T}^{(\mathbf{k})}) \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \mathbb{E}_{\mathbf{T}^{(\mathbf{k})}} \frac{\nabla_{p} P(\mathbf{T}^{(\mathbf{k})})}{P(\mathbf{T}^{(\mathbf{k})})} \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \sum_{\mathbf{T}^{(\mathbf{k})}} P(\mathbf{T}^{(\mathbf{k})}) \frac{\nabla_{p} P(\mathbf{T}^{(\mathbf{k})})}{P(\mathbf{T}^{(\mathbf{k})})} \\ &= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j}\neq \mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \sum_{\mathbf{T}^{(\mathbf{k})}} \nabla_{p} P(\mathbf{T}^{(\mathbf{k})}) \end{array}$$

$$= \mathbb{E}_{S} \frac{1}{I} \sum_{\mathbf{k}=1}^{I} \frac{1}{I-1} \sum_{\mathbf{j}\neq\mathbf{k}}^{I} \mathbb{E}_{\mathbf{T}^{(\mathbf{j})}} L_{B}(h([\mathbf{T}^{(\mathbf{j})}, \cdot]), S) \nabla_{\mathbf{p}} \sum_{\mathbf{T}^{(\mathbf{k})}} P(\mathbf{T}^{(\mathbf{k})})$$
(14)

Now, for any $i \in [n]$, we have:

$$\begin{split} \nabla_{\boldsymbol{p}_{i}} \sum_{\mathbf{T}^{(\mathbf{k})}} P(\mathbf{T}^{(\mathbf{k})})) &= \nabla_{\boldsymbol{p}_{i}} \sum_{t_{1}^{(\mathbf{k})} \in \mathcal{W}} \dots \sum_{t_{i}^{(\mathbf{k})} \in \mathcal{W}} \dots \sum_{t_{n}^{(\mathbf{k})} \in \mathcal{W}} P(t_{1}^{(\mathbf{k})}) \dots P(t_{i}^{(\mathbf{k})}) \dots P(t_{n}^{(\mathbf{k})}) \\ &= \nabla_{\boldsymbol{p}_{i}} \sum_{t_{i}^{(\mathbf{k})} \in \mathcal{W}} P(t_{i}^{(\mathbf{k})}) \left(\sum_{t_{1}^{(\mathbf{k})} \in \mathcal{W}} P(t_{1}^{(\mathbf{k})}) \left(\sum_{t_{2}^{(\mathbf{k})} \in \mathcal{W}} P(t_{2}^{(\mathbf{k})}) \left(\dots \left(\sum_{t_{i-1}^{(\mathbf{k})} \in \mathcal{W}} P(t_{n}^{(\mathbf{k})}) \right) \right) \right) \right) \right) \\ &= \nabla_{\boldsymbol{p}_{i}} \sum_{t_{i}^{(\mathbf{k})} \in \mathcal{W}} P(t_{i}^{(\mathbf{k})}) = \sum_{t_{i}^{(\mathbf{k})} \in \mathcal{W}} \nabla_{\boldsymbol{p}_{i}} P(t_{i}^{(\mathbf{k})}) = \sum_{t_{i}^{(\mathbf{k})} \in \mathcal{W}} \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \\ \leftarrow \text{ index } j \text{ s.t. the } j \text{-th word from } \mathcal{W} \text{ is } t_{i}^{(\mathbf{k})} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{bmatrix} \end{split}$$

Now, plugging the result above into 14, we obtain:

$$\mathbb{E}_{S}\mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}}\left[2\right] = \mathbb{E}_{S}\frac{1}{I}\sum_{\mathbf{k}=1}^{I}\frac{1}{I-1}\sum_{\mathbf{j}\neq\mathbf{k}}^{I}\mathbb{E}_{\mathbf{T}^{(\mathbf{j})}}L_{B}(h([\mathbf{T}^{(\mathbf{j})},\cdot]),S)\mathbf{1}_{|\mathcal{W}|}\mathbf{1}_{n}^{\top},$$

Continuing from 14, since the term in the sum in 14 is a constant (as for all j, the $\mathbf{T}_{j}^{(\mathbf{k})}$ are sampled i.i.d):

$$\mathbb{E}_{S}\mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}}\left[2\right] = \mathbb{E}_{S}\frac{I}{I}\frac{I-1}{I-1}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)\mathbf{1}_{|\mathcal{W}|}\mathbf{1}_{n}^{\top}$$
$$= \mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)\mathbf{1}_{|\mathcal{W}|}\mathbf{1}_{n}^{\top}$$
(15)

Therefore, plugging 11, 13 and 15 into 12, we obtain:

$$\mathbb{E}_{S}\mathbb{E}_{\{\mathbf{T}^{(\mathbf{k})}\}_{\mathbf{k}=1}^{I}}\boldsymbol{g}_{\boldsymbol{p}} = \mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)\mathbf{1}_{|\mathcal{W}|}\mathbf{1}_{n}^{\top} + \nabla_{\boldsymbol{p}}\mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)$$
$$- \mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)\mathbf{1}_{|\mathcal{W}|}\mathbf{1}_{n}^{\top}$$
$$= \nabla_{\boldsymbol{p}}\mathbb{E}_{S}\mathbb{E}_{\mathbf{T}}L_{B}(h([\mathbf{T},\cdot]),S)$$

B.2 PROOF OF LEMMA 2

Proof. Let $p \in C$. Let us denote by $\mathbb{E}_T := \mathbb{E}_{T \sim \mathcal{P}_p}$ for simplicity. We can express the gradient estimate as follows:

$$\nabla_{\boldsymbol{p}} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S) = \nabla_{\boldsymbol{p}} \mathbb{E}_{S} \mathbb{E}_{j_{1} \sim \operatorname{Cat}(\boldsymbol{p}_{1}), \dots, j_{n} \sim \operatorname{Cat}(\boldsymbol{p}_{n})} L_{B}(h([t_{j_{1}}, \dots, t_{j_{n}}, \cdot]), S))$$
$$= \mathbb{E}_{S} \nabla_{\boldsymbol{p}} \sum_{j_{1} \in [|\mathcal{W}|]} \dots \sum_{j_{n} \in [|\mathcal{W}|]} L_{B}(h([t_{j_{1}}, \dots, t_{j_{n}}, \cdot]), S) \Pi_{i=1}^{n} \boldsymbol{p}_{j_{i}, i})$$
$$= \mathbb{E}_{S} \sum_{j_{1} \in [|\mathcal{W}|]} \dots \sum_{j_{n} \in [|\mathcal{W}|]} L_{B}(h([t_{j_{1}}, \dots, t_{j_{n}}, \cdot]), S) \nabla_{\boldsymbol{p}} \Pi_{i=1}^{n} \boldsymbol{p}_{j_{i}, i})$$

