

Increasing Entropy to Boost Policy Gradient Performance on Personalization Tasks

Andrew Starnes
Lirio AI Research, Lirio LLC.
Knoxville, U.S.A.
astarnes@lirio.com

Anton Dereventsov
Lirio AI Research, Lirio LLC.
Knoxville, U.S.A.
adereventsov@lirio.com

Clayton Webster
Lirio AI Research, Lirio LLC.
Knoxville, U.S.A.
cwebster@lirio.com

Abstract—In this effort, we consider the impact of regularization on the diversity of actions taken by policies generated from reinforcement learning agents trained using a policy gradient. Policy gradient agents are prone to entropy collapse, which means certain actions are seldomly, if ever, selected. We augment the optimization objective function for the policy with terms constructed from various φ -divergences and Maximum Mean Discrepancy which encourages current policies to follow different state visitation and/or action choice distribution than previously computed policies. We provide numerical experiments using MNIST, CIFAR10, and Spotify datasets. The results demonstrate the advantage of diversity-promoting policy regularization and that its use on gradient-based approaches have significantly improved performance on a variety of personalization tasks. Furthermore, numerical evidence is given to show that policy regularization increases performance without losing accuracy.

Index Terms—personalization, entropy, regularization, reinforcement learning, discrepancy, divergence

I. INTRODUCTION

Recommendation system models (see, e.g., [1], [2]) have become critically important in retaining customers of industries such as retail, e-commerce, media apps, or even healthcare. Corporations like Netflix, Spotify, and Amazon, use sophisticated collaborative filtering and content-based recommendation systems for video, song, and/or product recommendations [3], [4], [5], [6], [7]. For a recent overview of recommender systems in the healthcare domain see, e.g., [8] and the references therein.

Conventional personalization focuses on personal, transactional, demographic, and possibly health-related information, such as an individual’s age, residential location, employment, purchases, medical history, etc. Additional applications of personalization include: web content personalization and layout customization [9], [10]; customer-centric interaction with healthcare providers [11], [12], [13], [14], [15], [16]; personalized medical treatments [17], [18]. One of the major challenges associated with personalization techniques is the time required to adapt and update such approaches to changes in individual behaviors, reactions, and choices.

Recently, reinforcement learning (RL) has been increasingly exploited in personalized recommendation systems that continually interact with users (see, e.g., [19] and the references therein). As opposed to traditional recommendation techniques, RL is a more complex and transformative approach that considers behavioral and real-time data produced as the result

of user action. Examples of this technique include online browsing behavior, communication history, in-app choices, and other engagement data. This allows for more individualized experiences like adding personalized engaging sections to the body of an email or sending push notifications at a time when the customer is typically active, which results in more customized communication and thus, ultimately, greater conversion.

One of the major challenges associated with personalized RL agents is that standard optimization techniques often stall or even fail to converge when applied to such complex problems. This results in highly localized policies having lower entropy which directly translates into very few actions taken by the agent throughout the training process. Improving the diversity of actions taken by the policy is critical to improving the performance of the RL agent on a variety personalization tasks [20].

The traditional approach for combating low-entropy models is to regularize the standard objective with an entropy (penalty) term, such that the optimal policy additionally aims to maximize its entropy at each visited state, see, e.g., [21], [22], [23] and the references therein. This is achieved by subtracting a weighted term for the entropy of the model’s prediction from the loss function, thereby encouraging a more entropic model. This is equivalent to adding the Kullback-Leibler (KL) divergence between the policy and the uniform distribution.

Comparing probability distributions is a fundamental component of many supervised, unsupervised, and RL problems. In the machine learning community, the first discrepancies that were introduced to compare two probability distributions are φ -divergences [24], with φ is a convex, lower semi-continuous function such that $\varphi(1) = 0$. Such divergences can be viewed as a weighted average (by φ) of the odds-ratio between the two measures. In particular, we compute the following

$$D_{\varphi}(\alpha\|\beta) = \mathbb{E}_{\beta}\left(\varphi\left(\frac{\alpha}{\beta}\right)\right). \quad (1)$$

The computational simplicity of φ -divergences has made them very popular; with the most widely used being the KL divergence (see, e.g., Table I and the work [25, Section 2]).

However, such approaches suffer from the major drawback of not metrizing weak-convergence, which is instrumental for discrepancies on measure, as it ensures that the losses

TABLE I: Definitions of φ -divergences.

