
Towards the Effect of Examples on In-Context Learning: A Theoretical Case Study

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In-context learning (ICL) has emerged as a powerful ability for large language
2 models (LLMs) to adapt to new tasks by leveraging a few (demonstration) examples.
3 Despite its effectiveness, the mechanism behind ICL remains underexplored. This
4 paper uses a Bayesian framework to investigate how ICL integrates pre-training
5 knowledge and examples for binary classification. In particular, we introduce a
6 probabilistic model extending from the Gaussian mixture model to exactly quantify
7 the impact of pre-training knowledge, label frequency, and label noise on the
8 prediction accuracy. Based on our analysis, when the pre-training knowledge
9 contradicts the knowledge in the examples, whether ICL prediction relies more on
10 the pre-training knowledge or the examples depends on the number of examples.
11 In addition, the label frequency and label noise of the examples both affect the
12 accuracy of the ICL prediction, where the minor class has a lower accuracy and
13 how the label error impacts the accuracy is determined by the specific error rate
14 of the two classes. Extensive simulations are conducted to verify the correctness
15 of the theoretical results, and real-data experiments also align with the theoretical
16 insights. Our work reveals the dual role of pre-training knowledge and examples in
17 ICL, offering a deeper understanding of LLMs' behaviors in classification tasks.

18 1 Introduction

19 Large language models (LLMs) have revolutionized various fields, such as GitHub Copilot for
20 software development, Microsoft 365 Copilot to embrace productivity, and medical applications such
21 as Med-Palm [1]. A particularly intriguing ability of LLMs is in-context learning, where LLMs can
22 adapt to new tasks only using a few examples at the inference stage without changing the model
23 parameters. As ICL enhances the predictive performance of LLMs, various existing literature attempts
24 to understand and quantify such a superiority [2, 3, 4].

25 During the ICL process, LLMs typically demonstrate two key abilities [5]: retrieving knowledge from
26 the pre-training data and learning from the examples in the prompt. Understanding how pre-training
27 knowledge and specific examples interact during the inference stage is crucial, especially given the
28 complex dynamics observed in practical applications. For instance, existing literature [6] conducts
29 various empirical evaluations to study ICL regarding the example size, demonstration order, prompt
30 templates, etc. Meanwhile, theoretical studies [7, 8, 9, 10, 11, 5] have explored the underlying
31 mechanisms of ICL from various perspectives, including Bayesian approaches and gradient descent,
32 primarily focusing on linear regression models.

33 However, the existing literature ¹ is insufficient to understand the behavior of ICL, especially for
34 classification tasks. First, previous works cannot draw a consensus on certain behaviors of ICL.
35 For example, in [6], it is empirically observed that injecting random noise to the example labels

¹Related works are discussed in section G.

36 does not hurt the ICL performance. They conjecture that the robustness of ICL against the noise is
37 because the pre-training knowledge dominates ICL. On the other hand, based on [5], when taking a
38 large number of examples (i.e., a large example size), the ICL will favor the knowledge provided
39 by the examples. However, there is no systematic understanding of the role of label quality, the
40 difference between pre-training and example knowledge, as well as the example size. Second, existing
41 theoretical frameworks may fail to explain the observed behaviors in classification tasks. For instance,
42 the balance of the example size in different classes matters in classification, while there is no such
43 concept in regression.

44 The above gaps drive the need for a theoretical exploration of how LLMs utilize pre-training knowl-
45 edge and specific examples in ICL in classification scenarios. In particular, we aim to answer:
46 **How do LLMs make predictions in classification tasks using their pre-training knowledge and**
47 **examples?**

48 This work aims to explore the above question by conducting an exact theoretical analysis in a binary
49 classification task. Our contributions are summarized as follows:

- 50 • We leverage Bayesian analysis to exactly quantify the ICL performance (measured by prediction
51 accuracy). When the example size is small, the pre-training knowledge will dominate ICL, and
52 when the example size is large, the ICL prediction mainly relies on the examples. Built upon this
53 finding, we further study the ICL performance under different scenarios such as label noise and
54 imbalanced examples mentioned above, and contradiction in pre-training and example knowledge.
55 Extensive simulations and real-data analysis are conducted to support our theoretical insights.
56 Technically, to perform the analysis, we assume all examples in the pre-training are selected
57 independently, and examine the posterior distribution of the parameters of the data generation
58 model with two distinct types of priors: one from the pre-training data and the other from the
59 examples. A central challenge in the analysis lies in the formulation and integration of these two
60 priors into a single coherent posterior for ICL prediction. Our result successfully accounts for
61 both the label distributions and the conditional output distribution within each class.
- 62 • When conducting simulations to verify the above theoretical insights, we surprisingly reveal
63 another counter-intuitive behavior when the examples are not selected independently: We fix
64 exactly 50% positive labels in each prompt in pre-training and provide only positive examples in
65 the test prompt, then the ICL prediction is a firm negative. We provide an intuitive explanation
66 and theoretical justification to explain this behavior. This finding can help understand how LLMs
67 consider dependency among tokens/sequences.

68 2 Classification Analysis via Bayesian

69 To analyze the ICL performance, we first introduce the model and data assumptions in Section
70 2.1, then derive the ICL accuracy under general situations in Section 2.2. We finally examine how
71 examples influence ICL under specific demonstration scenarios (Section 2.3).

72 2.1 Setups

73 To perform the exact analysis of the ICL prediction, in this subsection, we introduce the pre-training-
74 inference paradigm and impose some assumptions on the data generation distribution.

75 **Pre-training.** For the pre-training data, inspired by [2, 4, 5, 12], we form the prompts in the form
76 of $((x_1, y_1), (x_2, y_2), \dots, (x_k, y_k), (x_{query}, y_{query}))$ and the target is to predict the label y_{query} , i.e.,
77 performing ICL using the examples $\{(x_i, y_i)\}_{i \in [k]}$, where $x_i \in \mathbb{R}^m$ and $y_i \in \{-1, +1\}$. To simplify
78 the notation, we use (x, y) and $\{(x_i, y_i)\}$ to refer to (x_{query}, y_{query}) and $\{(x_i, y_i)\}_{i \in [k]}$ respectively
79 when no confusion arises. All (x_i, y_i) s and x in the same prompt are in general sampled from the
80 same distribution, and an exception considering label noise will be described in detail later in Section
81 2.3. In different prompts, the sampling distribution may vary. We assume that all examples in the
82 demonstration are independent, typically selected randomly from a prompt set [6]. A discussion
83 on the case where the examples are not independent is provided in Section E for a comprehensive
84 analysis.

85 Denote the pre-trained LLM as M . Without loss of generality, we assume that M learns the exact
86 distribution of the pre-training data and makes predictions based on the pre-training knowledge,
87 i.e., model output $M(x)$ follows the pre-training distribution. This assumption is supported by the
88 capabilities of LLMs and is commonly used in existing research [7, 13].

89 **Inference.** At the inference stage, we perform ICL of a test input x given the examples to predict
 90 its corresponding test label y , i.e. $\hat{y}_{ICL} = M(\{(x_i, y_i)\}, x)$. Our goal is to study the effect of the
 91 pre-training data and the examples on the distribution of \hat{y}_{ICL} . To simplify the notation, we denote
 92 $P(\hat{y}_{ICL} = s)$ as $P(y = s|x, \{(x_i, y_i)\}, M)$ for $s \in \{-1, +1\}$. Unless otherwise stated, $P(\cdot|\dots, M)$
 93 represents the meaning of “conditional on the pre-training knowledge”.

94 **Data generation process.** We mainly follow the idea of Bayesian inference to form the assumptions.
 95 Bayesian inference is a well-established theoretical method that has demonstrated its effectiveness in
 96 explaining the behavior of LLMs as shown in existing research [7, 8]. To perform Bayesian inference,
 97 we impose a prior distribution on the parameters in the data generation process and then use data and
 98 the prior distribution together to derive a posterior distribution of the parameters. The idea of prior
 99 distribution is widely used in uncertainty quantification in real applications such as various medical
 100 studies [14], and is justified by axioms of decision theory [15].

101 The following two assumptions are imposed in our main study. We consider the pre-training data,
 102 example, and test data to be in the same distribution family but with different parameters. Assumption
 103 1 describes how to generate a pair of (x, y) given a specific set of parameters, and Assumption 2
 104 explains how the parameters differ among datasets.

105 **Assumption 1** (Generate (x, y)). Assume $x \in \mathbb{R}^m$ and $y \in \{-1, +1\}$. Given parameters
 106 $(\theta_+, \theta_-, \pi, p_+, p_-)$, to generate (x, y) , y is first generated from a Bernoulli distribution with π ,
 107 i.e. $P(y = +1) = \pi$ and $P(y = -1) = 1 - \pi$, then x is generated from a class-wise input distribu-
 108 tion accordingly. Given $y = +1$, x follows a Gaussian distribution $N(\theta_+, \sigma_+^2 I)$ with probability p_+
 109 and sample from $N(\theta_-, \sigma_-^2 I)$ with probability $1 - p_+$; given $y = -1$, x is sampled from a Gaussian
 110 distribution $N(\theta_-, \sigma_-^2 I)$ with probability p_- and sample from $N(\theta_+, \sigma_+^2 I)$ with probability $1 - p_-$.
 111 In addition, the examples are independent with each other.

112 Assumption 1 follows the standard Gaussian mixture design for theoretical analysis in classification,
 113 e.g., [16, 17, 18]. We further consider “label noise”: When $y = +1$, the corresponding x can be from
 114 either of the two clusters. When $p_+ = p_- = 1$, it means that there is no label noise.

115 **Assumption 2** (Parameters). The parameter distributions for pre-training and the inference stage as
 116 as follows:

- 117 • Pre-train: $\theta_+ \sim N(\theta_M, \sigma_M^2 I_m)$, $\theta_- \sim N(-\theta_M, \sigma_M^2 I_m)$; $p_+ = p_- = 1$; $\pi \sim \text{Beta}(1, 1)$.
- 118 • Examples: $\theta_+ \sim N(\theta_+^e, \sigma_{e+}^2 I_m)$, $\theta_- \sim N(\theta_-^e, \sigma_{e-}^2 I_m)$; $p_+^e, p_-^e \in [0, 1]$, $\pi \in [0, 1]$.
- 119 • Test data (x, y) : $p_+^t = p_-^t = 1$. Examples and the test data in the same prompt share the same
 120 realization of (θ_+, θ_-) .

121 In pre-training, all (x_i, y_i) s and (x, y) in the same prompt are conditionally independent and share
 122 the same parameters. At the inference stage, the examples are conditionally independently sampled
 123 given the parameters and may incur label noise. For the test data (x, y) , while it shares the same
 124 (θ_+, θ_-) with the examples in the prompt, we do not further consider label noise in the test data.
 125 The proportion $P(y = +1)$ is not considered in the test data because the later accuracy analysis is
 126 performed on $y = +1$ and $y = -1$ separately.

127 Assumption 2 aligns with the common scenarios of ICL, i.e., the pre-training distribution and the
 128 example distribution at the inference stage can differ. In pre-training, we take $p_+ = p_- = 1$ to
 129 simplify the derivation. In this case, there is no label noise, and the misclassification of the Bayes
 130 classifier is only caused by the overlap of the two Gaussian clusters in the distribution. At the
 131 inference stage, the examples may have a distribution shift compared to the pre-training data, and we
 132 also consider potential label noise in the examples.

133 2.2 ICL Decision and Prediction Accuracy

134 To compute the ICL prediction accuracy, we first derive the posterior distribution of the parameters
 135 $(\theta_+, \theta_-, \pi)$ given the examples $\{(x_i, y_i)\}$ and the pre-training knowledge of θ_M , and then use
 136 $(\theta_+, \theta_-, \pi)$ to figure out the ICL accuracy.

