

PoRank: A Practical Framework for Learning to Rank Policies

Pengjie Gu¹, Mengchen Zhao^{2,†}, Xu He³ and Yi Cai² and Bo An^{1,4}

¹School of Computer Science and Engineering, Nanyang Technological University, Singapore

²School of Software Engineering, South China University of Technology, China

³Huawei Noah' Ark Lab

⁴Skywork AI, Singapore

{pengjie.gu, boan}@ntu.edu.sg, {zzmc, ycai}@scut.edu.cn, hexu27@huawei.com

Abstract

In many real-world scenarios, we need to select from a set of candidate policies before online deployment. Although existing off-policy evaluation (OPE) methods can be used to estimate the online performance, they suffer from high variance. Fortunately, we care only about the ranking of the candidate policies, rather than their exact online rewards. Based on this, we propose a novel framework PoRank for learning to rank policies. In practice, learning to rank policies faces two main challenges: 1) generalization over the huge policy space and 2) lack of supervision signals. To overcome the first challenge, PoRank uses a Policy Comparison Transformer (PCT) for learning cross-policy representations, which capture the core discrepancies between policies and generalizes well across the whole policy space. The second challenge arises because learning to rank requires online comparisons of policies as ground-truth labels, whereas deploying policies online might be highly expensive. To overcome this, PoRank adopts a crowdsourcing based learning-to-rank (LTR) framework, where a set of OPE algorithms are employed to provide weak comparison labels. Experimental results show that PoRank not only outperforms baselines when the ground-truth labels are provided, but also achieves competitive performance when the ground-truth labels are unavailable.

1 Introduction

¹ In many real-world scenarios such as trading [Zhang *et al.*, 2020], advertising [Cai *et al.*, 2023], autonomous vehicles [Shi *et al.*, 2021], and drug trials [Yang *et al.*, 2023], the task often involves selecting the most promising policy from a set of candidates prior to online deployment. This selection is critical, as online evaluation of each policy can be costly and potentially risky. The conventional approach for policy selection is to estimate the online performance of candidate policies via Off-Policy Evaluation (OPE) [Paduraru, 2013]. OPE

requires only off-policy data, which means that we can estimate the performance of a target policy using data generated by other policies.

Although OPE is a promising direction, existing OPE methods are still far from reliable in practice. For example, standard Inverse Propensity Scoring (IPS) based estimators such as importance sampling suffer from high variance due to the product of importance weights [Hanna *et al.*, 2019]. Direct Methods (DM) requires extra estimators of environmental dynamics or value functions, which are hard to learn when the observation data is high-dimensional or insufficient. Hybrid Methods (HM) such as doubly robust estimators combine IPS and DM [Jiang and Li, 2016], yet it often comes with additional hyperparameters that need to be carefully chosen.

Fortunately, in many real-world scenarios, we do not need to estimate the exact online performance of candidate policies. Instead, we only care about which policy would perform the best when deployed online. This inspires us to develop a policy ranker focusing on predicting the ranking of target policies regarding to their online performance. Learning such a policy ranker is similar to learning item rankers in recommender systems [Liu, 2009]. However, learning to rank policies faces two unique challenges. First, the policy space could be extremely large, even in simple environments. Therefore, a policy ranker with poor generalization ability would fail to rank unseen policies. Second, unless deployed online, we can hardly know the real performance of the policies, which means that we are lack of ground-truth labels for training the policy ranker. These challenges make it extremely hard to train useful policy rankers in real-world tasks.

In this paper, we propose a novel and practical framework called PoRank for learning to rank policies. PoRank composes of a Policy Comparison Transformer (PCT) and a learning-to-rank (LTR) module, where the PCT aims to learn compact representations of policy pairs and the LTR module focuses on predicting the order of input policies. Instead of directly encoding the raw trajectories to policy representations (as is done in existing works), PCT encodes different behaviors of two policies at the same states to represent a policy pair. In such a way, PCT successfully extracts the most useful features for policy ranking and significantly improves the generalization over the policy space. The LTR module is built upon the PCT. Conventionally, given two input policies, LTR requires a label to tell which policy performs bet-

[†] Corresponding author.

ter. If we can easily obtain such labels, for example in the game environments, we can directly train the policy ranker using supervised learning. However, in many industrial scenarios, obtaining such labels could be very expensive due to deployment cost. To resolve this, we propose a novel crowdsourcing based LTR module, where a set of OPE methods are employed to estimate the policy comparison labels. Specifically, these labels are constructed by comparing the estimated accumulated rewards of the target policies. In such a way, we can easily collect plenty of labels without deploying training policies online. In the experiments, we first demonstrate the advantage of the PCT architecture by comparing it with the state-of-the-art policy ranking network, where both networks are trained using ground-truth labels. Hence, we show that a simple crowdsourcing mechanism could compensate for the lack of ground-truth labels, which makes PoRank practical in many real-world scenarios.

2 Related Works

Off-Policy Evaluation/Selection/Ranking The goal of OPE is to precisely predict the online performance of target policies given trajectory data collected by some other behavior policies. Standard importance sampling approach suffers from exponential variance with respect to the time horizon [Li *et al.*, 2015; Jiang and Li, 2016]. Recent works such as Fitted-Q evaluation [Hoang *et al.*, 2019] and marginalized importance sampling [Liu *et al.*, 2018] achieve polynomial variance, yet they rely on additional function approximators. Direct methods avoid the large variance by learning the dynamic model or Q-function, which could be biased especially when the data is insufficient. Some works study the offline policy selection problem, yet their methods require running from scratch for each candidate policy [Zhang and Jiang, 2021; Mengjiao *et al.*, 2022]. By contrast, in offline policy ranking (OPR) problem we aim to develop a policy ranker that could directly rank policies. A recent work on OPR collects online performance of a set of policies and uses these labeled data to train the policy ranker [Jin *et al.*, 2022]. However, collecting such data might be extremely expensive in practice.

Learning from Crowds Crowdsourcing systems enable machine learners to collect labels of large datasets from crowds. One big issue with crowdsourcing is that the labels provided by crowds are often noisy [S. and Zhang, 2019]. To tackle this challenge, various probabilistic generative methods are proposed for statistical inference [Yuchen *et al.*, 2016; Tian and Zhu, 2015]. Another line of works use discriminative models that find the most likely label for each instance [Jing *et al.*, 2014; Jing *et al.*, 2015]. A recently work called Crowd Layer (CL) first describes an algorithm for jointly learning the target model and the reliability of workers [Filipe and Pereira, 2018]. CL proposes a simple yet efficient crowd layer that can train deep neural networks end-to-end directly from the noisy labels. In our work, we treat existing OPE methods as workers and adopt CL to process multiple labels due to its simplicity and effectiveness.

Policy Representations Compact but informative representations of policies not only benefit the policy learning process [Tang *et al.*, 2022], but also help with the policy trans-

fer among different tasks [Isac *et al.*, 2019; G. *et al.*, 2017]. A straightforward idea is to represent a policy by its network parameters, yet this leads to a very sparse representation space. Network Fingerprint [Harb *et al.*, 2020] proposes a differentiable representation that uses the concatenation of the vectors of actions outputted by the policy network on a set of probing states. Some recent works try to encode policy parameters as well as state-action pair data into a low-dimensional embedding space [Tang *et al.*, 2022; Jin *et al.*, 2022]. However, existing works focus on single policy representations, which fail to capture the relative discrepancies between different policies.

