

Do LLMs Play Dice? Exploring Probability Distribution Sampling in Large Language Models for Behavioral Simulation

Anonymous ACL submission

Abstract

With the rapid advancement of large language models (LLMs) for handling complex language tasks, an increasing number of studies are employing LLMs as agents to emulate the sequential decision-making processes of humans often represented as Markov decision-making processes (MDPs). The actions in MDPs adhere to specific probability distributions and require iterative sampling. This arouses curiosity regarding the capacity of LLM agents to comprehend probability distributions, thereby guiding the agent’s behavioral decision-making through probabilistic sampling and generating behavioral sequences. To answer the above question, we divide the problem into two main aspects: sequence simulation with known probability distribution and sequence simulation with unknown probability distribution. Our analysis indicates that LLM agents can understand probabilities, but they struggle with probability sampling. Their ability to perform probabilistic sampling can be improved to some extent by integrating coding tools, but this level of sampling precision still makes it difficult to simulate human behavior as agents.

1 Introduction

With the rapid development of artificial intelligence technology, large language models (LLMs), as a key component, have demonstrated powerful language understanding and generation capabilities (Zhao et al., 2023; Kojima et al., 2022). In addition to language generation, people have begun to explore the field of introducing LLMs into the decision-making process, such as simulating human sequential decision-making processes as an agent, and LLM agents have gradually shown their excellent potential in simulating human behaviors and decision-making (Li et al., 2023; Gao et al., 2023; Leng and Yuan, 2024; Brand et al., 2023; Horton, 2023). This brings great convenience to research in fields such as computer science, behav-

ioral science, psychology, and sociology. For example, the core conversational robots (Zhou et al., 2024b) and social robots of LLMs can more realistically simulate human speech and behavior. LLM agents have demonstrated promising results in various tasks simulating human behavior (Taylor et al., 2022; Yuan et al., 2022; Aher et al., 2023; He et al., 2024). However, whether LLM agents can effectively simulate human behavior sequences remains an open question and requires careful validation.

Human behavior can often be effectively modeled using Markov decision processes (MDPs). Numerous studies have applied this framework to simulate human behaviors such as planning (Callaway et al., 2017; Guojiang et al., 2007; Redding et al., 2012). MDPs consist of four key components: states, actions, transition probabilities, and rewards. Within this framework, an agent in a given state selects an action, leading to a transition to the next state based on defined probabilities, and receives associated rewards. The parallel between human behavior and MDPs stems from the tendency of decision-making processes to depend on current circumstances and available actions, rather than solely on past experiences. When confronted with a situation, humans typically make decisions based on their present state and the feasible choices, aiming to optimize future benefits or satisfaction. This mirrors the states, action selections, and rewards accrued through state transitions in MDPs.

In MDPs, sequences of actions follow a specific probability distribution. Similarly, human behavioral sequences theoretically adhere to probabilistic distributions as well. Considering the feasibility of using MDPs to model human behavior, an LLM-based agent should generate actions that match a probability distribution if it truly simulates human behavior. But can LLM agents do this? The most direct way to validate this is to explore their probability distribution sampling capabilities. Figure 1 illustrates our approach to assessing LLM’s under-

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

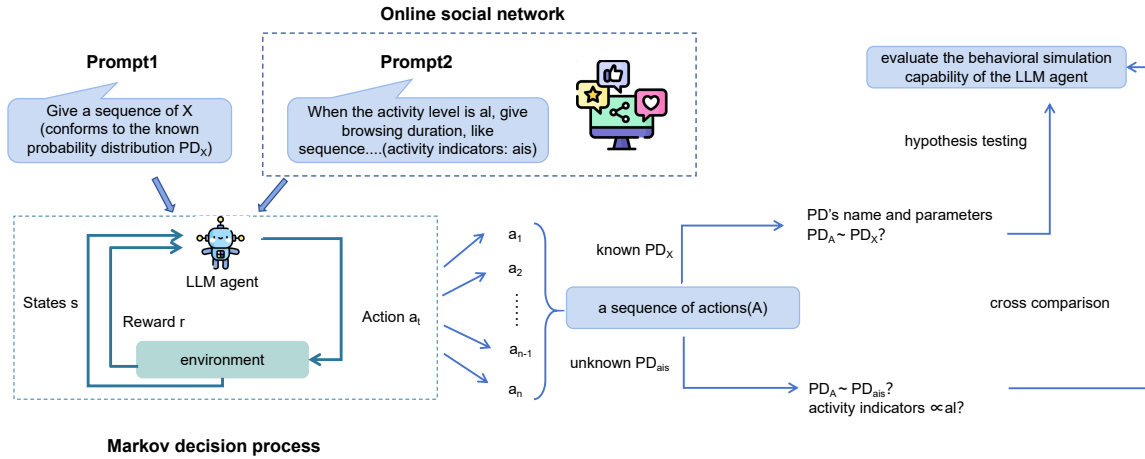


Figure 1: The decision-making process of the LLM agent is a MDP, and the generated action sequence A conforms to a certain probability distribution. We input *Prompt1* for the known probability distribution and *Prompt2* for the unknown probability distribution, analyze the probability distribution of A generated by the LLM agent, and finally evaluate the behavioral simulation capability of the LLM agent.

standing and sampling capabilities of probability distributions. We outline two tasks of varying complexity: simulating scenarios with known probability distributions and generating sequences with probability distribution is not clear. Our evaluation of multiple LLMs in these scenarios aims to assess the potential and limitations of LLM agents in simulating human decision-making processes.

In the first scenario, we propose questions about known probability distributions. Using these questions to prompt LLM agents for probability sampling, we analyze both the probability distributions and the generated sequences answered by the LLM agents. However, obtaining an accurate probability distribution is often challenging in practical scenarios. Therefore, we shift our focus to sequence generation in situations where the probability distribution is unclear. Online social networks are increasingly utilized in real-life scenarios (Sakaki et al., 2013; Panzarasa et al., 2009), prompting us to use these networks as contextual backgrounds. By observing varying activity levels corresponding to different behavior sequences generated by LLM agents, we indirectly assess the rationality of the sequences sampled. Ultimately, we instruct LLM agents to sample from probability distributions by generating code, aiming to evaluate whether programming tools can assist LLM agents in improving their sampling success rate.

Our analysis suggests that the probability sampling performance of LLM agents is not very good.

While LLM agents demonstrate some understanding of simple probability distributions, their success rates in sampling sequences adhering to known distributions are limited. Moreover, they struggle with unknown probability distributions. However, when coupled with programming tools, LLM agents can achieve higher success rates in sampling these distributions by invoking appropriate functions. Yet, their abilities to sample from complex probability distributions remain comparatively weak, and even with the integration of programming tools, enhancements in sampling capability are elusive. Consequently, prudent consideration is necessary before employing LLM directly for simulating human behavior.