137 **Posterior of parameters.** Our goal is to compute the posterior distribution of θ_+, θ_-, π given
 138 examples $(x_1, y_1), \dots, (x_k, y_k)$. Recall that in Assumption 2, $p_+ = p_- = 1$ in the pre-training stage,

139 and the p_+^e and p_-^e in the inference stage can be some values in $[0, 1]$ if label noise occurs. Denote
 140 $\#(y_i = +1)$ and $\#(y_i = -1)$ as the number of examples with positive/negative labels respectively
 141 (after possible flips if p_+^e or p_-^e is less than 1). The following lemma presents the posterior distribution
 142 of π, θ_+, θ_- .

Lemma 1. *Under Assumption 1 and Assumption 2, the posterior distribution of π, θ_+, θ_- satisfies*

$$P(\pi | \{(x_i, y_i)\}_{i \in [k]}, M) \propto \pi^{\#(y_i = +1)} (1 - \pi)^{\#(y_i = -1)},$$

$$\theta_+ \sim N \left(\frac{\sigma_+^2 \theta_M + \sigma_M^2 \sum_{y_i = +1} x_i}{\sigma_+^2 + \#(y_i = +1) \sigma_M^2}, \frac{\sigma_+^2 \sigma_M^2}{\sigma_+^2 + \#(y_i = +1) \sigma_M^2} I \right) \triangleq N(\hat{\theta}_+, \sigma_{\theta_+}^2 I),$$

and

$$\theta_- \sim N \left(\frac{\sigma_-^2 \sum_{y_i = -1} x_i - \sigma_M^2 \theta_M}{\sigma_-^2 + \#(y_i = -1) \sigma_M^2}, \frac{\sigma_-^2 \sigma_M^2}{\sigma_-^2 + \#(y_i = -1) \sigma_M^2} I \right) \triangleq N(\hat{\theta}_-, \sigma_{\theta_-}^2 I).$$

143 The proof of Lemma 1 can be found in Section B.1. In short, since the examples $\{(x_i, y_i)\}$ are given,
 144 we can directly write the likelihood for $(\pi, \theta_+, \theta_-)$ to derive the corresponding posterior distributions.

145 **ICL decision.** Given Lemma 1, denoting $z_k = (\#(y_i = -1) + 1) / (\#(y_i = +1) + 1)$, the following
 146 lemma shows the ICL decision boundary for the test data:

147 **Lemma 2.** *Under Assumption 1 and Assumption 2, the probability of $y = +1 / -1$ is as follows*

$$\begin{aligned} P(y = +1 | x, \{(x_i, y_i)\}, M) &= \frac{P(x, y = +1 | \{(x_i, y_i)\}, M)}{P(x, y = +1 | \{(x_i, y_i)\}, M) + P(x, y = -1 | \{(x_i, y_i)\}, M)} \\ &= \frac{\frac{\#(y_i = +1) + 1}{k+2} N(x)}{\frac{\#(y_i = +1) + 1}{k+2} N(x) + \frac{\#(y_i = -1) + 1}{k+2}}, \\ P(y = -1 | x, \{(x_i, y_i)\}, M) &= \frac{P(x, y = -1 | \{(x_i, y_i)\}, M)}{P(x, y = +1 | \{(x_i, y_i)\}, M) + P(x, y = -1 | \{(x_i, y_i)\}, M)} \\ &= \frac{\frac{\#(y_i = -1) + 1}{k+2}}{\frac{\#(y_i = +1) + 1}{k+2} N(x) + \frac{\#(y_i = -1) + 1}{k+2}}, \end{aligned}$$

where

$$N(x) = \left(\sqrt{\frac{\sigma_-^2 + \sigma_{\theta_-}^2}{\sigma_+^2 + \sigma_{\theta_+}^2}} \right)^m \exp \left[-\frac{(x - \hat{\theta}_+)^T (x - \hat{\theta}_+)}{2(\sigma_+^2 + \sigma_{\theta_+}^2)} + \frac{(x - \hat{\theta}_-)^T (x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\theta_-}^2)} \right].$$

148 The decision boundary is $\hat{y}_{ICL} = 1(f_{ICL}(x) > 0)$, where $f_{ICL}(x) = N(x) - z_k$.

149 The proof of Lemma 2 can be found in Section B.2. When $(\pi, \theta_+, \theta_-)$ are fixed, given y , x follows a
 150 Gaussian distribution. When integrating over all possible $(\pi, \theta_+, \theta_-)$, the marginal distribution of x
 151 given y still follows a Gaussian distribution. Hence, (x, y) marginally follows a Gaussian mixture
 152 distribution, and the decision boundary can be further obtained.

153 From Lemma 2, we can see how the pre-training distribution (θ_M, σ_M^2) and examples $\{(x_i, y_i)\}$
 154 impact $P(y = +1 | x, \{(x_i, y_i)\}, M)$ and $P(y = -1 | x, \{(x_i, y_i)\}, M)$, and further change the
 155 decision boundary correspondingly. The pre-training knowledge (θ_M, σ_M^2) and examples $\{(x_i, y_i)\}$
 156 first determines $(\hat{\theta}_+, \hat{\theta}_-, \sigma_{\theta_+}^2, \sigma_{\theta_-}^2)$, the latter of which further determines the decision boundary.
 157 More details about the interplay of pre-training and examples under different scenarios will be
 158 provided in Section 2.3. Besides, a larger π will result in higher weights of positive component in the
 159 conditional probability as shown in Lemma 2, and may lead to a higher probability of classifying test
 160 input as +1, as formally stated in Proposition 2 in Section 2.3.

161 **ICL Accuracy.** After obtaining the decision boundary from Lemma 2, we finally provide the general
 162 formula of the ICL prediction accuracy. In the following, we consider a simplified scenario and derive
 163 the exact accuracy of ICL in Theorem 1.

164 **Theorem 1.** Under Assumption 1 and Assumption 2, and further assume $\sigma_+^2 = \sigma_-^2 = \sigma^2$ and
 165 $k \rightarrow \infty$, we have the following probability of correct prediction for each class.

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) = 1 - \Phi \left(\frac{(\theta_+^e - \frac{1}{2}(\hat{\theta}_+ + \hat{\theta}_-))^T}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|_2} + \frac{\sigma^2 \log z_k}{\sqrt{\sigma^2 + \sigma_{e+}^2} \|\hat{\theta}_- - \hat{\theta}_+\|_2} \right),$$

$$P(\text{correct}|y = -1, \{(x_i, y_i)\}, M) = \Phi \left(\frac{(\theta_-^e - \frac{1}{2}(\hat{\theta}_+ + \hat{\theta}_-))^T}{\sqrt{\sigma^2 + \sigma_{e-}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|_2} + \frac{\sigma^2 \log z_k}{\sqrt{\sigma^2 + \sigma_{e+}^2} \|\hat{\theta}_- - \hat{\theta}_+\|_2} \right),$$

166 where $\Phi(\cdot)$ is the cumulative distribution function of a standard Gaussian distribution.

167 To prove Theorem 1, we first obtain the marginal distribution of $x|y$ given the example distribu-
 168 tion in Assumption 2 to remove internal parameters, denoted as $P(x|y = +1)$. Then the ICL
 169 performance, i.e., prediction accuracy, can be computed via $P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) =$
 170 $\int_{f_{ICL}(x) \geq 0} P(x|y = +1) dx$. Similar steps apply for $y = -1$. The detailed proof of Theorem 1 can
 171 be found in Section B.3.

172 Theorem 1 describes how the ICL performance is affected by the interplay between the pre-training
 173 knowledge and the examples. For example, in addition to the model parameters $(\theta_+^e, \sigma^2, \sigma_{e+}^2)$,
 174 $P(\text{correct}|y = +, \{(x_i, y_i)\}, M)$ is further determined by $\hat{\theta}_- - \hat{\theta}_+$ and $\hat{\theta}_- + \hat{\theta}_+$. One key insight
 175 is that these two terms are mixtures of examples and pre-training knowledge. The example size k , the
 176 variance of data σ^2 , as well as pre-training distribution σ_M^2 will also affect their exact formulas.

177 A final note is that, under Theorem 1, the decision boundary is a hyperplane, and one can integrate
 178 the above two probabilities. We also provide technical discussions when the decision boundary is not
 179 a hyperplane. In such a case, the boundary is a sphere, and the details can be found in Section A.

180 2.3 Different Demonstration Scenarios

181 In the following, we extend the above results to investigate how ICL is affected in specific situations.
 182 We consider contradicting knowledge, imbalanced examples, and label noise.

183 To simplify the analysis, we assume that in the inference stage, $\#(y_i = +1) = \pi$. Since $\#(y_i =$
 184 $+1)/k - \pi \rightarrow 0$ in k , the additional fluctuation in $\#(y_i = +1)$ does not affect the result.

185 **Contradicting knowledge.** In practical applications, it is possible that the examples exhibit different
 186 or even contradicting knowledge of the pre-training. To study this case, we compare $\theta_+^e = -\theta_M =$
 187 $-\theta_-^e$ and $\theta_+^e = \theta_M = -\theta_-^e$, i.e., the input distribution in examples is the opposite/same to that of
 188 pre-training distribution. The following result is obtained based on Theorem 1 in these scenarios:

189 **Proposition 1** (Contradicting knowledge). Assume the conditions of Theorem 1 hold, and also
 190 assume $\sigma_{e+}^2, \sigma_{e-}^2 \rightarrow 0$, and $\pi = 0.5$ at the inference stage. Then when $k\sigma_M^2 \ll \sigma^2$, i.e., insufficient
 191 example size,

$$P(\text{correct}|y = +1, \theta_+^e = \theta_M = -\theta_-^e) - P(\text{correct}|y = +1, \theta_+^e = -\theta_M = -\theta_-^e)$$

$$\rightarrow \Phi \left(\frac{\|\theta_M\|}{\sqrt{\sigma^2}} \right) - \Phi \left(-\frac{\|\theta_M\|}{\sqrt{\sigma^2}} \right) > 0.$$

192 When $k\sigma_M^2 \gg \sigma^2$, i.e., sufficient example size, both $P(\text{correct}|y = +1, \theta_+^e = -\theta_M = -\theta_-^e)$ and
 193 $P(\text{correct}|y = +1, \theta_+^e = \theta_M = -\theta_-^e)$ converges to $1 - \Phi \left(-\|\theta_M\|/\sqrt{\sigma^2} \right)$.

194 The accuracy of $y = -1$ exhibits a similar behavior.

195 The proof of Proposition 1 can be found in Section B.4. We mainly follow the result in Theorem 1
 196 and calculate the probabilities under the specific scenario.

197 There are two observations in Proposition 1. First, when there are not enough examples and the
 198 pre-training knowledge contradicts to the knowledge in the examples and the test data, there is an
 199 obvious drop in ICL performance compared to the case when the knowledge matches. Second, when
 200 there are enough examples, the knowledge from the examples will dominate, and ICL performance
 201 of contradicting knowledge converges to that of matching knowledge.

202 **Imbalanced examples.** In the following, we consider the case where the two classes are imbalanced
 203 at the inference stage, i.e. $\pi \neq 0.5$. In this case, the value of π will impact the ICL prediction.

Proposition 2 (Imbalanced examples). *Under Assumption 1 and Assumption 2, assume σ^2 and σ_M^2 are constants, $\pi k \rightarrow \infty$ and $(1 - \pi)k \rightarrow \infty$, then*

$$f_{ICL}(x) \rightarrow \exp \left[-\frac{(x - \hat{\theta}_+)^T (x - \hat{\theta}_+)}{2\sigma^2} + \frac{(x - \hat{\theta}_-)^T (x - \hat{\theta}_-)}{2\sigma^2} \right] - \frac{1 - \pi}{\pi}.$$

204 When $\pi \rightarrow 0$, $P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \rightarrow 0$.

205 The proof of Proposition 2 is in Section B.5. In short, we follow Lemma 2 to obtain the decision
206 boundary. Then, we repeat the steps of Theorem 1 to obtain the conclusion.

207 Based on Proposition 2, when the examples for one class are much fewer than the other class, ICL
208 performance for the minor class will significantly drop. In addition, the parameter π learned from
209 pre-training will be overlooked.