Kullback-Leibler	$D_{\text{KL}}(\alpha\ \beta) = \mathbb{E}_{\alpha} \left(\ln \left(\frac{\alpha(a)}{\beta(a)} \right) \right)$	$\varphi(x) = x \ln x$
Jensen-Shannon	$D_{\text{JS}}(\alpha\ \beta) = D_{\text{KL}}(\alpha\ \frac{1}{2}(\alpha+\beta)) + D_{\text{KL}}(\beta\ \frac{1}{2}(\alpha+\beta))$	$\varphi(x) = x \ln x - (1+x) \ln \left(\frac{1+x}{2} \right)$
Hellinger	$D_{H^2}(\alpha\ \beta) = \mathbb{E}_{\frac{1}{2}(\alpha+\beta)} \left((\sqrt{\alpha(a)} - \sqrt{\beta(a)})^2 \right)$	$\varphi(x) = (\sqrt{x} - 1)^2$
Total-Variation	$D_{\text{TV}}(\alpha\ \beta) = \sup_{A \in \mathcal{M}(\mathcal{A})} \alpha(A) - \beta(A) ^{\dagger}$	$\varphi(x) = \frac{1}{2} x - 1 $

$\dagger \mathcal{M}(\mathcal{A})$ denotes the α and β -measurable sets of \mathcal{A}

remain stable under small perturbations of the support of the measures. A class of discrepancies that satisfy this requirement are known as Maximum Mean Discrepancies (MMD) [26], which are a special case of integral probability metrics (IPM) [27]. Such approaches compare distributions without initially estimating their density functions. MMD is defined by the notion of representing distances between distributions as distances between *mean embeddings* of features, where the feature map is a kernel from a reproducing kernel Hilbert space (RKHS). This family of discrepancies presents the advantage of being efficiently computed from samples — both statistically since the estimates are robust with a small number of samples (reduced complexity) and also numerically as it can be computed in closed form.

In this work we augment the optimization objective function for the policy with various φ -divergence-based as well as MMD-based¹ term which encourages current policies to follow different state visitation and/or action choice distribution than previously computed policies. As such, by utilizing these more entropic variants of PG enables us to obtain a completely distinct set of policies.

Our main contributions are:

- formalization of φ -divergence-based as well as MMD-based regularization for personalized tasks in contextual bandit problems; and
- empirical demonstration of impact such regularization approaches have on RL.

A. Related work

The goal of this paper is to understand the impact that policy regularization has on an agent’s learning. It is often observed that policy gradient algorithms suffer from premature convergence to semi-deterministic, suboptimal policies. Avoiding this lack of diversity in actions is the motivation for adding entropy regularization to the REINFORCE algorithm [30], which is aptly named REINFORCE/MENT with MENT standing for Maximization of ENTropy. Using entropy regularization has also been found to improve agent performance (e.g., [31], [32]). While typical entropy regularization uses KL divergence between the policy and a uniform distribution over the actions, [33] uses KL divergence between the policy and the so-called default policy to improve performance. Bergmann divergence

¹MMD is the more popular IPM in machine learning applications, including, e.g., generative models (see [28], [29] and the references therein) due to the fact that it is applicable to a wide range of data types and distributions, computationally tractable even for high-dimensional data, and it is relatively robust to the curse of dimensionality [25].

is used in [34] to more safely train on-policy agents with off-policy data.

The work [35] presents diversity-driven approach for exploration, which can be easily combined with both off- and on-policy reinforcement learning algorithms. The authors show that by simply adding a distance measure regularization to the loss function, the proposed methodology significantly enhances an agent’s exploratory behavior. Similarly, the effort [36] presents an MMD-based approach for identifying a collection of near-optimal policies with significantly different distributions of trajectories.

Soft policy optimization was introduced in [22], [23]. These works show that the impact of entropy regularization goes beyond providing the agent with extra exploration, but also serves as more stable training process by avoiding a collapse onto a select set of actions. Similar to our work, empirical results also show the robustness of these approaches when compared with standard optimization algorithms.