915 Therefore:

917 $\nabla_{\boldsymbol{p}_{k,l}} \mathbb{E}_S \mathbb{E}_{\mathbf{T}} L_B(h([\mathbf{T},\cdot]),S) = \mathbb{E}_S \sum_{j_1 \in [|\mathcal{W}|]} \dots \sum_{j_n \in [|\mathcal{W}|]} L_B(h([t_{j_1},\dots,t_{j_n},\cdot]),S) \nabla_{\boldsymbol{p}_{k,l}} \prod_{i=1}^n \boldsymbol{p}_{j_i,i}$ 918 Now, we have:

$$\nabla_{\boldsymbol{p}_{k,l}} \Pi_{i=1}^{n} \boldsymbol{p}_{j_{i},i} = \begin{cases} \Pi_{i=1,i\neq l}^{n} \boldsymbol{p}_{j_{i},i} \text{ if } j_{l} = k\\ 0 \text{ otherwise} \end{cases}$$

Therefore:

$$\nabla_{\boldsymbol{p}_{k,l}} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T}, \cdot]), S)$$

$$= \mathbb{E}_{S} \sum_{j_{1} \in [|\mathcal{W}|]} \dots \sum_{j_{n} \in [|\mathcal{W}|]} L_{B}(h([t_{j_{1}}, \dots, t_{j_{n}}, \cdot]), S) \nabla_{\boldsymbol{p}_{k,l}} \Pi_{i=1}^{n} \boldsymbol{p}_{j_{i},i}$$

$$= \mathbb{E}_{S} \sum_{j_{1} \in [|\mathcal{W}|]} \dots \sum_{j_{l-1} \in [|\mathcal{W}|]} \sum_{j_{l+1} \in [|\mathcal{W}|]} \dots \sum_{j_{n} \in [|\mathcal{W}|]} L_{B}(h([t_{j_{1}}, \dots, t_{l-1}, t_{k}, t_{l+1}, \dots, t_{j_{n}}, \cdot]), S) \Pi_{i=1, i \neq l}^{n} \boldsymbol{p}_{j_{i},i}$$

Note that the last expression above can be expressed as an expectation, in the case where each column of p defines a probability distribution:

$$\nabla_{\boldsymbol{p}_{k,l}} \mathbb{E}_{S} \mathbb{E}_{\mathbf{T}} L_{B}(h([\mathbf{T},\cdot]),S) = \mathbb{E}_{S} \mathbb{E}_{j_{1},\dots,j_{l-1},j_{l+1}\dots,j_{n}} L_{B}(h([t_{j_{1}},\dots,t_{l-1},t_{k},t_{l+1},\dots,t_{j_{n}},\cdot]),S)$$

Similarly, we can compute the Hessian of such cost function:

$$\frac{\partial^2}{\partial \boldsymbol{p}_{k,l} \partial \boldsymbol{p}_{m,q}} \mathbb{E}_S \mathbb{E}_{\mathbf{T}} L_B(h([\mathbf{T}, \cdot]), S) \\
= \mathbb{E}_S \frac{\partial}{\partial \boldsymbol{p}_{m,q}} \left[\sum_{j_1 \in [|\mathcal{W}|]} \dots \sum_{j_{l-1} \in [|\mathcal{W}|]} \sum_{j_{l+1} \in [|\mathcal{W}|]} \dots \sum_{j_n \in [|\mathcal{W}|]} L_B(h([t_{j_1}, \dots, t_{l-1}, t_k, t_{l+1}, \dots, t_{j_n}, \cdot]), S) \Pi_{i=1, i \neq l}^n \boldsymbol{p}_{j_i, i} \right] \\
= \mathbb{E}_S \left[\sum_{j_1 \in [|\mathcal{W}|]} \dots \sum_{j_{l-1} \in [|\mathcal{W}|]} \sum_{j_{l+1} \in [|\mathcal{W}|]} \dots \sum_{j_n \in [|\mathcal{W}|]} L_B(h([t_{j_1}, \dots, t_{j_n}, \cdot]), S) \frac{\partial}{\partial \boldsymbol{p}_{m,q}} \Pi_{i=1, i \neq l}^n \boldsymbol{p}_{j_i, i} \right]$$

Now, similarly as before, we have:

$$\frac{\partial}{\partial \boldsymbol{p}_{m,q}} \prod_{i=1, i \neq l}^{n} \boldsymbol{p}_{j_i, i} = \begin{cases} \prod_{i=1, i \neq l, i \neq k}^{n} \boldsymbol{p}_{j_i, i} \text{ if } j_q = m \text{ and } q \neq l \\ 0 \text{ otherwise} \end{cases}$$

Therefore, the last expression above can be expressed as an expectation, if each column of p defines a probability distribution:

$$\begin{split} &\frac{\partial^2}{\partial p_{k,l}\partial p_{m,n}} \mathbb{E}_S \mathbb{E}_{\mathbf{T}} L_B(h([\mathbf{T},\cdot]),S) \\ &= \begin{cases} \mathbb{E}_S \mathbb{E}_{j_1,\dots,j_{l-1},j_{l+1}\dots,j_q-1}, L_B(h([t_{j_1},\dots,t_{j_l-1},t_k,t_{j_l+1},\dots,t_{j_q-1},t_m,t_{j_q+1},\dots,t_{j_n},\cdot],S)) \text{ if } q \neq l \\ 0 \text{ otherwise.} \end{cases} \end{split}$$