210 **Label noise.** It is common that there exist label noises in the examples for ICL. For example, an
211 example x_i sampled from $N(\theta_-^e, \sigma_{e-}^2)$ may be labeled as $+1$. Therefore, we change the value of p_+^e
212 and p_-^e in the examples to see how these changes affect the ICL performance, the result of which is
213 summarized as follows:

214 **Proposition 3** (Label noise). *Under the conditions of Theorem 1, assume σ^2 and σ_M^2 are constants,
215 $\theta_M = \theta_+^e$ and $\theta_M = -\theta_-^e$, and $\sigma_{e+}^2, \sigma_{e-}^2 \rightarrow 0$. Also assume $\pi = 0.5$ at the inference stage. When
216 $1 - p_+^e - p_-^e < 0$, and $k \rightarrow \infty$, $P(\text{correct}|y = +1, p_+^e, p_-^e)$ increases in p_+^e , and $P(\text{correct}|y =$
217 $-1, p_+^e, p_-^e)$ increases in p_-^e .*

218 The proof of Proposition 3 can be found in Section B.6 and is a direct extension from Theorem 1.

219 In Proposition 3, recall that $p_+^e = p_-^e = 1$ implies no random flip on the example labels. Intuitively,
220 when keeping $p_-^e = 1$ and decreasing p_+^e , the positive class becomes a mixture of two Gaussian
221 distributions. In this case, $\hat{\theta}_+$ is closer to zero, and the decision boundary will shift towards $-\theta_M$.
222 Therefore, it is more likely that ICL predicts a negative label for x , which aligns with the change in
223 $P(\text{correct}|y = +1, p_+^e, p_-^e)$ and $P(\text{correct}|y = -1, p_+^e, p_-^e)$ in Proposition 3. When $1 - p_+^e - p_-^e =$
224 0, the decision boundary set $\{f_{ICL}(x) > 0\}$ will degenerate to either \emptyset or full space. Therefore, in
225 these special cases, the positive accuracy and negative accuracy will be either (0,1), (1,0), or (0.5,0.5).

226 Due to the page limit, we postpone all the simulations and real-data experiments to Appendix D. The
227 simulation results for the next section can also be found in Appendix E

228 2.4 Mean Reversion

229 When the fraction of positive and negative is fixed in the pre-training, we notice an interesting
230 phenomenon ‘‘Mean Reversion’’.

231 **Theorem 2 (Mean Reversion, informal version of Theorem 3).** *Let $frac$ denote the fraction of $+1$
232 among the set of labels in the pre-training set. Under some mild conditions, assume in each prompt
233 in pre-training, $frac$ is always a fixed π , then in the testing prompt: (1) If $\#(y_i = +1)/k < \pi$, then
234 the prediction of x is $+1$. (2) If $\#(y_i = +1)/k > \pi$, then the prediction of x is -1 .*

235 We direct the reader into Appendix B.7 for the formal statement and detailed proof. Theorem 4
236 indicates that the conditional probability of y is determined by the fraction of labels within the
237 pre-training set and the examples during inference, in addition to the inputs. A direct corollary is
238 that when the fraction of $y_i = +1$ is fixed as 0.5 during the pre-training, and all y_i are negative in the
239 inference stage, the prediction for y is always positive.

240 3 Conclusion

241 In this paper, we analyze the behavior of ICL in a binary classification model. We study the ICL
242 performance under different scenarios, including contradicting knowledge, imbalanced examples, and
243 label noise. In addition to the above analysis in which we assume examples are independently chosen
244 in pre-training, we also find out a counter-intuitive phenomenon when the examples are selected in
245 a dependent way. When fixing the number of positive labels and negative labels in the prompt, the
246 ICL prediction behaves in a mean-reversion manner. We believe that our observations and theoretical
247 results can provide deep insights into understanding ICL. A future direction could be to relax the
248 conditions in this paper and consider more general data distributions.

249 References

- 250 [1] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung,
251 Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models
252 encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- 253 [2] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models
254 in-context. *arXiv preprint arXiv:2306.09927*, 2023.
- 255 [3] Yingqian Cui, Jie Ren, Pengfei He, Jiliang Tang, and Yue Xing. Superiority of multi-head
256 attention in in-context linear regression. *arXiv preprint arXiv:2401.17426*, 2024.
- 257 [4] Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. *arXiv*
258 *preprint arXiv:2310.05249*, 2023.
- 259 [5] Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning. *arXiv preprint*
260 *arXiv:2402.18819*, 2024.
- 261 [6] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and
262 Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning
263 work? *arXiv preprint arXiv:2202.12837*, 2022.
- 264 [7] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of
265 in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 266 [8] Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large
267 language models are implicitly topic models: Explaining and finding good demonstrations for
268 in-context learning. In *Workshop on Efficient Systems for Foundation Models@ ICML2023*,
269 2023.
- 270 [9] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers
271 learn in-context? a case study of simple function classes. *Advances in Neural Information*
272 *Processing Systems*, 35:30583–30598, 2022.
- 273 [10] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What
274 learning algorithm is in-context learning? investigations with linear models. *arXiv preprint*
275 *arXiv:2211.15661*, 2022.
- 276 [11] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander
277 Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by
278 gradient descent. In *International Conference on Machine Learning*, pages 35151–35174.
279 PMLR, 2023.
- 280 [12] Chi Han, Ziqi Wang, Han Zhao, and Heng Ji. Explaining emergent in-context learning as kernel
281 regression. 2023.
- 282 [13] Kabir Ahuja, Madhur Panwar, and Navin Goyal. In-context learning through the bayesian prism.
283 *arXiv preprint arXiv:2306.04891*, 2023.
- 284 [14] Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce C
285 Ho, Carl Yang, and May Dongmei Wang. Ehragent: Code empowers large language models for
286 few-shot complex tabular reasoning on electronic health records. In *ICLR 2024 Workshop on*
287 *Large Language Model (LLM) Agents*, 2024.
- 288 [15] James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business
289 Media, 2013.
- 290 [16] Chen Dan, Yuting Wei, and Pradeep Ravikumar. Sharp statistical guarantees for adversarially
291 robust gaussian classification. In *International Conference on Machine Learning*, pages 2345–
292 2355. PMLR, 2020.
- 293 [17] Hossein Taheri, Ramtin Pedarsani, and Christos Thrampoulidis. Asymptotic behavior of
294 adversarial training in binary linear classification. In *2022 IEEE International Symposium on*
295 *Information Theory (ISIT)*, pages 127–132. IEEE, 2022.

- 296 [18] Ke Wang and Christos Thrampoulidis. Benign overfitting in binary classification of gaussian
297 mixtures. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and*
298 *Signal Processing (ICASSP)*, pages 4030–4034. IEEE, 2021.
- 299 [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
300 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
301 *processing systems*, 30, 2017.
- 302 [20] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
303 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 304 [21] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning.
305 In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48,
306 2009.
- 307 [22] Jeffrey L Elman. Learning and development in neural networks: The importance of starting
308 small. *Cognition*, 48(1):71–99, 1993.
- 309 [23] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien,
310 Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward
311 Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In
312 *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- 313 [24] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
314 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
315 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 316 [25] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
317 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*
318 *preprint arXiv:1804.07461*, 2018.
- 319 [26] Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large
320 language models are latent variable models: Explaining and finding good demonstrations for
321 in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 322 [27] Yiming Zhang, Shi Feng, and Chenhao Tan. Active example selection for in-context learning.
323 *arXiv preprint arXiv:2211.04486*, 2022.
- 324 [28] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
325 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
326 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 327 [29] Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. *arXiv*
328 *preprint arXiv:2310.15916*, 2023.
- 329 [30] Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau.
330 Function vectors in large language models. *arXiv preprint arXiv:2310.15213*, 2023.
- 331 [31] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen.
332 What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
- 333 [32] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically
334 ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv*
335 *preprint arXiv:2104.08786*, 2021.
- 336 [33] Zhiyong Wu, Yaoliang Wang, Jiacheng Ye, and Lingpeng Kong. Self-adaptive in-context
337 learning: An information compression perspective for in-context example selection and ordering.
338 *arXiv preprint arXiv:2212.10375*, 2022.
- 339 [34] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use:
340 Improving few-shot performance of language models. In *International conference on machine*
341 *learning*, pages 12697–12706. PMLR, 2021.

- 342 [35] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
343 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
344 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 345 [36] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can
346 gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers.
347 *arXiv preprint arXiv:2212.10559*, 2022.
- 348 [37] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to imple-
349 ment preconditioned gradient descent for in-context learning. *Advances in Neural Information*
350 *Processing Systems*, 36, 2024.
- 351 [38] Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is
352 provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint*
353 *arXiv:2307.03576*, 2023.
- 354 [39] Hong Jun Jeon, Jason D Lee, Qi Lei, and Benjamin Van Roy. An information-theoretic analysis
355 of in-context learning. *arXiv preprint arXiv:2401.15530*, 2024.
- 356 [40] Jane Pan. *What in-context learning “learns” in-context: Disentangling task recognition and*
357 *task learning*. PhD thesis, Princeton University, 2023.

358 The structure of the appendix is as follows. In Section A, we provide the discussion when the
 359 decision boundary in Lemma 2 is not a hyperplane. Section B collects the proof for all lemmas and
 360 theorems in the main content. Section C describes the simulation setups, and Section F includes
 361 additional experiment results corresponding to Section D and Section E for both independent and
 362 dependent scenarios. Section D provides detailed experiments and results, including simulations and
 363 real-data experiments. Section E provides a detailed discussion on the mean reversion phenomenon
 364 both theoretically and empirically. Section G discusses the related works.

365 A Technical Discussion when $\sigma_+^2 + \sigma_{\hat{\theta}_+}^2 \neq \sigma_-^2 + \sigma_{\hat{\theta}_-}^2$

General scenario In Theorem 1 and the propositions in Section 2.3, our assumptions aim to simplify the analysis so that the decision boundary is a hyperplane. When the examples are provided such that $\sigma_{\hat{\theta}_+}^2 \neq \sigma_{\hat{\theta}_-}^2$, the general intuition holds, but the decision boundary changes from a hyperplane into a sphere. To be specific, based on Lemma 2, we have

$$N(x) = \left(\sqrt{\frac{\sigma_-^2 + \sigma_{\hat{\theta}_-}^2}{\sigma_+^2 + \sigma_{\hat{\theta}_+}^2}} \right)^m \exp \left[-\frac{(x - \hat{\theta}_+)^T(x - \hat{\theta}_+)}{2(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)} + \frac{(x - \hat{\theta}_-)^T(x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)} \right].$$

366 Then for some constant a ,

$$\begin{aligned} N(x) &> a \\ \Leftrightarrow -\frac{(x - \hat{\theta}_+)^T(x - \hat{\theta}_+)}{2(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)} + \frac{(x - \hat{\theta}_-)^T(x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)} &> \log(a) - m \log \left(\sqrt{\frac{\sigma_-^2 + \sigma_{\hat{\theta}_-}^2}{\sigma_+^2 + \sigma_{\hat{\theta}_+}^2}} \right), \end{aligned}$$

where the decision boundary

$$-\frac{(x - \hat{\theta}_+)^T(x - \hat{\theta}_+)}{2(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)} + \frac{(x - \hat{\theta}_-)^T(x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)} = \log(a) - m \log \left(\sqrt{\frac{\sigma_-^2 + \sigma_{\hat{\theta}_-}^2}{\sigma_+^2 + \sigma_{\hat{\theta}_+}^2}} \right)$$

367 is a sphere.

368 In such a case, in Theorem 1, for

$$P(\text{correct} | y = +1) = \int_{f_{ICL}(x) \geq 0} P(x | y = +1) dx,$$

369 instead of directly using Φ to represent the probability, we use the noncentral Chi-square distribution
 370 to write the probability. The following is the definition of noncentral Chi-square distribution.