3 Problem Statement

Markov Decision Process We consider the underlying environment as a Markov decision process (MDP) and define an MDP as a finite-horizon tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{P}, \mathcal{R}, \gamma)$. Here, \mathcal{S} is the state space, and \mathcal{A} is the action space. \mathcal{T} is the length of time horizon. \mathcal{P} and \mathcal{R} are the transition function and the reward function, respectively. $\mathcal{P}(s_{t+1}|s_t, a_t)$ represents the probability of transitioning from state s_t to state $s_{t+1} \in \mathcal{S}$ when the agent takes action $a_t \in \mathcal{A}$ under state $s_t \in \mathcal{S}$ and $\mathcal{R}(s_t, a_t)$ represents the immediate reward the agent receives. The expected return of a policy π can be computed by $\mathbb{E}_{\mathcal{P}} \left[\sum_{t=1}^{\mathcal{T}} [\gamma^t \mathcal{R}(s_t, \pi(s_t))] \right]$, where $\gamma \in (0, 1]$ is the discount factor.

Ranking Policies without Online Deployment The goal of OPE is to estimate the expected return of a policy π without deploying it online, given an offline dataset $\mathcal{D} = \{\tau^k\}_{k=1}^N$, where $\tau^k = (s_0, a_0, r_0, \dots, s_{\mathcal{T}}, a_{\mathcal{T}}, r_{\mathcal{T}})^k$ are trajectories generated by some behavior policies. OPE is usually used for model selection: We are required to select the most promising policy from a candidate set of available policies before actual deployment. Take recommender systems as example, we can easily obtain a set of candidate policies by adjusting the training data or the hyperparameters of the model. However, we often need to select very few policies from the candidates for online test, since a bad policy would harm the user experience. Therefore, we care more about the ranking of the candidate policies, instead of their exact expected reward. We formally define the off-policy ranking problem as follows.

Definition 1 (Offline Policy Ranking, OPR). Given a set of trajectory data $\mathcal{D} = \{\tau^k\}_{k=1}^N$ generated by some behavior policies and a set of target policies $\Pi = \{\pi^i\}_{i=1}^M$, the OPR problem is to seek a ranking of the target policies that aligns with their online expected accumulated rewards.

Intuitively, OPR is a simpler problem than OPE because we do not care about the exact online performance of the candidate policies. While OPE as well as off-policy selection (OPS) can also be used to rank candidate policies, they often need to inputrun an estimation procedure from scratch for each candidate policy, thus leading to poor efficiency. In this work, we aim to learn an universal policy ranker that could directly return the ranking results for any input candidate policies.

4 Approach

In this section, we will elaborate our PoRank framework. PoRank borrows the idea from learning-to-rank literature and explicitly solves two unique challenges in ranking policies. Firstly, given the vastness of the policy space, we need an effective policy representation scheme so that the ranking module can accurately comprehend and compare different policies. Secondly, although we can easily generate many policies for training, we can hardly know their actual online performance, therefore we are lack of ground-truth labels for training the ranker. Overall, PoRank consists of a Policy Comparison Transformer (PCT) and a Learning-to-Rank (LTR) module, as shown in Figure 2(a). We will introduce each of them in the following subsections.

4.1 Learning Cross-policy Representations

Cross-policy Representation A policy can be considered as a conditional distribution over actions given the state. Therefore, a policy can be naturally represented by a set of state-action pairs where the actions are sampled from the policy. However, such a naive policy representation could be inefficient since the number of state-action pairs can be extremely large. Previous works address this issue by extracting high-level features from the state-action pairs using deep neural networks [Jin *et al.*, 2022]. Although these representations reflect the features of single policies, they fail to capture the discrepancies of different policies at some crucial states.

To this end, we aim to learn cross-policy representations by comparing two policies’ decisions at the same set of states. Formally, given a set of states $\{s_1, \dots, s_K\}$ and two policies π^i, π^j , we can construct the following sequence of state-action pairs by taking actions at these states:

$$\xi^{i>j} = \left\{ (s_1, a_1^i), (s_1, a_1^j), \dots, (s_K, a_K^i), (s_K, a_K^j) \right\}, \quad (1)$$

where $a^i \sim \pi^i(\cdot|s)$, and $a^j \sim \pi^j(\cdot|s)$. We denote by $\chi^{i>j} = g(\xi^{i>j}) \in \mathbb{R}^n$ the cross-policy representation where g is a function that maps $\xi^{i>j}$ to an n -dimensional representation space. Figure 2(a) shows the computation of cross-policy representations. By contrast, Figure 2(b) shows the computation of single-policy representations, which is adopted in [Jin *et al.*, 2022]. Intuitively, cross-policy representations focus on encoding the discrepancies between policies, while single-policy representations only encode features of single policies. Thus, cross-policy representations are more effective for downstream learning-to-rank tasks.

Policy Comparison Transformer (PCT) Transformers are proved to be effective for learning dependencies between different positions in sequences. Prior works has employed transformers to extract features from trajectory sequences [Lili *et al.*, 2021; Michael *et al.*, 2021]. However, existing transformer architectures fail to capture the differences of two policies’ decisions. In our work, we propose the PCT architecture to learn cross-policy representations. Unlike previous works where the positional encodings indicate the positions of state-action pairs in a trajectory, PCT uses positional encoding to distinguish the relative order of two policies. In

this way, the learned cross-policy representation $\chi^{i>j}$ can be directly used to predict whether π^i performs better than π^j .

Figure 1 shows the construction of input tokens. We first sample K states from \mathcal{D} and then use a linear encoder f to map the K state-action pairs into $2K$ tokens:

$$x_k^i = f(s_k, a_k^i), \quad x_k^j = f(s_k, a_k^j), \quad k = 1, \dots, K \quad (2)$$

where i and j represent the indexes of two policies. In order to represent the relative order of π^i and π^j , we introduce two one-hot positional encodings $e_+ = [1, 0]$ and $e_- = [0, 1]$, where e_+ indicates the policy ranked higher and e_- indicates the policy ranked lower. We also use an aggregation token e_0 , which is a learnable vector for aggregating the information from the other $2K$ tokens [Zhu *et al.*, 2021]. The final inputs that indicate π^i ranked higher than π^j can be represented as:

$$z^{i>j} = \left[e_0, x_1^i + e_+, x_1^j + e_-, \dots, x_K^i + e_+, x_K^j + e_- \right] \quad (3)$$

This construction of inputs has two advantages. First, the two policies share the same set of states, thus their discrepancies are naturally represented by the different actions taken at these states. Second, we can easily get a mirrored representation $z^{j>i}$ by simply exchange the positional encoding e_+ and e_- used in $z^{i>j}$.

