
Categorical Decision Mamba : On Tractable Binning in Sequence Model based Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently, sequence modeling methods have been applied to solve the problem of
2 off-policy reinforcement learning. One notable example is the work on Decision
3 Mamba, incorporating Mamba block into the Decision-Transformer-type neural
4 network architecture. In this work, we begin our exploration with the latest sequen-
5 tial decision-making model, leveraging its strengths as a foundation for further
6 development. We propose a theoretical measure of alignment on the policy of the
7 agent with the human expert, known as Expected Agent Alignment Error (EA2E).
8 Furthermore, we provide a complete theoretical proof that reducing the Wasserstein-
9 1 distance between distributions of the present model (agent) and the target model
10 (agent) effectively aligns the agent’s policy with the potential expert’s. Building
11 upon theoretical results, we propose Categorical Decision Mamba (CDMamba),
12 which originates from Decision Mamba (DMamba). The core improvements of
13 CDMamba involve utilizing histograms of categorical distributions as inputs to the
14 Mamba model, minimizing the Wasserstein-1 distance between distributions, which
15 ultimately yields a trained model with aligned policy and superior performance.

16 1 Introduction

17 Offline Reinforcement Learning [1, 2, 3, 4] has been a promising approach for training agents that
18 does not necessitate online experience in an environment, which is advantageous when online ex-
19 perience is expensive or when offline experience is abundant. In the past few years, Transformers
20 [5] have shown impressive results across a number of problem domains in Natural Language Pro-
21 cessing [6, 7, 8] and Computer Vision [9, 10]. Inspired by these recent successes, formulating offline
22 reinforcement learning as a sequence modeling problem [1, 2, 3, 4, 11] has become a novel idea for
23 solution, where the Transformer model predicts the next element in a sequence of states, actions,
24 rewards, and then tackled the issue with techniques similar to those employed in large language
25 models. The attention mechanism in Transformer does have produced several impressive results
26 [1, 11]; however, investigating alternative mechanisms to further enhance model performance remains
27 an open and intriguing research question [3].

28 Recently, state space sequence models (SSMs) have gained popularity as efficient and effective
29 building blocks for constructing deep networks, achieving great performance in analyzing continuous
30 long-sequence data [12, 13]. In particular, structured state space sequence models (S4) have been
31 effective in various applications [12, 14]. Mamba [15] enhances S4 [12] by incorporating a selective
32 mechanism, allowing the model to selectively focus on input-dependent, relevant information. This
33 improvement, combined with hardware-aware implementation, enables Mamba to outperform Trans-
34 formers on dense modalities, such as language and genomics. Leveraging the numerous advantages
35 of the Mamba architecture and significant success achieved in various domains [16, 17, 18, 19],
36 **Decision Mamba** [3] investigate the integration of the Mamba framework as a new architectural

37 choice within the Decision Transformer. Empirical study in the work Ota [3] shows that Decision
 38 Mamba is competitive to existing DT-type models, suggesting the effectiveness of Mamba framework
 39 for RL tasks.

40 In work [20], the exploration problem is reformulated as a State Marginal Matching (SMM) issue, in
 41 which a target state distribution is given, and a policy is learned to make the state marginal distribution
 42 match this target distribution. Inspired by such notion, Furuta et al. [2] proposed the Categorical
 43 Decision Transformer (CDT) [2], taking histograms of categorical distribution (i.e. discrete approxi-
 44 mations of feature distributions; \mathcal{B} -dim vector) as the inputs of the transformer. With obtaining the
 45 desired information statistics for all trajectories, they are fed to the Categorical Decision Transformer
 46 during training/test time. While such approach is effective for reinforcement learning in multiple
 47 tasks, it is crucial to consider the limitations of training a black-box model. Specifically, for an agent
 48 generated by such large model, we cannot help but question: What about the alignment of the trained
 49 agent? Furthermore, does the reduction in Wasserstein-1 distance between distributions of existing
 50 model (agent) and target model (agent) effectively aligns the agent’s policy with the potential human
 51 expert’s?