Definition 1 (Noncentral Chi-square distribution²). *Let (X_1, X_2, \dots, X_k) be k independent and normally distributed random variables with means μ_i and unit variances. Then the random variable*

$$\sum_{i=1}^k X_i^2$$

is distributed according to the noncentral chi-square distribution. It has two parameters: k which specifies the number of degrees of freedom and λ which is related to the mean of X_i s by

$$\lambda = \sum_{i=1}^k \mu_i^2.$$

371 **Case when $\sigma_+^2 = \sigma_-^2$** When $\sigma_+^2 = \sigma_-^2$, asymptotically, when $k \rightarrow \infty$, the difference between $\sigma_{\hat{\theta}_-}^2$
 372 and $\sigma_{\hat{\theta}_+}^2$ does not hurt the decision boundary. To explain this, based on the formula of $\sigma_{\hat{\theta}_-}^2$ and $\sigma_{\hat{\theta}_+}^2$,
 373 both of them are in $O(1/k)$, which quickly diminishes to zero in k . On the other hand, for the other
 374 terms in the decision boundary, e.g., $\hat{\theta}_+$ and $\hat{\theta}_-$ in $N(x)$, they converge to their expectation in a rate
 375 of $O(1/\sqrt{k})$. As a result, the effect of $\sigma_{\hat{\theta}_-}^2$ and $\sigma_{\hat{\theta}_+}^2$ are negligible compared to the other quantities in
 376 the decision boundary formula.

²https://en.wikipedia.org/wiki/Noncentral_chi-squared_distribution

377 **B Proofs**

378 **B.1 Proof of Lemma 1**

379 *Proof of Lemma 1.* Given the prior distribution of θ_-, θ_+, π and the data $\{(x_i, y_i)\}$, the likelihood
380 becomes

$$\begin{aligned}
& P(\theta_+, \theta_-, \pi | (x_1, y_1), \dots, (x_k, y_k), M) \\
& \propto P((x_1, y_1), \dots, (x_k, y_k) | \theta_+, \theta_-, \pi, M) P(\theta_+, \theta_-, \pi | M) \\
& = P(\theta_+, \theta_-, \pi) \prod_{i=1}^k P((x_i, y_i) | \theta_+, \theta_-, \pi) \quad (\text{Omit } M \text{ for simplicity}) \\
& = P(\theta_+, \theta_-, \pi) \\
& \quad \cdot \prod_{y_i=+1} \pi \frac{1}{(\sqrt{2\pi\sigma_+^2})^m} \exp\left[-\frac{1}{2\sigma_+^2}(x_i - \theta_+)^T(x_i - \theta_+)\right] \\
& \quad \cdot \prod_{y_i=-1} (1 - \pi) \frac{1}{(\sqrt{2\pi\sigma_+^2})^m} \exp\left[-\frac{1}{2\sigma_+^2}(x_i - \theta_+)^T(x_i - \theta_+)\right] \\
& = [\pi^{\#(y_i=+1)}(1 - \pi)^{\#(y_i=-1)}] P(\theta_+, \theta_-) \prod_{y_i=+1} \pi \frac{1}{(\sqrt{2\pi\sigma_+^2})^m} \exp\left[-\frac{1}{2\sigma_+^2}(x_i - \theta_+)^T(x_i - \theta_+)\right] \\
& \quad \cdot \prod_{y_i=-1} (1 - \pi) \frac{1}{(\sqrt{2\pi\sigma_+^2})^m} \exp\left[-\frac{1}{2\sigma_+^2}(x_i - \theta_+)^T(x_i - \theta_+)\right]
\end{aligned} \tag{1}$$

Posterior of π . Since all parameters are independent, we can obtain the posterior distribution of π as

$$P(\pi | \{(x_i, y_i)\}, M) \propto P(\pi | M) \pi^{\#(y_i=+1)} (1 - \pi)^{\#(y_i=-1)} \propto \pi^{\#(y_i=+1)} (1 - \pi)^{\#(y_i=-1)}.$$

381 Therefore, the posterior of π is *Beta*($\#(y_i = +1) + 1, \#(y_i = -1)k + 1$).

Posterior of θ_+, θ_- . The likelihood of θ_+ satisfies

$$\begin{aligned}
P(\theta_+ | \{(x_i, y_i)\} | M) & \propto P(\theta_+ | M) \prod_{y_i=+1} \pi \frac{1}{(\sqrt{2\pi\sigma_+^2})^m} \exp\left[-\frac{1}{2\sigma_+^2}(x_i - \theta_+)^T(x_i - \theta_+)\right] \\
& \propto \exp\left[-\frac{1}{2\sigma_M^2}(\theta_+ - \theta_M)^T(\theta_+ - \theta_M) - \frac{1}{2\sigma_+^2} \sum_{y_i=+1} (x_i - \theta_+)^T(x_i - \theta_+)\right],
\end{aligned}$$

which means that the posterior of θ_+ follows a Gaussian distribution, i.e.

$$\theta_+ \sim N\left(\frac{\sigma_+^2 \theta_M + \sigma_M^2 \sum_{y_i=+1} x_i}{\sigma_+^2 + \#(y_i = +1)\sigma_M^2}, \frac{\sigma_+^2 \sigma_M^2}{\sigma_+^2 + \#(y_i = +1)\sigma_M^2} I\right) = N(\hat{\theta}_+, \sigma_{\hat{\theta}_+}^2 I).$$

Similarly, the posterior of θ_- follows

$$\theta_- \sim N\left(\frac{\sigma_M^2 \sum_{y_i=-1} x_i - \sigma_-^2 \theta_M}{\sigma_-^2 + \#(y_i = -1)\sigma_M^2}, \frac{\sigma_M^2 \sigma_-^2}{\sigma_-^2 + \#(y_i = -1)\sigma_M^2} I\right) = N(\hat{\theta}_-, \sigma_{\hat{\theta}_-}^2 I).$$

382 □

383 **B.2 Proof of Lemma 2**

384 *Proof of Lemma 2.* Given $p_+ = p_- = 1$ and the posterior of θ_+, θ_-, π , we obtain that

$$P(x, y = +1 | \{(x_i, y_i)\}, M)$$

$$\begin{aligned}
&= \int_{\pi} \int_{\theta_+, \theta_-} \pi \frac{1}{\left(\sqrt{2\pi\sigma_+^2}\right)^m} \exp\left[-\frac{1}{2\sigma_+^2}(x - \theta_+)^T(x - \theta_+)\right] \\
&\quad \cdot P(\pi|\{(x_i, y_i)\}, M) P(\theta_+|\{(x_i, y_i)\}, M) P(\theta_-|\{(x_i, y_i)\}, M) d\pi d\theta_+ d\theta_- \\
&= \int_{\pi} \int_{\theta_+} \pi \frac{1}{\left(\sqrt{2\pi\sigma_+^2}\right)^m} \exp\left[-\frac{1}{2\sigma_+^2}(x - \theta_+)^T(x - \theta_+)\right] \\
&\quad \cdot P(\pi|\{(x_i, y_i)\}, M) P(\theta_+|\{(x_i, y_i)\}, M) d\pi d\theta_+ \\
&= \int_{\pi} \int_{\theta_+} \pi \frac{1}{\left(\sqrt{2\pi\sigma_+^2}\right)^m \left(\sqrt{2\pi\sigma_{\hat{\theta}_+}^2}\right)^m} \exp\left[-\frac{1}{2\sigma_+^2}(x - \theta_+)^T(x - \theta_+)\right] \\
&\quad \cdot P(\pi|\{(x_i, y_i)\}, M) \exp\left[-\frac{1}{2\sigma_{\hat{\theta}_+}^2}(\theta_+ - \hat{\theta}_+)^T(\theta_+ - \hat{\theta}_+)\right] d\pi d\theta_+ \\
&= \frac{\#(y_i = +1) + 1}{k + 2} \int_{\theta_+} \frac{1}{\left(\sqrt{2\pi\sigma_+^2}\right)^m \left(\sqrt{2\pi\sigma_{\hat{\theta}_+}^2}\right)^m} \exp\left[-\frac{1}{2\sigma_+^2}(x - \theta_+)^T(x - \theta_+)\right] \\
&\quad \cdot \exp\left[-\frac{1}{2\sigma_{\hat{\theta}_+}^2}(\theta_+ - \hat{\theta}_+)^T(\theta_+ - \hat{\theta}_+)\right] d\theta_+ \\
&= \frac{\#(y_i = +1) + 1}{k + 2} \frac{1}{\left(\sqrt{2\pi(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)}\right)^m} \exp\left[-\frac{1}{2(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)}(x - \hat{\theta}_+)^T(x - \hat{\theta}_+)\right].
\end{aligned}$$

Similarly, we have

$$P(x, y = -1|\{(x_i, y_i)\}, M) = \frac{\#(y_i = -1) + 1}{k + 2} \frac{1}{\left(\sqrt{2\pi(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)}\right)^m} \exp\left[-\frac{(x - \hat{\theta}_-)^T(x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)}\right].$$

385 Then we can obtain the predicted probability and decision boundary.

$$\begin{aligned}
P(y = +1|x, \{(x_i, y_i)\}, M) &= \frac{P(x, y = +1|\{(x_i, y_i)\}, M)}{P(x, y = +1|\{(x_i, y_i)\}, M) + P(x, y = -1|\{(x_i, y_i)\}, M)} \\
&= \frac{\frac{\#(y_i = +1) + 1}{k + 2} N(x)}{\frac{\#(y_i = +1) + 1}{k + 2} N(x) + \frac{\#(y_i = -1) + 1}{k + 2}}, \\
P(y = -1|x, \{(x_i, y_i)\}, M) &= \frac{P(x, y = -1|\{(x_i, y_i)\}, M)}{P(x, y = +1|\{(x_i, y_i)\}, M) + P(x, y = -1|\{(x_i, y_i)\}, M)} \\
&= \frac{\frac{\#(y_i = -1) + 1}{k + 2}}{\frac{\#(y_i = +1) + 1}{k + 2} N(x) + \frac{\#(y_i = -1) + 1}{k + 2}},
\end{aligned}$$

where

$$N(x) = \left(\sqrt{\frac{\sigma_+^2 + \sigma_{\hat{\theta}_-}^2}{\sigma_+^2 + \sigma_{\hat{\theta}_+}^2}}\right)^m \exp\left[-\frac{(x - \hat{\theta}_+)^T(x - \hat{\theta}_+)}{2(\sigma_+^2 + \sigma_{\hat{\theta}_+}^2)} + \frac{(x - \hat{\theta}_-)^T(x - \hat{\theta}_-)}{2(\sigma_-^2 + \sigma_{\hat{\theta}_-}^2)}\right].$$

386 Finally, the decision boundary is $\hat{y}_{ICL} = 1(f_{ICL}(x) > 0)$, where $f_{ICL}(x) = N(x) - \frac{\#(y_i = -1) + 1}{\#(y_i = +1) + 1}$.
387 \square

388 B.3 Proof of Theorem 1

389 In the following, we first provide the posterior of the parameters in a simplified scenario in Lemma 3,
390 and then use this result in Theorem 1 to derive the exact accuracy of ICL.

391 **Lemma 3.** Under the conditions of Theorem 1, when taking $k \rightarrow \infty$, we can simplify $\hat{\theta}_+$, $\hat{\theta}_-$, $\sigma_{\hat{\theta}_+}^2$,
 392 $\sigma_{\hat{\theta}_-}^2$, $N(x)$, $f_{ICL}(x)$ as follows:

$$\begin{aligned}\hat{\theta}_+ &= \frac{\sigma^2 \theta_M + \sigma_M^2 \sum_{y_i=+1} x_i}{\sigma^2 + \#(y_i = +1) \sigma_M^2}, \quad \hat{\theta}_- = \frac{\sigma_M^2 \sum_{y_i=-1} x_i - \sigma^2 \theta_M}{\sigma^2 + \#(y_i = -1) \sigma_M^2}, \\ \sigma_{\hat{\theta}_+}^2 &= \frac{\sigma^2 \sigma_M^2}{\sigma^2 + \#(y_i = +1) \sigma_M^2} \rightarrow 0, \quad \sigma_{\hat{\theta}_-}^2 = \frac{\sigma^2 \sigma_M^2}{\sigma^2 + \#(y_i = -1) \sigma_M^2} \rightarrow 0, \\ N(x) &= \exp \left[\frac{1}{\sigma^2} (\hat{\theta}_+ - \hat{\theta}_-)^T x - \frac{1}{2\sigma^2} (\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-) \right], \quad f_{ICL}(x) = N(x) - z_k.\end{aligned}$$

393 As mentioned in Section A, since $\sigma_{\hat{\theta}_+}^2$ and $\sigma_{\hat{\theta}_-}^2$ are negligible compared to other terms in $N(x)$, we
 394 remove them from f_{ICL} .