Our main contributions are:

- Explore the LLM agent’s ability to understand and sample probability distributions from both the perspectives of known and unknown probability distributions.
- A method of using programming tools to assist sampling is proposed to improve the effect of the LLM agent in the probability distribution sampling process and is effective in simple probability distributions.
- Our experimental results reveal the limitations of LLM agents in terms of probability distributions when simulating human behavior.

2 Background

2.1 MDPs of Human Behavior

The modeling of human behavior in MDPs could include the following aspects (Tan and Cheng, 2009):

Human behavior can be understood as transitions between different states of state space (S), which may include emotional states, social identities, and more. In each state, humans can choose from a series of possible actions or decisions of action space (A). These actions lead to state changes with certain probabilities, represented by the state transition probability P . Human behavior is driven by rewards and punishments, described by the reward function R . To maximize benefits, humans adopt various strategies, captured by a policy π , which maps the probability of taking action a in state s .

Humans typically have long-term goals, such as maximizing life satisfaction, career achievement, or health. The objective in an MDP is to find a policy π that maximizes the expected cumulative reward. A common objective function is the total discounted return G_t :

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \quad (1)$$

where t represents current time, γ ($0 \leq \gamma \leq 1$) is the discount factor, representing the present value of future rewards.

Based on the above discussion, human behavior can usually be regarded as conforming to a certain probability distribution. Given characteristics, there is a defined probability distribution (Viti et al., 2015). Some simple behaviors conform to common probability distributions (Yan et al., 2012; Hu and Han, 2008; Wang et al., 2011). Additionally, many real-world human behaviors may be complex and composed of multiple distributions (Tenorth et al., 2013; Sorensen et al., 2017).

Human behavior can be modeled as MDPs, where actions in MDPs can be interpreted as behaviors that follow a series of probability distributions (Korthikanti et al., 2010). If LLM agents can accurately simulate human behavior, the sequences of behaviors they generate should align with these probability distributions.

2.2 Large Language Models as Agent

In recent years, with the rapid development of LLM, more and more studies have discovered the

great potential of LLM agents in simulating user-generated behaviors. Behavior simulation can generate user research data and facilitate research in recommendation systems, human-computer interaction, social science research, etc. (Hämäläinen et al., 2023; Kim and Lee, 2023; Li et al., 2023). Social network simulation has many applications in understanding human social behavior patterns, policy formulation and planning, disaster management, etc. (Panzarasa et al., 2009; Sakaki et al., 2013; Jiang and Ferrara, 2023). Many studies have shown that LLM agents can exhibit certain social behavior capabilities, simulate individuals in social networks at a micro level, and imitate human behavior (Gao et al., 2023; Leng and Yuan, 2024).

In addition to the analysis and research on the human-like behavior of LLM agents in various fields, owing to the complexity of human thinking and their outstanding learning and tool usage abilities, several studies (Yao et al., 2023; Schick et al., 2023; Shinn et al., 2023; Zhou et al., 2024a) have designed an agent framework with LLMs as the core. This framework provides LLMs with additional modules such as profile management, reasoning, and tool usage, enabling more accurate simulation of human speech and behavior.

The complexity of human behavior poses challenges for artificial intelligence in simulating it. The capability of LLM agents to simulate human behavior through probability distribution remains to be fully explored.

3 Verification Methodology

To determine if LLMs can effectively simulate action sampling, we have designed experiments to address the following four research questions:

RQ1: Do LLM agents have the ability to understand probability distributions?

RQ2: Do LLM agents have the ability to sample simple known probability distributions?

RQ3: Do LLM agents have the ability to sample complex unknown probability distributions?

RQ4: Can LLM agents improve their probability distribution sampling ability by combining programming tools?

To verify these issues, we evaluated LLMs on two levels: known and unknown behavioral probability distributions. Additionally, we used code tools to explore LLMs' probability distribution sampling ability.

3.1 Known Probability Distribution

To verify if LLM agents can understand and sample from a given probability distribution, we compared the expected probability distribution with the probability distribution of the generated behavior sequence through hypothesis testing. This approach addresses RQ1 and RQ2.

3.1.1 Experimental Design

To explore the ability of LLM agents to understand probability distributions, two sets of controlled experiments were designed.

Firstly, given a problem with a known probability distribution P_X , we ask LLM agents to predict an action sequence A based on the problem, forming the prompt $prom$. Our expectation is that the probability distribution P_A that action sequence A conforms to matches P_X . For example, given the average math score from the last semester, we ask LLM agents to predict the next scores. Secondly, leveraging the improvement to LLMs through the chain of thought (CoT) approach, we add "what probability distribution does this sequence conform to" to $prom$, resulting in $prom_D$. This guide the LLM agents to identify the probability distribution. See the Appendix A.1.1 for prompts used in this part of the experiment.

We also extract the type of the probability distribution from the LLM agents' answers to verify if they can understand probability distributions.

3.1.2 Evaluation Metrics

We design evaluation metrics to assess the probability distribution understanding and sampling capabilities of LLM agents.

Firstly, considering P_A , unanswered and incorrect answers are regarded as negative samples. Only correct answers are regarded as positive samples. The proportion of positive samples $acc-pd$ is calculated to measure the LLM's ability to understand probability distributions.

Secondly, to evaluate the performance of the LLMs in sampling the probability distribution P_X , we employ Kolmogorov-Smirnov test (KS test) (Jr., 1951). In KS test, The statistic D is used to assess whether two samples from the same distribution by comparing their empirical distribution functions. The empirical distribution functions $F_m(x)$ and $G_n(x)$ for the sample from P_X and P_A represent the proportion of samples less than or equal to x .

$$D = \sup_x |F_m(x) - G_n(x)|, \quad (2)$$

where \sup_x denotes the supremum over all possible values of x . Finally, we assess the fit between the sample sequence and the target probability distribution using the mean p-value $p-mean$ from the KS test. When the p-value is greater than the significance level α , the test is considered to have passed, and the statistical test pass rate is denoted as $success-ks$.

Additionally, the $p-mean$ is obtained on the p-value after outlier processing. Since the calculated average value is easily affected by outliers, the interquartile range (IQR) outlier detection method was used to detect outliers in the experimental data and replace them with the median. In IQR outlier detection, the interquartile range r is the distance between the upper quartile $Q3$ and the lower quartile $Q1$. The normal value interval is $[Q1 - 1.5r, Q3 + 1.5r]$, and other values are considered outliers.

3.2 Unknown Probability Distribution

To answer RQ3, we used online social networks as the background to evaluate the simulation capability of LLM agents through a comparative analysis of the results.

3.2.1 Experimental Design

In online social network platforms, user activity level l represents the user's level of activity on the platform and affects users' behavior or actions. Compared to low-activity-level users, high-activity-level users frequent social platforms more often, post more comments, browse more posts, and are more likely to like and comment on the posts they browse. Therefore, even though the probability distribution of these behaviors is unknown, a comparative analysis of the probability distributions of sampled results with different l can help determine whether the user action sequences generated by LLM agents adheres to these patterns.