52 Inspired by the innovative frameworks in sequential decision-making, such as Decision-Mamba [3],
 53 we would like to reevaluating the role of SMM method in sequence-based decision-making. There-
 54 fore, in this paper, we propose **Categorical Decision Mamba: further taking the histograms of**
 55 **categorical distribution as the input of Mamba, minimizing the distance between distributions**
 56 **and finally get a trained model with excellent and aligned policy.**

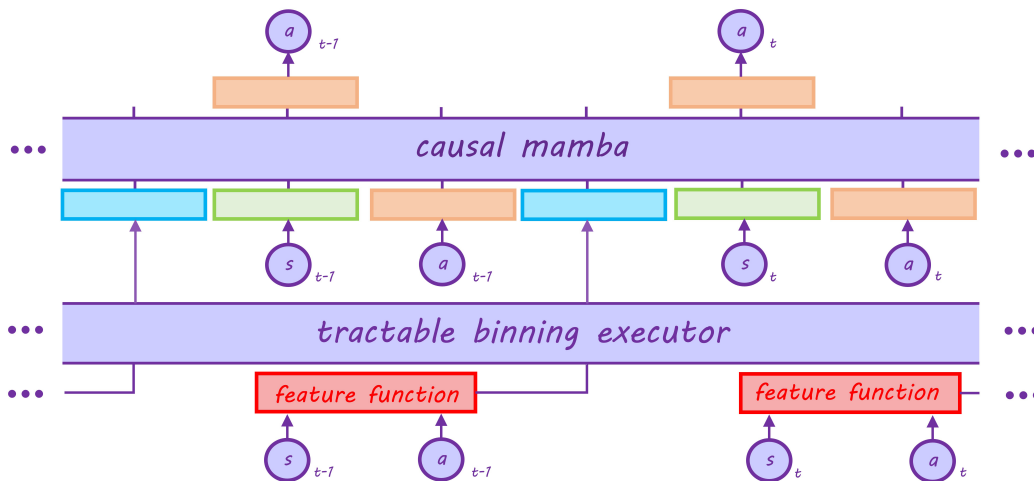


Figure 1: Categorical Decision Mamba (CDMamba) architecture

57

58 To summary, the contributions of this work are as follows:

- 59 • We propose a theoretical measure of alignment towards the policy of an agent embedded in
 60 a large model within the background of state marginal matching.
- 61 • We provide a theoretical proof that reducing Wasserstein-1 distance between distributions of
 62 existing and target models aligns the agent’s policy with the potential human expert’s.
- 63 • We present an enhanced approach for Decision Mamba (DMamba) termed as **Categorical**
 64 **Decision Mamba (CDMamba).**

65 2 Methodology

66 2.1 Expected Agent Alignment Error (EA2E)

67 Offline reinforcement learning (RL) [1, 2, 3, 4] has been a promising approach for training agents
 68 that does not necessitate online experience in an environment, which is advantageous when online
 69 experience is expensive or when offline experience is abundant. Notwithstanding, it also gives rise

70 to several other concerns. One notable issue is that the policy learned by an agent may deviate
 71 significantly from the underlying human strategy, yet still achieve a higher return, leading to the
 72 misconception that a good policy has been learned. To facilitate a comprehensive evaluation of the
 73 learned policy, it is essential to define a metric that explains the confidence of the model (agent) in its
 74 strategy. Similar to ECE [21], we define the Expected Agent Alignment Error (EA2E).

75 **Definition 2.1.** Defining $U = \mathbb{I}[\text{dis}(\pi_1, \pi_2) \leq \delta]$, where $\mathbb{I}(\cdot)$ is the indicator function, $U = 1$
 76 denotes that the distance between policy π_1 and policy π_2 is smaller than the threshold δ . Hence, for
 77 the model (agent), a perfect alignment (full confidence) can be expressed as:

$$\mathbb{P}(U = 1 | \hat{s} = s) = s \quad \forall s \in [0, 1] \quad (1)$$

78 **Definition 2.2.** Expected Agent Alignment Error (EA2E): Define the misalignment of the model
 79 (agent) by computing the expectation of alignment error over predicted confidence \hat{s} :

$$\text{EA2E} = \mathbb{E}_{\hat{s}} [|\mathbb{P}(U = 1 | \hat{s} = s) - s|] \quad (2)$$

80 2.2 Theoretical Analysis

81 In this section, we theoretically delve deeper into the benefits of mitigating distribution discrepancies,
 82 demonstrating the superiority of our approach.