We adopt a widely used transformer architecture as our encoder [Dosovitskiy *et al.*, 2021]. It contains L alternating layers of multi-head self-attention (MSA) and multi-layer perceptron (MLP) blocks. Layernorm (LN) and residual connections are applied to the outputs of each block. For brevity, we rewrite the inputs in Equation 3 as $\mathbf{z}^{(0)}$. And the computations at each block can be represented as:

$$\begin{aligned} \hat{\mathbf{z}}^l &= \text{MSA}(\text{LN}(\mathbf{z}^{(l-1)})) + \mathbf{z}^{(l-1)} & l = 1, \dots, L \\ \mathbf{z}^l &= \text{MLP}(\text{LN}(\hat{\mathbf{z}}^{(l-1)})) + \hat{\mathbf{z}}^{(l-1)} & l = 1, \dots, L \\ \chi^{i>j} &= \text{LN}(\mathbf{z}^L). \end{aligned} \quad (4)$$

The final cross-policy representation $\chi^{i>j}$ is the corresponding outputs of the aggregation token e_0 taken from \mathbf{z}^L . Note that $\chi^{i>j}$ changes to $\chi^{j>i}$ when we exchange the positional encodings e_+ and e_- , but they are permutation invariant to the order of inputted state-action pairs.

4.2 Learning to Rank Policies

In this section, we will introduce how to train the PCT in two cases regarding to the existence of ground-truth ranking labels of policies. First, we show that the policy ranking problem can be reduced to a binary classification problem since our cross-policy representations can be directly used to predict the ranking of two policies. Second, we will introduce a learning paradigm where multiple OPE methods are modeled as label providers. We will also show how to train the PCT leveraging the inaccurate labels provided by the workers.

Reducing OPR to Binary Classification We first consider the case when there is a training set $\Pi = \{(\pi^i, R^i)\}_{i=1}^T$ consisting of T deployed policies π^i as well as their real expected accumulated rewards R^i . In this case, we can directly construct binary labels by comparing the performance of the two

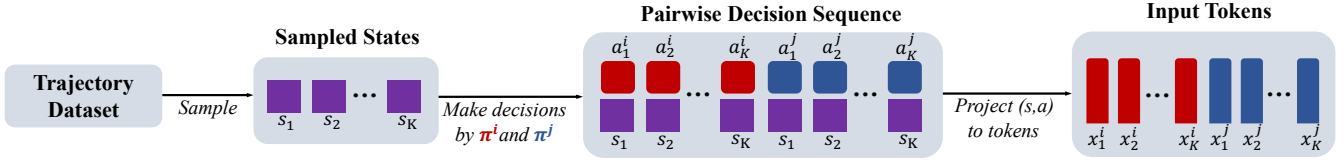


Figure 1: Generation of input tokens.

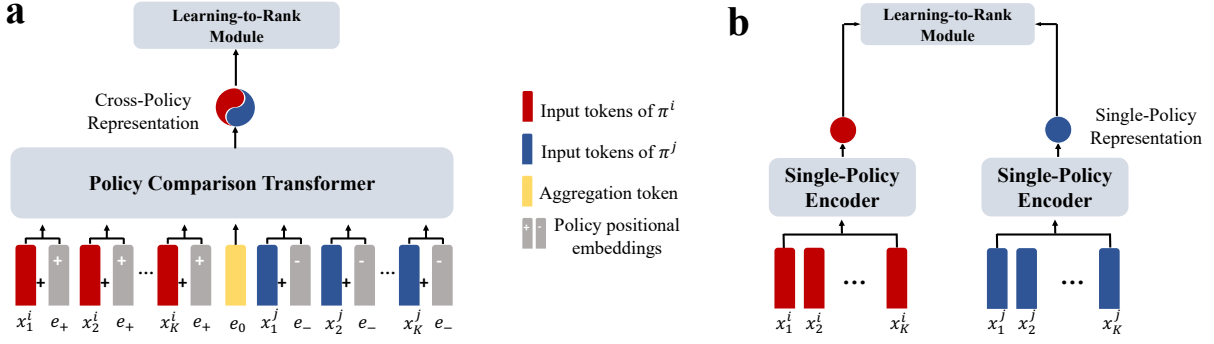


Figure 2: Comparison between Cross-Policy Representation (left) and Single-Policy Representation (right).

294 policies. We use an indicator $\mathbb{1}_{R^i > R^j}$ to represent the label
 295 of a pair of policies (π^i, π^j) . The PCT can be trained by min-
 296 imizing the following binary cross entropy loss:

$$\mathcal{L}_{sup} = - \mathbb{E}_{\pi^i, \pi^j \sim \Pi} \left[\left(\mathbb{1}_{R^i > R^j} \right) \cdot \log \left(\hat{y}_{i>j} \right) + \left(\mathbb{1}_{R^i \leq R^j} \right) \cdot \log \left(1 - \hat{y}_{i>j} \right) \right], \quad (5)$$

297 where $\hat{y}_{i>j} = \text{sigmoid}(\phi(\chi^{i>j}))$ represents the predicted
 298 probability that π^i performs better than π^j . ϕ is a function
 299 that projects $\chi^{i>j}$ to a real number. The final ranking of test
 300 policies is based on their scores computed by:

$$\text{score}_i = \frac{1}{N} \sum_{j \neq i} \hat{y}_{i>j}, i = 1, \dots, N, \quad (6)$$

301 which can be interpreted as the expected probability that π^i
 302 performs better than other test policies [Rodrigo *et al.*, 2019].

303 Figure 3 shows the framework of the Learning-to-rank
 304 module. In the presence of ground-truth label $y_{i>j}$, the LTR
 305 module is simplified to Situation 1. Otherwise, we will rely
 306 on OPE workers to provide supervision signals. Situation 2
 307 will be illustrated below.

308 **Learning from OPE Workers** Supervised training is effi-
 309 cient when the dataset $\Pi = \{(\pi^i, R^i)\}_{i=1}^T$ contains enough
 310 policies. Unfortunately, collecting such training data can be
 311 extremely expensive in many real applications. Meanwhile,
 312 we note that although existing OPE methods are not robust
 313 enough, they actually provide candidate solutions to the OPR
 314 problem. To this end, we borrow ideas from crowdsourc-
 315 ing domain as an alternative way to approach the OPR prob-
 316 lem. Specifically, suppose that there exists a set of OPE al-
 317 gorithms estimating the policy performance, we can employ

318 them as crowd workers to generate possibly inaccurate labels
 319 and make use of these labels to train our models. The intu-
 320 ition is that the inaccurate labels generated by OPE workers
 321 are implicitly conditioned on the ground-truth performance
 322 of policies. If we can take advantage of these labels and learn
 323 their relationships with the ground-truth labels, our prediction
 324 $\hat{y}_{i>j}$ would be more close to the ground-truth labels.

325 In the framework of PoRank, we adopt Crowd Layer (CL,
 326 [Filipe and Pereira, 2018]) as our backend for learning from
 327 crowd labels. CL is able to automatically distinguish the good
 328 from the unreliable workers and capture their individual bi-
 329 ases in many other domains, such as image annotation [Guan
 330 *et al.*, 2018; Li *et al.*, 2022] and music genre classification
 331 [Rodrigues *et al.*, 2013]. In addition, CL is naturally com-
 332 patible with deep learning approaches since it simply adds a
 333 crowd layer to the deep neural networks and can be trained in
 334 an end-to-end way. As shown in Figure 2, we add CL to the
 335 top of our predicted probability $\hat{y}_{i>j}$. During training, CL ad-
 336 justs the gradients coming from these noisy labels according
 337 to its reliability by scaling them and adjusting their bias. The
 338 adjusted gradients are then backpropagated to PCT according
 339 to the chain rule.