II. BACKGROUND

We consider a contextual bandit environment [37], with a continuous state (context) space $\mathcal{S} \subset \mathbb{R}^m$, a discrete action space $\mathcal{A} = \{1, 2, \dots, n\}$ consisting of n available actions, and a reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. Using this conventional setting for recommendation and personalization tasks [38], [39], we define the reward \mathcal{J} of the policy π , which is given by the expectation return under the policy, i.e.,

$$\mathcal{J}(\pi) = \mathbb{E} \left[r(s, a) \mid s \sim \mathcal{S}, a \sim \pi(s) \right], \quad (2)$$

where $\pi(s)$ denotes the action probability distribution as state s . In this setting, traditional approaches for reinforcement learning aim to find a policy π that maximizes the reward function \mathcal{J} . However, to promote a more entropic model, we augment this optimization functional with various φ -divergence-based as well as an MMD-based regularization function \mathcal{R} . In other words, without loss of generality, our goal is to find a policy π that solves the following regularized optimization problem, namely:

$$\max_{\theta \in \mathbb{R}^d} \mathcal{J}(\pi_{\theta}) + \lambda \mathcal{R}(\pi_{\theta}), \quad (3)$$

where $\theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$ is a d -dimensional parameter that represents, e.g., the weights of a neural network, and $\lambda \in \mathbb{R}$ is a regularization (penalty) parameter. In what follows, we detail the construction of the regularized problem (3) and solutions will be sought for the policy gradient technique.

TABLE II: φ -divergences up to an additive constant in a finite action space and their gradients.

φ -Divergence		Definition	Gradient (wrt θ)
Kullback-Leibler	$D_{\text{KL}}(\pi_\theta(s) u)$	$\sum_{a=1}^n \pi_\theta(a s) \ln \pi_\theta(a s)$	$\sum_{a=1}^n (1 + \ln \pi_\theta(a s)) \nabla \pi_\theta(a s)$
Jensen-Shannon	$D_{\text{JS}}(\pi_\theta(s) u)$	$D_{\text{KL}}(\pi_\theta(s) u) - D_{\text{KL}}(\frac{1}{2}(\pi_\theta(s) + u) u)$	$\sum_{a=1}^n \ln \left(\frac{2\pi_\theta(a s)}{\pi_\theta(a s) + u(a)} \right) \nabla \pi_\theta(a s)$
Hellinger	$D_{H^2}(\pi_\theta(s) u)$	$\sum_{a=1}^n (\pi_\theta(a s) - 2\sqrt{u(a)\pi_\theta(a s)})$	$\sum_{a=1}^n \left(1 - \sqrt{\frac{u(a)}{\pi_\theta(a s)}} \right) \nabla \pi_\theta(a s)$
Total-Variation	$D_{\text{TV}}(\pi_\theta(s) u)$	$\max_{a=1, \dots, n} \pi_\theta(a s) - u(a) $	$\text{sgn}(\pi_\theta(a^* s) - u(a^*)) \nabla \pi_\theta(a^* s)^\ddagger$

$^\ddagger a^*$ is the action maximizing $|\pi_\theta(a|s) - u(a)|$

A. Relative entropy

The distribution of the agent’s policy π is often critical in practical applications as it directly translates to the actions the agent is taking throughout the training process. A conventional way to quantify the policy distribution is by computing its entropy $H(\pi)$ given by

$$H(\pi) = \mathbb{E} \left[\sum_{a \in \mathcal{A}} \pi(a|s) \log \pi(a|s) \mid s \sim \mathcal{S} \right]. \quad (4)$$

Entropy quantifies the amount of uncertainty involved in the value of a random variable or the outcome of a random process. In RL, entropy indicates how distributed the policy is, with more localized policies having lower entropy, which is known to lead to undesirable results, discussed in, e.g., [40].

B. Policy gradient methods

Policy gradient (PG) makes use of gradients to iteratively optimize a policy $\pi_\theta(s, a)$, parameterized by $\theta \in \mathbb{R}^d$. In order to maximize \mathcal{J} , we apply the Policy Gradient Theorem (see 13.2 of [41] for example), which shows

$$\begin{aligned} \nabla_\theta \mathcal{J}(\pi_\theta) &= \sum_{a \in \mathcal{A}} r(s, a) \pi_\theta(a|s) \nabla \ln \pi_\theta(a|s) \\ &= \mathbb{E}_{a \sim \pi_\theta(s)} (r(s, a) \nabla \ln \pi_\theta(a|s)). \end{aligned} \quad (5)$$

We use a Monte Carlo approximation of this expectation in order to estimate the gradient, denoted $\nabla_\theta \mathcal{J}_{\text{PG}}$ and given by (10).