83 **Theorem 1.** Suppose $\pi_{\theta_1}(a|s)$ and $\pi_{\theta_2}(a|s)$ are two policies of the model (agent), $\rho^{\pi_{\theta_1}}(s, a)$ and
 84 $\rho^{\pi_{\theta_2}}(s, a)$ are the state-action marginal distributions of two agents (models), then:

$$\text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) \leq \mathbb{E} [4 \cdot \text{TV}(\rho^{\pi_{\theta_1}}(s, a), \rho^{\pi_{\theta_2}}(s, a))] \quad (3)$$

85 , where $\text{TV}(\cdot)$ is the total variation distance.

86 **Proposition 1.** Suppose $\pi(a|s)$ is the potential policy of human expert , $\pi_{\theta}(a|s)$ is the policy of
 87 model (agent), $\rho^{\pi}(s, a)$ is the state-action marginal distribution of potential human expert agent and
 88 $\rho^{\pi_{\theta}}(s, a)$ is the state-action marginal distribution of current agent (model), then:

$$\text{EA2E}(\pi) - \text{EA2E}(\pi_{\theta}) \leq \mathbb{E} \left[\frac{4}{d_{\min}} \cdot W_1(\rho^{\pi}(s, a), \rho^{\pi_{\theta}}(s, a)) \right] \quad (4)$$

89 , where $W_1(\rho^{\pi}(s, a), \rho^{\pi_{\theta}}(s, a))$ is the Wasserstein-1 distance between $\rho^{\pi}(s, a)$ and $\rho^{\pi_{\theta}}(s, a)$; let
 90 $\rho^{\pi}(s, a) \in \mu$ and $\rho^{\pi_{\theta}}(s, a) \in \nu$, setting $\Omega = \text{supp}(\mu) \cup \text{supp}(\nu)$, $d_{\min} = \inf_{\rho^{\pi} \neq \rho^{\pi_{\theta}} \in \Omega} \|\rho^{\pi} - \rho^{\pi_{\theta}}\|$.

91 In Theorem 1, it is shown that the differences of EA2E is bounded by total variation distance of the
 92 policies. And in Proposition 1, we can observe that the deviation of EA2E between the potential
 93 human expert and the current model (agent) is bounded by Wasserstein-1 distance between the
 94 distribution of the expert agent and the current agent (model). Detailed proofs are available in
 95 Appendix A and Appendix B.

96 2.3 Categorical Decision Mamba

97 Recently, the Mamba framework have been introduced [15], known as an sequence modeling
 98 framework that leverages a selective structured state space model to achieve efficient and effective
 99 performance. Decision Mamba is a novel approach that replaces traditional self-attention with Mamba
 100 block. Empirical verification has shown that such modification can improve the model’s capacity of
 101 capturing complex dependencies in sequential decision-making tasks, thereby potentially enhancing
 102 its decision-making capabilities in diverse and challenging environments. Based on such preliminary,
 103 we construct categorical approximations of continuous distributions by leveraging the discretization
 104 of feature spaces as a substitute for return-to-go (RTG). The architecture of Categorical Decision
 105 Mamba (CDMamba) is shown in Figure 1. From the architecture, we can see that CDMamba takes
 106 binnings of distribution (rewards or state dimensions like xyz-velocities), states and actions as input
 107 and actions of future as output. And furthermore, we also evaluate the model with Wasserstein-1
 108 distance between categorical distributions of features, in order to demonstrate the effectiveness of
 109 state-feature distribution matching.

110 3 Experiments

111 We conduct the experiments on the MuJoCo tasks (Halfcheetah, Hopper and Walker2d) from the
 112 widely-used D4RL [22] benchmarks. Firstly, we sort all the trajectories by their cumulative rewards.
 113 For comparison with CDT [2], we similarly hold out five best trajectories and five 50 percentile
 114 trajectories as a test set (10 trajectories in total), and use the rest as a train set. We report the results
 115 averaged over 20 rollouts every 3 random seeds in Table 1. We select DT [1], RvS [23], DS4 [24],
 116 DC [25] and DMamba [3] as our baselines. The results show that CDMamba is competitive to
 117 DMamba and existing DT-type models, suggesting an effectiveness of CDMamba architecture for
 118 RL tasks. Furthermore, we evaluate the Wasserstein-1 distance between categorical distributions of
 119 features in Table 2. The practical computation of Wasserstein-1 distance is conducted by package
 120 in [26]. Distance datas of DT and CDT are from work [8]. The results show that CDMamba
 121 matches approximate distribution much better than CDT and DT. For reproducibility, we provide the
 122 hyperparameter of experiments in Table 3 (Appendix D). And for intuitive understanding, we provide
 typical visualizations in Figure 2 (Appendix C).