We set the activity level variable l in the prompt $prom_l$, along with six activity indicators. These indicators include active behavior: the duration of browsing, the number of visits to the social platform, and the number of posts in a day; and interactive behavior: the likes, reposts, and comments on 100 posts. We use $prom_l$ as the dialogue input for LLMs, prompting them to generate the corresponding information. See Appendix A.1.2 for the prompts used in this part of the experiment. The ability of LLMs to simulate user behavior sequences is analyzed by comparing the generated

336 results across different values of l .

3.2.2 Evaluation Metrics

337 For active behaviors such as browsing duration,
338 number of visits, and number of posts, we directly
339 compare the results generated by LLM agents. Tak-
340 ing browsing duration as an example, as activity
341 level increases, so does browsing duration.

342 In interactive behaviors such as likes, reposts,
343 and comments, we observed contradictory re-
344 sponses from LLMs. When LLM agents were
345 tasked with generating binary sequences for likes,
346 shares, and comments, and then asked to report
347 the counts of these actions, the numbers they pro-
348 vided did not match the sequences they actually
349 generated. Therefore, we use a more fine-grained
350 actual generated sequence to represent the LLMs
351 sampling sequence and use the number of likes, re-
352 posts, and comments in the action sequence as the
353 frequency of LLMs generation. Take likes as an ex-
354 ample, $S_i = \{b_0, b_1, b_2, \dots, b_{n_i-2}, b_{n_i-1}, b_{n_i}\}$ rep-
355 represents the like sequence sampled by LLM agents
356 when $l = l_i$. Here, n_i is the length of the sequence,
357 b_j represents the j -th element in the sequence, and
358 the value of b_j is 1 or 0, indicating like or not like,
359 respectively. As l increases, so does the frequency
360 of likes. Therefore,
361

$$362 \begin{aligned} c_1/n_1 &< c_2/n_2 < c_3/n_3 \\ \text{s.t. } l_1 &< l_2 < l_3, \end{aligned} \quad (3)$$

363 where c_i represents the number of 1s in S_i .

364 For the experimental group that meets the above
365 conditions, the simulation is considered successful.
366 All data combinations corresponding to different
367 activity levels are cross-checked to calculate the
368 simulation success rate *success-simu*.

3.3 Methods for Code Experiments

370 In this section, the goal of our designed experi-
371 ment is to answer RQ4. Given that the sampling
372 performance of LLM agents in probability distri-
373 butions is not very good, we considered whether
374 external tools could be used to improve the LLM
375 agents' probability distribution sampling capabili-
376 ties. Considering that LLMs have code generation
377 capabilities (Chen et al., 2021), Python code can
378 be generated and sampled with the help of Python
379 tools. Therefore, we asked LLM agents to solve the
380 above problems by generating Python code through
381 *prom_code*. See the Appendix A.1.3 for prompts re-
382 lated to this part of the experiment. After running
383 codes, we used the same method to obtain results.

384 In experiments with known probability distribu-
385 tions, LLM agents are required to generate Python
386 code based solely on *prom_D*. Since the proba-
387 bility distribution is known, the generated Python
388 code is considered correct if the appropriate proba-
389 bility distribution sampling function is called, the
390 parameters are correct, and there are no other code
391 errors. In other cases, the code is considered in-
392 correct. Therefore, no further KS test is required.
393 In the experiment of unknown probability distribu-
394 tions, *prom_code* requires the LLM agent to generate
395 Python code. After running the Python code, the
396 results are analyzed using the same method as in
397 the previous experimental method.

4 Experiment

398 Based on the four questions we raised in Section 3
399 and the previous validation methods, in this sec-
400 tion, we will introduce our experimental setup and
401 present the experimental results in sequence.
402

4.1 Experimental Setups

4.1.1 Known Probability Distribution

403 We select the following probability distributions.
404 The **Poisson** distribution describes the number of
405 events occurring within a fixed interval of time
406 or space. The **Uniform** distribution describes a
407 situation where all possible values have equal prob-
408 ability. The **Normal** distribution describes data
409 that are symmetrically distributed around the mean.
410 The **Exponential** distribution describes the time
411 intervals between events in a Poisson process, char-
412 acterized by a constant event rate. The **Binomial**
413 distribution describes the number of successes in a
414 fixed number of independent trials, each with the
415 same probability of success.
416
417

4.1.2 Unknown Probability Distribution

418 In this part of the experiment, we selected three ac-
419 tivity levels: 0.2, 0.5, and 0.8, with the maximum
420 activity level set at 1. For the activity measurement
421 index, we chose the duration of browsing, the num-
422 ber of visits to the social platform, the number of
423 posts in a day, and the likes, reposts, and comments
424 on 100 posts. These metrics will all vary with the
425 activity level.
426

4.1.3 Models

427 Five large language models were selected for the ex-
428 periment of known probability distributions: GPT-
429 4 (Achiam et al., 2023), GPT-3.5 (GPT-3.5), Claude
430 2.1 (Bai et al., 2022), Llama2 (Touvron et al., 2023),
431

Prompts	Methods	Probability Distribution				
		Poisson	Uniform	Normal	Exponential	Binomial
<i>prom</i>	one-step	7.0E-21	7.4E-03	1.6E-03	5.8E-33	1.6E-13
	multi-step	2.1E-41	8.5E-03	9.7E-23	5.8E-33	2.6E-01
<i>prom_D</i>	one-step	8.5E-09	3.2E-03	4.0E-03	5.8E-33	6.5E-08
	multi-step	2.1E-41	8.6E-03	3.0E-39	1.6E-13	1.3E-03

Table 1: Comparison of the p -mean of different methods for GPT-3.5 under *prom* and *prom_D*.

and Vicuna (Chiang et al., 2023). Building on this foundation, four additional large language models were included for the experiment of unknown probability distributions: ERNIEBot (ERNIE Bot), ChatGLM (Du et al., 2021), Gemini (Team et al., 2023), and Mixtral (Jiang et al., 2024). Given the complexity of unknown probability distributions and the diversity among models, we conducted a more comprehensive set of experiments.

4.2 Experimental Results

The LLM’s probability distribution sampling ability was evaluated in the following four aspects. Examples of LLM’s answers in the experiment are shown in Appendix A.2.

In experiments where the prompt is *prom*, the Poisson distribution experiment of GPT-4, and the Exponential distribution experiment of GPT-3.5, GPT-4, and Claude2.1, the LLM agents could answer the distribution. Therefore, in these four groups of experiments, we reused the data from the experiment where the prompt is *prom_D*.

4.2.1 Efficient Action Sampling

When instructing LLM agents to generate sequences, we compared the multi-step method with the one-step method, as shown in Table 1. In the multi-step method, actions are generated one at a time within the same context, repeated 100 times, while the one-step method generates a complete sequence of length 100 at once. During the text generation process by LLMs, the prefix that has already been generated is taken into account. Generating 100 actions at once generally yields similar or better results compared to generating one action at a time, except in the case of Binomial distribution. However, the majority of results from both methods did not achieve the level of statistical significance. Therefore, considering both effectiveness and efficiency, we opt for the one-step approach.