340 Formally, assume that there are W workers of OPE meth-
 341 ods. For each worker w_m , its estimation about the expected
 342 return of π^i is denoted as R_m^i . The goal of CL is to train a
 343 mapping function $\hat{y}_{i>j}^m = \zeta^m(\hat{y}_{i>j})$ to predict the noisy bi-
 344 nary label generated by worker w_m : $y_{i,j}^m = \mathbb{1}_{R_m^i > R_m^j}$. The
 345 overall objective can be written as:

$$\mathcal{L}_{CL} = - \mathbb{E}_{\substack{m=1, \dots, W \\ \pi^i, \pi^j \sim \Pi}} \left[y_{i,j}^m \cdot \log \left(\hat{y}_{i>j}^m \right) + (1 - y_{i,j}^m) \cdot \log \left(1 - \hat{y}_{i>j}^m \right) \right]. \quad (7)$$

346 The complete training procedures of PoRank is summa-

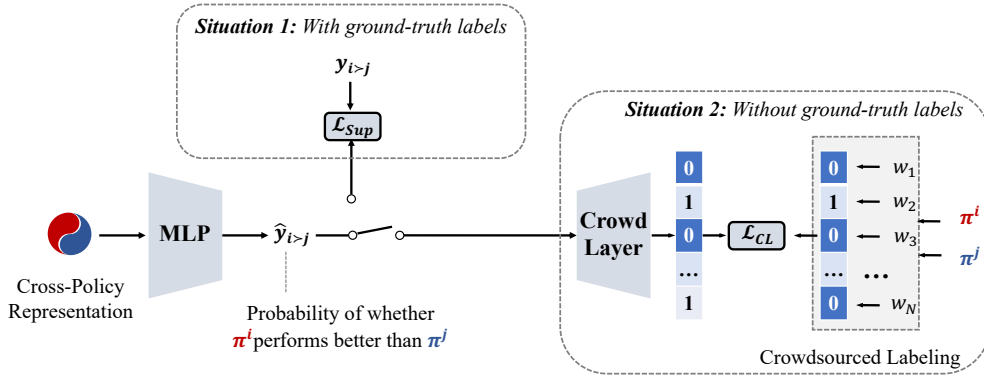


Figure 3: Framework of the Learning-to-Rank Module.

347 rized in Appendix. In practice, to reduce the computational
 348 cost brought by CL, we set ζ^m as a linear projection followed
 349 by a sigmoid function. Therefore, the number of additional
 350 parameters only grows linearly with the number of workers.
 351 Note that the CL is only available during training since we
 352 still use $\hat{y}_{i>j}$ to generate the predicted ranking of policies.

353 5 Experiments

354 In this section, we compare PoRank with widely-used OPE
 355 methods on various tasks. We present ablation studies with
 356 respect to PCT and CL, which are the main components of
 357 PoRank. **Note that** hyper-parameter selection and more ex-
 358 perimental results can be found in Appendix.

359 5.1 Experimental settings

360 **Trajectory Set** We evaluate PoRank and all baseline
 361 OPE methods on D4RL dataset consisting of various trajec-
 362 tory sets [Fu *et al.*, 2020]. Overall, we adopt trajec-
 363 tory sets collected from 2 environments of Mujoco games:
 364 HalfCheetah-v2 and Walker2d-v2. Besides, there
 365 are 3 different types of trajectory sets for each environment:
 366 expert, full-replay and medium. The difference be-
 367 tween them is that the behavioral policies collecting these 3
 368 types of trajectories show different performance. And these
 369 behavioral policies are trained by the Soft Actor-Critic (SAC)
 370 algorithm online [Haarnoja *et al.*, 2018].

371 **Policy Set** To evaluate the abilities of all methods to cor-
 372 rectly rank policies. We use the policy set released by [Jin
 373 *et al.*, 2022]. For each trajectory set mentioned above, there
 374 are 2 types of policy sets (referred as Set I and Set II)
 375 in which the expected return of policies are evenly spaced
 376 in the performance range of them. As mentioned in [Jin
 377 *et al.*, 2022], Set I and Set II aim to simulate two kinds
 378 of OPE cases. The policies contained in Set I are trained
 379 by offline RL algorithms (CQL [Kumar *et al.*, 2020], BEAR
 380 [Kumar *et al.*, 2019], CRR [Wang *et al.*, 2020]) and show
 381 diverse behavioral performance. This is aligned with prac-
 382 tical cases where the sources of policies are diverse and un-
 383 known. Set II contains policies trained by SAC, which is
 384 the same as the algorithm of training the behavioral policies.
 385 Therefore, Set II corresponds to the practical OPE cases

386 in which the target policies share many common properties
 387 with the policies generating the trajectory data.

388 **Baselines**² We compare PoRank with six state-of-the-
 389 art baselines. **i)** Fitted Q-Evaluation (FQE [Hoang *et al.*,
 390 2019]). It is a value-based OPE method, which learns a neu-
 391 ral network to approximate the Q-function of the evaluated
 392 policy by temporal difference learning on the trajectory set.
 393 **ii)** Model-based estimation (MB [Paduraru, 2013]). It learns
 394 the dynamics model of environment, and estimates the ex-
 395 pected return of evaluated policies by computing their av-
 396 erage returns of Monte-Carlo rollouts in the model environ-
 397 ment. **iii)** Weighted importance sampling (IW [Mahmood *et al.*,
 398 2014]). It leverages weighted importance sampling to cor-
 399 rect the weight of the reward, regarding the collected trajec-
 400 tory data distribution to the data distribution of the evaluated
 401 policy. **iv)** DualDICE [Nachum *et al.*, 2019]. It also aims to
 402 achieve distribution correction yet without directly using im-
 403 portance sampling. It learns an estimation of the state-action
 404 stationary distribution for achieving distribution correction.
 405 **v)** Doubly Robust (DR [Jiang and Li, 2016]). It utilizes an
 406 unbiased estimator that leverages an estimated environment
 407 model to decrease the variance of the unbiased estimates pro-
 408 duced by importance sampling techniques.

409 **Worker Set** In all experiments, we use 15 models trained
 410 by 5 OPE methods mentioned above (IW, MB, DR,
 411 DualDICE, and FQE, each trained 3 models with different
 412 seeds) as our OPE workers. We present the architecture de-
 413 tails for them in Appendix.

414 **Evaluation Metrics** We evaluate all models accord-
 415 ing two widely-used metrics. **i)** Spearman’s Rank
 416 Correlation. It is the Pearson correlation between the
 417 ground truth rank sequence and the evaluated rank sequence
 418 of the evaluated policies. So, a higher rank correlation indi-
 419 cates a better policy ranker. **ii)** Normalized Regret@k.
 420 It is the normalized difference between the actual value of
 421 the best policy in the policy set, and the actual value of the
 422 best policy in the estimated top-k set. Mathematically, it can

²We leverage a popular implementation of OPE algo-
 rithms: [https://github.com/google-research/google-research/tree/
 master/policy_eval](https://github.com/google-research/google-research/tree/master/policy_eval). It contains the first 5 baselines used in our paper

Table 1: Comparing PoRank with other OPE baselines.