Table 1: The offline results of CDMamba and baselines in MuJoCo domain. We abbreviate dataset names as follows: ‘medium’ as ‘m’, ‘medium-replay’ as ‘m-r’, ‘medium-expert’ as ‘m-e’.

Dataset	DT	RvS	DS4	DC	DMamba	CDMamba
halfcheetah-m	42.6	41.6	42.5	43.0	42.8±0.08	43.2±0.06
hopper-m	68.4	60.2	54.2	92.5	83.5±12.5	70.7±1.66
walker2d-m	75.5	71.7	78.0	79.2	78.2±0.6	79.2±0.12
halfcheetah-m-r	37.0	38.0	15.2	41.3	39.6±0.1	39.9±0.05
hopper-m-r	85.6	73.5	49.6	94.2	82.6±4.6	88.8±2.17
walker2d-m-r	71.2	60.6	69.0	76.6	70.9±4.3	78.2±2.88
halfcheetah-m-e	88.8	92.2	92.7	93.0	91.9±0.6	88.4±1.89
hopper-m-e	109.6	101.7	110.8	110.4	111.1±0.3	110.6±0.53
walker2d-m-e	109.3	106.0	105.7	109.6	108.3±0.5	108.6±0.39

123

Table 2: Quantitative evaluation of reward distribution matching via measuring Wasserstein-1 distance between the rollout and target distributions.

Method	Halfcheetah		Hopper		Walker2d	
	m	m-e	m	m-e	m	m-e
DT	1.039±1.548	0.846±1.134	0.091±0.035	0.159±0.111	0.626±0.495	0.341±0.452
CDT	1.002±1.458	0.838±1.054	0.064±0.017	0.111±0.077	0.114±0.037	0.105±0.030
CDMamba	0.166±0.044	0.469±0.229	0.081±0.025	0.134±0.026	0.215±0.014	0.028±0.019

124 4 Discussion

125 In this work, we provide a novel perspective on the agent policy alignment: decreasing Wasserstein-1
 126 distance between distributions of existing and target models aligns the agent’s policy with that
 127 of potential human experts. Building upon the theoretical foundations, we introduce CDMamba,
 128 a novel approach that synergistically integrates the superior performance of Mamba in handling
 129 sequential problems with strengths of distributional matching for agent alignment. Hence, this work
 130 makes theoretical contributions to the field of agent alignment in reinforcement learning; and further
 131 validates the Mamba architecture’s robust representational capabilities in sequence model-based
 132 reinforcement learning tasks, offering valuable insights for future research in this area.

133 In future work, we would like to further explore the agent policy alignment in reinforcement learning
 134 based on sequence models, foundation models and large language models; and further study practical
 135 methods of agent policy alignment. Moreover, we also would like to study additional implementations
 136 that better effectively leverage Mamba’s or other sequence modeling architectures’ advantages in the
 137 sequence model based reinforcement learning.