395 *Proof of Theorem 1.* We can compute the average ICL accuracy when test samples are sampled from
 396 the example distribution. We first derive the marginal distribution for test input x . For any fixed θ_+
 397 and θ_- , following the data generation assumption in Section 2.1, we have

$$\begin{aligned}P(x|y = +1) &= \frac{1}{(\sqrt{2\pi}\sigma^2)^m} \exp \left[-\frac{(x - \theta_+)^T (x - \theta_+)}{2\sigma^2} \right] \\ P(x|y = -1) &= \frac{1}{(\sqrt{2\pi}\sigma^2)^m} \exp \left[-\frac{(x - \theta_-)^T (x - \theta_-)}{2\sigma^2} \right].\end{aligned}$$

398 As a result, when integrating over all possible θ_+ and θ_- , it becomes

$$\begin{aligned}&P(x|y = +1) \\ &= \int_{\theta_+, \theta_-} P(x|y = +1, \theta_+, \theta_-) P_M(\theta_+) P_M(\theta_-) d\theta_+ d\theta_- \\ &= \int_{\theta_+, \theta_-} \frac{1}{(\sqrt{2\pi}\sigma^2)^m} \exp \left[-\frac{1}{2\sigma^2} (x - \theta_+)^T (x - \theta_+) \right] \\ &\quad \cdot \frac{1}{(\sqrt{2\pi}\sigma_{e+}^2)^m} \exp \left[-\frac{1}{2\sigma_{e+}^2} (\theta_+ - \theta_+^e)^T (\theta_+ - \theta_+^e) \right] \\ &\quad \cdot \frac{1}{(\sqrt{2\pi}\sigma_{e-}^2)^m} \exp \left[-\frac{1}{2\sigma_{e-}^2} (\theta_- - \theta_-^e)^T (\theta_- - \theta_-^e) \right] d\theta_+ d\theta_- \\ &= \int_{\theta_+} \frac{1}{(\sqrt{2\pi}\sigma^2)^m} \exp \left[-\frac{1}{2\sigma^2} (x - \theta_+)^T (x - \theta_+) \right] \\ &\quad \cdot \frac{1}{(\sqrt{2\pi}\sigma_{e+}^2)^m} \exp \left[-\frac{1}{2\sigma_{e+}^2} (\theta_+ - \theta_+^e)^T (\theta_+ - \theta_+^e) \right] d\theta_+ \\ &= \frac{1}{(\sqrt{2\pi}(\sigma^2 + \sigma_{e+}^2))^m} \exp \left[-\frac{1}{2(\sigma^2 + \sigma_{e+}^2)} (x - \theta_+^e)^T (x - \theta_+^e) \right].\end{aligned}$$

399 Similarly,

$$P(x|y = -1) = \frac{1}{(\sqrt{2\pi}(\sigma^2 + \sigma_{e-}^2))^m} \exp \left[-\frac{1}{2(\sigma^2 + \sigma_{e-}^2)} (x - \theta_-^e)^T (x - \theta_-^e) \right].$$

400 Then, we compute the probability of correct prediction for each class respectively.

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) = \int_{f_{ICL}(x) \geq 0} P(x|y = +1) dx$$

Based on Lemma 2, we know that

$$\{x : f_{ICL}(x) \geq 0\} = \left\{ x : (\hat{\theta}_+ - \hat{\theta}_-)^T x - \frac{\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-}{2} \geq \sigma^2 \log z_k \right\}.$$

401 Let $z = (\hat{\theta}_- - \hat{\theta}_+)^T x - \frac{\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-}{2} - \sigma^2 \log z_k$, then we have

$$z|y = +1, \{(x_i, y_i)\}, M \sim N \left((\hat{\theta}_+ - \hat{\theta}_-)^T \theta_+^e - \frac{\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-}{2} - \sigma^2 \log z_k, \|\hat{\theta}_- - \hat{\theta}_+\|_2^2 (\sigma^2 + \sigma_{e+}^2) \right),$$

402 which is still a Gaussian distribution. Therefore, we have

$$\begin{aligned} P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) &= \int_{z \geq 0} P(z|y = +1) dz \\ &= \left(1 - \Phi \left(-\frac{(\hat{\theta}_+ - \hat{\theta}_-)^T \theta_+^e - \frac{\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-}{2} - \sigma^2 \log z_k}{\|\hat{\theta}_- - \hat{\theta}_+\|_2 \sqrt{\sigma^2 + \sigma_{e+}^2}} \right) \right) \\ &= \left(1 - \Phi \left(\frac{\theta_+^e - (\hat{\theta}_+ + \hat{\theta}_-)/2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|} + \frac{\sigma^2 \log z_k}{\sqrt{\sigma^2 + \sigma_{e+}^2} \|\hat{\theta}_- - \hat{\theta}_+\|} \right) \right), \end{aligned}$$

403 and

$$\begin{aligned} P(\text{correct}|y = -1, \{(x_i, y_i)\}, M) &= \int_{z < 0} P(z|y = -1) dz \\ &= \Phi \left(-\frac{(\hat{\theta}_+ - \hat{\theta}_-)^T \theta_-^e - \frac{\hat{\theta}_+^T \hat{\theta}_+ - \hat{\theta}_-^T \hat{\theta}_-}{2} - \sigma^2 \log z_k}{\|\hat{\theta}_- - \hat{\theta}_+\|_2 \sqrt{(\sigma^2 + \sigma_{e-}^2)}} \right) \\ &= \Phi \left(\frac{(\theta_-^e - (\hat{\theta}_+ + \hat{\theta}_-)/2)^T}{\sqrt{\sigma^2 + \sigma_{e-}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|} + \frac{\sigma^2 \log z_k}{\sqrt{\sigma^2 + \sigma_{e-}^2} \|\hat{\theta}_- - \hat{\theta}_+\|} \right). \end{aligned}$$

404

□

405 **B.4 Proof of Proposition 1**

406 *Proof of Proposition 1.* Denote

$$\bar{x}_+ = \frac{1}{\#(y_i = +1)} \sum_{y_i = +1} x_i, \bar{x}_- = \frac{1}{\#(y_i = -1)} \sum_{y_i = -1} x_i, \bar{x} = \frac{1}{k} \sum_{i=1}^k x_i.$$

From Assumption 1 and Assumption 2, we know that

$$\bar{x}_- \sim N \left(\theta_-^e, \frac{\sigma_{e+}^2 + \sigma^2}{\#(y_i = +1)} I \right), \bar{x}_+ \sim N \left(\theta_+^e, \frac{\sigma_{e-}^2 + \sigma^2}{\#(y_i = -1)} I \right),$$

which implies that

$$\bar{x} \sim N \left(\frac{\#(y_i = -1)}{k} \theta_-^e + \frac{\#(y_i = +1)}{k} \theta_+^e, \left(\frac{\#(y_i = +1)}{k} \sigma_{e+}^2 + \frac{\#(y_i = -1)}{k} \sigma_{e-}^2 \right) I \right).$$

407 We rewrite \bar{x}_+ and \bar{x}_- via introducing zero-mean variables:

$$\bar{z}_+ = \bar{x}_+ - \theta_+^e, \bar{z}_- = \bar{x}_- - \theta_-^e.$$

408 When $\sigma_{e+}^2, \sigma_{e-}^2 \rightarrow 0$, the mean of \bar{z}_+ , \bar{z}_- , and \bar{x} are always zero.

409 Then, when $\theta_+^e + \theta_-^e = 0$, we have $\mathbb{E}\bar{x} = 0$, and

$$\begin{aligned} & P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \\ &= 1 - \Phi \left(\frac{\theta_+^e - (\hat{\theta}_+ + \hat{\theta}_-)/2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|} + \frac{\sigma^2 \log z_k}{\sqrt{\sigma^2 + \sigma_{e+}^2} \|\hat{\theta}_- - \hat{\theta}_+\|} \right) \\ &= 1 - \Phi \left(\left(\frac{\theta_+^e - \frac{k\sigma_M^2 \bar{x}}{2\sigma^2 + k\sigma_M^2}}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right)^T \frac{0.5k\sigma_M^2[\theta_-^e - \theta_+^e + (\bar{z}_- - \bar{z}_+)] - 2\sigma^2\theta_M}{\|0.5k\sigma_M^2[\theta_-^e - \theta_+^e + (\bar{z}_- - \bar{z}_+)] - 2\sigma^2\theta_M\|_2} \right). \end{aligned}$$

410 Now we compare the case of contradicted knowledge and matched knowledge.

411 **Contradict knowledge.** When $\theta_-^e = \theta_M = -\theta_+^e$, we have

$$\begin{aligned} & P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \\ &= 1 - \Phi \left(\left(\frac{-\theta_M - \frac{k\sigma_M^2 \bar{x}}{2\sigma^2 + k\sigma_M^2}}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right)^T \frac{0.5k\sigma_M^2(\bar{z}_- - \bar{z}_+) + (k\sigma_M^2 - 2\sigma^2)\theta_M}{\|0.5k\sigma_M^2(\bar{z}_- - \bar{z}_+) + (k\sigma_M^2 - 2\sigma^2)\theta_M\|_2} \right). \end{aligned}$$

412 When $k\sigma_M^2 \ll \sigma^2$, we have

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \rightarrow \left(1 - \Phi \left(\frac{\|\theta_M\|_2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right) \right).$$

413 When $k\sigma_M^2 \gg \sigma^2$, we have

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \rightarrow \left(1 - \Phi \left(\frac{-\|\theta_M\|_2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right) \right).$$

414 **Matched knowledge.** When $\theta_-^e = -\theta_M = -\theta_+^e$, we have

$$\begin{aligned} & P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \\ &= 1 - \Phi \left(\left(\frac{\theta_M - \frac{k\sigma_M^2 \bar{x}}{2\sigma^2 + k\sigma_M^2}}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right)^T \frac{0.5k\sigma_M^2(\bar{z}_- - \bar{z}_+) - (k\sigma_M^2 + 2\sigma^2)\theta_M}{\|0.5k\sigma_M^2(\bar{z}_- - \bar{z}_+) - (k\sigma_M^2 + 2\sigma^2)\theta_M\|_2} \right). \end{aligned}$$

415 When $k\sigma_M^2 \ll \sigma^2$, we have

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) = \left(1 - \Phi \left(\frac{-\|\theta_M\|_2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right) \right).$$

416 When $k\sigma_M^2 \gg \sigma^2$, we have

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) \rightarrow \left(1 - \Phi \left(\frac{-\|\theta_M\|_2}{\sqrt{\sigma^2 + \sigma_{e+}^2}} \right) \right).$$

417

□

418 B.5 Proof of Proposition 2

Proof of Proposition 2. When $k \rightarrow \infty$, we have $\hat{\theta}_+ \rightarrow \theta_M$ and $\hat{\theta}_- \rightarrow -\theta_M$, as well as $\sigma_{\hat{\theta}_+}^2 \rightarrow 0$, $\sigma_{\hat{\theta}_-}^2 \rightarrow 0$. In this case, the ICL decision boundary is still a hyperplane, and we can use the result in Lemma 2 and simplify the decision function into

$$f_{ICL}(x) = \exp \left[-\frac{(x - \hat{\theta}_+)^T (x - \hat{\theta}_+)}{2\sigma^2} + \frac{(x - \hat{\theta}_-)^T (x - \hat{\theta}_-)}{2\sigma^2} \right] - \frac{1 - \pi}{\pi}.$$

