ZeroDL: Zero-shot Distribution Learning for Text Clustering via Large Language Models

Anonymous ACL submission

Abstract

The advancements in large language models (LLMs) have brought significant progress in NLP tasks. However, if a task cannot be fully described in prompts, the models could fail to carry out the task. In this paper, we propose a simple yet effective method to contextualize a task toward a LLM. The method utilizes (1) open-ended zero-shot inference from the entire dataset, (2) aggregating the inference results, and (3) finally incorporate the aggregated metainformation for the actual task. We show the effectiveness in text clustering tasks, empowering LLMs to perform text-to-text-based clustering and leading to improvements on several datasets. Lastly, we explore the generated class labels for clustering, showing how the LLM understands the task through data.

1 Introduction

001

004

011

013

017

021

037

Large language models (LLMs) have demonstrated impressive performances on various downstream tasks (Devlin et al., 2019; Radford et al., 2019). These also exhibit the ability to understand the context of input text, known as in-context learning (ICL) (Brown et al., 2020; OpenAI, 2023). ICL allows to leverage LLMs for specific tasks without further extensive training. However, effective use of ICL hinges on well-designed prompts.

While prompts with few-shot examples demonstrably improve performance, they can easily overfit a model to the examples (Perez et al. (2021); Mizrahi et al. (2023); *inter alia*). This led to a growing interest in zero-shot learning, which reduces the need for intricate few-shot selection. Recent advancements in zero-shot learning involve incorporating more sophisticated use of prompt structures, such as Chain-of-Thought (Wei et al., 2022), zero-shot reasoning (Kojima et al., 2022), and models trained to follow instructions (Ouyang et al., 2022; Chung et al., 2024). However, how to design prompts for target tasks remains challenging.



Figure 1: Illustration of proposed method ZeroDL. While in-context learning (ICL) relies on examples (D) tailored to specific tasks, ZeroDL aggregates all the outputs from these zero-shot inferences (\hat{m}) , resulting in meta-level information (\hat{M}) . This information is then used by the LLM to generate its final predictions.

Motivated by the core principle of ICL– providing task and data contexts *within* prompts– we propose an approach to construct more effective zero-shot prompts by understanding how LLMs describe datasets *across* prompting outputs. As illustrated in Figure 1, **Zero**-shot **D**istribution Learning (ZeroDL) aims to learn data distributions through zero-shot inferences. The method comprises two key components: open-ended zero-shot inference and output aggregation. Zero-shot prompts are then constructed with the generated meta information, and used for actual task. This method takes advantage of the self-generated frame of LLMs to successfully carry out a given task.

We exemplify the effectiveness of ZeroDL on text clustering tasks where complete task descriptions cannot be provided because of an absence of ground-truth class labels. Furthermore, we show how LLMs understand the task compared with human-labeled classes. In addition, our method works in a text-to-text format unlike traditional clustering algorithms, allowing clustering with specific context. For instance, "I love this movie" and "I hate this movie" express opposite sentiment but

064

041

042

043

065

081

100

102

104

105

106

107

108 109

110

Distribution Learning that leverages zero-shot inferences to generate meta-level information

about the data distribution by aggregating open-ended inference outputs from datasets. • ZeroDL allows models to perform text-based clustering, empowering them to handle data

belong to the same cluster of movie reviews¹.

Our contributions in this paper are as follows:

• We propose a novel approach called **Zero**-shot

with specific context, which offers advantages over embedding-based clustering methods.

Related Works 2

While there have been a few attempts to leverage LLMs for clustering tasks, majority of existing approaches rely on traditional methods like K-Means clustering based on LLM-generated embeddings (Petukhova et al., 2024; BehnamGhader et al., 2024). In contrast, our approach leverages text-level prompting, which injects specific viewpoints into LLMs, enabling it to perform more targeted and contextualized clustering. While Huang and He (2024) proposed a related clustering framework for LLMs, their approach relies on few-shot learning using human labels, which deviates significantly from traditional clustering settings.

The importance of appropriate ground-truth labels extends beyond clustering and permeates ICL. While Min et al. (2022) observed cases where the input-label correspondence does not play significant roles, Yoo et al. (2022) argued that the impact heavily depends on target tasks and experiment settings. We believe that this work would serve as a reference that generate the appropriate class information automatically by LLM itself.

Proposed Method: ZeroDL 3

Stage 1: Open-Ended Inference. We begin by designing a prompt for zero-shot classification. This prompt intentionally avoids any detailed information about the task, minimizing the risk of overfitting. Based on the idea, we opt for the simplest prompt format:

Text: [text]\n\nClassify the text to the best [type_of_task] class.

where [text] is the input data and [type_of_task] provides view of the task. In the experiment, it can be either sentiment or topic. Leveraging this prompt, we perform model inferences on the input data (D). This process generates open-ended class predictions, denoted as $\hat{m_1} \cdots \hat{m_N}$.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

Stage 2: Aggregation. The open-ended predictions lack constraints, leading to potentially inconsistent output formats. For instance, the model might predict "positive" and "non-positive" classes, while the ground-truth is "positive" and "negative". The predictions could even be entire sentences. To address the inconsistencies in the open-ended predictions, we employ aggregation strategy.