Therefore, using Assumption 1, we have that:

$$\frac{\partial^2}{\partial \boldsymbol{p}_{k,l} \partial \boldsymbol{p}_{m,n}} \mathbb{E}_S \mathbb{E}_{\mathbf{T}} L_B(h([\mathbf{T},\cdot]),S) \leq C$$

which implies, with $H(\mathbf{p})$ denoting the Hessian of $\mathbb{E}_S \mathbb{E}_{\mathbf{T}} L_B(h([\mathbf{T}, \cdot]), S)$ with respect to \mathbf{p} :

$$\begin{split} ||H(p)||_{F} \leq \sqrt{n}|W|C^{2} = \sqrt{n}|W|C \\ \text{And therefore we have :} \\ ||H(p)||_{F} \leq ||H(p)||_{F} \leq \sqrt{n}|W|C, \\ \text{using (2.3.7) in Golub & Van Loan (2013). We can now use Lemma 1.2.2 in Nesterov et al. (2018) to relate such bound on the Hessian to the smoothness constant of $\mathbb{E}_{S}\mathbb{P}_{T}L_{H}(h([T, \cdot]), S). \\ \\ \text{B.3. PROOF OF BOUNDED VARIANCE (REMARK 2) \\ Proof. Consider the following constraints set: $C = \{p \in \mathbb{R}^{|W| \times n} : \forall i \in [n], ||p_{i}||_{1} = 1, \forall j \in [|W|], \nu \leq p_{j,i} \leq 1\}, \text{ for some } \nu \in (0, 1]. \text{ For simplicity, we denote } L(\mathbf{T}^{(k)}) := \\ L_{H}(h([\mathbf{T}^{(k)}, \cdot]), S), \text{ and denote } e_{i} = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ \vdots \\ 0 \end{bmatrix} = -i. \text{ For any } i \in [n], \text{ we have:} \\ \text{Var}(g_{p_{i}}) = \text{Var}\left(L_{avg}\mathbf{1}_{|W|} + \frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right) \\ = \text{Var}\left(L_{avg}\mathbf{1}_{|W|} + \frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right) \\ \leq \mathbb{E}\left\|L_{avg}\mathbf{1}_{|W|}\right\|^{2} + \mathbb{E}\left\|\frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ = \mathbb{E}\|L_{avg}\mathbf{1}_{|W|}\|^{2} + \mathbb{E}\left\|\frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\mathbb{E}\langle L_{avg}\mathbf{1}_{|W|}\right|^{2} + \frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\mathbb{E}\langle L_{avg}\mathbf{1}_{|W|}\right\|^{2} + \frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\mathbb{E}\langle L_{avg}\mathbf{1}_{|W|}\right\|^{2} + \frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\frac{1}{I-1}\sum_{k}\left[\mathbb{E}L_{avg}\left(\frac{L(\mathbf{T}^{(k)}) - L_{avg}}{P(t_{i}^{(k)})}\right)(1_{|W|}, e_{j}^{(v)})\right] \\ = \|W\|_{E_{T}}\|L(\mathbf{T})\|^{2} + \frac{|W|(I-1)}{I}\|\mathbb{E}_{T}L(\mathbf{T})\|^{2} + \mathbb{E}\left\|\left|\frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\frac{1}{I-1}\sum_{k}\left[\mathbb{E}L_{avg}\left(\frac{L(\mathbf{T}^{(k)}) - L_{avg}}{P(t_{i}^{(k)})}\right)\right] \\ = \|W\|_{E_{T}}\|L(\mathbf{T})\|^{2} + \frac{|W|(I-1)}{I}\|\mathbb{E}_{T}L(\mathbf{T})\|^{2} + \mathbb{E}\left\|\left|\frac{1}{I-1}\sum_{k}\frac{e_{j_{i}(k)}}{P(t_{i}^{(k)})}(L(\mathbf{T}^{(k)}) - L_{avg})\right\|^{2} \\ + 2\frac{1}{I-1}\sum_{k}\left[\mathbb{E}L_{avg}\left(\frac{L(\mathbf{T}^{(k)}) - L$$$$

$$\begin{split} &= \frac{|\mathcal{W}|}{I} \mathbb{E}_{\mathbf{T}} \|L(\mathbf{T})\|^{2} + \frac{|\mathcal{W}|(I-1)}{I} \|\mathbb{E}_{\mathbf{T}} L(\mathbf{T})\|^{2} + \mathbb{E} \left\| \frac{1}{I-1} \Sigma_{k} \frac{e_{j_{k}^{(k)}}}{P(t_{k}^{(k)})} (L(\mathbf{T}^{(k)}) - L_{avg}) \right\|^{2} \\ &+ 2 \frac{1}{I-1} \Sigma_{k} \left[\mathbb{E}_{t_{k}^{(k)}, \dots, t_{k-1}^{(k)}, t_{k+1}^{(k)}, \dots, t_{k}^{(k)}} \sum_{t_{k}^{(k)} \in \mathcal{W}} L_{avg} \left(L(\mathbf{T}^{(k)}) - L_{avg} \right) \right] \\ &\leq \frac{|\mathcal{W}|}{I} C^{2} + \frac{|\mathcal{W}|(I-1)}{I} C^{2} + \mathbb{E} \left\| \frac{1}{I-1} \Sigma_{k} \frac{e_{j_{k}^{(k)}}}{P(t_{k}^{(k)})} (L(\mathbf{T}^{(k)}) - L_{avg}) \right\|^{2} + 2 \frac{I}{I-1} |\mathcal{W}| 2C^{2} \\ &+ \frac{1}{(I-1)^{2}} \left(\sum_{k} \mathbb{E} \left\| \frac{e_{j_{k}^{(k)}}}{P(t_{k}^{(k)})} (L(\mathbf{T}^{(k)}) - L_{avg}) \right\|^{2} \\ &+ \sum_{\mathbf{1}} \sum_{\mathbf{m}, \mathbf{m} \neq \mathbf{I}} \mathbb{E} \left\{ \frac{e_{j_{k}^{(k)}}}{P(t_{k}^{(1)})} (L(\mathbf{T}^{(1)}) - L_{avg}), \frac{e_{j_{k}^{(m)}}}{P(t_{k}^{(m)})} (L(\mathbf{T}^{(m)}) - L_{avg}) \right\} \right\} + 2 \frac{I}{I-1} |\mathcal{W}| 2C^{2} \\ &\leq \frac{|\mathcal{W}|}{I} C^{2} + \frac{|\mathcal{W}|(I-1)}{I} C^{2} + \frac{1}{(I-1)^{2}} \left(\sum_{\mathbf{k}} \mathbb{E} \left(\frac{1}{P(t_{k}^{(m)})} (L(\mathbf{T}^{(m)}) - L_{avg}) \right) \right)^{2} \\ &+ \mathbb{E} \sum_{\mathbf{1}} \sum_{\mathbf{m}, \mathbf{m} \neq \mathbf{I}} \left\| \frac{e_{j_{k}^{(1)}}}{P(t_{k}^{(1)})} (L(\mathbf{T}^{(1)}) - L_{avg}) \right\| \left\| \frac{e_{j_{k}^{(m)}}}{P(t_{k}^{(m)})} (L(\mathbf{T}^{(m)}) - L_{avg}) \right\| \right) + 2 \frac{I}{I-1} |\mathcal{W}| 2C^{2} \\ &= \frac{|\mathcal{W}|}{I} C^{2} + \frac{|\mathcal{W}|(I-1)}{I} C^{2} + \frac{1}{(I-1)^{2}} \left(\sum_{\mathbf{k}} \mathbb{E} \left\{ t_{j_{k}^{(1)}} \right\}_{j=1, t_{k}^{(1)}} \mathbb{E} \left\{ t_$$