III. DIVERSITY PROMOTING POLICY REGULARIZATION

In this section we develop all the necessary machinery to improve the existing policy gradient method by including an additional diversity-promoting term, thus, resulting in more entropic approaches.

A. φ -divergence regularization

The traditional approach for combating low-entropy models is to augment the standard objective with an entropy (penalty) term, such that the optimal policy additionally aims to maximize its entropy at each visited state [23]. This is achieved by adding a weighted term that measures the diversity of the model’s prediction from the loss function, thereby encouraging a more entropic model. One way to accomplish this is by adding any of the φ -divergences in Table I calculated between the policy

π_θ and the uniform distribution. Table II provides simplified definitions of the φ -divergences in the case of a finite action space, where we use

$$u \sim \text{Unif}(1, 2, \dots, n) \quad (6)$$

as the uniform distribution here but can be replaced by any other distribution of the actions. Therefore, using (3) and given a regularization constant $\lambda \in \mathbb{R}$, our goal is to find a policy π that maximizes a new objective function, namely:

$$\max_{\theta \in \mathbb{R}^d} \mathcal{J}(\pi_\theta) + \lambda D_\varphi(\pi_\theta || u). \quad (7)$$

B. MMD regularization

In addition to φ -divergence regularization, defined in Table II, we also propose to exploit a family of discrepancies known as Maximum Mean Discrepancy (MMD). Given a RKHS \mathcal{H} with kernel k , MMD between two probability measures α and β is given by

$$\begin{aligned} \text{MMD}_k^2(\alpha, \beta) &:= \left(\sup_{\{f: \|f\|_{\mathcal{H}} \leq 1\}} |\mathbb{E}_\alpha f(x) - \mathbb{E}_\beta f(y)| \right)^2 \\ &= \mathbb{E}_{\alpha \otimes \alpha} k(x, x') + \mathbb{E}_{\beta \otimes \beta} k(y, y') - 2\mathbb{E}_{\alpha \otimes \beta} k(x, y). \end{aligned} \quad (8)$$

This family of discrepancies presents the advantage of being efficiently estimated from samples of the measures — both statistically since the estimates are robust with a small number of samples (reduced complexity) and also numerically, as (8) can be computed in closed form. Therefore, using (3) and given a regularization constant $\lambda \in \mathbb{R}$, our goal is to find a policy π that maximizes a new objective function, namely:

$$\max_{\theta \in \mathbb{R}^d} \mathcal{J}(\pi_\theta) + \lambda \text{MMD}_k^2(\pi_\theta, u). \quad (9)$$

There are many choices for k (or equivalently \mathcal{H}). For our examples in this paper, we use the Gaussian kernel, $k(x, y) = \exp(-\|x - y\|^2)$, where $\|x - y\| = \mathbb{1}(x = y)$ for $x, y \in \mathcal{A}$. We do this because each of the examples focuses on correctly labeling and the arithmetic difference between two labels is not meaningful.

C. Diversity-promoting policy gradients

We will use $\theta \in \mathbb{R}^d$ to denote the parameters of a neural network that takes as input $s \in \mathcal{S}$ and outputs a probability distribution over \mathcal{A} with the policy output mapping $\mathcal{Z} : \mathcal{S} \rightarrow$

\mathbb{R}^n . For model parameters θ , the action selection distribution as a particular state, $s \in \mathcal{S}$, is denoted by $\pi_\theta(s)$.

The standard gradient loss estimate for policy gradient is given by

$$\nabla \mathcal{J}_{\text{PG}}(\pi_\theta(s)) = r(s, a) \nabla \pi_\theta(a|s). \quad (10)$$

The gradients of each of the φ -divergences can be found in Table II. Lastly, the gradient of MMD in the contextual bandit setting is

$$\begin{aligned} \nabla_\theta \text{MMD}_k^2(\pi_\theta(s), u) &= 2\mathbb{E}_{\pi_\theta \otimes \pi_\theta \otimes u} \left((k(a, a') - k(a, a^*)) \nabla_\theta \ln \pi_\theta(a|s) \right) \\ &= \sum_{a \in \mathcal{A}} c_{\theta, s, a}(a', a^*) \nabla_\theta \pi_\theta(a|s), \end{aligned} \quad (11)$$

where

$$c_{\theta, s, a}(a', a^*) = \sum_{a', a^*} (k(a, a') - k(a, a^*)) u(a^*) \pi_\theta(a'|s).$$

The gradient update from $\nabla_\theta \mathcal{J}_{\text{PG}}$ only depends on the gradient based on the action that was selected. Three of the φ -divergences, KL, Jensen-Shannon, and Hellinger, as well as MMD have gradients that are weighted sums of the gradients over all of the actions, not just the selected action. On the other hand, for Total-Variation the gradient only depends on the action whose likelihood is furthest away from the policy u , given by (6).