Before the aggregation, we count the frequency of each predictions and sort it. After that, the predictions which frequency is only 1 are dropped in order to filter out extraordinary predictions and save computation. Next, we iteratively construct subsets of the predictions by removing the least frequent predictions one by one. This process results in a list of subsets, denoted as [$\{\hat{m}_1, \hat{m}_2, \cdots, \hat{m}_U\}$, $\{\hat{m}_1, \hat{m}_2, \cdots, \hat{m}_{U-1}\}, \cdots, \{\hat{m}_1\}\}$ where U denotes the number of remained predictions. We input all the subset to provide more weights to the frequently occurred predictions. The prompt for aggregation is as follows:

[type_of_task] List:	
- $\hat{m_1} = \hat{m_2} = \dots = \hat{m_U}$ # prompt (sub)set	
<pre>\n\nAggregate the [type_of_task] List into</pre>	
[NUM_CLUSTER_CLASS] classes.	

where [NUM_CLUSTER_CLASS] is pre-defined number of classes to cluster.

However, the model outputs often still lack coherence, especially in generating exact number of classes. To address this, we select the LLM outputs where the number of aggregated classes matches the pre-defined number of cluster classes. We then use the most frequent class pairs aggregated from prompt (sub)sets as meta-information M. This information represents the model's understanding of potential views over the entire data.

Stage 3: Leveraging Meta-Information We incorporate the aggregated meta-information (\hat{M}) into the original prompt² to enhance the model's prediction capabilities:

<pre>Text: [text]\n\nClass description:</pre>
- Class 0: $\hat{M_0}$
- Class 1: \hat{M}_1
- Class \cdots : \hat{M}_{\cdots}
- Class $C: \hat{M_C}$
\n\nBased on the class description, classify
the text to the best [type_of_task] class.

By incorporating the meta-information (M) into

²The order of Text: and Class description: can be reversed. The placeholders - Class n: serve as templates to parse.

Refer to Table 10 in the Appendix for potential risk of absence of the perspectives.

Model	Method	IMDB	SST-2	SST-5	YRev	AGNew	DBp(F)DBp(H	3)Yah(F)Yah(B)	Macro	Micro
mistral -7b-inst	ZeroDL(C-T) ZeroDL(T-C)	85.3 92.2	$\frac{82.4}{79.4}$	43.4 36.0	$\frac{51.6}{48.2}$	$\frac{75.4}{81.0}$	64.0 62.2	73.3 78.4	<u>47.2</u> 51.8	$\frac{69.4}{73.4}$	$\frac{65.8}{67.0}$	65.1 66.4
	Gold(C-T) Gold(T-C)	87.5 91.7	77.0 82.5	41.8 43.3	50.8 51.8	60.7 82.7	74.2 84.1	85.2 82.7	40.8 50.9	62.5 73.8	64.5 71.5	68.1 72.7
	llm2vec+KMeans	s 62.1	55.2	30.3	56.0	84.4	96.1	70.8	48.0	46.2	61.0	67.5
Llama-3 -8b-inst	ZeroDL(C-T) ZeroDL(T-C)	$\frac{89.8}{66.7}$	$\frac{74.2}{79.1}$	$\frac{42.6}{41.1}$	37.6 49.2	$\frac{69.0}{80.4}$	65.9 73.2	52.5 63.7	$\frac{46.8}{49.1}$	48.2 48.6	$\frac{58.5}{61.2}$	56.4 60.8
	Gold(C-T) Gold(T-C)	66.5 93.2	73.7 84.2	42.1 46.0	48.9 53.5	55.9 70.2	71.3 74.6	75.5 79.2	35.2 49.3	54.7 75.4	58.2 69.5	61.1 70.3
	llm2vec+KMeans	s 53.2	53.2	29.6	54.8	85.1	94.5	72.7	47.4	50.4	60.1	66.2
Llama-3 -70b-inst	ZeroDL(C-T) ZeroDL(T-C)	94.3 95.0	90.0 89.4	42.7 41.5	47.7 49.3	83.3 72.6	55.3 47.0	35.4 <u>36.1</u>	51.2 53.2	82.0 81.4	64.7 62.8	57.5 56.2
	Gold(C-T) Gold(T-C)	94.4 95.2	88.9 89.4	53.0 52.3	56.6 58.2	83.2 81.1	94.8 95.8	95.8 91.0	65.3 63.8	82.4 82.6	79.4 78.8	80.7 80.3

Table 1: The performance of ZeroDL for text clustering. C-T denotes the prompt order with class information then input text. T-C is the reversed. Bold means the best accuracy and underline the outperforming cases than ground-truth (i.e., Gold) class label setting. We present more baselines in Table 6 of the Appendix.

the prompt, we enable the LLM to perform conditioned classification within the clustering context. The generated meta-information might not directly correspond to the ground-truth classes, potentially leading to variations in performance.