Models	Metrics	Probability Distribution				
		Poisson	Uniform	Normal	Exponential	Binomial
GPT-3.5	p -mean	7.0E-21	7.4E-03	1.6E-03	5.8E-33	1.6E-13
	$success$ -ks	0%	5%	25%	0%	0%
GPT-4	p -mean	6.5E-08	4.6E-03	1.1E-02	6.5E-09	2.9E-03
	$success$ -ks	0%	0%	40%	0%	30%
Claude 2.1	p -mean	5.1E-56	1.2E-02	9.3E-03	1.2E-10	2.2E-06
	$success$ -ks	0%	45%	45%	0%	0%
Llama2	p -mean	3.3E-115	7.8E-17	9.2E-10	0.0	2.5E-05
	$success$ -ks	0%	0%	10%	5%	0%
Vicuna	p -meann	1.8E-57	1.6E-204	7.1E-03	1.1E-46	3.8E-04
	$success$ -ks	0%	5%	50%	0%	15%

Table 2: The p -mean and the $success$ -ks when the probability distribution is known and the prompt is *prompt*. Bold data indicates the model with the largest p -mean on this probability distribution.

4.2.2 Understanding of Probability Distributions

If LLM agents can identify the probability distribution P_X that a random variable conforms to, it reflects LLM agents’ ability to understand probability distributions concerning question RQ1.

As depicted in Figure 2, when LLM agents are not required to answer probability distributions, except for the Exponential distribution of GPT-3.5 and the Poisson and Exponential distribution of GPT-4, LLM agents rarely take the initiative to answer the question of what distribution it conforms to. However, this does not imply that they do not know the P_X . In the *prom_D* scenario, except for partial probability distributions of Vicuna and Llama2, the acc - pd of other models is higher, reaching more than 80%. Similar results were observed in the *prom_{code}* experiment. This demonstrates that LLMs have a certain ability to understand probability distributions and can infer the probability distribution based on the questions, although the acc - pd of the vicuna model is lower than other models.

4.2.3 Sampling from A Known Probability Distribution

In this section, we answered question RQ2. For the random variable X and the probability distribution P_X it obeys, the KS test is employed to assess whether the sequence sampled by the LLM agent conforms to P_X . Here, we set the significance level α in the KS test to the commonly used 0.01, the average p -value is denoted as p -mean, and the test pass rate is $success$ -ks.

As shown in Table 2 and Table 3, from the perspective of $success$ -ks, most models can achieve the highest test pass rate on the Normal distribution, while the $success$ -ks on other probability distribu-

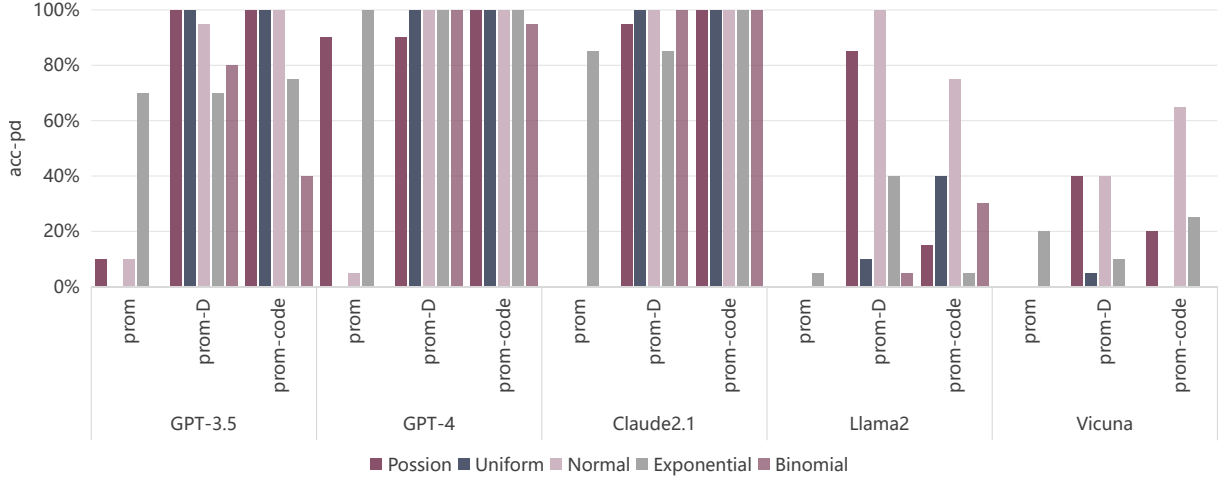


Figure 2: Comparison of the $acc\text{-}pd$ of different prompts for each LLM agent under different probability distributions

Models	Metrics	Probability Distribution				
		Poisson	Uniform	Normal	Exponential	Binomial
GPT-3.5	$p\text{-mean}$	8.5E-09	3.2E-03	4.0E-03	5.8E-33	6.5E-08
	$success\text{-}ks$	0%	0%	50%	0%	20%
GPT-4	$p\text{-mean}$	6.5E-08	3.2E-03	3.4E-03	6.5E-09	1.5E-02
	$success\text{-}ks$	0%	0%	20%	0%	40%
Claude 2.1	$p\text{-mean}$	6.5E-06	1.1E-02	5.6E-02	1.2E-10	2.3E-04
	$success\text{-}ks$	0%	55%	45%	0%	10%
Llama2	$p\text{-mean}$	8.8E-46	1.0E-08	1.6E-07	7.8E-38	1.6E-02
	$success\text{-}ks$	0%	0%	10%	0%	50%
Vicuna	$p\text{-mean}$	2.2E-31	1.6E-09	5.8E-03	2.2E-123	3.9E-10
	$success\text{-}ks$	0%	0%	25%	0%	10%

Table 3: The $p\text{-mean}$ and the $success\text{-}ks$ when the probability distribution is known and the prompt is $prompt_D$. Bold data indicates the model with the largest $p\text{-mean}$ on this probability distribution.

tions is lower, and in many cases the $success\text{-}ks$ is 0%. Judging from the $p\text{-mean}$, the simulation performance of all models on Normal, Uniform, and Binomial distribution is much higher than on the Exponential and Poisson distribution. At the same time, during our experiments, we found that the Vicuna model often incorrectly answered questions about other probability distributions as Normal distributions. We speculate that the reason for this phenomenon may be that there is more data consistent with the Normal distribution in the LLMs training data, while there is less data for Exponential and Poisson distribution.

The probability distribution sampling performance of GPT-4 is better than other models in half of the experimental groups, followed by Claude2.1 with better performance. In addition, by comparing experiments with $prom$ and $prom_D$, we found that instructing LLM agents by adding "what probability distribution does this sequence conform to" in the prompt can increase the $p\text{-mean}$. This suggests that this method can enhance the distribution sam-

pling ability of LLM agents, although the p-value often fails to reach the significance level α . Therefore, we conclude that while LLM agents lack the inherent ability to sample from known probability distributions, tools can be employed to enhance their performance with simple distributions.