419 Then, following the same definition of z as Theorem 1, we can compute the ICL accuracy:

$$P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) = \int_{z \geq \log(\frac{1-\pi}{\pi})} P(z|y = +1) dz.$$

420 Since $\log \frac{1-\pi}{\pi}$ goes to ∞ when $\pi \rightarrow 0$, $P(\text{correct}|y = +1) \rightarrow 0$. \square

421 B.6 Proof of Proposition 3

422 *Proof of Proposition 3.* Following the definition of $\hat{\theta}_+$ and $\hat{\theta}_-$, when taking $k \rightarrow \infty$ and $\theta_+^e = \theta_M =$
423 $-\theta_-^e$, we obtain

$$\begin{aligned} \hat{\theta}_- + \hat{\theta}_+ &\rightarrow (1 + p_+^e - p_-^e)\theta_+^e + (1 - p_+^e + p_-^e)\theta_-^e = 2(p_+^e - p_-^e)\theta_M, \\ \hat{\theta}_- - \hat{\theta}_+ &\rightarrow (1 - p_+^e - p_-^e)(\theta_+^e - \theta_-^e) = 2(1 - p_+^e - p_-^e)\theta_M. \end{aligned}$$

424 Taking the above into the formula of $P(\text{correct}|y = +1)$ in Theorem 1, we obtain

$$\begin{aligned} P(\text{correct}|y = +1, \{(x_i, y_i)\}, M) &= 1 - \Phi \left(\frac{\theta_+^e - (\hat{\theta}_+ + \hat{\theta}_-)/2}{\sqrt{(\sigma^2 + \sigma_{e+}^2)}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|} + \frac{(\sigma^2 + \sigma_\theta^2) \log z_k}{\|\hat{\theta}_- - \hat{\theta}_+\| \sqrt{(\sigma^2 + \sigma_{e+}^2)}} \right) \\ &= 1 - \Phi \left(C_1(1 - p_+^e + p_-^e)\|\theta_M\| \text{sign}(1 - p_+^e - p_-^e) + C_2 \frac{\log z_k}{|1 - p_+^e - p_-^e|} \right), \end{aligned}$$

425 and

$$\begin{aligned} P(\text{correct}|y = -1, \{(x_i, y_i)\}, M) &= \Phi \left(\frac{(\theta_-^e - (\hat{\theta}_+ + \hat{\theta}_-)/2)^T}{\sqrt{\sigma^2 + \sigma_{e-}^2}} \frac{\hat{\theta}_- - \hat{\theta}_+}{\|\hat{\theta}_- - \hat{\theta}_+\|_2} + \frac{(\sigma^2 + \sigma_\theta^2) \log z_k}{\|\hat{\theta}_- - \hat{\theta}_+\| \sqrt{(\sigma^2 + \sigma_{e+}^2)}} \right) \\ &= \Phi \left(-C_1(1 + p_+^e - p_-^e) \text{sign}(1 - p_+^e - p_-^e) + C_2 \frac{\log z_k}{|1 - p_+^e - p_-^e|} \right). \end{aligned}$$

426 where $C_1 = 2 \frac{\|\theta_M\|}{\sqrt{(\sigma^2 + \sigma_{e+}^2)}} > 0$, $C_2 = \frac{\sigma^2}{\sqrt{(\sigma^2 + \sigma_{e+}^2)} \|\theta_M\|} > 0$. \square

427 B.7 Formal statement and proof of Theorem 3

428 **Theorem 3** (Dependent examples). Assume $\{(x_i, y_i)\}$ are not independent and are considered as a
429 sequence of inputs. Let frac denote the fraction of 1 among the set of labels in the pre-training set.
430 Assume frac approximately³ follows Beta(α, β) with $\alpha, \beta \rightarrow \infty$ and $\alpha/\beta \rightarrow \pi/(1 - \pi)$ for some
431 constant $\pi \in (0, 1)$, i.e., frac is around π with probability tending to 1. Further, assume that all y_i s
432 and y have an equal chance of being positive in pre-training. Then when $P(x|y = +1, \{(x_i, y_i)\}, M)$
433 and $P(x|y = -1, \{(x_i, y_i)\}, M)$ are both bounded and bounded away from zero, the following holds:

$$P(y = +1|x, \{(x_i, y_i)\}, M) \rightarrow \begin{cases} 1 & \text{if } \frac{\#(y_i=1)}{k+1} < \frac{\lfloor \pi(k+1) \rfloor - 1}{k+1} \\ 0 & \text{if } \frac{\#(y_i=1)}{k+1} > \frac{\lceil \pi(k+1) \rceil + 1}{k+1} \end{cases}.$$

434 *Proof of Theorem 3.* During pre-training, since there are k examples and one test sample in the
435 prompt, the fraction of positive labels can only take values in the form of $i/(k+1)$ for $i = 0, \dots, k$,
436 rather than a continuous variable in $[0, 1]$. As a result, to connect the distribution of frac with the
437 Beta distribution, we assume $k+1$ is odd and denote B as a random variable following Beta(α, β).
438 Then we set the following:

$$P(\text{frac} = i/(k+1)|M) = \begin{cases} P\left(B < \frac{1}{2(k+1)}\right) & i = 0 \\ P\left(\frac{2i-1}{2(k+1)} \leq B < \frac{2i+1}{2(k+1)}\right) & 0 < i < k+1 \\ P\left(B \geq \frac{2k+1}{2(k+1)}\right) & i = k+1 \end{cases}.$$

³The fraction number given a finite number of examples follows a discrete distribution. Here we approximate it to the Beta distribution and focus on the intuition. Details can be found in the proof.

439 In addition, we assume that all y_i s and y have equal chance of being positive.

440 Based on our assumption, when LLM learns from the pre-training data, it can exactly learn the
441 distribution of $frac$, and use the likelihood to make a decision when receiving the testing data.

442 In the testing stage, when receiving $\{(x_i, y_i)\}_{i \in [k]}$ and x . From the definition of conditional probab-
443 ility, we know that

$$P(y = +1|x, \{(x_i, y_i)\}, M) = \frac{P((x, y = +1)|\{(x_i, y_i)\}, M)}{P((x, y = +1)|\{(x_i, y_i)\}, M) + P((x, y = -1)|\{(x_i, y_i)\}, M)},$$

444 where

$$P((x, y = +1)|\{(x_i, y_i)\}, M) = P(x|y = +1, \{(x_i, y_i)\}, M)P(y = +1|\{(x_i, y_i)\}, M).$$

445 From the above, we need to figure out the following quantity:

$$\begin{aligned} & P((x, y = +1)|\{(x_i, y_i)\}, M) \\ &= P(x|y = +1, (x_i, y_i), M)P(y = +1|\{(x_i, y_i)\}, M) \\ &= P(x|y = +1, \{(x_i, y_i)\}, M)P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right), \end{aligned}$$

446 where

$$\begin{aligned} & P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right) \\ &= \frac{P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1}, \{(x_i, y_i)\} \middle| M\right)}{P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1}, \{(x_i, y_i)\} \middle| M\right) + P\left(frac = \frac{\#(y_i = +1)}{k + 1}, \{(x_i, y_i)\} \middle| M\right)}. \quad (2) \end{aligned}$$

To calculate $P(frac = (1 + \#(y_i = +1))/(k + 1), \{(x_i, y_i)\} \middle| M)$, when $frac = (1 + \#(y_i = +1))/(k + 1)$, it means that there are $1 + \#(y_i = +1)$ examples (and the query) which have a positive label. Given a total of $k + 1$ data, there are $C_{k+1}^{1 + \#(y_i = +1)}$ different combinations. As a result, for a fixed $\{(x_i, y_i)\}$, we have

$$P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1}, \{(x_i, y_i)\} \middle| M\right) = \frac{1}{C_{k+1}^{1 + \#(y_i = +1)}} P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| M\right).$$

Similarly, we obtain that

$$P\left(frac = \frac{\#(y_i = +1)}{k + 1}, \{(x_i, y_i)\} \middle| M\right) = \frac{1}{C_{k+1}^{\#(y_i = +1)}} P\left(frac = \frac{\#(y_i = +1)}{k + 1} \middle| M\right).$$

447 Taking the above into (2), it becomes

$$\begin{aligned} & P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right) \\ &= \frac{P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| M\right) C_{k+1}^{\#(y_i = +1)}}{P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| M\right) C_{k+1}^{\#(y_i = +1)} + P\left(frac = \frac{\#(y_i = +1)}{k + 1} \middle| M\right) C_{k+1}^{1 + \#(y_i = +1)}} \\ &= \frac{1}{1 + \frac{P\left(frac = \frac{\#(y_i = +1)}{k + 1} \middle| M\right) k - \#(y_i = +1) + 1}{P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| M\right) 1 + \#(y_i = +1)}}. \end{aligned}$$

448 To further look at the exact value of $P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right)$, we need to figure out

449 $P(frac = i/(k + 1) \middle| M)$ using the Beta distribution. Recall that the probability density function f
450 of $\text{Beta}(\alpha, \beta)$ satisfies

$$f(u) = \frac{u^{\alpha-1}(1-u)^{\beta-1}}{B(\alpha, \beta)}$$

451 for Beta function $B(\alpha, \beta)$. Recall that we assume that $\alpha/\beta = \pi/(1-\pi)$ for some $\pi \in (0, 1)$, and both
 452 $\alpha, \beta \rightarrow \infty$. When $u < (\alpha - 1)/(\alpha + \beta - 2) \approx \pi$, we have f is an increasing function in u , otherwise
 453 f is decreasing. This implies that the largest probability of $frac$ may be taken from $P(frac =$
 454 $(\lfloor \pi(k+1) \rfloor + 1)/(k+1)$), $P(frac = \lfloor \pi(k+1) \rfloor / (k+1))$ or $P(frac = (\lfloor \pi(k+1) \rfloor - 1)/(k+1))$.
 455 When $frac < (\lfloor \pi(k+1) \rfloor - 1)/(k+1)$, when α and β are large enough, one can obtain that

$$\frac{P\left(frac = \frac{\#(y_i=+1)}{k+1} | M\right)}{P\left(frac = \frac{1+\#(y_i=+1)}{k+1} | M\right)} \frac{k - \#(y_i = +1) + 1}{1 + \#(y_i = +1)} \rightarrow 0,$$

456 which implies that

$$P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right) = \frac{1}{1 + \frac{P(frac = \frac{\#(y_i=+1)}{k+1} | M)}{P(frac = \frac{1+\#(y_i=+1)}{k+1} | M)} \frac{k - \#(y_i=+1) + 1}{1 + \#(y_i=+1)}} \rightarrow 1.$$

457 Similarly, when $frac > (\lceil \pi(k+1) \rceil + 1)/(k+1)$,

$$P\left(frac = \frac{1 + \#(y_i = +1)}{k + 1} \middle| \{(x_i, y_i)\}, M\right) = \frac{1}{1 + \frac{P(frac = \frac{\#(y_i=+1)}{k+1} | M)}{P(frac = \frac{1+\#(y_i=+1)}{k+1} | M)} \frac{k - \#(y_i=+1) + 1}{1 + \#(y_i=+1)}} \rightarrow 0.$$