4 Experiments

152

153

154

155

156

157

159

160

161

162

163

164

165

166

167

168

169

171

172

173

174

176

178

179

180

181

Our experiments are conducted on Models. mistral-7b-instruct-v0.2 (Jiang et al., 2023) due to its impressive performance even with a relatively small parameter size compared with other LLMs. We additionally report results obtained with the Llama-3 family models (Meta, 2024) ran with vLLM library (Kwon et al., 2023). Greedy decoding is employed for generation. We follow the official instruction formatting guidelines. The results are averaged from 5 runs. We also tested llm2vec (BehnamGhader et al., 2024) with KMeans approach where applicable. The method represents a baseline using (not exactly but) the same backbone³ computed by embedding-level, like encoders. Other baselines with different backbones are presented in Table 6 of the Appendix⁴. Setting. Our ZeroDL method is designed to consider all data distributions and converge (or aggregate) them into several classes. We believe that clustering datasets often contain too many classes relative to the number of data points. Consequently, we opted for text classification datasets, which typically provide a larger number of data instances per class. Furthermore, the availability of groundtruth labels allows for a direct comparison between generated and actual labels, enabling a robust qualitative analysis of our method's effectiveness.

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201

202

203

204

205

206

207

209

Datasets and Evaluation. We use 6 text classification datasets, as mentioned in Setting: IMDB (Maas et al., 2011), SST-2, SST-5 (Socher et al., 2013), YelpReivews (Zhang et al., 2015) for sentiment classification and AGNews, DBpedia (Lehmann et al., 2015), YahooAnswers (Chang et al., 2008) for topic classification. Further details are presented in Table 4 of the Appendix.

We evaluate the model performance using accuracy. To determine the predicted class, we leverage the Class n: anchor tokens within the LLM outputs. LLMs might not directly predict the same classes as the ground-truth labels, for example, Class 0 means Positive in the prediction while Negative is labeled as Class 0 in the ground-truth. We thus test all possible mapping combinations and report the best performed mapping. However, this results in a total of factorial of C potential mappings, so we split datasets with more than 7 classes (i.e., DBpedia and YahooAnswers⁵) into 2 subsets Front (F) and Back (B) to avoid out-of-memory⁶.

5 Results

Table 1 reports the performance of our method. Interestingly, ZeroDL achieves even better performance than models provided with ground-truth

³llm2vec involves additional training steps.

⁴In the Appendix, Table 5 provides a detailed explanation of why some approaches are not considered fair baselines.

⁵After removing the 3 smallest size of classes: CarsAnd-Transportation, SocialScience, Sports

⁶The total computational cost is factorial to the number of classes (e.g., 7! = 5040 for 7 classes, 14! = 87 billion for 14). To avoid excessively small or large class sizes, we balanced the number of classes to 7.

Data	Method	ClassLabels
SST-5	ZeroDL	 Neutral Sentiment: This class includes all the sentiment labels that express a neutral sentiment towards the movie or documentary. Examples include () Negative Sentiment: This class includes all the sentiment labels that express a negative or sad emotion towards the movie or documentary. Examples include () Ambiguous Sentiment: This class includes all the sentiment labels that do not clearly express a positive or negative emotion towards the movie or documentary. Examples include () Mixed Sentiment: This class includes all the sentiment labels that express a mixed sentiment towards the movie or documentary. Examples include () Mixed Sentiment: This class includes all the sentiment labels that express a mixed sentiment towards the movie or documentary. Examples include () Positive Sentiment: This class includes all the sentiment labels that express a positive emotion towards the movie or documentary. Examples include ()
	Gold	Very Negative, Negative, Neutral, Positive, Very Positive
AGNews	ZeroDL	 International Relations and Politics: This class includes topics related to international relations, diplomacy, Middle East politics, terrorism, nuclear politics, and elections. Sports and Entertainment: This class includes topics related to sports, tennis, golf, basketball, baseball, cricket, and entertainment. Business and Economy: This class includes topics related to the economy, finance, stocks, mergers and acquisitions, retail, real estate, and labor markets. Technology and Science: This class includes topics related to technology, computing, internet, cybersecurity, space exploration, and science.
	Gold	World, Sports, Business, Sci/Tech

Table 2: The example of generated class labels in 5-class sentiment classification (SST-5), and topic classification (AGNews). ZeroDL can generates alternative class labels and its description. Additional examples are in Table 9.

	S2(C-T)	S2(T-C)	S5(C-T)	S5(T-C)
RandToken	51.6	59.5	28.6	28.6
AutoL(Best)	86.8	84.9	41.8	38.5
AutoL(Worst)	82.5	80.7	41.7	38.6
Gold	77.0	82.5	41.8	43.3
ZeroDL(Mistral) ZeroDL(GPT3.5)	82.4 73.3	79.4 79.1	43.4 42.8	36.0 42.1

Table 3: The performance with various class labels. The best AutoL in SST-2 are [Wonderful, Bad] and the worst labels are [Irresistible, Pathetic]. In SST-5, [Terrible, Better, Good, Extraordinary, Unforgettable] are the best while [Aw-ful, Better, Hilarious, Perfect, Wonderful] are the worst.

(i.e., Gold) class labels on several datasets. This suggests that ZeroDL might uncover richer or more nuanced class structures within the data compared with the pre-defined labels. Our method demonstrates comparable performance to K-Means clustering that utilizes LLM embeddings, outperforming on tasks with relatively smaller datasets. ZeroDL achieves this by flexible zero-shot prompting without any modification. Table 10 in the Appendix shows the importance of the constraints utilized by ZeroDL with detailed examples.