		HalfCheetah-v2				Walker2d-v2			
		Rank Correlation \uparrow		Regret @1 \downarrow		Rank Correlation \uparrow		Regret @1 \downarrow	
		Set I	Set II	Set I	Set II	Set I	Set II	Set I	Set II
<i>Expert</i>	PoRank (Ours)	0.65 \pm 0.10	0.35 \pm 0.03	0.00 \pm 0.00	0.32 \pm 0.02	0.85 \pm 0.12	0.83 \pm 0.08	0.01 \pm 0.01	0.02 \pm 0.03
	FQE	-0.53 \pm 0.14	0.31 \pm 0.10	0.23 \pm 0.06	0.60 \pm 0.07	0.62 \pm 0.20	-0.08 \pm 0.03	0.04 \pm 0.07	0.02 \pm 0.07
	DualDICE	0.47 \pm 0.18	0.27 \pm 0.06	0.38 \pm 0.02	0.22 \pm 0.04	0.52 \pm 0.16	0.31 \pm 0.14	0.43 \pm 0.07	0.28 \pm 0.03
	MB	0.39 \pm 0.04	0.24 \pm 0.08	0.34 \pm 0.05	0.18 \pm 0.02	0.46 \pm 0.11	0.29 \pm 0.12	0.39 \pm 0.03	0.23 \pm 0.08
	IW	0.18 \pm 0.03	0.09 \pm 0.02	0.23 \pm 0.04	0.02 \pm 0.01	0.24 \pm 0.07	0.06 \pm 0.03	0.28 \pm 0.02	0.14 \pm 0.07
	DR	-0.12 \pm 0.03	-0.18 \pm 0.04	0.13 \pm 0.06	0.08 \pm 0.01	0.14 \pm 0.09	0.12 \pm 0.06	0.19 \pm 0.03	0.04 \pm 0.02
<i>Full-replay</i>	PoRank (Ours)	0.72 \pm 0.19	0.34 \pm 0.08	0.01 \pm 0.02	0.31 \pm 0.01	0.82 \pm 0.04	0.81 \pm 0.17	0.09 \pm 0.06	0.21 \pm 0.09
	FQE	0.24 \pm 0.18	0.52 \pm 0.01	0.36 \pm 0.09	0.03 \pm 0.02	0.71 \pm 0.13	-0.19 \pm 0.19	0.06 \pm 0.04	0.48 \pm 0.02
	DualDICE	-0.57 \pm 0.18	0.06 \pm 0.09	0.53 \pm 0.08	0.27 \pm 0.03	-0.26 \pm 0.14	-0.24 \pm 0.18	0.42 \pm 0.01	0.28 \pm 0.06
	MB	0.23 \pm 0.04	-0.19 \pm 0.02	0.17 \pm 0.06	0.34 \pm 0.07	0.63 \pm 0.13	0.71 \pm 0.03	0.08 \pm 0.01	0.14 \pm 0.09
	IW	-0.24 \pm 0.01	-0.31 \pm 0.04	0.42 \pm 0.07	0.46 \pm 0.02	0.11 \pm 0.09	-0.65 \pm 0.08	0.09 \pm 0.01	0.43 \pm 0.08
	DR	0.03 \pm 0.04	0.18 \pm 0.03	0.34 \pm 0.08	0.09 \pm 0.02	0.31 \pm 0.06	0.02 \pm 0.07	0.06 \pm 0.04	0.18 \pm 0.05
<i>Medium</i>	PoRank (Ours)	0.82 \pm 0.13	0.81 \pm 0.04	0.02 \pm 0.01	0.03 \pm 0.04	0.82 \pm 0.11	0.72 \pm 0.13	0.13 \pm 0.02	0.08 \pm 0.03
	FQE	0.48 \pm 0.12	-0.12 \pm 0.09	0.05 \pm 0.03	0.12 \pm 0.08	0.71 \pm 0.14	0.72 \pm 0.19	0.08 \pm 0.04	0.13 \pm 0.01
	DualDICE	-0.42 \pm 0.19	-0.28 \pm 0.11	0.32 \pm 0.07	0.13 \pm 0.02	0.43 \pm 0.17	0.23 \pm 0.14	0.03 \pm 0.08	0.28 \pm 0.06
	MB	0.12 \pm 0.03	-0.19 \pm 0.14	0.23 \pm 0.02	0.12 \pm 0.01	0.47 \pm 0.11	0.02 \pm 0.10	0.18 \pm 0.03	0.13 \pm 0.02
	IW	-0.53 \pm 0.02	-0.78 \pm 0.04	0.63 \pm 0.01	0.57 \pm 0.06	0.27 \pm 0.09	0.67 \pm 0.03	0.38 \pm 0.01	0.27 \pm 0.09
	DR	0.63 \pm 0.01	0.17 \pm 0.04	0.03 \pm 0.08	0.28 \pm 0.09	0.07 \pm 0.05	0.37 \pm 0.02	0.28 \pm 0.04	0.38 \pm 0.08

423 be computed by $regret@k = \frac{V_{max} - V_{topk}}{V_{max} - V_{min}}$, where V_{max} and
424 V_{min} is the expected return of the best and the worse policies,
425 respectively, in the entire set, while V_{topk} is the estimated top
426 k policies. So, a lower regret value indicates a better policy
427 ranker.

428 5.2 Experimental Results

429 **Comparison with Other OPE Baselines** We evaluated Po-
430 Rank against five baseline methods in two environments. Fig-
431 ure 1 presents the rank correlation and regret@1 values of the
432 estimated rank sequences generated by each model.

433 Our results demonstrate that PoRank consistently outper-
434 forms the baseline methods. Specifically, PoRank achieved
435 the highest rank correlations in 11 of the 12 policy sets and
436 the lowest regret@1 values in 7 sets. This indicates PoRank’s
437 capability to deliver robust and effective performance across
438 diverse policy scenarios.

439 Notably, PoRank excels in learning from noisy labels gen-
440 erated by other OPE workers. In our experiments, the OPE
441 workers, which generate these noisy labels for PoRank, were
442 directly sampled from the trained models of the baselines.
443 For instance, in the Set I of the expert trajectory set in
444 HalfCheetah-v2, while all OPE baselines exhibited poor per-
445 formance, PoRank achieved a high rank correlation of about
446 0.65—despite relying on these low-quality labels. We at-
447 tribute this resilience to two key factors: **i)** The Policy
448 Comparison Transformer (PCT) in PoRank effectively mit-
449 igates the biases of inaccurate workers. **ii)** The crowd layer
450 adeptly distinguishes reliable OPE workers from unreliable
451 ones and adjusts for their individual biases. In conclusion,
452 PoRank demonstrates highly effective and robust OPR results
453 across a variety of policy sets. It effectively reduces the bi-
454 ases induced by OPE workers, thereby outperforming these
455 workers significantly.