References

- 138
- 139 [1] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter
140 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning
141 via sequence modeling, 2021.
- 142 [2] Hiroki Furuta, Yutaka Matsuo, and Shixiang Shane Gu. Generalized decision transformer for
143 offline hindsight information matching, 2022.
- 144 [3] Toshihiro Ota. Decision mamba: Reinforcement learning via sequence modeling with selective
145 state spaces, 2024.
- 146 [4] Hao Liu and Pieter Abbeel. Emergent agentic transformer from chain of hindsight experience,
147 2023.
- 148 [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
149 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- 150 [6] OpenAI et.al. Gpt-4 technical report, 2024.
- 151 [7] Hugo Touvron and Louis Martin et.al. Llama 2: Open foundation and fine-tuned chat models,
152 2023.
- 153 [8] Gemini Team et.al. Gemini: A family of highly capable multimodal models, 2024.
- 154 [9] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan,
155 and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys*, 54(10s):1–41,
156 January 2022. ISSN 1557-7341. doi: 10.1145/3505244.
- 157 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
158 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
159 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
160 recognition at scale, 2021.
- 161 [11] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big
162 sequence modeling problem. In *Advances in Neural Information Processing Systems*, 2021.
- 163 [12] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
164 state spaces, 2022.
- 165 [13] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré.
166 Combining recurrent, convolutional, and continuous-time models with linear state-space layers,
167 2021.
- 168 [14] Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It’s raw! audio generation with
169 state-space models, 2022.
- 170 [15] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces,
171 2023.
- 172 [16] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical
173 image segmentation, 2024.
- 174 [17] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and
175 Yunfan Liu. Vmamba: Visual state space model, 2024.
- 176 [18] Lianghai Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang.
177 Vision mamba: Efficient visual representation learning with bidirectional state space model,
178 2024.
- 179 [19] Zhaohu Xing, Tian Ye, Yijun Yang, Guang Liu, and Lei Zhu. Segmamba: Long-range sequential
180 modeling mamba for 3d medical image segmentation, 2024.
- 181 [20] Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Ruslan
182 Salakhutdinov. Efficient exploration via state marginal matching, 2020.

- 183 [21] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural
184 networks, 2017.
- 185 [22] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
186 deep data-driven reinforcement learning, 2021.
- 187 [23] Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential
188 for offline rl via supervised learning?, 2022.
- 189 [24] Shmuel Bar-David, Itamar Zimmerman, Eliya Nachmani, and Lior Wolf. Decision s4: Efficient
190 sequence-based rl via state spaces layers, 2023.
- 191 [25] Jeonghye Kim, Suyoung Lee, Woojun Kim, and Youngchul Sung. Decision convformer: Local
192 filtering in metaformer is sufficient for decision making, 2023.
- 193 [26] Patrick T. Komiske, Eric M. Metodiev, and Jesse Thaler. The hidden geometry of particle
194 collisions. *Journal of High Energy Physics*, 2020(7), July 2020. ISSN 1029-8479. doi:
195 10.1007/jhep07(2020)006.
- 196 [27] Victor M. Panaretos and Yoav Zemel. Statistical aspects of wasserstein distances. *Annual
197 Review of Statistics and Its Application*, 6(1):405–431, March 2019. ISSN 2326-831X. doi:
198 10.1146/annurev-statistics-030718-104938.

199 **A Proof of Theorem 1**

200 Before the prove of the main theorem, let's consider two lemmas.

201 **Lemma 2.** *Suppose a, b, c are three vectors, then we have:*

$$\langle |a-b|, 2b \rangle - \langle |a-c|, 2c \rangle \leq \langle b+c+|a-b|+|a-c|, |b-c| \rangle$$

202 *Proof.* Consider the following inequation:

$$\langle |a-c|, b-c-|b-c| \rangle \leq \langle |a-b|, c-b+|c-b| \rangle \quad (5)$$

203 According to the nature of the absolute value, it is obvious that $b-c-|b-c| \leq 0$ and that $c-b+|c-b| \geq 0$, which shows the constancy of (5) is obvious.

205 Transform (5), we can get:

$$\langle |a-c|, b \rangle + \langle |a-b|, b-|b-c| \rangle \leq \langle |a-c|, c+|b-c| \rangle + \langle c, |a-b| \rangle$$

206 Based on the absolute value inequality, we can get:

$$|a-b|-|b-c| \leq |a-c|; |a-c|+|b-c| \geq |a-b|$$

207 Hence:

$$\begin{aligned} \langle |a-b|-|b-c|, b \rangle + \langle |a-b|, b-|b-c| \rangle &\leq \langle |a-c|, b \rangle + \langle |a-b|, b-|b-c| \rangle \\ &\leq \langle |a-c|, c+|b-c| \rangle + \langle c, |a-b| \rangle \\ &\leq \langle |a-c|, c+|b-c| \rangle + \langle c, |a-c|+|b-c| \rangle \end{aligned}$$

208 Combine some items of the same kind:

$$\langle |a-b|-|b-c|, b \rangle + \langle |a-b|, b-|b-c| \rangle \leq \langle |a-c|+|b-c|, c \rangle + \langle |a-c|, c+|b-c| \rangle$$