210

211

212

213

214

215

216

217

218

219

221

222

224

227

229

Table 3 investigates the significance of class labels in text clustering tasks. We explore the performance of mistral-7b-instruct-v0.2 using class labels suggested by AutoL (Gao et al., 2021) designed for prompt-based model fine-tuning. These labels represent a curated selection. We also investigate the class labels generated by gpt-3.5-turbo. The results demonstrate that performances are generally higher when using manually selected labels compared with the original dataset labels. This suggests that carefully chosen labels can significantly im-

prove clustering outcomes. In the context, ZeroDL performs particularly well on SST-5, implying that the critical role of class labels in text clustering and the potential of ZeroDL to capitalize on informative class labels automatically.

232

233

234

235

236

238

239

240

241

242

243

244

245

246

247

248

249

251

252

253

254

255

257

258

259

260

261

262

Table 2 shows the examples of class labels generated by ZeroDL. Notably, these labels provide richer explanations for the classifications compared with the original ground-truth labels⁷. Besides, ZeroDL might uncover newly emerged classes based on the data, such as "Ambiguous Sentiment" and "Mixed Sentiment" (see more examples in Table 8 of the Appendix).

ZeroDL involves a trade-off in computational cost. We thus measure the performance changes based on the amount of input data used (see Table 7 in the Appendix). The results indicate that using only 10% of the data yields plausible performance but it leads to inconsistency in the model performances, considering increased standard deviation.

6 Conclusion

We introduce ZeroDL, a novel approach to contextualize LLM tasks for a given LLM. ZeroDL employs open-ended zero-shot inference and output aggregation to learn data distributions. We demonstrate its effectiveness, showing competitive performances against embedding-based clustering methods and superior performance than ground-truth labels in some cases. Beyond its clustering capabilities, ZeroDL offers the generation of informative class labels that provide deeper insights into LLMs.

⁷The generation of class description might depend on datasets and LLMs; Llama-3 usually does not make it.

263

272

273

274

276

278

279

290

292

293

296

297

299

305

306

308

309

7 Limitations

Prompt Dependency and Heuristics. ZeroDL
relies on carefully designed prompts to guide LLMs
towards effective clustering. While our focus was
on using simple and intuitive prompts, prompt selection can potentially influence the model's behavior and introduce biases. Future work could explore
more sophisticated prompt engineering techniques
to further enhance ZeroDL's performance.

Experiments with Diverse LLMs and Prompts. While we acknowledge the computational limitations (and price) of ZeroDL, investigating its behavior with a wider range of LLMs (including commercial models like GPT-4, Claude, and Gemini) and prompt templates could provide valuable insights into the generalizability and robustness of the approach.

Lower than State-of-the-Art Performance. Achieving state-of-the-art performance is not the sole focus of ZeroDL but it offers a valuable framework in understanding data distributions with zeroshot inference via LLMs. By addressing the limitations mentioned above, ZeroDL has the potential to become a powerful and versatile tool not only for text clustering but data exploration.

Expensive Computational Cost in Inferences. Although ZeroDL have an alternative approach to reduce computational burden through data sampling, effectively sampling data to generate appropriate class labels remains a challenge. The technique within the framework presents a valuable future direction.

References

- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. 2024. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Ming-Wei Chang, Lev-Arie Ratinov, Dan Roth, and Vivek Srikumar. 2008. Importance of semantic representation: Dataless classification. In *AAAI*, volume 2, pages 830–835.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yanjun Gao, Ting-Hao Huang, and Rebecca J. Passonneau. 2021. Learning clause representation from dependency-anchor graph for connective prediction. In Proceedings of the Fifteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-15), pages 54–66, Mexico City, Mexico. Association for Computational Linguistics.
- Chen Huang and Guoxiu He. 2024. Text clustering as classification with llms. *arXiv preprint arXiv:2410.00927*.
- AQ Jiang, A Sablayrolles, A Mensch, C Bamford, DS Chaplot, D de las Casas, F Bressand, G Lengyel, G Lample, L Saulnier, et al. 2023. Mistral 7b (2023). *arXiv preprint arXiv:2310.06825*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles.*
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 6(2):167–195.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. https://ai.meta.com/blog/meta-llama-3/.

Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe,

Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-

moyer. 2022. Rethinking the role of demonstrations:

What makes in-context learning work? In Proceed-

ings of the 2022 Conference on Empirical Methods in

Natural Language Processing, pages 11048–11064,

Abu Dhabi, United Arab Emirates. Association for

Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2023. State

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in neural in-

formation processing systems, 35:27730–27744.

Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021.

Alina Petukhova, Joao P Matos-Carvalho, and Nuno Fachada. 2024. Text clustering with llm embeddings.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan,

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on

Empirical Methods in Natural Language Processing

and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages

3982-3992, Hong Kong, China. Association for Com-

Richard Socher, Alex Perelygin, Jean Wu, Jason

Chuang, Christopher D Manning, Andrew Y Ng, and

Christopher Potts. 2013. Recursive deep models for

semantic compositionality over a sentiment treebank.