Where (a) and (b) follow from the fact that for some random variable X: $Var(X) = \mathbb{E}||X - \mathbb{E}X||^2 = \mathbb{E}[||X||^2] - ||\mathbb{E}[X]||^2 \le \mathbb{E}[||X||^2]$, (c) follows from Assumption 1 (which implies that $|L(\mathbf{T})| \le C$ for all \mathbf{T} and also consequently that $|L_{avg}| \le C$), and (d) follows from the Cauchy-Schwarz inequality. Therefore, for any $i \in [n]$, $Var(g_{p_i})$ is indeed bounded, and consequently, the final gradient estimator is also bounded.

B.4 PROOFS FOR THEOREM 1

Lemma 3. Let $\{p\}_{t=1}^{T}$ be generated by mDP-DPG. Let $\eta_t \in (0, 1]$ and $\gamma \in (0, \frac{1}{2L\eta_t}]$, we have

$$\mathcal{L}(\boldsymbol{p}_{t+1}) \leq \mathcal{L}(\boldsymbol{p}_t) + \eta_t \gamma \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 - \frac{\eta_t}{2\gamma} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2.$$
(16)

Proof. Recall that $\mathcal{L}(\boldsymbol{p}) := \mathbb{E}_{S} \mathbb{E}_{\mathbf{T} \sim \mathcal{P}_{\boldsymbol{p}}} L_{B}(h([\mathbf{T}, \cdot]), S)$. According to Lemma 2, $\mathcal{L}(\boldsymbol{p})$ is L-smooth. Then we have

$$\begin{split} \mathcal{L}(\boldsymbol{p}_{t+1}) \leq & \mathcal{L}(\boldsymbol{p}_t) + \langle \nabla \mathcal{L}(\boldsymbol{p}_t), \boldsymbol{p}_{t+1} - \boldsymbol{p}_t \rangle + \frac{L}{2} \|\boldsymbol{p}_{t+1} - \boldsymbol{p}_t\|^2 \\ \leq & \mathcal{L}(\boldsymbol{p}_t) + \eta_t \left\langle \nabla \mathcal{L}(\boldsymbol{p}_t), \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t \right\rangle + \frac{L\eta_t^2}{2} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 \\ \leq & \mathcal{L}(\boldsymbol{p}_t) + \eta_t \left\langle \nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t, \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t \right\rangle + \eta_t \left\langle \boldsymbol{m}_t, \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t \right\rangle + \frac{L\eta_t^2}{2} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 \end{split}$$

Since $\tilde{p}_{t+1} = proj_{\mathcal{C}}(p_t - \gamma m_t) = \arg \min_{p \in \mathcal{C}} \frac{1}{2} \|p - (p_t - \gamma m_t)\|^2$, we have $\forall p \in \mathcal{C}$, $\langle \tilde{p}_{t+1} - (p_t - \gamma m_t), p - \tilde{p}_{t+1} \rangle \geq 0$. Set $p = p_t$, we have

$$\langle oldsymbol{m}_t, oldsymbol{ ilde{p}}_{t+1} - oldsymbol{p}_t
angle \leq -rac{1}{\gamma} \|oldsymbol{ ilde{p}}_{t+1} - oldsymbol{p}_t\|^2$$

Thus we have

$$\begin{aligned} \mathcal{L}(p_{t+1}) \\ \mathcal{L}(p_{t+1}) \\ \leq \mathcal{L}(p_t) + \eta_t \langle \nabla \mathcal{L}(p_t) - m_t, \tilde{p}_{t+1} - p_t \rangle + \eta_t \langle m_t, \tilde{p}_{t+1} - p_t \rangle + \frac{L\eta_t^2}{2} \|\tilde{p}_{t+1} - p_t\|^2 \\ \leq \mathcal{L}(p_t) + \eta_t \gamma \|\nabla \mathcal{L}(p_t) - m_t\|^2 + \frac{\eta_t}{4\gamma} \|\tilde{p}_{t+1} - p_t\|^2 - \frac{\eta_t}{\gamma} \|\tilde{p}_{t+1} - p_t\|^2 + \frac{L\eta_t^2}{2} \|\tilde{p}_{t+1} - p_t\|^2 \\ \leq \mathcal{L}(p_t) + \eta_t \gamma \|\nabla \mathcal{L}(p_t) - m_t\|^2 - \frac{\eta_t}{2\gamma} \|\tilde{p}_{t+1} - p_t\|^2 - (\frac{\eta_t}{4\gamma} - \frac{L\eta_t^2}{2}) \|\tilde{p}_{t+1} - p_t\|^2 \\ \leq \mathcal{L}(p_t) + \eta_t \gamma \|\nabla \mathcal{L}(p_t) - m_t\|^2 - \frac{\eta_t}{2\gamma} \|\tilde{p}_{t+1} - p_t\|^2. \end{aligned}$$
Where the last inequality holds due to $0 < \gamma < \frac{1}{2L\eta_t}. \end{aligned}$

Lemma 4.