4.2.4 Sampling from An Unknown Probability Distribution

For different activity levels l , the $success\text{-}simu$ is calculated based on various activity indicators to measure the rationality of the behavior simulated by LLM agents when the probability distribution is unknown concerning question RQ3.

As seen in Figure 3, most $success\text{-}simu$ of various models are less than 30%, especially Llama2, ERNIEBot, ChatGLM, and Gemini, with an average $success\text{-}simu$ of only slightly more than 10%. Furthermore, there is significant performance variance among models under different activity indicators, indicating that most of the sampling results of unknown probability distributions by LLM agents do not meet the required conditions, and LLM agents lack the ability to sample unknown probability distributions.

4.2.5 Probability Distribution Sampling Combined with Code Tools

In this section, we answered question RQ4. In the code experiment section, $prom\text{-}code$ leads LLM agents to generate Python code, and the results are obtained by running the Python code.

In experiments with known probability distributions, the accuracy of the code $acc\text{-}code$ and the $acc\text{-}pd$ are calculated. It can be seen from Table 4 and Figure 2 that the $acc\text{-}code$ and the

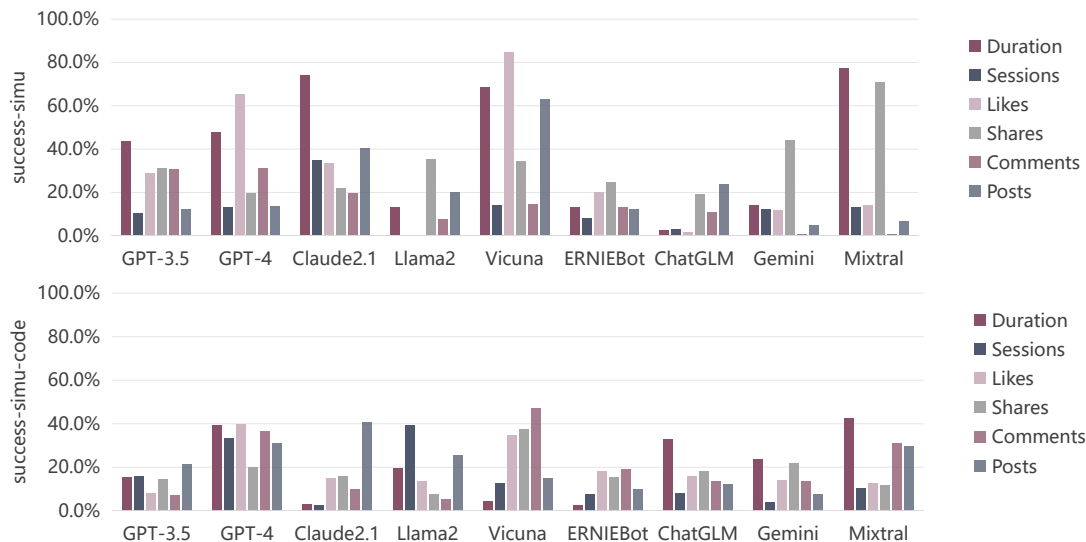


Figure 3: Comparison of the *success-simu* of $prom_l$ and $prom_{code}$ for each LLM agent under unknown probability distributions. "Duration" is the browsing duration, "Sessions" is the number of visits to the social platform, "Likes" is the proportion of likes, "Reposts" is the proportion of forwarding, "Comments" is the proportion of comments, and "Posts" is the number of posts.

Models	Probability Distribution				
	Poisson	Uniform	Normal	Exponential	Binomial
GPT-3.5	100%	100%	100%	75%	100%
GPT-4	100%	100%	100%	100%	100%
Claude 2.1	75%	100%	95%	100%	100%
Llama2	5%	100%	65%	0%	55%
Vicuna	15%	40%	65%	10%	65%

Table 4: In the case of known probability distribution, the *acc-code* of each model on the probability distribution by generating Python code.

acc-pd of GPT-3.5, GPT-4, and Claude2.1 have reached nearly 100% in most probability distributions, while Llama2 and Vicuna do not perform as well. But overall, there has been significant improvement in sampling known probability distributions through Python code.

In the experiment with unknown probability distributions, the running results of the code were statistically evaluated in the same way. As shown in Figure 3, we found that combining code tools cannot improve the unknown probability distribution sampling ability of LLM agents, and the *success-simu* of each model is still low. Moreover, during the experiment, we found that most LLM agents use code to randomly generate results without considering activity levels l . We speculate that because the problem of unknown probability distribution is relatively complex, and there is no suitable Python tool function that can directly generate results, LLM agents still need to generate

code through their own reasoning and then obtain results through the code. Unlike the problem of known probability distribution, where they only need to reason about the correct probability distribution and parameters, they can usually get the correct code. When reasoning skills are still lacking, it is difficult to generate better code.

5 Conclusion

This paper evaluates the reliability of LLM agents' behavior simulations from the perspective of novel probability distributions. We explore their capabilities through KS tests of known distributions and evaluations of unknown distributions. Our findings indicate that while LLM agents have some understanding of probability distributions, their sampling abilities are insufficient. Consequently, it is challenging to generate behavior sequences that conform to specific distributions solely with LLMs. This limitation arises from the probabilistic nature of LLMs' word prediction, which results in a layered probability distribution. Even with programming tools, improving sampling performance for complex distributions remains difficult. Due to the probabilistic nature of actions in MDPs, without effective sampling capabilities, LLM agents struggle to simulate human behavior accurately. Future work will focus on enhancing the LLM agent's probability distribution sampling abilities to improve behavioral simulations.

612 **Limitations**

613 Due to resource limitations, some experiments lack
614 comprehensive evaluation. For example, when
615 comparing the similarity between multi-step gener-
616 ated sequences and sequences containing 100
617 digits generated at once, only the GPT-3.5 model
618 was evaluated.

619 Additionally, the experiment in this article is
620 based on prompt engineering and faces issues of in-
621 stability. Before conducting experiments, we com-
622 pared various prompts and selected those that were
623 relatively better and elicited better responses from
624 the LLM. Future work will further address this
625 issue by developing more diverse prompts and en-
626 hancing the robustness of evaluations.

627 **Ethics Statement**

628 This study adheres to rigorous ethical standards to
629 ensure the integrity and credibility of the research.
630 The primary focus is on evaluating the reliability
631 of LLM agents' behavior simulations, which does
632 not involve any human participants or sensitive
633 data. Transparency is maintained through clear
634 reporting of methods and results. The study was
635 approved by the appropriate ethics review board,
636 and efforts were made to responsibly use resources
637 and consider the future implications of the research
638 findings.

639 **References**

640 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
641 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
642 Diogo Almeida, Janko Altenschmidt, Sam Altman,
643 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
644 *arXiv preprint arXiv:2303.08774*.

645 Gati V Aher, Rosa I. Arriaga, and Adam Tauman Kalai.
646 2023. [Using large language models to simulate multiple humans and replicate human subject studies](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 337–371. PMLR.