In Proceedings of the 2013 conference on empiri-

cal methods in natural language processing, pages

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten

Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,

et al. 2022. Chain-of-thought prompting elicits rea-

soning in large language models. Advances in neural

information processing systems, 35:24824–24837.

Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyun-

Dario Amodei, Ilva Sutskever, et al. 2019. Language models are unsupervised multitask learners. OpenAI

arXiv preprint arXiv:2403.15112.

True few-shot learning with language models. Advances in neural information processing systems,

of what art? a call for multi-prompt llm evaluation.

Computational Linguistics.

arXiv preprint arXiv:2401.00595.

OpenAI. 2023. Gpt-4 technical report.

34:11054-11070.

blog, 1(8):9.

putational Linguistics.

- 376 377
- 378
- 379

- 386
- 387
- 390

- 394 395
- 396
- 400 401 402

403 404

- 405
- 406 407
- 408
- 409 410
- 411 412

413 414

415

416 417

418

419

soo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee, and Taeuk Kim. 2022. Ground-truth labels matter: A 420

1631-1642.

deeper look into input-label demonstrations. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 2422-2437, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

421

422

423

424

425

426

427

428

429

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In Advances in neural information processing systems, pages 649-657.

	#Train	#Valid	#Test	#Class	Pre-defined ClassTitle
IMDB	25,000	3,750	25,000	2	Negative, Positive
SST-2	9,645	1,101	2,210	2	Negative, Positive
SST-5	9,645	1,101	2,210	5	Very Negative, Negative, Neutral, Positive, Very Positive
YelpReviews	650,000	97,500	49,999	5	Very Negative, Negative, Neutral, Positive, Very Positive
AGNews	120,000	18,000	7,600	4	World, Sports, Business, Sci/Tech
DBpedia(F)	280,000	41,939	35,000	7	Company, EducationalInstitution, Artist, Athlete, OfficeHolder, MeanOfTransportation, Building
DBpedia(B)	280,000	42,213	35,000	7	NaturalPlace, Village, Animal, Plant, Album, Film, WrittenWork
Yahoo(F)	59,518	8,879	10,489	7	ArtsAndHumanities, BeautyAndStyle, BusinessAndFinance, ComputersAndInternet, ConsumerElectronics, EducationAn- dReference, EntertainmentAndMusic
Yahoo(B)	59,493	8,973	10,514	7	FoodAndDrink, GamesAndRecreation, Health, HomeAndGar- den, Pets, PregnancyAndParenting, SocietyAndCulture

Table 4: Data statistics used in the experiments.

Method	Base Models	Further Trained?	Predefined Labels are Used?	Justification
ZeroDL	Decoder-based LLM, text-level	No	No	-
ZeroDL(with Gold label)	Decoder-based LLM, text-level	No	Yes	Baseline with human labels
llm2vec + K-Means (BehnamGhader et al., 2024)	Decoder-based LLM, embedding- level	Yes	No	Baseline with embedding-level approach (but further trained)
TF-IDF+K-Means	- , embedding-level	No	No	Traditional baseline but too weak
SBERT + K-Means (Reimers and Gurevych, 2019)	Encoders, embedding-level	Yes	No	Different backbone and further trained

Table 5: Justification for appropriate baselines and other methods. We believe a comparison of GT cases and llm2vec, both of which utilize the same LLM, would be more appropriate.

Model	Method	IMDB	SST-2	SST-5	YRev	AGNew	DBp(F)DBp(B)Yah(F)Yah(B)	Macro	Micro
TF-IDF	KMeans	52.1	53.0	26.2	35.2	43.8	55.9	62.8	44.7	46.5	46.7	48.8
SBERT	KMeans	65.0	54.0	33.9	30.3	82.7	95.4	93.9	67.2	68.3	65.6	67.5
mistral -7b-inst	ZeroDL(C-T) ZeroDL(T-C)	85.3 92.2	$\frac{82.4}{79.4}$	43.4 36.0	$\frac{51.6}{48.2}$	$\frac{75.4}{81.0}$	64.0 62.2	73.3 78.4	<u>47.2</u> <u>51.8</u>	$\frac{69.4}{73.4}$	$\frac{65.8}{67.0}$	65.1 66.4
	Gold(C-T) Gold(T-C)	87.5 91.7	77.0 82.5	41.8 43.3	50.8 51.8	60.7 82.7	74.2 84.1	85.2 82.7	40.8 50.9	62.5 73.8	64.5 71.5	68.1 72.7
	llm2vec+KMeans	62.1	55.2	30.3	56.0	84.4	96.1	70.8	48.0	46.2	61.0	67.5
Llama-3 -8b-inst	ZeroDL(C-T) ZeroDL(T-C)	$\frac{89.8}{66.7}$	$\frac{74.2}{79.1}$	$\frac{42.6}{41.1}$	37.6 49.2	$\frac{69.0}{80.4}$	65.9 73.2	52.5 63.7	$\frac{46.8}{49.1}$	48.2 48.6	$\frac{58.5}{61.2}$	56.4 60.8
	Gold(C-T) Gold(T-C)	66.5 93.2	73.7 84.2	42.1 46.0	48.9 53.5	55.9 70.2	71.3 74.6	75.5 79.2	35.2 49.3	54.7 75.4	58.2 69.5	61.1 70.3
	llm2vec+KMeans	53.2	53.2	29.6	54.8	85.1	94.5	72.7	47.4	50.4	60.1	66.2
Llama-3 -70b-inst	ZeroDL(C-T) ZeroDL(T-C)	94.3 95.0	90.0 89.4	42.7 41.5	47.7 49.3	83.3 72.6	55.3 47.0	35.4 <u>36.1</u>	51.2 53.2	82.0 81.4	64.7 62.8	57.5 56.2
	Gold(C-T) Gold(T-C)	94.4 95.2	88.9 89.4	53.0 52.3	56.6 58.2	83.2 81.1	94.8 95.8	95.8 91.0	65.3 63.8	82.4 82.6	79.4 78.8	80.7 80.3