11.9 a 02 - 2

$$\mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{m}_{t+1}\|^{2}$$

$$\leq (1 - \theta_{t})^{2} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_{t}) - \boldsymbol{m}_{t}\|^{2} + 2(1 - \theta_{t})^{2} L^{2} \eta_{t}^{2} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_{t}\|^{2} + 2\theta_{t}^{2} \sigma^{2}.$$
(18)

Proof. According to the update rule of m_{t+1} , we have

 $\mathbb{E} \| \nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{m}_{t+1} \|^2$ $= \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{g}_{\boldsymbol{p}_{t+1}, S_{t+1}} - (1 - \theta_{t+1})(\boldsymbol{m}_t - \boldsymbol{g}_{\boldsymbol{p}_t, S_{t+1}})\|^2$ $= \mathbb{E} \| (1 - \theta_{t+1}) (\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t) + \theta_t (\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{g}_{\boldsymbol{p}_{t+1}, S_{t+1}}) \|$ $+ (1 - \theta_{t+1}) (\nabla \mathcal{L}(p_{t+1}) - \nabla \mathcal{L}(p_t) - (g_{p_{t+1},S_{t+1}} - g_{p_t,S_{t+1}})) \|^2$ $\leq (1-\theta_t)^2 \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2(1-\theta_t)^2 \|\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \nabla \mathcal{L}(\boldsymbol{p}_t) - (\boldsymbol{g}_{\boldsymbol{p}_{t+1},S_{t+1}} - \boldsymbol{g}_{\boldsymbol{p}_t,S_{t+1}})\|^2$ + $2\theta_t^2 \mathbb{E} \| \nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{g}_{\boldsymbol{p}_{t+1},S_{t+1}} \|^2$ $\leq (1-\theta_t)^2 \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2(1-\theta_t)^2 \|\boldsymbol{g}_{\boldsymbol{p}_{t+1},S_{t+1}} - \boldsymbol{g}_{\boldsymbol{p}_t,S_{t+1}}\|^2 + 2\theta_t^2 \sigma^2$ $\leq (1-\theta_t)^2 \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2(1-\theta_t)^2 L^2 \|\boldsymbol{p}_{t+1} - \boldsymbol{p}_t\|^2 + 2\theta_t^2 \sigma^2$ $= (1 - \theta_t)^2 \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2(1 - \theta_t)^2 L^2 \eta_t^2 \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + 2\theta_t^2 \sigma^2.$

(17)

Let $\eta_t = \frac{k}{(\xi+t)^{1/3}}$, we have

$$\begin{aligned} \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} &= \frac{1}{k} ((\xi+t)^{1/3} - (\xi+t-1)^{1/3}) \le \frac{1}{3k(\xi+t-1)^{2/3}} \le \frac{1}{3k(\xi/2+t)^{2/3}} \\ &\le \frac{2^{2/3}}{3k(\xi+t)^{2/3}} = \frac{2^{2/3}}{3k^3} \eta_t^2 \le \frac{2}{3k^3} \eta_t, \end{aligned}$$

 $\leq (\frac{(1-\theta_t)^2}{\eta_t} - \frac{1}{\eta_{t-1}}) \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2(1-\theta_t)^2 L^2 \eta_t \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + \frac{2\theta_t^2 \sigma^2}{\eta_t}$

 $\leq (\frac{1-\theta_t}{n_t} - \frac{1}{n_{t-1}})\mathbb{E}\|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2L^2\eta_t\|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + \frac{2\theta_t^2\sigma^2}{n_t}$

 $= (\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - c\eta_t) \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + 2L^2 \eta_t \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + \frac{2\theta_t^2 \sigma^2}{m}.$

$$\leq \frac{}{3k(\xi +$$

where the first inequality is due to $(x+1)^{1/3} \le x^{1/3} + \frac{1}{3x^{2/3}}$ and the second inequality is due to $\xi > 2$. Let $c \geq \frac{2}{3k^3} + \frac{5}{4}$, then we have

 $\frac{1}{n_t} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{m}_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2$

Then we define the Lyapunov function $R_t = \mathbb{E}[\mathcal{L}(p_t) + \frac{\gamma}{n_{t-1}} \|\nabla \mathcal{L}(p_t) - m_t\|^2]$. Then we have

$$\begin{aligned} R_{t+1} - R_t = & \mathbb{E}[\mathcal{L}(\boldsymbol{p}_{t+1}) - \mathcal{L}(\boldsymbol{p}_t)] + \frac{1}{\eta_t} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_{t+1}) - \boldsymbol{m}_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 \\ \leq & (\eta_t \gamma - \frac{5\eta_t \gamma}{4}) \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 - \frac{\eta_t}{2\gamma} \mathbb{E} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + 2L^2 \eta_t \gamma \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + \frac{2\theta_t^2 \sigma^2 \gamma}{\eta_t} \\ \leq & - \frac{\eta_t \gamma}{4} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 - \frac{\eta_t}{4\gamma} \mathbb{E} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 + \frac{2\theta_t^2 \sigma^2 \gamma}{\eta_t}, \end{aligned}$$

where the last inequality is due to $\gamma \leq \frac{1}{2\sqrt{2L}}$. Rearranging the above inequality, we have

$$\frac{\eta_t \gamma}{4} \mathbb{E} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + \frac{\eta_t}{4\gamma} \mathbb{E} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2 \le R_t - R_{t+1} + \frac{2\theta_t^2 \sigma^2 \gamma}{\eta_t}$$

Taking average over timesteps $t = 1, \ldots, T$, we have

$$\begin{aligned} \frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\left[\frac{\eta_t\gamma}{4}\|\nabla\mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + \frac{\eta_t}{4\gamma}\|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2\right] \\ \leq \frac{\mathcal{L}(\boldsymbol{p}_1) - \mathcal{L}^*}{T} + \frac{\gamma\|\nabla\mathcal{L}(\boldsymbol{p}_1) - \boldsymbol{m}_1\|^2}{T\eta_0} + \sum_{t=1}^{T}\frac{2\theta_t^2\sigma^2\gamma}{T\eta_t} \leq \frac{\mathcal{L}(\boldsymbol{p}_1) - \mathcal{L}^*}{T} + \frac{\gamma\sigma^2}{T\eta_0} + \sum_{t=1}^{T}\frac{2\theta_t^2\sigma^2\gamma}{T\eta_t} \\ = \frac{\mathcal{L}(\boldsymbol{p}_1) - \mathcal{L}^*}{T} + \frac{\gamma\xi^{1/3}\sigma^2}{kT} + \sum_{t=1}^{T}\frac{2c^2\eta_t^3\sigma^2\gamma}{T}.\end{aligned}$$