456 **Comparing with Other Label Aggregation Methods** We
457 propose to use PCT and crowd layer to aggregate the infor-
458 mation of multiple OPE workers. There are also other simple

459 mechanism to aggregate worker information. For example,
460 ranking policies by the average score of workers (denoted by
Average) or ranking policies by the number of worker votes
461 (denoted by **Major Voting**).

Table 2: Performance with different label aggregation methods.

Environment	Avg. Score	Major Voting	Ours	Ours with RA
HalfCheetah-expert	-0.34	-0.27	0.71	0.72
HalfCheetah-full-replay	0.24	0.31	0.74	0.73
HalfCheetah-medium	0.32	0.57	0.81	0.80
Walker2d-expert	0.53	0.23	0.85	0.83
Walker2d-full-replay	0.41	0.31	0.82	0.75
Walker2d-medium	0.21	0.29	0.80	0.87

Table 3: Performance with different batch size of state-action pairs.

Batch size	Avg. Rank Correlation
8	0.21
16	0.27
32	0.37
64	0.39
128	0.56
256	0.65
512	0.64
1024	0.66
2048	0.65

462 Actually, the superiority of CL over these baselines has been
463 demonstrated in [Filipe and Pereira, 2018]. However, it is
464 still valuable to reproduce this superiority in the context of
465 OPR. Specifically, we report the rank correlations in the fol-
466 lowing Table 2. We can see from the first four columns that
467 our method dominates these two baselines in all of the six
468 environments. Note that the framework of PoRank can also
469 combine with more advanced crowdsourcing methods other
470 than CL. On the other hand, the line of works in rank aggre-
471 gation, focusing on aggregating a set of pairwise comparisons
472

Table 4: **Ablations.** Comparison of policy ranking performance among PoRank (PCT with CL), and its ablation models PCT with GL, SOPR-T with GL, and SOPR-T with CL across different policy sets. Models with PCT architectures consistently achieves higher scores than models with SOPR-T architectures, underscoring the advantage of PCT in generating meaningful cross-policy representations. Notably, PCT with CL and SOPR-T with CL show competitive performance against their GL counterparts, despite lacking additional supervised information from deployed policies, affirming the effectiveness of the CL in learning from OPE-generated noisy labels.

		HalfCheetah-v2				Walker2d-v2			
		Rank Correlation \uparrow		Regret @1 \downarrow		Rank Correlation \uparrow		Regret @1 \downarrow	
		Set I	Set II	Set I	Set II	Set I	Set II	Set I	Set II
<i>Expert</i>	PoRank (PCT w/ CL)	0.65 \pm 0.10	0.35 \pm 0.03	0.00 \pm 0.00	0.32 \pm 0.02	0.85 \pm 0.12	0.83 \pm 0.08	0.01 \pm 0.01	0.02 \pm 0.03
	PCT w/ GL	0.43 \pm 0.14	0.63 \pm 0.10	0.13 \pm 0.06	0.05 \pm 0.07	0.82 \pm 0.20	0.78 \pm 0.03	0.04 \pm 0.07	0.05 \pm 0.07
	SOPR-T w/ CL	-0.47 \pm 0.18	0.47 \pm 0.06	0.58 \pm 0.02	0.52 \pm 0.04	0.72 \pm 0.16	0.71 \pm 0.14	0.13 \pm 0.07	0.18 \pm 0.03
	SOPR-T w/ GL	-0.29 \pm 0.04	0.70 \pm 0.08	0.34 \pm 0.05	0.18 \pm 0.02	0.56 \pm 0.11	0.80 \pm 0.12	0.19 \pm 0.03	0.03 \pm 0.08
<i>Full-replay</i>	PoRank (PCT w/ CL)	0.72 \pm 0.19	0.34 \pm 0.08	0.01 \pm 0.02	0.31 \pm 0.01	0.82 \pm 0.04	0.81 \pm 0.17	0.09 \pm 0.06	0.21 \pm 0.09
	PCT w/ GL	0.57 \pm 0.18	0.37 \pm 0.06	0.08 \pm 0.02	0.32 \pm 0.04	0.72 \pm 0.16	0.81 \pm 0.14	0.03 \pm 0.07	0.18 \pm 0.05
	SOPR-T w/ CL	-0.37 \pm 0.18	0.36 \pm 0.06	0.58 \pm 0.02	0.52 \pm 0.04	0.74 \pm 0.16	0.71 \pm 0.11	0.13 \pm 0.07	0.19 \pm 0.03
	SOPR-T w/ GL	-0.29 \pm 0.04	0.24 \pm 0.28	0.34 \pm 0.05	0.08 \pm 0.02	0.66 \pm 0.11	0.79 \pm 0.12	0.19 \pm 0.03	0.23 \pm 0.08
<i>Medium</i>	PoRank (PCT w/ CL)	0.82 \pm 0.13	0.81 \pm 0.04	0.02 \pm 0.01	0.03 \pm 0.04	0.82 \pm 0.11	0.72 \pm 0.13	0.13 \pm 0.02	0.08 \pm 0.03
	PCT w/ GL	0.73 \pm 0.14	0.88 \pm 0.10	0.03 \pm 0.01	0.00 \pm 0.00	0.72 \pm 0.20	0.88 \pm 0.03	0.14 \pm 0.07	0.02 \pm 0.07
	SOPR-T w/ CL	-0.27 \pm 0.18	0.47 \pm 0.06	0.58 \pm 0.02	0.02 \pm 0.04	0.72 \pm 0.16	0.81 \pm 0.14	0.14 \pm 0.07	0.08 \pm 0.03
	SOPR-T w/ GL	0.59 \pm 0.04	0.84 \pm 0.08	0.14 \pm 0.05	0.08 \pm 0.02	0.76 \pm 0.11	0.89 \pm 0.12	0.19 \pm 0.03	0.03 \pm 0.08

473 into a ranking list. Therefore, we use a more recent and simple RA method [Maystre and Grossglauser, 2017] (denoted by **RA**) to replace Equation 6 in our work. From the last two columns we can see that this method indeed further improves PoRank in some environments. However, this does not contradicts to our main contribution: modeling the OPR problem from the perspective of crowdsourcing.

480 **Selection of the Batch Size of State-action Pairs** In the training phase, the batch size of state-action pairs feeded into the Transformer is an important hyper-parameter in our model. It play the role in balancing the computational cost and the performance. We chose the number 256 as the batch size. This choice is supported by the experimental results reported in Table 3, which show the averaged rank correlations of our model with the batch size growing. We can find that when the batch size is larger than 256, the performance of our model tends to be stable.