209 Thus,

$$2\langle |a-b|, b \rangle - 2\langle |a-c|, c \rangle \leq \langle b, |b-c| \rangle + \langle c, |b-c| \rangle + \langle |a-b|, |b-c| \rangle + \langle |a-c|, |b-c| \rangle$$

210 Then,

$$\langle |a-b|, 2b \rangle - \langle |a-c|, 2c \rangle \leq \langle b+c+|a-b|+|a-c|, |b-c| \rangle$$

211 , which completes the proof. □

212 **Remark A.1.** *Suppose a, b, c are three vectors, then we have:*

$$\langle |a-b|, b \rangle - \langle |a-c|, c \rangle \leq \left\langle \frac{b+c+|a-b|+|a-c|}{2}, |b-c| \right\rangle$$

213 **Lemma 3** (Holder's inequation). *Set $p > 1, 1/p + 1/q = 1$, if $a_1, a_2 \dots a_n$ and $b_1, b_2 \dots b_n$ is nonnegative, then we have:*

$$\sum_{i=1}^n a_i b_i \leq \left(\sum_{i=1}^n a_i^p \right)^{\frac{1}{p}} \left(\sum_{i=1}^n b_i^q \right)^{\frac{1}{q}}$$

215 **Remark A.2.** *a and b are two vectors, and each of their terms is nonnegative. Then, we can get:*

$$\langle a, b \rangle \leq \langle \|a\|_1, \|b\|_\infty \rangle$$

216 , where $\|a\|_1$ represents the L_1 -norm of vector a and $\|b\|_\infty$ represents the L_∞ -norm of vector b .

217 *Proof.* Setting p as ∞ and q as 1 , then according to:

$$\sum_{i=1}^n a_i b_i \leq \left(\sum_{i=1}^n b_i^p \right)^{\frac{1}{p}} \left(\sum_{i=1}^n a_i^q \right)^{\frac{1}{q}}$$

218 We can have:

$$\sum_{i=1}^n a_i b_i \leq \|b\|_\infty \cdot \|a\|_1$$

219 Thus,

$$\langle a, b \rangle \leq \langle \|a\|_1, \|b\|_\infty \rangle$$

220 , which completes the proof □

221 Now, let's start considering the proof of **Theorem 1** as follows:

222 *Proof.* For any state-action marginal distribution, we set that $\hat{s} = \rho^{\pi_\theta}(s, a)$. We know that $\rho^{\pi_\theta}(s, a)$
 223 $= \rho^{\pi_\theta}(s) \cdot \pi_\theta(a|s)$, and given policy π , the action is sampled with the policy $\pi_\theta(a|s)$. Hence, we can
 224 get:

$$\begin{aligned} \text{EA2E}(\pi_\theta) &= \mathbb{E}_A \left[\mathbb{E}_{\hat{s}|A} [\mathbb{P}(U = 1 | \hat{s} = s, A = a) - s] \right] \\ &= \mathbb{E}_A \left[\mathbb{E}_{\rho^{\pi_\theta}(s, a)} [\mathbb{P}(U = 1 | \hat{s} = \rho^{\pi_\theta}(s, a), A = a) - \rho^{\pi_\theta}(s, a)] \right] \\ &= \mathbb{E}_A \left[\sum \rho^{\pi_\theta}(s, a) \cdot |\mathbb{P}(U = 1 | \hat{s} = \rho^{\pi_\theta}(s, a), A = a) - \rho^{\pi_\theta}(s, a)| \right] \end{aligned}$$

225 Setting ρ^π be the conditional distribution $\mathbb{P}(U = 1 | \hat{s} = \rho^{\pi_\theta}(s, a), A = a)$. Thus,

$$\text{EA2E}(\pi_\theta) = \mathbb{E}_A [\langle |\rho^\pi - \rho^{\pi_\theta}(s, a)|, \rho^{\pi_\theta}(s, a) \rangle]$$