458 Finally, we put $P\left(frac = \frac{1+\#(y_i=+1)}{k+1} \middle| \{(x_i, y_i)\}, M\right)$ into $P(y = +1|x, \{(x_i, y_i), M\})$:

$$\begin{aligned} & P(y = +1|x, \{(x_i, y_i)\}, M) \\ &= \frac{P((x, y = +1)|\{(x_i, y_i)\}, M)}{P((x, y = +1)|\{(x_i, y_i)\}, M) + P((x, y = -1)|\{(x_i, y_i)\}, M)} \\ &= \frac{P(x|y = +1, \{(x_i, y_i)\}, M)P(frac = \frac{1+\#(y_i=+1)}{k+1} | \{(x_i, y_i)\}, M)}{P(x|y = +1, \{(x_i, y_i)\}, M)P(frac = \frac{1+\#(y_i=+1)}{k+1} | \{(x_i, y_i)\}, M) + P(x|y = -1, \{(x_i, y_i)\}, M)P(frac = \frac{\#(y_i=+1)}{k+1} | \{(x_i, y_i)\}, M)}. \end{aligned}$$

459 When $P(x|y = +1, \{(x_i, y_i), M\})$ and $P(x|y = -1, \{(x_i, y_i), M\})$ are both bounded and bounded
 460 away from zero, we have

$$P(y = +1|x, \{(x_i, y_i)\}, M) \rightarrow \begin{cases} 1 & \text{if } \frac{\#(y_i=+1)}{k+1} < \frac{\lfloor \pi(k+1) \rfloor - 1}{k+1} \\ 0 & \text{if } \frac{\#(y_i=+1)}{k+1} > \frac{\lceil \pi(k+1) \rceil + 1}{k+1} \end{cases},$$

461 which completes the proof. \square

462 C Simulation setups

463 In this section, we provide details of experimental setups for simulation.

464 **Model structure.** We pre-train a decoder-only Transformer [19] from the GPT-2 [20] family. This
465 model has 12 layers, 8 attention heads, and a 256-dimensional embedding space. The model input
466 takes the form of $(x_1, y_1, x_2, y_2, \dots)$. In our training, $x_i \in \mathbb{R}^5$ and $y_i \in \{1, -1\}$. We map y_i to the
467 same dimension of x_i by appending zeros. Then the whole prompt will be projected into the latent
468 embedding space of the Transformer through a (learnable) MLP layer. Another (learnable) MLP
469 layer is used to project embeddings back to scalars in the output.

470 **Pre-training.** We train the model using a cross-entropy loss function for binary classification. We
471 sample a batch of random prompts at each training step and update the model through a gradient
472 update. We train with a batch size of 64 and for 50k steps. This training is done from scratch, that is,
473 we do not fine-tune a pre-trained language model, nor do we train on actual text. Following previous
474 work [9], we also use curriculum learning [21, 22]. In particular, we start with a shorter length of
475 prompts (10 input-output pairs) and increase the length by 2 every 2000 training steps. For the other
476 hyperparameters, e.g., learning rate, we use the default values as in [9].

477 **Pre-training data.** We follow the data generation model in Section 2.1. We first select label
478 $y \in \{+1, -1\}$ with probability π (positive probability). Then for inputs with positive labels, we
479 first sample a mean value θ_+ from a Gaussian distribution $N(\theta_M, \sigma_M^2 I)$, and then sample data x
480 from Gaussian distribution $N(\theta_+, \sigma^2 I)$; similar for the inputs with label -1 , we sample θ_- from
481 $N(-\theta_M, \sigma_M^2 I)$, and sample x from $N(\theta_-, \sigma^2 I)$. Specifically, we let $\theta_M = 0.5\mathbf{1}$, $\sigma_M^2 = \sigma^2 = 1$.
482 During the pre-training, to ensure the transformer can learn the population information rather than
483 overfitting a particular set of data, we sample a new pair of (θ_+, θ_-) for each iteration and generate
484 corresponding sample pair (x_i, y_i) .

485 **Computation resources.** Both simulations and real-world experiments are running on a server with
486 8 Nvidia RTX A6000 GPU (48G GPU memory each) and 32 AMD EPYC 7302 16-Core Processors.

487 D Experiments

488 In this section, we empirically verify the analysis in Section 2. In summary, both simulation and
489 real-data experiments are consistent Section 2⁴.

490 D.1 Simulation

491 To set up the experiment, we pre-train a decoder-only Transformer [19] from the GPT-2 [20] family.
492 We follow Section 2.1 to construct the pre-training data and follow [9] to perform next-token
493 prediction to estimate all y_i s and y_{query} . During the pre-training, we sample a new pair of (θ_+, θ_-)
494 for each iteration and generate corresponding demonstration examples. A detailed setting for
495 simulation can be found in Appendix C.

496 We implement the scenarios as in Section 2.3:

497 **Contradicting knowledge** We pre-train the model with $\theta_M = 0.5\mathbf{1}_5$, $\sigma_M^2 = 1$, $\sigma^2 = 1$. During
498 the inference stage, we generate examples and test data with $\theta_+^e = -0.5\mathbf{1}_5$, $\theta_-^e = 0.5\mathbf{1}_5$ which
499 contradicts the pre-training distribution. We let $\sigma_{e+}^2 = \sigma_{e-}^2 = 1$, and test on various $\sigma^2 \in \{1, 2, 4\}$.
500 The results for $\sigma^2 = 2, 4$ are postponed to Section F.1. We compute the ICL accuracy for each
501 class when k increases. For comparison, we also generate examples with matched knowledge
502 ($\theta_+^e = 0.5\mathbf{1}_5 = -\theta_-^e$) and examine the ICL accuracy.

503 The results can be found in Figures 1, where the X-axis represents the number of examples k , and the
504 Y-axis is the ICL accuracy. The red dash denotes σ^2/σ_M^2 based on Proposition 1 and $\sigma^2/\sigma_M^2 = 1$
505 in our simulation. There are two observations. First, when $k \leq \sigma^2/\sigma_M^2$, the ICL performance of
506 contradicting knowledge is worse than that of matching knowledge, verifying that the transformer
507 heavily relies on the pre-training knowledge when there are limited examples. Second, when k

⁴Code is available in <https://anonymous.4open.science/r/ICL-understanding-classification-DC1C>

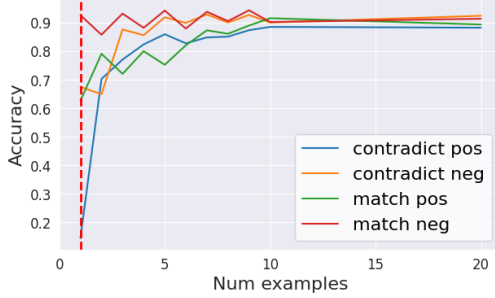


Figure 1: Contradicting knowledge when $\sigma^2 = 1$

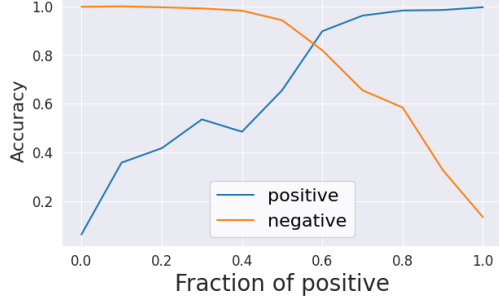


Figure 2: Imbalanced examples

508 increases, the ICL performance for contradicting knowledge increases to around 87% when $k = 20$,
 509 indicating that the knowledge from the examples will dominate when k is large.

510 **Imbalanced examples** We pre-train the model with $\theta_M = 0.5\mathbf{1}_5, \sigma_M^2 = 1, \sigma^2 = 1$. During
 511 the inference stage, we generate examples and the test data with $\theta_+^e = 0.5\mathbf{1}_5, \theta_-^e = -0.5\mathbf{1}_5$ and
 512 $\sigma_{e+}^2 = \sigma_{e-}^2 = 1, \sigma^2 = 1$. We test with various fraction $\pi \in \{0, 0.1, 0.2, \dots, 0.9, 1.0\}$. In Figure 2,
 513 the X-axis represents the fraction of positive examples among all examples in the demonstration, i.e.,
 514 π . We can observe that when the fraction of positive is increasing, ICL accuracy for positive inputs is
 515 increasing and finally reaches 100%, while ICL accuracy for negative inputs is decreasing, which is
 516 the same as described in Proposition 2.

517 **Label noise** We pre-train the model with $\theta_M = 0.5\mathbf{1}_5, \sigma_M^2 = 1, \sigma^2 = 1$. During the inference
 518 stage, we generate examples and the test data with $\theta_+^e = 0.5\mathbf{1}_5, \theta_-^e = -0.5\mathbf{1}_5$, and fix $\sigma_{e+}^2 = \sigma_{e-}^2 =$
 519 $1, \sigma^2 = 0.01$. The example size k is 100, and π is 0.5. For the examples in each class, we randomly
 520 flip their label with probability $1 - p_+^e, 1 - p_-^e$ respectively, and test ICL accuracy for each class for
 521 $1 - p_+^e, 1 - p_-^e \in \{1, 0.9, 0.8, \dots, 0.1, 0\}$. The ICL accuracies for each class and the overall result
 522 are summarized in Figures 3.

523 The phenomenon in these heatmaps is consistent with our conclusion in Proposition 3. Take the
 524 ICL accuracy of the positive class as an example (the first figure in Figure 3); we observe that when
 525 the flipping probability in the negative class is fixed, smaller flipping probability (higher p_+^e) in the
 526 positive class usually leads to higher accuracy in the positive class. Moreover, the diagonal from the
 527 bottom left to the upper right represents cases when $p_+^e + p_-^e = 1$ and it is obvious that the positive
 528 accuracy is either approximately 0, 0.5, or 1. Similar observations can be found for the negative class
 529 as well (the middle panel of Figure3). In terms of the overall accuracy, only when both p_+^e and p_-^e
 530 are close to 1, the overall accuracy is greater than 80%.

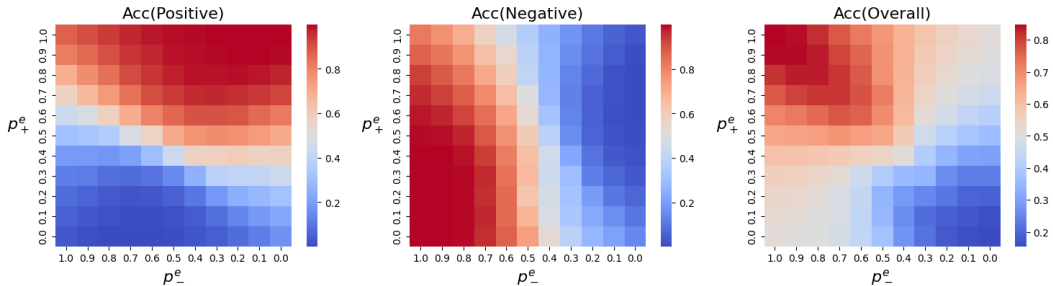


Figure 3: Simulation results on positive and negative accuracy facing label noises.

531 **D.2 Real-Data Experiment**

532 In this subsection, we conduct experiments on real datasets to show that the theoretical insights in
 533 Section 2 also align with the practical scenarios.

534 We consider two popular pre-trained LLMs, Pythia-6.9B [23] and Llama2-7B [24]. We test on a
 535 sentiment analysis dataset, SST2 dataset [25], which is also a binary classification task (labeled as
 536 “positive” and “negative”). During the inference, we randomly select k samples from the training set
 537 as examples and compute the ICL accuracy for each class. We repeat the process 10 times and record
 538 the average accuracy. If not specified, $k = 50$.

539 **Label noise** Similar to the simulation, we also randomly flip the label of examples from positive
 540 and negative classes and the correct probability is p_+^e, p_-^e respectively. We let $p_+^e, p_-^e \in$
 541 $\{0.0, 0.1, \dots, 0.9, 1.0\}$ and record ICL accuracy in each class. Results are shown in Figure 4, 5, 6,
 542 7. It can be consistently observed that when p_-^e is fixed, larger p_+^e leads to higher accuracy in the
 543 positive class (Figure 4 and 6); when p_+^e is fixed, larger p_-^e leads to higher accuracy in the negative
 class (Figure 5 and 7). This observation is consistent with our analysis in Proposition 3.