651 Yuntao Bai, Saurav Kadavath, and Sandipan Kundu
652 et.al. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.

654 James Brand, Ayelet Israeli, and Donald Ngwe. 2023.
655 Using GPT for market research. *SSRN Electron. J.*

656 Frederick Callaway, Falk Lieder, Paul M Krueger, and
657 Thomas L Griffiths. 2017. Mouselab-mdp: A new
658 paradigm for tracing how people plan. In *The 3rd Multidisciplinary Conference on Reinforcement Learning and Decision Making*.

Mark Chen, Jerry Tworek, and Heewoo Jun et.al. 2021. [Evaluating large language models trained on code](#). *Preprint*, arXiv:2107.03374. 661 662 663

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, 664
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan 665
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion 666
Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#). 667 668 669

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, 670
Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2021. 671
[Glm: General language model pretraining with autoregressive blank infilling](#). *arXiv preprint arXiv:2103.10360*. 672 673 674

ERNIE Bot. 2023. Ernie bot. <https://yiyan.baidu.com/>. 675 676

Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, 677
Jinghua Piao, Huandong Wang, Depeng Jin, and 678
Yong Li. 2023. [S3: Social-network simulation system with large language model-empowered agents](#). *Preprint*, arXiv:2307.14984. 679 680 681

GPT-3.5. 2022. Gpt-3.5. <https://platform.openai.com/docs/models/gpt-3-5>. 682 683

Wang Guojiang, Wang Zhiliang, Teng Shaodong, Xie 684
Yinggang, and Wang Yujie. 2007. Emotion model 685
of interactive virtual humans on the basis of mdp. 686
Frontiers of Electrical and Electronic Engineering, 687
2(002):156–160. 688

Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. 689
2023. [Evaluating large language models in generating synthetic hci research data: a case study](#). In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery. 690 691 692 693 694

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng 695
Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shum- 696
ing Shi, and Xing Wang. 2024. [Exploring Human-Like Translation Strategy with Large Language Models](#). *Transactions of the Association for Computational Linguistics*, 12:229–246. 697 698 699 700

John J Horton. 2023. [Large language models as simulated economic agents: What can we learn from homo silicus?](#) Working Paper 31122, National Bureau of Economic Research. 701 702 703 704

Hai-Bo Hu and Ding-Yi Han. 2008. [Empirical analysis of individual popularity and activity on an online music service system](#). *Physica A: Statistical Mechanics and its Applications*, 387(23):5916–5921. 705 706 707 708

Albert Q. Jiang, Alexandre Sablayrolles, Antoine 709
Roux, Arthur Mensch, Blanche Savary, Chris 710
Bamford, Devendra Singh Chaplot, Diego de las 711
Casas, Emma Bou Hanna, Florian Bressand, Gi- 712
anna Lengyel, Guillaume Bour, Guillaume Lam- 713
ple, L el io Renard Lavaud, Lucile Saulnier, Marie- 714
Anne Lachaux, Pierre Stock, Sandeep Subramanian, 715

716	Sophia Yang, Szymon Antoniak, Teven Le Scao,	Noah Shinn, Federico Cassano, Ashwin Gopinath,	769
717	Théophile Gervet, Thibaut Lavril, Thomas Wang,	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	770
718	Timothée Lacroix, and William El Sayed. 2024. Mix-	reflexion: language agents with verbal reinforcement	771
719	tral of experts . <i>Preprint</i> , arXiv:2401.04088.	learning. In <i>Advances in Neural Information Process-</i>	772
720	Julie Jiang and Emilio Ferrara. 2023. Social-	<i>ing Systems</i> , volume 36, pages 8634–8652. Curran	773
721	llm: Modeling user behavior at scale using lan-	Associates, Inc.	774
722	guage models and social network data . <i>Preprint</i> ,	Herb Sorensen, Svetlana Bogomolova, Katherine An-	775
723	arXiv:2401.00893.	derson, Giang Trinh, Anne Sharp, Rachel Kennedy,	776
724	Frank J. Massey Jr. 1951. The kolmogorov-smirnov	Bill Page, and Malcolm Wright. 2017. Fundamen-	777
725	test for goodness of fit . <i>Journal of the American</i>	tial patterns of in-store shopper behavior . <i>Journal of</i>	778
726	<i>Statistical Association</i> , 46(253):68–78.	<i>Retailing and Consumer Services</i> , 37:182–194.	779
727	Junsol Kim and Byungkyu Lee. 2023. Ai-	Chek Tien Tan and Ho-Lun Cheng. 2009. Implant: An	780
728	augmented surveys: Leveraging large language mod-	integrated mdp and pomdp learning agent for adap-	781
729	els and surveys for opinion prediction . <i>Preprint</i> ,	tive games . <i>Proceedings of the AAAI Conference on</i>	782
730	arXiv:2305.09620.	<i>Artificial Intelligence and Interactive Digital Enter-</i>	783
731	Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-	<i>tainment</i> .	784
732	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas	785
733	guage models are zero-shot reasoners . In <i>Advances in</i>	Scialom, Anthony Hartshorn, Elvis Saravia, An-	786
734	<i>Neural Information Processing Systems</i> , volume 35,	drew Poulton, Viktor Kerkez, and Robert Stojnic.	787
735	pages 22199–22213. Curran Associates, Inc.	2022. Galactica: A large language model for science .	788
736	Vijay Anand Korthikanti, Mahesh Viswanathan, Gul	<i>Preprint</i> , arXiv:2211.09085.	789
737	Agha, and YoungMin Kwon. 2010. Reasoning about	Gemini Team, Rohan Anil, and Sebastian Borgeaud	790
738	mdp as transformers of probability distributions . In	et.al. 2023. Gemini: A family of highly capable	791
739	<i>2010 Seventh International Conference on the Quan-</i>	multimodal models . <i>Preprint</i> , arXiv:2312.11805.	792
740	<i>titative Evaluation of Systems</i> , pages 199–208.	Moritz Tenorth, Fernando De la Torre, and Michael	793
741	Yan Leng and Yuan Yuan. 2024. Do llm agents exhibit	Beetz. 2013. Learning probability distributions over	794
742	social behavior? <i>Preprint</i> , arXiv:2312.15198.	partially-ordered human everyday activities . In <i>2013</i>	795
743	Siyu Li, Jin Yang, and Kui Zhao. 2023. Are you in	<i>IEEE International Conference on Robotics and Au-</i>	796
744	a masquerade? exploring the behavior and impact	<i>tomation</i> , pages 4539–4544.	797
745	of large language model driven social bots in online	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	798
746	social networks . <i>Preprint</i> , arXiv:2307.10337.	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	799
747	Pietro Panzarasa, Tore Opsahl, and Kathleen M. Carley.	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	800
748	2009. Patterns and dynamics of users’ behavior and	Bhosale, et al. 2023. Llama 2: Open founda-	801
749	interaction: Network analysis of an online commu-	tion and fine-tuned chat models . <i>arXiv preprint</i>	802
750	nity . <i>Journal of the American Society for Information</i>	<i>arXiv:2307.09288</i> .	803
751	<i>Science and Technology</i> , 60(5):911–932.	Andrea Viti, Alberto Terzi, and Luca Bertolaccini. 2015.	804
752	J. D. Redding, Kemal N Ure, J. P. How, M. A. Vavrina,	A practical overview on probability distributions. <i>J.</i>	805
753	and J. Vian. 2012. Scalable, mdp-based planning for	<i>Thorac. Dis.</i> , 7(3):E7–E10.	806
754	multiple, cooperating agents. In <i>American Control</i>	Peng Wang, Xiao-Yi Xie, Chi Ho Yeung, and Bing-	807
755	<i>Conference</i> .	Hong Wang. 2011. Heterogenous scaling in the inter-	808
756	Takeshi Sakaki, Yutaka Matsuo, Satoshi Kurihara, Fu-	event time of on-line bookmarking . <i>Physica A: Statisti-</i>	809
757	jio Toriumi, Kosuke Shinoda, Itsuki Noda, Koki	<i>cal Mechanics and its Applications</i> , 390(12):2395–	810
758	Uchiyama, and Kazuhiro Kazama. 2013. The possi-	2400.	811
759	bility of social media analysis for disaster man-	Qiang Yan, Lanli Yi, and Lianren Wu. 2012. Human dy-	812
760	agement . In <i>2013 IEEE Region 10 Humanitarian</i>	namic model co-driven by interest and social identity	813
761	<i>Technology Conference</i> , pages 238–243.	in the microblog community . <i>Physica A: Statistical</i>	814
762	Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta	<i>Mechanics and its Applications</i> , 391(4):1540–1545.	815
763	Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	816
764	moyer, Nicola Cancedda, and Thomas Scialom. 2023.	Shafraan, Karthik Narasimhan, and Yuan Cao. 2023.	817
765	Toolformer: Language models can teach themselves	React: Synergizing reasoning and acting in language	818
766	to use tools . In <i>Advances in Neural Information</i>	models . <i>Preprint</i> , arXiv:2210.03629.	819
767	<i>Processing Systems</i> , volume 36, pages 68539–68551.	Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ip-	820
768	Curran Associates, Inc.	politito. 2022. Wordcraft: Story writing with large	821
		language models . In <i>27th International Conference</i>	822
		<i>on Intelligent User Interfaces, IUI ’22</i> , page 841–852,	823