226 , where $|\rho^\pi - \rho^{\pi_\theta}(s, a)|$ and $\rho^{\pi_\theta}(s, a)$ are vectors and $\langle |\rho^\pi - \rho^{\pi_\theta}(s, a)|, \rho^{\pi_\theta}(s, a) \rangle$ represents the inner
 227 product of $|\rho^\pi - \rho^{\pi_\theta}(s, a)|$ and $\rho^{\pi_\theta}(s, a)$. Further, let's compare the EA2E of two models (agents)
 228 $\theta_1, \theta_2 \in \Theta$:

$$\text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) = \mathbb{E}_A [\langle |\rho^\pi - \rho^{\pi_{\theta_1}}(s, a)|, \rho^{\pi_{\theta_1}}(s, a) \rangle - \langle |\rho^\pi - \rho^{\pi_{\theta_2}}(s, a)|, \rho^{\pi_{\theta_2}}(s, a) \rangle]$$

229 According to the Lemma 2, we can get:

$$\begin{aligned} \text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) &= \mathbb{E}_A [\langle |\rho^\pi - \rho^{\pi_{\theta_1}}(s, a)|, \rho^{\pi_{\theta_1}}(s, a) \rangle - \langle |\rho^\pi - \rho^{\pi_{\theta_2}}(s, a)|, \rho^{\pi_{\theta_2}}(s, a) \rangle] \\ &\leq \mathbb{E}_A \left[\left\langle \frac{|\rho^{\pi_{\theta_1}}(s, a) + \rho^{\pi_{\theta_2}}(s, a)| + |\rho^\pi - \rho^{\pi_{\theta_1}}(s, a)| + |\rho^\pi - \rho^{\pi_{\theta_2}}(s, a)|}{2}, |\rho^{\pi_{\theta_1}}(s, a) - \rho^{\pi_{\theta_2}}(s, a)| \right\rangle \right] \end{aligned}$$

230 According to Lemma 3(Holder's inequality), we have that:

$$\begin{aligned} \text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) &\leq \mathbb{E}_A \left[\left\langle \frac{|\rho^{\pi_{\theta_1}} + \rho^{\pi_{\theta_2}}| + |\rho^\pi - \rho^{\pi_{\theta_1}}| + |\rho^\pi - \rho^{\pi_{\theta_2}}|}{2}, |\rho^{\pi_{\theta_1}} - \rho^{\pi_{\theta_2}}| \right\rangle \right] \\ &\leq \mathbb{E}_A \left[\|\rho^{\pi_{\theta_1}} - \rho^{\pi_{\theta_2}}\|_1 \cdot \left\| \frac{|\rho^{\pi_{\theta_1}} + \rho^{\pi_{\theta_2}}| + |\rho^\pi - \rho^{\pi_{\theta_1}}| + |\rho^\pi - \rho^{\pi_{\theta_2}}|}{2} \right\|_\infty \right] \end{aligned}$$

231 Setting

$$m(\pi_{\theta_1}, \pi_{\theta_2}, \pi) = \left\| \frac{|\rho^{\pi_{\theta_1}} + \rho^{\pi_{\theta_2}}| + |\rho^\pi - \rho^{\pi_{\theta_1}}| + |\rho^\pi - \rho^{\pi_{\theta_2}}|}{2} \right\|_\infty$$

232 For the sake that each term of the distributions $\rho^{\pi_{\theta_1}}, \rho^{\pi_{\theta_2}}$ and ρ^π are bounded in $[0, 1]$, hence it is
 233 evident that $m(\pi_{\theta_1}, \pi_{\theta_2}, \pi) \leq 2$. Therefore, we can get:

$$\text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) \leq \mathbb{E}_A [2 \cdot \|\rho^{\pi_{\theta_1}}(s, a), \rho^{\pi_{\theta_2}}(s, a)\|_1] = \mathbb{E} [4 \cdot \text{TV}(\rho^{\pi_{\theta_1}}(s, a), \rho^{\pi_{\theta_2}}(s, a))]$$

234 , which completes the prove.

235 □

236 B Proof of Proposition 1

237 According to Theorem 1, we have that:

$$\text{EA2E}(\pi_{\theta_1}) - \text{EA2E}(\pi_{\theta_2}) \leq \mathbb{E} [4 \cdot \text{TV}(\rho^{\pi_{\theta_1}}(s, a), \rho^{\pi_{\theta_2}}(s, a))]$$