1188 Dividing both sides with $\gamma \eta_T$, we have

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\left[\frac{1}{4} \|\nabla \mathcal{L}(\boldsymbol{p}_t) - \boldsymbol{m}_t\|^2 + \frac{1}{4\gamma^2} \|\tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_t\|^2\right]$$

$$\leq \frac{\mathcal{L}(\boldsymbol{p}_{1}) - \mathcal{L}^{*}}{T\eta_{T}\gamma} + \frac{\xi^{1/3}\sigma^{2}}{kT\eta_{T}} + \sum_{t=1}^{T} \frac{2c^{2}\eta_{t}^{3}\sigma^{2}}{T\eta_{T}} \leq \frac{\mathcal{L}(\boldsymbol{p}_{1}) - \mathcal{L}^{*}}{T\eta_{T}\gamma} + \frac{\xi^{1/3}\sigma^{2}}{kT\eta_{T}} + \frac{2c^{2}\sigma^{2}}{T\eta_{T}}\int_{1}^{T} \frac{k^{3}}{\xi + t}dt$$

$$\mathcal{L}(\boldsymbol{p}_{1}) - \mathcal{L}^{*} = \xi^{1/3}\sigma^{2} - 2c^{2}\sigma^{2}k^{3}, \quad (z = z)$$

$$\leq \frac{\mathcal{L}(\mathbf{p}_{1}) - \mathcal{L}}{T\eta_{T}\gamma} + \frac{\zeta}{kT\eta_{T}} + \frac{1}{T\eta_{T}} \log(\xi + T)$$

$$= \frac{\mathcal{L}(\mathbf{p}_{1}) - \mathcal{L}^{*}}{T\gamma k} (\xi + T)^{\frac{1}{3}} + \frac{\xi^{1/3}\sigma^{2}}{kT} (\xi + T)^{\frac{1}{3}} + \frac{2c^{2}\sigma^{2}k^{3}}{T} (\xi + T)^{\frac{1}{3}} \log(\xi + T)$$

$$\leq \frac{M}{T} (\xi + T)^{\frac{1}{3}},$$

where $M = \frac{\mathcal{L}(p_1) - \mathcal{L}^*}{\gamma k} + \frac{\xi^{1/3} \sigma^2}{k} + 2c^2 \sigma^2 k^3 \log(\xi + T)$. Using Jensen's inequality, we have

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[\frac{1}{2} \| \nabla \mathcal{L}(\boldsymbol{p}_{t}) - \boldsymbol{m}_{t} \| + \frac{1}{2\gamma} \| \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_{t} \| \right] \\ & \frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[\frac{1}{2} \| \nabla \mathcal{L}(\boldsymbol{p}_{t}) - \boldsymbol{m}_{t} \| + \frac{1}{2\gamma} \| \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_{t} \| ^{2} \right] \\ & \frac{1}{212} \\ & 2 \sum_{t=1}^{T} \mathbb{E} \left[\frac{1}{4} \| \nabla \mathcal{L}(\boldsymbol{p}_{t}) - \boldsymbol{m}_{t} \|^{2} + \frac{1}{4\gamma^{2}} \| \tilde{\boldsymbol{p}}_{t+1} - \boldsymbol{p}_{t} \|^{2} \right] \right)^{\frac{1}{2}} \\ & \frac{1}{213} \\ & \frac{1}{214} \\ & 2 \sum_{t=1}^{\sqrt{2M}} (\xi + T)^{1/6} \leq \frac{\sqrt{2M}}{T^{1/2}} (\xi^{1/6} + T^{1/6}), \end{aligned}$$

1218 where the first inequality is due to $x + y \le (2x^2 + 2y^2)^{1/2}$ and the last inequality is due to $(x + y)^{1/6} \le x^{1/6} + y^{1/6}$. Then we have

$$\begin{aligned} \frac{1221}{1222} & \frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \| G_{\mathcal{C}}(\boldsymbol{p}_{t}, \gamma) \| = \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \| \boldsymbol{p}_{t} - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \\ = \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \| \boldsymbol{p}_{t} - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) + \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \\ = \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \| \boldsymbol{p}_{t} - \tilde{\boldsymbol{p}}_{t+1} + \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \\ = \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \| \boldsymbol{p}_{t} - \tilde{\boldsymbol{p}}_{t+1} + \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \\ \leq \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \| \boldsymbol{p}_{t} - \tilde{\boldsymbol{p}}_{t+1} \| + \| \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \\ \leq \frac{1}{T} \sum_{t=1}^{T} \frac{1}{\gamma} \mathbb{E} \left[\| \boldsymbol{p}_{t} - \tilde{\boldsymbol{p}}_{t+1} \| + \| \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \boldsymbol{m}_{t}) - \operatorname{proj}_{\mathcal{C}}(\boldsymbol{p}_{t} - \gamma \nabla \mathcal{L}(\boldsymbol{p}_{t})) \| \right] \\ \leq \frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[\| \nabla \mathcal{L}(\boldsymbol{p}_{t}) - \boldsymbol{m}_{t} \| + \frac{1}{\gamma} \| \boldsymbol{p}_{t} - \tilde{\boldsymbol{p}}_{t+1} \| \right] \\ \leq \frac{2\sqrt{2M}}{T^{1/2}} \xi^{1/6} + \frac{2\sqrt{2M}}{T^{1/3}}, \end{aligned}$$

where the first inequality is due to Triangle inequality, the second inequality is due to the non-expansivity of convex projection.

1242 C IMPLEMENTATION DETAILS

1244 C.1 MANUAL TEMPLATES

Table 4: Input templates, and output label words used in RoBERTa-large. $\langle S \rangle$ represents the sentences in the dataset. [MASK] represents the mask token.

1249	Task	Dataset	Input Template	Output Label Words
1250	Peol World Datasets	Amazon Book	$\langle S \rangle$ It was [MASK].	positive, negative
1251	Real-world Datasets	Amazon Electronics	$\langle S \rangle$ It was [MASK].	positive, negative
1252	GLUE Datasets	CoLA	$\langle S \rangle$ correct? [MASK].	no, yes
1253	OLUE Datasets	MRPC	$\langle S_1 \rangle \langle S_2 \rangle$ entailment? [MASK].	no, yes

Table 5: Input templates, and output label words used in GPT2-XL and Llama3. $\langle S \rangle$ represents the sentences in the dataset.

Task	Dataset	Input Template	Output Label Words
Peol World Datasets	Amazon Book	$\langle S \rangle$ It was	positive, negative
Real-world Datasets	Amazon Electronics	$\langle S \rangle$ It was	positive, negative
CLUE Detecto	CoLA	$\langle S \rangle$ correct?	no, yes
GLUE Datasets	MRPC	$\langle S_1 \rangle \langle S_2 \rangle$ entailment?	no, yes

C.2 HYPERPARAMETERS

Table 6: Main hyperparameters used in our algorithm.