Table 6: The performance of ZeroDL for text clustering compared with more methods for references. C-T denotes the prompt order with class information then input text. T-C is the reversed. Bold means the best accuracy and underline the outperforming cases than ground-truth (i.e., Gold) class label setting. Note that llm2vec for Llama3-70b is not publicly available at the moment. Refer to Table 5 why other baselines (especially SBERT) cannot be a fair comparison.

[C-T]	IMDB	SST-2	SST-5	YRev	AGNew	DBp(F)	DBp(L)	Yah(F)	Yah(L)	Macro	Micro
All (std)	85.3 3.40	82.4 3.67	43.4 2.93	51.6 1.35	75.4 2.23	64.0 3.17	73.3 3.82	47.2 1.02	69.4 2.65	65.8	65.1
1% (std)	75.2 23.01	70.7 13.47	38.3 6.17	51.1 1.30	-	48.9 10.51	37.9 22.91	23.7 1.82	17.2 2.50	45.4	47.8
5% (std)	92.9 0.30	82.4 2.34	45.8 2.78	52.1 1.64	-	59.2 6.77	58.5 7.37	35.2 9.03	51.7 14.15	59.7	60.1
10% (std)	91.9 0.82	82.0 3.17	35.0 4.45	50.3 3.47	-	67.2 6.35	68.6 3.27	37.3 2.47	57.1 6.24	61.2	63.5
[T-C]	IMDB	SST-2	SST-5	YRev	AGNew	DBp(F)	DBp(L)	Yah(F)	Yah(L)	Macro	Micro
[T-C] All (std)	IMDB 92.2 2.34	SST-2 79.4 5.12	SST-5 36.0 3.84	YRev 48.2 3.22	AGNew 81.0 1.81	DBp(F) 62.2 4.65	DBp(L) 78.4 7.14	Yah(F) 51.8 1.07	Yah(L) 73.4 1.91	Macro 67.0	Micro 66.4
[T-C] All (std) 1% (std)	IMDB 92.2 2.34 76.1 23.88	SST-2 79.4 5.12 78.3 6.01	SST-5 36.0 3.84 39.4 5.17	YRev 48.2 3.22 50.0 2.48	AGNew 81.0 1.81	DBp(F) 62.2 4.65 49.3 13.03	DBp(L) 78.4 7.14 34.3 19.56	Yah(F) 51.8 1.07 21.0 3.67	Yah(L) 73.4 1.91 17.2 2.50	Macro 67.0 45.7	Micro 66.4 46.9
[T-C] All (std) 1% (std) 5% (std)	IMDB 92.2 2.34 76.1 23.88 93.3 0.93	SST-2 79.4 5.12 78.3 6.01 81.0 2.59	SST-5 36.0 3.84 39.4 5.17 41.7 4 09	YRev 48.2 3.22 50.0 2.48 48.4 4.98	AGNew 81.0 1.81 - - -	DBp(F) 62.2 4.65 49.3 13.03 55.3 1.84	DBp(L) 78.4 7.14 34.3 19.56 67.9 4.28	Yah(F) 51.8 1.07 21.0 3.67 56.3 6.09	Yah(L) 73.4 1.91 17.2 2.50 55.7 12.68	Macro 67.0 45.7 62.5	Micro 66.4 46.9 61.7

Table 7: ZeroDL performances according to the number of training data. ZeroDL fails to find exact number of clusters in AGNews dataset.

Predicted	Gold	Example
	VeryNeg	It 's hard not to feel you 've just watched a feature-length video game with some really heavy back story .
	Neg	But it pays a price for its intricate intellectual gamesmanship.
Mixed Sentiment: ()	Noutral	The appearance of Treebeard and Gollum 's expanded role will either
	Neutral	have you loving what you 're seeing, or rolling your eyes .
	Dec	An utterly compelling ' who wrote it ' in which the reputation of
	Pos	the most famous author who ever lived comes into question .
	VeryPos	a roller-coaster ride of a movie
	VeryNeg	It 's difficult to say whether The Tuxedo is more boring or embarrassing
		- I'm prepared to call it a draw.
		Like most Bond outings in recent years,
Ambimum Continuet. ()	Neg	some of the stunts are so outlandish that they border on being cartoonlike .
Ambiguous Sentiment: ()	Neutral	Effective but too-tepid biopic
	Pos	But he somehow pulls it off.
	VamuDaa	Emerges as something rare, an issue movie that 's so honest
	veryPOS	and keenly observed that it does n't feel like one .