238 In work [27], the following conclusion has been given:

239 **Lemma 4.** Setting X, Y are finitely discrete random variables, and they are bounded; $X \in \mu$ and
 240 $Y \in \nu, \Omega = \text{supp}(\mu) \cup \text{supp}(\nu); d_{\min} = \inf_{x \neq y \in \Omega} \|x - y\|$, then we have :

$$\text{TV}(X, Y) \leq \frac{1}{d_{\min}} \cdot W_1(X, Y)$$

241 Hence, let's consider the proof of **Proposition 1**.

242 *Proof.* According to Theorem 1, we have:

$$\text{EA2E}(\pi) - \text{EA2E}(\pi_\theta) \leq \mathbb{E}[4 \cdot \text{TV}(\rho^\pi(s, a), \rho^{\pi_\theta}(s, a))]$$

243 According to Lemma 4, setting $\rho^\pi(s, a) \in \mu, \rho^{\pi_\theta}(s, a) \in \nu, \Omega = \text{supp}(\mu) \cup \text{supp}(\nu), d_{\min} =$
 244 $\inf_{\rho^\pi \neq \rho^{\pi_\theta} \in \Omega} \|\rho^\pi - \rho^{\pi_\theta}\|$, we can derive that:

$$\begin{aligned} \text{EA2E}(\pi) - \text{EA2E}(\pi_\theta) &\leq \mathbb{E}[4 \cdot \text{TV}(\rho^\pi(s, a), \rho^{\pi_\theta}(s, a))] \\ &\leq \mathbb{E}\left[\frac{4}{d_{\min}} \cdot W_1(\rho^\pi(s, a), \rho^{\pi_\theta}(s, a))\right] \end{aligned}$$

245 ,which completes the proof. □

246 C Visualizations

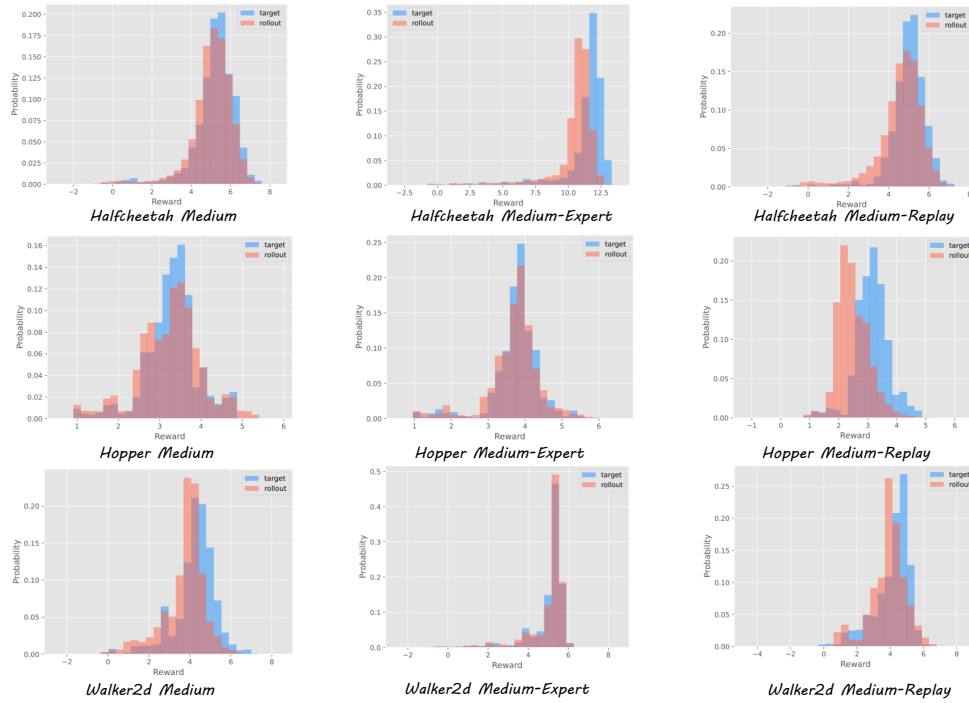


Figure 2: Visualizations of the distributional matching results on offline datasets.

247 **D Experimental Hyperparameters**

Table 3: Hyperparameters of CDMamba on the D4RL datasets.

Hyperparameter	Value
Number of Layers	3
Batch Size	64
Context Length K	20
Embedding Dimension	128
Distribution Dimension	30
Number of bins for categorical distribution	31
Learning Rate	1×10^{-4}
Learning Rate Decay	Linear warmup for first 100k training steps
Grad Norm Clip	0.25
Weight Decay	1×10^{-4}