490 **Ablations** We conducted ablation studies to evaluate the importance of each component in our framework, using the same policy sets as the primary experiments. The results are depicted in Table 4. In our framework, termed PoRank, we integrate two key components: the Policy Comparison Transformer (PCT) for generating cross-policy representations, and the Crowd Layer (CL) that aggregates information from OPE workers. This integration allows PoRank to rank policies without requiring additional ground truth labels. We compared PCT with CL (PoRank) against three different ablations: **i)** PCT with GL: This model uses our PCT architecture but discards the CL, relying on extra ground truth labels (GL) for training. It is trained using additional sets of deployed policies released by [Jin *et al.*, 2022]. Without CL, it lacks the capability to learn from OPE workers. **ii)** SOPR-T with GL: [Jin *et al.*, 2022], a transformer-based model designed to learn individual policy representations. Unlike PCT, SOPR-T does not focus on capturing the differences between policies at the decision level. Like PCT with GL, it also requires extra ground truth labels for training. **iii)** SOPR-T

510 with CL: This variation of SOPR-T attempts to learn from OPE workers by incorporating the CL, similar to PoRank, but does not require additional supervised information from deployed policies.

514 As shown in Table 4, PCT with CL outperforms SOPR-T with CL on 8 policy sets, while PCT with GL surpasses SOPR-T with GL on 8 policy sets. This suggests that our cross-policy representation, which aims to discern the nuanced decision differences between policies, exhibits stronger representational power compared to single-policy representations used by SOPR-T. Moreover, PCT with CL (PoRank) and SOPR-T with CL both demonstrate competitive performance against their counterparts PCT with GL and SOPR-T with GL, despite the absence of additional supervised information from deployed policies. This underscores the effectiveness of using the CL to learn from OPE-generated noisy labels when extra policies providing supervised labels are not readily available.

6 Conclusions

529 In this study, we introduced PoRank, a novel framework designed to learn robust off-policy rankers from a set of unreliable off-policy estimators. Unlike existing approaches that require the deployment of policies online for gathering supervision signals, PoRank innovatively leverages labels generated by existing Off-Policy Evaluation (OPE) methods. This feature significantly reduces the training cost of the ranker. Our theoretical analysis elucidates the relationships between a worker’s bias, variance, and overall quality. We further contribute a unique Policy Comparison Transformer (PCT) architecture, developed to discern the relative discrepancies between policies through effective cross-policy representations. Empirical results confirm PoRank’s superior performance over baseline models across diverse tasks and policy sets. Importantly, PoRank demonstrates excellent generalizability across various policy sets. Our ablation studies further validate the effectiveness of each individual component within the PoRank framework.

547 Acknowledgments

548 This research is supported by the National Research Founda-
549 tion, Singapore under its Industry Alignment Fund – Pre-
550 positioning (IAF-PP) Funding Initiative. Any opinions, find-
551 ings and conclusions or recommendations expressed in this
552 material are those of the author(s) and do not reflect the views
553 of National Research Foundation, Singapore.

554 References

555 [Cai *et al.*, 2023] Tianchi Cai, Jiyan Jiang, Wenpeng Zhang,
556 Shiji Zhou, Xierui Song, Li Yu, Lihong Gu, Xiaodong
557 Zeng, Jinjie Gu, and Guannan Zhang. Marketing bud-
558 get allocation with offline constrained deep reinforcement
559 learning. In *Proceedings of the Sixteenth ACM Interna-*
560 *tional Conference on Web Search and Data Mining*, pages
561 186–194, 2023.

562 [Dosovitskiy *et al.*, 2021] Alexey Dosovitskiy, Lucas Beyer,
563 Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
564 Thomas Unterthiner, Mostafa Dehghani, Matthias Min-
565 derer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and
566 Neil Houlsby. An image is worth 16x16 words: Trans-
567 formers for image recognition at scale. In *International*
568 *Conference on Learning Representations*, 2021.

569 [Filipe and Pereira, 2018] Rodrigues Filipe and Francisco
570 Pereira. Deep learning from crowds. In *Proceedings*
571 *of the AAAI Conference on Artificial Intelligence*, page
572 1611–1618, 2018.

573 [Fu *et al.*, 2020] Justin Fu, Aviral Kumar, Ofir Nachum,
574 George Tucker, and Sergey Levine. D4RL: Datasets for
575 deep data-driven reinforcement learning. *arXiv*, 2020.

576 [G. *et al.*, 2017] Bellemare Marc G., Will Dabney, and Rémi
577 Munos. A distributional perspective on reinforcement
578 learning. In *International Conference on Machine Learn-*
579 *ing*, pages 449–458, 2017.

580 [Guan *et al.*, 2018] Melody Guan, Varun Gulshan, Andrew
581 Dai, and Geoffrey Hinton. Who said what: Modeling in-
582 dividual labelers improves classification. In *Proceedings*
583 *of the AAAI Conference on Artificial Intelligence*, page
584 3109–3118, 2018.

585 [Haarnoja *et al.*, 2018] Tuomas Haarnoja, Aurick Zhou,
586 Pieter Abbeel, and Sergey Levine. Soft actor-critic:
587 Off-policy maximum entropy deep reinforcement learning
588 with a stochastic actor. In *International Conference on*
589 *Machine Learning*, pages 1861–1870, 2018.

590 [Hanna *et al.*, 2019] Josiah Hanna, Niekum Scott, and Peter
591 Stone. Importance sampling policy evaluation with an es-
592 timated behavior policy. In *International Conference on*
593 *Machine Learning*, pages 2605–2613, 2019.

594 [Harb *et al.*, 2020] Jean Harb, Tom Schaul, Doina Precup,
595 and Pierre-Luc Bacon. Policy evaluation networks. *arXiv*,
596 2020.

597 [Hoang *et al.*, 2019] Le Hoang, Cameron Voloshin, and
598 Yisong Yue. Batch policy learning under constraints.
599 In *International Conference on Machine Learning*, pages
600 3703–3712, 2019.

[Isac *et al.*, 2019] Arnekvist Isac, Danica Kragic, and Jo-
601 hannes A. Stork. VPE: Variational policy embedding for
602 transfer reinforcement learning. In *ICRA*, pages 36–42,
603 2019. 604

[Jiang and Li, 2016] Nan Jiang and Lihong Li. Doubly ro-
605 bust off-policy value evaluation for reinforcement learn-
606 ing. In *International Conference on Machine Learning*,
607 pages 652–661, 2016. 608

[Jin *et al.*, 2022] Yue Jin, Yue Zhang, Tao Qin, Xudong
609 Zhang, Jian Yuan, Houqiang Li, and Tie-Yan Liu. Super-
610 vised off-policy ranking. In *International Conference on*
611 *Machine Learning*, pages 10323–10339, 2022. 612

[Jing *et al.*, 2014] Zhang Jing, Xindong Wu, and Victor S.
613 Sheng. Imbalanced multiple noisy labeling. *Transac-*
614 *tions on Knowledge and Data Engineering*, 27(2):489–
615 503, 2014. 616

[Jing *et al.*, 2015] Zhang Jing, Victor S. Sheng, Jian Wu, and
617 Xindong Wu. Multi-class ground truth inference in crowd-
618 sourcing with clustering. *Transactions on Knowledge and*
619 *Data Engineering*, 28(4):1080–1085, 2015. 620

[Kumar *et al.*, 2019] Aviral Kumar, Justin Fu, George
621 Tucker, and Sergey Levine. Stabilizing off-policy Q-
622 Learning via bootstrapping error reduction. In *Advances*
623 *in Neural Information Processing Systems*, pages 11761–
624 11771, 2019. 625