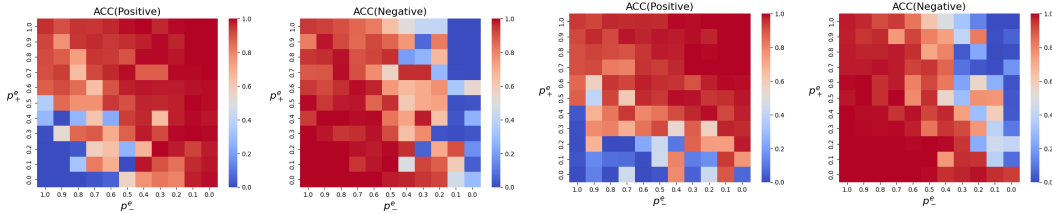


Figure 4: Pythia-6.9B: Figure 5: Pythia-6.9B: Figure 6: Llama2-7B: Figure 7: Llama2-7B:
 ICL acc, positive class. ICL acc, negative class. ICL acc, positive class. ICL acc, negative class.

544

545 **Imbalanced examples** We also conduct experiments when the fraction π of examples from the positive
 546 class is not 0.5. Specifically, we test with $\pi \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$.
 547 As depicted in Figure 8 and 9, when the number of positive examples increases, the accuracy of the
 548 positive class increases as that of the negative class decreases, which also supports our analysis in
 Proposition 2.

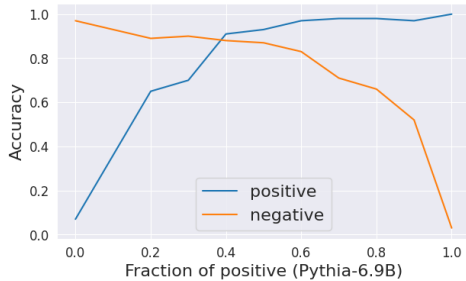


Figure 8: Imbalance examples, Pythia-6.9B.

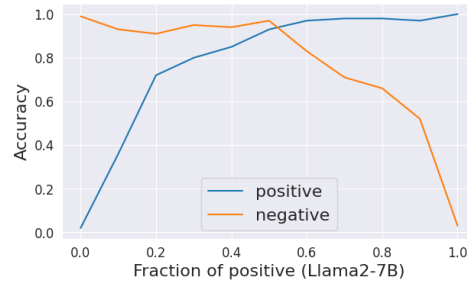


Figure 9: Imbalance examples, Llama2-7B.

549

550 E Mean Reversion in ICL with Dependent Examples

551 The previous analysis and experiments provide a comprehensive understanding of the effect of
 552 pre-training and examples under the independent-example scenario. However, it is also common
 553 when examples are sampled dependently, especially when examples are strategically selected to
 554 serve a specific objective, such as ensuring a balanced representation of 50% positive and 50%
 555 negative examples to prevent dataset imbalance [26, 27]. Surprisingly, we discover a counter-intuitive
 556 phenomenon, named as “**mean reversion**”, under this scenario. We first empirically illustrate this
 557 phenomenon and then provide a theoretical analysis.

558 **Empirical illustration of Mean Reversion** We follow a similar procedure as introduced in Section
 559 D, while the difference is that during the pre-training stage, we fix the fraction of positive labels
 560 (examples + test data) to be exactly 0.5. This differs from the independent case since the fraction π

561 may fluctuate around 0.5 instead of strictly equal to 0.5. During the inference, in-context examples are
 562 sampled from both classes but are all labeled as positive or negative. We test with various fractions of
 563 positive examples in the prompt to see how the prediction for the test input x is affected. Results are
 564 shown in Figure 10 11.

565 In Figure 10, all examples are labeled as positive and the X-axis reflects the fraction of examples
 566 truly from the positive class; while in 11, all examples are labeled as negative and the X-axis reflects
 567 the fraction of examples truly from the negative class. There are four lines in the two figures. The
 568 “Positive”/“Negative” refers to the probability of the prediction being positive/negative when the
 569 correct label is positive/negative. The “Total pos”/“Total neg” represents the marginal probability of
 570 the prediction being positive/negative.

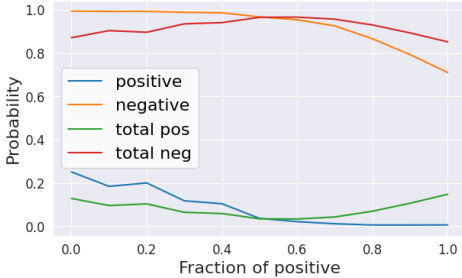


Figure 10: All y_i s are positive.

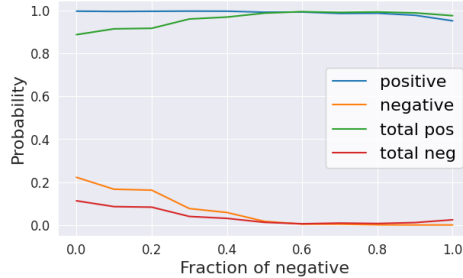


Figure 11: All y_i s are negative.

571 As shown in Figure 10, we can see that regardless of the change in the fraction of examples’ true
 572 classes (X-axis), we always obtain a low positive rate for the true positive test data, and the overall
 573 positive rate is low. This contradicts the independent case in Figure 2. Intuitively, since LLM learns
 574 that the fraction of positive labels is exactly 0.5 in the pre-training, at the inference stage, the joint
 575 label distribution of examples and test input appears to converge to 0.5-0.5⁵.

576 In the following, we provide a rigorous theoretical analysis to explain this phenomenon:

577 **Theorem 4 (Mean Reversion, informal version of Theorem 3).** *Let $frac$ denote the fraction of +1*
 578 *among the set of labels in the pre-training set. Under some mild conditions, assume in each prompt*
 579 *in pre-training, $frac$ is always a fixed π , then in the testing prompt: (1) If $\#(y_i = +1)/k < \pi$, then*
 580 *the prediction of x is +1. (2) If $\#(y_i = +1)/k > \pi$, then the prediction of x is -1.*

581 We direct the reader into Appendix B.7 for the formal statement and detailed proof. In short, when cal-
 582 culating $P((x, y = +1)|\{(x_i, y_i)\}, M) = P(x|y = +1, \{(x_i, y_i)\}, M)P(y = +1|\{(x_i, y_i)\}, M)$,
 583 since the examples are not independent, we need to follow the dependency among y_i s and
 584 y to determine $P(y = +1|\{(x_i, y_i)\}, M)$, which is determined by the relationship between
 585 $\#(y_i = +1)/(k + 1)$ in the testing data and $frac$ in pre-training.

586 Theorem 4 indicates that the conditional probability of y is determined by the fraction of labels within
 587 the pre-training set and the examples during inference, in addition to the inputs. A direct corollary is
 588 that when the fraction of $y_i = +1$ is fixed as 0.5 during the pre-training, and all y_i are negative in
 589 the inference stage, the prediction for y is always positive. This is consistent with the observation in
 590 Figure 10, 11.

591 F Additional Experiment Results

592 F.1 Simulation for Independent Exampels

593 Figure 12, 13, 14 represents additional results corresponding to the contradict knowledge setting in
 594 Figure 1. The observations are similar to Figure 1.

595 F.2 Mean Reversion

596 We further pre-train GPT models with different fractions ($frac = 0.2, 0.5, 0.8$) and test the posterior
 597 distribution of labels when the fraction of positive labels within examples varies. We do not add noise

⁵This is similar to the “mean reversion” in certain stochastic differential equations (SDEs) where the variable tends to move toward a long-term average over time, thus we also name our observation as “mean reversion”.

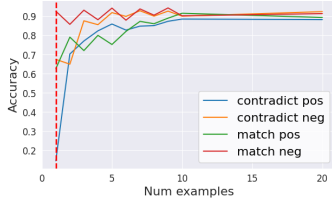


Figure 12: $\sigma^2 = 1$

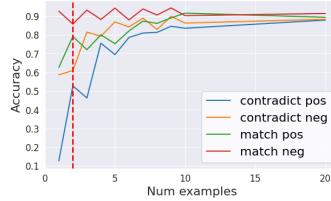


Figure 13: $\sigma^2 = 2$

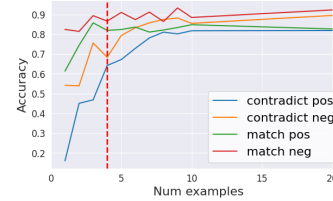


Figure 14: $\sigma^2 = 4$

598 and keep other settings unchanged. We observe a dramatic change around 0.2,0.5,0.8 respectively in
 599 Figure 15, and these figures directly verify our results: in Theorem 3, the cutting point are 0.2, 0.5, 0.8
 600 respectively in the three settings.

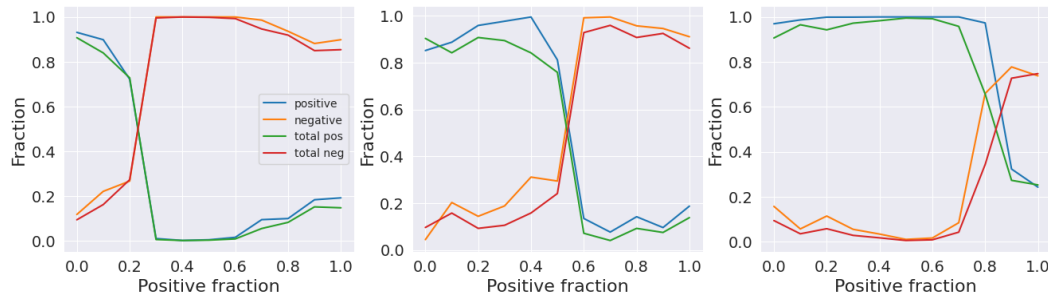


Figure 15: Pre-training with a fraction of positive 0.2.

601 G Related Works

602 **Empirical findings of ICL.** There are many empirical studies working on understanding ICL. [28]
 603 first reveals that LLMs can learn from examples, and refers to it as in-context learning (ICL). Later,
 604 to investigate the properties of ICL, [9] empirically shows that a transformer-based model can learn
 605 linear functions in context. [29, 30] find that transformer models can encode input-output relationships
 606 in the hidden space of attention layers. [6] observes that the key respects of the demonstration are
 607 label space, distribution of input, and format of the prompt. They also notice that randomly replacing
 608 labels barely hurts the performance when the example size is not large. [31, 32, 33] reveals the
 609 importance of examples, including orders and templates. More works are proposed to select examples
 610 [32, 26, 27] or design prompts [34, 35] to improve the ICL performance.

611 **Theoretical understanding of ICL.** To theoretically understand ICL, one popular line of research is
 612 to treat the ICL process as an implicit gradient descent procedure on examples. [2] shows that a single
 613 linear self-attention layer trained by gradient flow results in a competitive prediction error with the
 614 best linear predictor during ICL. [10, 11, 36] shows that one attention layer can be exactly constructed
 615 to perform gradient descent. [37, 38] further prove that under some conditions, a transformer with
 616 one or more attention layers trained on noisy linear regression task minimizing the pre-training loss
 617 will implement gradient descent algorithm on examples. [12] show that ICL can asymptotically
 618 converge to kernel regression as the number of examples increases.

619 Another line of research focuses on Bayesian inference. [7] first leverages a Hidden Markov Model
 620 to represent the pre-training data and prove that a transformer trained on such data exhibits the ICL
 621 ability. [39] introduces an information-theoretic tool to show how ICL error decays in the number
 622 and length of examples. [5] introduces a probabilistic model to understand two modes of ICL, i.e.,
 623 task learning and task retrieval [40], on the linear regression tasks.

624 However, these analyses primarily focus on regression tasks with continuous outputs, and lack precise
 625 quantification for classification scenarios. Furthermore, they often overlook scenarios where the
 626 distribution of in-context examples diverges from pre-training data, such as cases of label noise,
 627 imbalanced examples, or contradictory information. Our work addresses these limitations by focusing

628 on classification problems, providing exact quantification of example effects on predictions, and
629 offering insights into the impact of label noise, imbalanced examples, and contradictory knowledge
630 on in-context predictions.