Table 8: Example of data predicted to newly generated classes in SST-5. Randomly selected.

Data	Method	ClassLabels
		Negative Sentiment: The list also includes various expressions of negative sentiment towards movies, films,
	ZeroDI	shows, and documentaries. Some examples include ()
IMDB	LCIODL	Positive Sentiment: The list includes various expressions of positive sentiment towards movies, films, shows,
	Cald	and documentaries. Some examples include ()
	GOID	Positive Sentiment: All the sentiment labels that express a positive sentiment towards the movie film document
		tary, or subject. For example, ()
SST-2	ZeroDL	Negative or Neutral Sentiment: All the sentiment labels that do not express a positive sentiment towards the
		movie, film, documentary, or subject. For example, ()
	Gold	Negative, Positive
		Negative Sentiment: This class includes sentences expressing negative sentiments towards a place, food, or
		Very Positive Sentiment: This class includes sentences expressing highly positive sentiments towards a place
		food, or experience. Examples include ()
	ZeroDI	Mixed Sentiment: This class includes sentences expressing mixed sentiments towards a place, food, or experi-
Yelp	LCIODE	ence. Examples include ()
Rev		Neutral Sentiment: This class includes sentences expressing neutral sentiments towards a place, food, or experience. Examples include (
		Positive Sentiment: This class includes sentences expressing positive sentiments towards a place, food, or
		experience. Examples include ()
	Gold	Very Negative, Negative, Neutral, Positive, Very Positive
		Aviation and Transportation: This class includes topics related to aviation, aerospace technology, military
		Business and Economy. This class includes topics related to business finance, industries, companies, and
		economics.
		Sports and Biographies: This class includes topics related to sports, athletes, and their biographies.
		Politics and Government: This class includes topics related to politics, government, elections, and specific
DBp	ZeroDL	political parties. Education: This class includes topics related to education, universities, schools, and specific educational
(F)		institutions.
		History and Architecture: This class includes topics related to history, architecture, historic sites, castles, and
		landmarks.
		Art and Entertainment: This class includes topics related to art, music, entertainment, and specific artists or record labels
	Gold	Company, EducationalInstitution, Artist, Athlete, OfficeHolder, MeanOfTransportation, Building
		Science and Technology: This class includes topics related to paleontology, geology, volcanology, space
		exploration, and academic journals.
	ZeroDI	Music and Entertainment : This class includes topics related to music, album releases, jazz music, heavy metal
		Geography and Hydrology: This class includes topics related to geography, hydrology, rivers, water bodies.
		and water resources.
		Botany and Plant Sciences: This class includes topics related to botany, horticulture, plant taxonomy, plant
DBp	LCIODE	conservation, and endangered species.
(B)		and academic publications
		Zoology and Entomology : This class includes topics related to zoology, entomology, moths, butterflies, fish
		species, and arachnids.
		Film and Television: This class includes topics related to film, cinema, movies, movie reviews, Bollywood, and
	Gold	NaturalPlace Village Animal Plant Album Film WrittenWork
	Gora	Personal Finance and Economics: Topics related to personal finance, credit scores, debt management, taxes,
		and economics.
		Health and Wellness : Topics related to health, medicine, fitness, nutrition, and wellness.
		Technology and Computing: Topics related to computers, technology, software, internet, telecommunications
17 - I-	ZeroDL	and mobile phones.
ran (F)		Miscellaneous: Topics that do not fit neatly into any of the above categories, such as philosophy, religion,
(-)		science, and humor.
		Education and Careers : Topics related to raunon, clothing, makeup, cosmetics, nair care, and beauty.
		and employment.
	Gold	ArtsAndHumanities, BeautyAndStyle, BusinessAndFinance, ComputersAndInternet, ConsumerElectron-
		ics, EducationAndReference, EntertainmnentAndMusic
		rood and Cooking: This class includes topics related to various cuisines, recipes, food items, and cooking techniques
		Home Improvement and DIY : This class includes topics related to home repair, renovation, decorating,
		gardening, and DIY projects.
		Miscellaneous: This class includes topics that do not fit neatly into any of the above categories, such as politics,
	ZaroDI	education, art, and entertainment.
Yah	LEIODL	technology, and internet culture.
(B)		Health and Wellness: This class includes topics related to physical and mental health, nutrition, dieting, weight
		loss, fitness, exercise, and medical conditions.
		Ammais and rets: This class includes topics related to various religions, theology, philosophy, and Religion and Philosophy. This class includes topics related to various religions, theology, philosophy, and
		spirituality.
	Gold	FoodAndDrink, GamesAndRecreation, Health, HomeAndGarden, Pets, PregnancyAndParenting, Soci-
	3010	etyAndCulture

Table 9: Additional examples of generated class labels. Class title is marked as bold and its description is colored as gray.