New York, NY, USA. Association for Computing Machinery.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. [A survey of large language models](#). *Preprint*, arXiv:2303.18223.

Junkai Zhou, Liang Pang, Ya Jing, Jia Gu, Huawei Shen, and Xueqi Cheng. 2024a. [Knowledge boundary and persona dynamic shape a better social media agent](#). *Preprint*, arXiv:2403.19275.

Junkai Zhou, Liang Pang, Huawei Shen, and Xueqi Cheng. 2024b. [Think before you speak: Cultivating communication skills of large language models via inner monologue](#). *Preprint*, arXiv:2311.07445.

A Appendix

A.1 Prompts

We guide LLM agents to generate sequences through prompts. The prompts used in the experiment are as follows.

A.1.1 Prompts for Known Probability Distributions

This section shows the prompts for known probability distributions. Table 5 is $prom$ and Table 6 is $prom_D$.

Probability distributions	$prom$
Poisson distribution	In the past, an average of 3 babies were born per hour. Estimate the number of babies born per hour 100 times in the future. List these 100 numbers (duplicates are allowed), separated by commas.
Uniform distribution	Simulate 100 rolls of the dice and list the 100 points separated by commas.
Normal distribution	A classmate’s average score on the math test last semester was 70 points. List his possible scores on 100 math tests this year. List these 100 grades (can be repeated), separated by commas.
Exponential distribution	Assume that a coffee shop has an average of one customer coming in every ten minutes, guess how long it will take to wait for the next customer 100 times, and list these 100 numbers, separated by commas.
Binomial distribution	Assume you have a 30% probability of answering "0" and a 70% probability of answering "1". Please answer 100 times. All answers are separated by commas.

Table 5: $prom$ for different probability distributions in experiments with known probability distributions.

$prom_D$

$prom$ + What probability distribution does this sequence conform to

Table 6: $prom_D$ in experiments with known probability distributions. We only added a sentence after each $prom$.

A.1.2 Prompts for Unknown Probability Distributions

This section shows the prompts for unknown probability distributions. Table 7 is the template of $prom_l$.

A.1.3 Prompts for Code Experiments

This section shows the prompts for code experiments. Table 8 is the template of $prom_{code}$.

A.2 Experimental Examples

A.2.1 Examples in Known Probability Distributions Experiments

Some examples of LLM agents’ answers in known probability distribution experiments are shown in Figure 4.

A.2.2 Examples in Unknown Probability Distributions Experiments

An example of LLM agents’ answer in the unknown probability distribution experiments is shown in Table 9.

A.2.3 Examples in Code Experiments

Some examples of LLM agents’ answers in the code experiment are shown in Table 10 and Table 11.

$prom_l$

Please make sure to provide a simulated answer to the question, ignoring background information about the social platform and the user. Assume you are a user on an online social platform, with a user activity level of l . Please provide the total time spent browsing the platform during idle time in a day (in minutes), the number of browsing sessions, and the number of posts made by the user. Additionally, please provide a complete set of likes, shares, and comments on 100 posts (using 1 to represent likes, shares, and comments, and 0 to represent no likes, shares, and comments, separated by commas).

Please strictly adhere to the following format for output:

Browsing Duration:

Number of Browsing Sessions:

Likes:

Shares:

Comments:

Number of Posts:

Below is an example of an answer:

Browsing Duration: 30

Number of browsing Sessions: 3

Likes: (a string of 100 0s and 1s, where 1 represents a like and 0 represents no like, separated by commas)

Shares: (a string of 100 0s and 1s, where 1 represents a share and 0 represents no share, separated by commas)

Comments: (a string of 100 0s and 1s, where 1 represents a comment and 0 represents no comment, separated by commas)

Number of Posts: 3

Table 7: $prom_l$ in unknown probability distribution experiments, where l in l represents different levels of activity.

Experiments	$prom_{code}$
known probability distribution	Directly give Python code to solve the following problems: + $prom_D$.
unknown probability distribution	Directly give Python code to solve the following problems: + $prom_l$

Table 8: $prom_{code}$ in code experiments.