Hyperparameter	RoBERTa-large	GPT2-XL	Llama3
query limit	32000	3200	1600
train batch size	32	32	16
prompt length	$\{50, 20\}$	$\{50, 20\}$	$\{50, 20\}$
step size	1e-3	1e-3	1e-3

D ADDITIONAL EXPERIMENT RESULTS

1277 D.1 EXPERIMENT RESULTS ON MRPC

We conduct experiments on MRPC using the same experiment setups as on CoLA and observe similar phenomena as those on CoLA. The results are shown in Table 7.

Table 7: Comparison of AUC scores (mean±std.) on constructed imbalanced scenarios of MRPC.
We conduct three groups of experiments on pre-trained RoBERTa-large, GPT2-XL, and Llama3 with a prompt length of 20. The best results are highlighted in **bold**.

1285	Imbalanced Ratio	Method	RoBERTa-large	GPT2-XL	Llama3
1286		Manual Prompt	$.4764 \pm .0855$	$.4556 \pm .0834$.5264±.0531
1287		BBT	$.5236 \pm .0497$	$.4986 \pm .0024$	$.5000 \pm .0000$
1288	$\tau = 20$	GAP3	$.5000 \pm .0000$	$.4972 \pm .0024$	$.4708 \pm .0331$
1289		BDPL	$.4639 \pm .1660$	$.4917 {\pm} .0072$	$.4972 {\pm} .0048$
1290		mDP-DPG (ours)	.5292±.0573	.5278±.0808	$.5083 \pm .0144$
1291		Manual Prompt	.4700±.1386	$.4433 \pm .1444$.5206±.1275
1292		BBT	.5400±.0173	$.4972 \pm .0024$	$.5250 \pm .0433$
1293	$\tau = 50$	GAP3	$.5000 \pm .0000$	$.5517 \pm .1721$.4717±.0407
1204		BDPL	$.3050 \pm .0976$	$.5667 \pm .1241$	$.5000 \pm .0000$
1205		mDP-DPG (ours)	.5767±.0580	.5983±.1234	.5317±.0548

1296 D.2 OVERCOMING THE LIMITATION OF BINARY CLASSIFICATION

Model

RoBERTa-Large

GPT2-XL

Llama3

Although our current use of AUC loss is specific to binary classification, however, it could also be generalized to the multi-class classification dataset by using the micro averaging. That is, for each class, calculate the AUC loss for that class against all others and sum the losses and average them over all classes. Additionally, we conduct experiments on the multi-class datasets MNLI and SNLI and compare our methods with 3 representative baselines. The results are presented in the Table 8. It can be observed that our methods also maintain optimal performance on multi-class datasets.

Method Manual Prompt

BBT

BDPL

mDP-DPG(ours)

Manual Prompt

BBT

BDPL

mDP-DPG(ours)

Manual Prompt

BBT

BDPL

mDP-DPG(ours)

Table 8: Comparison of AUC values on 2 multi-class datasets

MNLI (len=50)

 $.4636 \pm .0340$

 $.4589 \pm .0271$

 $.4670 \pm .0462$

.4814±.0538

 $.4718 \pm .0360$

.4761±.0413

.4518±.0503

.4823±.0382

 $.4036 \pm .0106$

 $.5025 \pm .0034$

 $.4946 \pm .0124$

.5079±.0284

SNLI (len=50)

 $.5290 \pm .0458$

 $.5845 \pm .0470$

.5697±.0189

.5897±.0296

 $.5150 \pm .0434$

.5177±.0113

 $.5104 \pm .0189$

.5243±.0197

 $.5478 \pm .0122$

.4612±.0556

 $.6034 \pm .0233$

.6171±.0298

1304 1305

1306

307

1309 1310

1311

1313 1314

1315 1316

1317 1318

1319

1321

1320 D.3 BALANCED SCENARIO

Although experiments in the paper have demonstrated that our methods outperform baseline methods in imbalanced scenarios, their effectiveness in balanced settings is equally important. If our methods were to suffer from performance collapse in balanced scenarios, their utility would be compromised. To verify the performance of our methods in the 16-shot setting, we have conducted additional experiments, with the results provided in Table 9. In balanced scenarios, the performance of our methods is similar to that of various baselines in most cases, with a few instances where they even surpass all baselines.

Table 9: Comparison of ACC on 3 datasets in the balanced scenario with prompt length of 20.

Model	Method	CoLA	Book	Elec
	BBT	.5717±.0159	$.9364 {\pm} .0008$.9149 ±.0040
RoBERTa-Large	GAP3	$.5254 \pm .0739$	$.9019 {\pm} .0308$.8519±.1271
	BDPL	.4851±.0515	$.9349 \pm .0016$.8794±.0229
	mDP-DPG(ours)	.5762±.0605	.9384±.0006	.9112±.0075
	BBT	.3321±.0244	$.6873 \pm .0397$	$.5423 \pm .0671$
GPT2 YI	GAP3	$.4624 \pm .0774$.8257±.1216	.6312±.3736
UI 12-AL	BDPL	.6497±.0221	$.7527 {\pm} .0283$.6741±.0941
	mDP-DPG(ours)	.6142±.0339	$.8034 \pm .0217$.7777±.0608

1339 1340

1341 1342 D.4 TRAINING EFFICIENCY

On the one hand, both mDP-DPG and BDPL use variance-reduced policy gradient estimators. On the other hand, since mDP-DPG requires sampling pairs of examples and involves additional forward passes through the black-box model. To mitigate this computational cost, we employ smaller mini-batch size, which reduces the number of forward passes while achieving comparable experimental results. Furthermore, although pairwise sampling necessitates multiple forward passes, the computational burden is still much lower when compared backpropagation. To visually compare the training efficiency between BDPL and our methods, we provide Figure 4 showing the progression of the current best AUC on the development set across epochs.



We provide the prompts learned by our methods in Table 10 and 11, along with some correctly predicted examples. Our prompts, like those in Diao et al. (2022), are sequences of discrete words without explicit natural language semantics. Additionally, from the black-box optimization perspective, we prefer to consider the prompts as tunable parameters of LLM, and we can adapt the model to downstream tasks at a lower cost by optimizing prompts.