[Kumar *et al.*, 2020] Aviral Kumar, Aurick Zhou, George
626 Tucker, and Sergey Levine. Conservative Q-Learning for
627 offline reinforcement learning. In *Advances in Neural In-*
628 *formation Processing Systems*, pages 1179–1191, 2020. 629

[Li *et al.*, 2015] Lihong Li, Rémi Munos, and Csaba
630 Szepesvári. Toward minimax off-policy value estimation.
631 In *Artificial Intelligence and Statistics*, pages 608–616,
632 2015. 633

[Li *et al.*, 2022] Junbing Li, Changqing Zhang, Joey Tianyi
634 Zhou, Huazhu Fu, Shuyin Xia, and Qinghua Hu. Deep-
635 lift: Deep label-specific feature learning for image anno-
636 tation. *IEEE Transactions on Cybernetics*, 52(8):7732–
637 7741, 2022. 638

[Lili *et al.*, 2021] Chen Lili, Lu Kevin, Rajeswaran Aravind,
639 Lee Kimin, Grover Aditya, Laskin Michael, Abbeel Pieter,
640 Srinivas Aravind, and Mordatch Igor. Decision trans-
641 former: Reinforcement learning via sequence modeling.
642 In *Advances in Neural Information Processing Systems*,
643 pages 15084–15097, 2021. 644

[Liu *et al.*, 2018] Qiang Liu, Lihong Li, Ziyang Tang, and
645 Dengyong Zhou. Breaking the curse of horizon: Infinite-
646 horizon off-policy estimation. In *Advances in Neural In-*
647 *formation Processing Systems*, page 5361–5371, 2018. 648

[Liu, 2009] Tie-Yan Liu. Learning to rank for information
649 retrieval. *Information Retrieval*, 3(3):225–331, 2009. 650

[Mahmood *et al.*, 2014] A. Rupam Mahmood, Hado P van
651 Hasselt, and Richard S Sutton. Weighted importance sam-
652 pling for off-policy learning with linear function approx-
653 imation. In *Advances in Neural Information Processing*
654 *Systems*, 2014. 655

- 656 [Maystre and Grossglauser, 2017] Lucas Maystre and
657 Matthias Grossglauser. Just sort it! a simple and effective
658 approach to active preference learning. In *International*
659 *Conference on Machine Learning*, pages 2344–2353.
660 PMLR, 2017.
- 661 [Mengjiao *et al.*, 2022] Yang Mengjiao, Bo Dai, Ofir
662 Nachum, George Tucker, and Dale Schuurmans. Offline
663 policy selection under uncertainty. In *International*
664 *Conference on Artificial Intelligence and Statistics*, pages
665 4376–4396, 2022.
- 666 [Michael *et al.*, 2021] Janner Michael, Qiyang Li, and
667 Sergey Levine. Offline reinforcement learning as one big
668 sequence modeling problem. In *Advances in Neural Infor-*
669 *mation Processing Systems*, pages 1273–1286, 2021.
- 670 [Nachum *et al.*, 2019] Ofir Nachum, Yinlam Chow, Bo Dai,
671 and Lihong Li. DualDICE: Behavior-agnostic estimation
672 of discounted stationary distribution corrections. In *Ad-*
673 *vances in Neural Information Processing Systems*, 2019.
- 674 [Paduraru, 2013] Cosmin Paduraru. *Off-policy evaluation in*
675 *Markov decision processes*. PhD thesis, McGill University
676 Libraries, 2013.
- 677 [Rodrigo *et al.*, 2019] Nogueira Rodrigo, Wei Yang,
678 Kyunghyun Cho, and Jimmy Lin. Multi-stage document
679 ranking with bert. *arXiv preprint arXiv:1910.14424*,
680 2019.
- 681 [Rodrigues *et al.*, 2013] Filipe Rodrigues, Francisco Pereira,
682 and Bernardete Ribeiro. Learning from multiple annota-
683 tors: Distinguishing good from random labelers. *Pattern*
684 *Recognition Letters*, 34(12):1428–1436, 2013.
- 685 [S. and Zhang, 2019] Sheng Victor S. and Jing Zhang. Ma-
686 chine learning with crowdsourcing: A brief summary of
687 the past research and future directions. In *Proceedings*
688 *of the AAAI Conference on Artificial Intelligence*, pages
689 9837–9843, 2019.
- 690 [Shi *et al.*, 2021] Tianyu Shi, Dong Chen, Kaian Chen, and
691 Zhaojian Li. Offline reinforcement learning for au-
692 tonomous driving with safety and exploration enhance-
693 ment. *arXiv preprint arXiv:2110.07067*, 2021.
- 694 [Tang *et al.*, 2022] Hongyao Tang, Zhaopeng Meng, Jianye
695 Hao, Chen Chen, Daniel Graves, Dong Li, Changmin
696 Yu, Hangyu Mao, Wulong Liu, Yaodong Yang, Wenyuan
697 Tao, and Li Wang. What about inputting policy in value
698 function: Policy representation and policy-extended value
699 function approximator. In *Proceedings of the AAAI Con-*
700 *ference on Artificial Intelligence*, pages 8441–8449, 2022.
- 701 [Tian and Zhu, 2015] Tian Tian and Jun Zhu. Max-margin
702 majority voting for learning from crowds. In *Advances in*
703 *Neural Information Processing Systems*, page 1621–1629,
704 2015.
- 705 [Wang *et al.*, 2020] Ziyu Wang, Alexander Novikov, Konrad
706 Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E
707 Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre,
708 Nicolas Heess, and Nando de Freitas. Critic regularized
709 regression. In *Advances in Neural Information Processing*
710 *Systems*, pages 7768–7778, 2020.
- [Yang *et al.*, 2023] Chan-Yun Yang, Chamani Shiranthika, 711
Chung-Yih Wang, Kuo-Wei Chen, and Sagara Sumath- 712
ipala. Reinforcement learning strategies in cancer 713
chemotherapy treatments: A review. *Computer Methods* 714
and Programs in Biomedicine, 229:107280, 2023. 715
- [Yuchen *et al.*, 2016] Zhang Yuchen, Xi Chen, Dengyong 716
Zhou, and Michael I. Jordan. Spectral methods meet em: 717
A provably optimal algorithm for crowdsourcing. *JMLR*, 718
17(1):3537–3580, 2016. 719
- [Zhang and Jiang, 2021] Siyuan Zhang and Nan Jiang. To- 720
wards hyperparameter-free policy selection for offline re- 721
inforcement learning. In *Advances in Neural Information* 722
Processing Systems, 2021. 723
- [Zhang *et al.*, 2020] Zihao Zhang, Stefan Zohren, and 724
Roberts Stephen. Deep reinforcement learning for trading. 725
The Journal of Financial Data Science, 2020. 726
- [Zhu *et al.*, 2021] Chen Zhu, Wei Ping, Chaowei Xiao, Mo- 727
hammad Shoeybi, Tom Goldstein, Anima Anandkumar, 728
and Bryan Catanzaro. Long-short transformer: Efficient 729
transformers for language and vision. In *Advances in Neu-* 730
ral Information Processing Systems, pages 17723–17736, 731
2021. 732