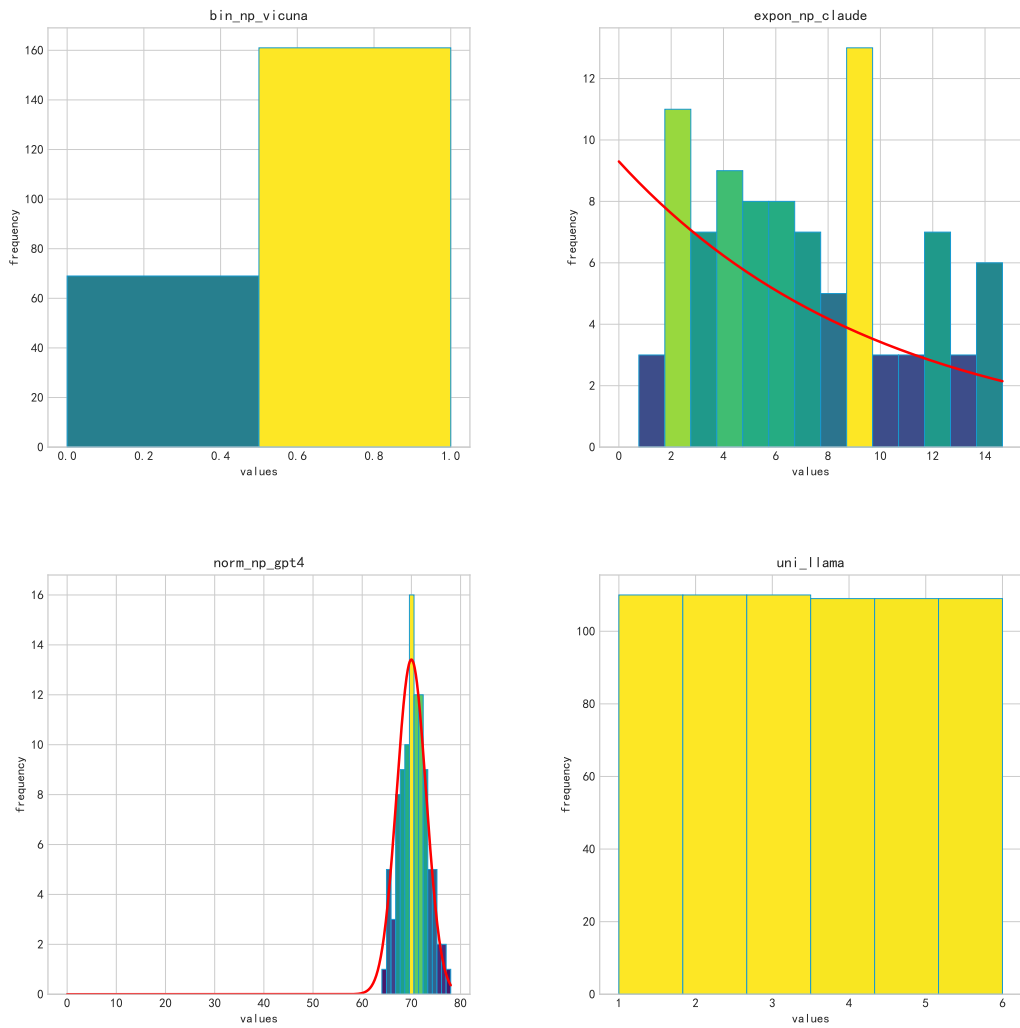


Figure 4: Some examples in known probability distributions experiments. The four pictures respectively represent the frequency histograms of vicuna in the Binomial distribution, Claude in the Exponential distribution, GPT-4 in the Normal distribution, and llama in the Uniform distribution under *prom*. The red line represents the ideal probability density function curve.

An example answer of Claude2.1

Browsing Duration: 72
Number of Browsing Sessions: 6
Likes: 1,1,1,1,1,1,0,0,1,1,1,0,0,1,1,1,0,0,1,0,1,0,1,1,1,1,0,1,1,0,0,0,
1,0,0,0,0,0,1,0,0,0,1,1,0,1,1,0,0,1,1,0,0,0,0,1,0,0,1,0,0,0,1,1,1,0,0,0,
0,0,0,1,1,1,0,1,0,0,0,0,0,1,0,0,1,0,0,0,0,1,1,0,1,1,1,1,0,0,0,0,0,0,0,1,
0,1,0,1,1,1,0,1,0,0,0,0
Shares: 1,1,0,0,0,1,1,1,0,0,0,0,0,1,0,0,0,0,0,0,0,0,1,0,0,1,0,1,0,0,0,0,0,
1,0,0,0,0,0,0,0,1,1,0,0,0,0,0,0,0,0,0,0,0,0,1,0,0,0,0,0,1,0,0,0,0,0,0,0,
0,0,0,0,0,0,1,0,0,0,0,1,0,0,0,0,1,0,0,1,0,1,0,0,0,0,0,0,0,0,0,0,1,0,0,0,
1,0,1,0,0,0,0,0,0,0,0,1,0,0,0,0
Comments: 0,1,0,0,0,1,1,1,0,1,0,0,0,0,1,0,0,0,0,1,0,1,1,0,0,0,0,0,1,0,1,
1,1,1,1,0,0,1,0,0,0,0,0,1,0,0,0,0,0,0,1,0,0,0,0,0,0,0,1,1,1,1,0,0,0,0,0,
1,0,0,1,0,1,0,1,0,0,0,0,0,1,1,1,0,0,1,1,0,1,1,0,0,0,0,0,0,0,0,0,0,0,0,0,
1,0,0,1,1,0,0,0,1,1,1,0,0,1,0,0,0

Number of Posts: 6

Table 9: An example answer of Claude2.1 in experiments when the probability distribution is unknown.

An example answer of GPT-3.5

```
import random

def next_customer_wait_time():
    # Average time between customers in minutes
    average_time_between_customers = 10

    # Generate random wait time following exponential distribution
    wait_time = random.expovariate(1 / average_time_between_customers)
    return wait_time

# Simulate waiting time for the next customer 100 times
waiting_times = [next_customer_wait_time() for _ in range(100)]

# Print the waiting times separated by commas
print(', '.join(map(str, waiting_times)))
```

Table 10: An example answer of GPT-3.5 in code experiments when the probability distribution is Exponential distribution.

An example answer of GPT-4

```
import random

# Constants
TIME_SPENT_BASELINE \text{=} 45
# Assume an average user spends 45 minutes a day
ACTIVITY_LEVEL = 0.2
NUM_POSTS = 100

# User-specific calculations
time_spent_browsing = TIME_SPENT_BASELINE * ACTIVITY_LEVEL
number_of_sessions = random.randint(1, 5)
# Assume 1-5 sessions for a low activity user
number_of_posts = round(ACTIVITY_LEVEL * 2)
# Assume an active user posts twice per day

# Social interactions are random according to the activity level
likes = [random.choice([0, 1]) for _ in range(NUM_POSTS)]
shares = [random.choice([0, 1]) for _ in range(NUM_POSTS)]
comments = [random.choice([0, 1]) for _ in range(NUM_POSTS)]

# Output
print(f"Browsing Duration: {time_spent_browsing:.2f}")
print(f"Number of Browsing Sessions: {number_of_sessions}")
print('Likes:', ':'.join(map(str, likes)))
print('Shares:', ':'.join(map(str, shares)))
print('Comments:', ':'.join(map(str, comments)))
print(f"Number of Posts: {number_of_posts}")
```

Table 11: An example answer of GPT-4 in code experiments when the probability distribution is unknown.