Data	TaskType	Generated Class Labels (Descriptions are omitted)
	sentiment	Negative, Very Positive, Highly Negative, Mixed, Neutral-Positive, Positive, Highly Positive, Neutral-Negative, Extremely Negative, Neutral; [Total 10]
IMDB	topic	Film and Television Influence on Business, Film and Television Influence on Politics, Film and Television Production Companies, Film and Television Influence on Art, Film and Television Productions, Film Analysis or Review, Film and Television Influence on Society and Culture, Film and Television Influence on Technology, Film and Television Recommendations, Film and Television Awards, Film and Television, Film and Television Genres, Film and Television Influence on Society, Film and Television Marketing, Movie Reviews or Film Criticism, Film and Television Critics, Film and Television Technologies, Film and Television Festivals, Film or Media Criticism, Film and Television History, Film and Television Technology Trends, Film and Television Education, Film and Television Influence on Entertainment, Film and Television Influence on Education, Personal Opinion or Review, Film and Television Industry, Film and Television Distribution; [Total 27]
SST	sentiment	The text expresses a negative sentiment towards the subject being described, The text expresses a negative or sad sentiment towards the film, The text expresses a negative or sad sentiment, The text has a negative tone, The text expresses a negative or cautionary sentiment, The text expresses a negative or slightly negative sentiment, Negative Sentiment, The text expresses a negative or slightly negative sentiment towards the film, The text expresses a negative sentiment towards the film, The text expresses a negative sentiment, Negative sentiment, The text expresses a negative sentiment, to text expresses a negative sentiment towards the film, The text expresses a negative sentiment, The text expresses a negative sentiment towards the film, The text expresses a negative sentiment towards the film, The text expresses a negative sentiment towards the movie being described, The text has a negative sentiment; [Total 11]
	topic	Literature and Writing, Media and Entertainment, Mental Health and Psychology, Travel and Adventure, Film and Television, Family and Relationships, Science and Technology, Food and Cooking, Performing Arts, Sports, Education and Learning, Business and Finance, Greetings and Open-Ended Texts, Entertainment; [Total 14]
	sentiment	Positive Sentiment, Highly Negative Sentiment, Neutral to Positive Sentiment, Negative Sentiment, Mixed Sentiment, Neutral Sentiment, Very Positive Sentiment, Neutral to Negative Sentiment; [Total: 8]
Yelp Rev topic		Automotive, Environment of Nature, Legal of Law, Customer Service of Business, Mightine of Entertainment, Education or Training, Sports or Fitness, Discrimination or Racism, Religion or Spirituality, Technology or Gadgets, Health or Medical, Travel or Tourism, Education or Learning, Beauty or Personal Care, Customer Reviews or Testimonials, Politics or Government, Food or Cooking, Food or Beverage, Shopping or Retail, Home Improvement or Construction, Science or Technology, Real Estate or Housing, Food Safety or Food Poisoning, Entertainment or Leisure, Arts or Culture, Personal Care or Beauty, Business or Economy, Personal Experiences, Dining Experience or Food Review: [Total 29]
AG News	sentiment	Negative: 12 expressions that contain negative sentiment, Positive: 25 expressions that contain positive sentiment, Negative: 12, Sentiment not clear: 3, Mixed: 11 expressions that contain a mix of positive and negative sentiment, Positive: 25, Neutral: 33, Neutral: 33 expressions that do not contain any clear positive or negative sentiment, Sentiment not clear: 3 expressions that do not provide enough context to determine a clear sentiment, Mixed: 11; [Total 10]
	topic	Same Result with Table 2
DBp (F)	topic	Neutral, Positive, Negative; [Total 3] Baseball, Aviation, Sports and Biographies, History, Aircraft, Soccer or Football, People and Biographies, Ice Hockey, Education, American Football, General, Aircraft Design, Music or Entertainment, Higher Education or Universities, Healthcare or Hospitals, Football or Soccer; [Total 15]
	sentiment	Ambiguous, Neutral, Positive, Romantic, Negative, Mixed, Objective; [Total 6]
DBp (B)	topic	Botany or Biology (specifically, Plant Science or Taxonomy), Botany or Endangered Species or Conservation Biology, Botany or Horticulture, Botany or Algae, Botany or Brazilian Flora, Botany or Cacti, Botany or Mexican Flora, Botany or Tillandsia species, Botany or Orchids, Botany or Aquatic Plants, Botany or Plant Science, Botany or Tropical Plants, Botany or Hawaiian Flora, Botany or Plant Taxonomy, Botany or Palm Trees [Total 15]
	sentiment	Positive, Neutral, Mixed, Informational, Negative [Total 5]
Yah (F)	topic	Miscellaneous (for topics that do not fit neatly into any specific category), Education and Careers, Makeup and Beauty, Gaming and Technology, Fashion and Clothing, Telecommunications, Housing and Real Estate, Music and Entertainment, Pop Culture and Entertainment, Personal Finance and Credit Scores; [Total 10]
	sentiment	Neutral; [Total 1]
Yah (B)	topic	Health and Medical Concerns, Pop Culture and Entertainment, Mental Health and Psychology, Literature and Writing, Business and Finance, Food and Cooking, Video Games and Technology, Art and Creativity, Philosophy and Ethics, Education and Learning, Humor and Satire, Sports and Fitness, Religion and Theology, Home Improvement and DIY, Science and Technology, Travel and Adventure, Pets and Animals, Politics and Society, Pregnancy and Reproductive Health; [Total 19]

Table 10: The examples of generated class labels when no constraints are given; we experiment with wrong task type and unlimited the number of clusters.