1381 1382

1383

D.6 REAL-WORLD APPLIACTION

To show performance of our method in real-world applications, we add three additional represen-1384 tative imbalance datasets: **BB** (Burfoot and Baldwin) Burfoot & Baldwin (2009), originally de-1385 veloped for satire news detection; Job Scams and SMS datasets Boumber et al. (2024), derived 1386 for fraudulent job postings and spam message detection, respectively. The BB dataset consists of 1387 4,000 true news articles and 233 satire articles. Its challenge lies in satire articles mimicking the 1388 tone and style of true news while incorporating exaggerated or absurd content, requiring semantic 1389 understanding and background knowledge for accurate classification. The Job Scams dataset, de-1390 rived from the Employment Scam Aegean Dataset, includes 14,295 cleaned job advertisements, of 1391 which 599 are fraudulent postings, presenting a significant class imbalance. Fraudulent postings 1392 often use fake job positions to deceive applicants, typically featuring short and structured texts. The 1393 SMS dataset contains 6,574 messages, of which 1,274 are spam or phishing. These deceptive mes-1394 sages are typically brief and generic promotional content, whereas genuine messages reflect more 1395 personalized communication. The diversity of these datasets ensures the broad applicability of the experiments. 1396

We compare the performance of five methods—Manual Prompt, BDPL, APE (Zhou et al., 2022),
EvoPrompt (Guo et al., 2023), and mDP-DPG—evaluated under prompt lengths of 5 and 20 using
true black-box LLM GPT-4. For APE and EvoPrompt, the final prompts are generated based on the
manual prompt pool and therefore do not have a fixed prompt length. The results are shown in Table 12. The experimental results demonstrate that mDP-DPG outperforms other methods across all
datasets. On the BB dataset, mDP-DPG achieves an AUC of 0.5972 (length = 5), significantly surpassing BDPL and Manual Prompt. Despite the class imbalance challenge on the Job Scams dataset,
mDP-DPG achieves an AUC of 0.5307 (length = 20), outperforming BDPL's 0.5024. Overall, the re-

INICUIUU	Dataset	Prompt+Sentence	Prediction	Label
	CoLA	Sandy was trying to work out which stu-	no	ves
nDP-DPG		dents would be able to solve a certain prob-		5.5%
		lem, but she wouldn't tell us which one.		
		This never in had for of her It his if They then with some not her set the second the se	yes	
		than with some not know think would I hat		
		students would be able to solve a certain		
		problem but she wouldn't tell us which		
		one.		
	Book	Not impressed! This is just another gluten	positive	negative
		Tree COOKDOOK - albeit with some great		U
		recipes. I don't find fecipes that contain agave pectar to be "sugar" free. What's		
		the purpose - anyone can use other sweet-		
		eners than plain old sugar While the		
		recipes might be sugar-free they are def-		
		initely NOT without sweeteners. Sorry I		
		bought the bookperhaps there are better		
		books using NO sweeteners that satisfy the		
		sweet tooth of a gluten-free person???		
		is their books at is one It through good	negative	
		can really he I so will just It never so will		
		Not impressed! This is just another gluten		
		free cookbook - albeit with some great		
		recipes. I don't find recipes that contain		
		the purpose - anyone can use other sweet-		
		eners than plain old sugar While the		
		recipes might be sugar-free they are def-		
		initely NOT without sweeteners. Sorry I		
		bought the bookperhaps there are better		
		books using NO sweeteners that satisfy the		
		sweet tooth of a gluten-free person???		
Tat	ole 11: Lea	rned prompts on GPT2-XL and correctly class	ified examples	
Tat Method	ole 11: Lear Dataset	rned prompts on GPT2-XL and correctly class Prompt+Sentence	ified examples Prediction	Label
Tat Method	ole 11: Lean Dataset CoLA	Prompt+Sentence Jerry attempted to blow up the Pentagon.	ified examples Prediction no	Label yes
Tat Method nDP-DPG	De 11: Lean Dataset CoLA	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What	ified examples Prediction no yes	Label yes
Tal Method nDP-DPG	De 11: Lear Dataset CoLA	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the	ified examples Prediction no yes	Label yes
Tał Method nDP-DPG	Die 11: Lean Dataset CoLA	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon	ified examples Prediction no yes	Label yes
Tal Method nDP-DPG	Dataset CoLA	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon.	ified examples Prediction no yes positive	Label yes
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds. I thought this was going to be	ified examples Prediction no yes positive	Label yes negative
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec-	ified examples Prediction no yes positive	Label yes negative
Tal Method mDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year	ified examples Prediction no yes positive	Label yes negative
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on	ified examples Prediction no yes positive	Label yes negative
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one.	ified examples Prediction no yes positive	Label yes negative
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a	ified examples Prediction no yes positive negative	Label yes negative
Tal Method mDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a most characters not up she their a my book	ified examples Prediction no yes positive negative	Label yes negative
Tal Method mDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a most characters not up she their a my book people only most I had read the 1 year to	ified examples Prediction no yes positive negative	Label yes negative
Tal Method mDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a most characters not up she their a my book people only most I had read the 1 year to an organized life, by Ms Leeds, I thought	ified examples Prediction no yes positive negative	Label yes negative
Tal Method mDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a most characters not up she their a my book people only most I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the ex-	ified examples Prediction no yes positive negative	Label yes negative
Tal Method nDP-DPG	Dataset CoLA Book	Prompt+Sentence Jerry attempted to blow up the Pentagon. which did not by her from people What him do eat were people will on It a more was made Jerry attempted to blow up the Pentagon. I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the exact same book!! I rec- ommend buying the newer version: 1 year to an organized life, don't waste money on this one. well into there most love reading into a most characters not up she their a my book people only most I had read the 1 year to an organized life, by Ms Leeds, I thought this was going to be different. It's the ex- act same book!! I recommend buying the	ified examples Prediction no yes positive negative	Label yes

Length	Method	BB ($\tau = 20$)	Job Scams ($\tau = 20$)	SMS ($\tau = 5$)
-	Manual Prompt	$.4333 {\pm} .0191$	$.5098 \pm .0546$	$.5252 \pm .0131$
≤ 20	APE	$.4968 \pm .0398$.5000±.0042	$.5268 \pm .0184$
≤ 20	EvoPrompt	$.5002 \pm .0373$.4972±.0064	$.5271 \pm .0231$
5	BDPL	$.4486 \pm .0146$.4987±.0235	$.5200 \pm .0114$
	mDP-DPG (ours)	.5972±.1428	.5314±.0386	.5272±.0090
20	BDPL	$.4403 \pm .0244$	$.5024 \pm .0306$	$.5135 \pm .0064$
	mDP-DPG (ours)	.5236±.1545	.5307±.0043	.5278±.0119

Table 12: Comparison of AUC values across different datasets using GPT-4.

sults validate the adaptability of the mDP-DPG method, particularly in complex or class-imbalanced tasks where it exhibited remarkable advantages.