When π_θ is found using the softmax function, we can further expand all of the above gradients by

$$\nabla \pi_\theta(a|s) = \pi_\theta(a|s) [\mathbb{1}(a = a') - \pi(a'|s)]_{a'=1}^n \times \nabla \mathcal{Z}(s).$$

Using the gradient information given by (10), optimal solutions to the φ -divergence-based diversity-promoting objective, given by (7), as well optimal solutions to the MMD-based diversity-promoting objective, given by (9), can then be solved with standard gradient-based methods.

IV. NUMERICAL EXAMPLES

In this section we conduct numerical experiments comparing performance of the RL agents with policy regularization methods described in Section III. Specifically, we consider the following agents:

- 1) pg: the default policy gradient agent without any regularization to act as a baseline;
- 2) pg_ent: pg-agent with entropy regularization;
- 3) pg_mmd: pg-agent with MMD regularization;
- 4) pg_js: pg-agent with Jensen-Shannon regularization;
- 5) pg_hl: pg-agent with Hellinger regularization; and
- 6) pg_tv: pg-agent with total variation regularization.

We chose these algorithms so that we could easily identify the impact that the regularizers have in the absence of additional constraints imposed by other algorithms such as TRPO [42] or PPO [43].

The agents and regularizer losses are manually implemented in TensorFlow and the network training is performed with Adam optimizer with the default hyperparameters. For all of

our algorithm configurations, we use a batch size of 100. An agent policy is parameterized by a 2-layer feed-forward neural network with 32 nodes on each layer. For each regularized agent we perform a hyperparameter search to determine the appropriate value of the regularization coefficient.

We deploy the agents on various personalization tasks that are given by contextual bandit environments. For each agent and environment we report the following metrics, computed over the test set: the agent reward, the policy entropy, and the action selection histogram. For the simplicity of presentation, the histograms are sorted to emphasize the agent’s action distribution over the test set.

The presented examples are performed using Python3.8 with Tensorflow 2.12 on personal laptops. The source code reproducing the given experiments is available at <https://github.com/acstarnes/wain23-policy-regularization>.

A. MNIST Environment

We use MNIST dataset² to create a contextual bandit environment, as done in [44], [45], [46], [47]. In this formulation the images act as observations and the labels act as the actions that agents can take. The reward for selecting the correct label is 1 and $-1/9$ for an incorrect classification. Defining the reward function this way means that the expected return for the uniformly random policy is 0 and for the optimal policy is 1. The agent reward, policy entropy, and action selection histograms for the various regularizers on MNIST environment are shown in Figure 1.

We note that all regularized agents solve the environment and demonstrate a comparable performance while outperforming the baseline agent. From the action selection histogram we observe that the unregularized policy gradient agent only selects 7 out of 10 actions, which results in the agent reward value plateauing at about 0.6. In contrast, all of the regularized agents maintain diverse action selection throughout the training process and achieve reward values close to 1, which indicates fully learning the environment.

B. CIFAR10 Environment

As in the previous example, we use CIFAR10 dataset³ to create a contextual bandit environment. The agent reward, policy entropy, and action selection histograms for the various regularizers on CIFAR10 environment are shown in Figure 2.

Unlike the previous example, the agents fail to fully solve the environment in this case due to the increased complexity of CIFAR10 dataset. In fact, even the most successful agents achieve reward values of only about 0.3, which roughly equates to a 35% classification accuracy on CIFAR10 dataset. Such poor performance is due to the constrained network architecture and the contextual bandit formulation of the problem. Nonetheless, the advantage of regularized agents is evident, both from policy reward and entropy perspectives. In particular, we observe that most of the regularized agents are able to maintain diverse

²<http://yann.lecun.com/exdb/mnist/>

³<https://www.cs.toronto.edu/~kriz/cifar.html>

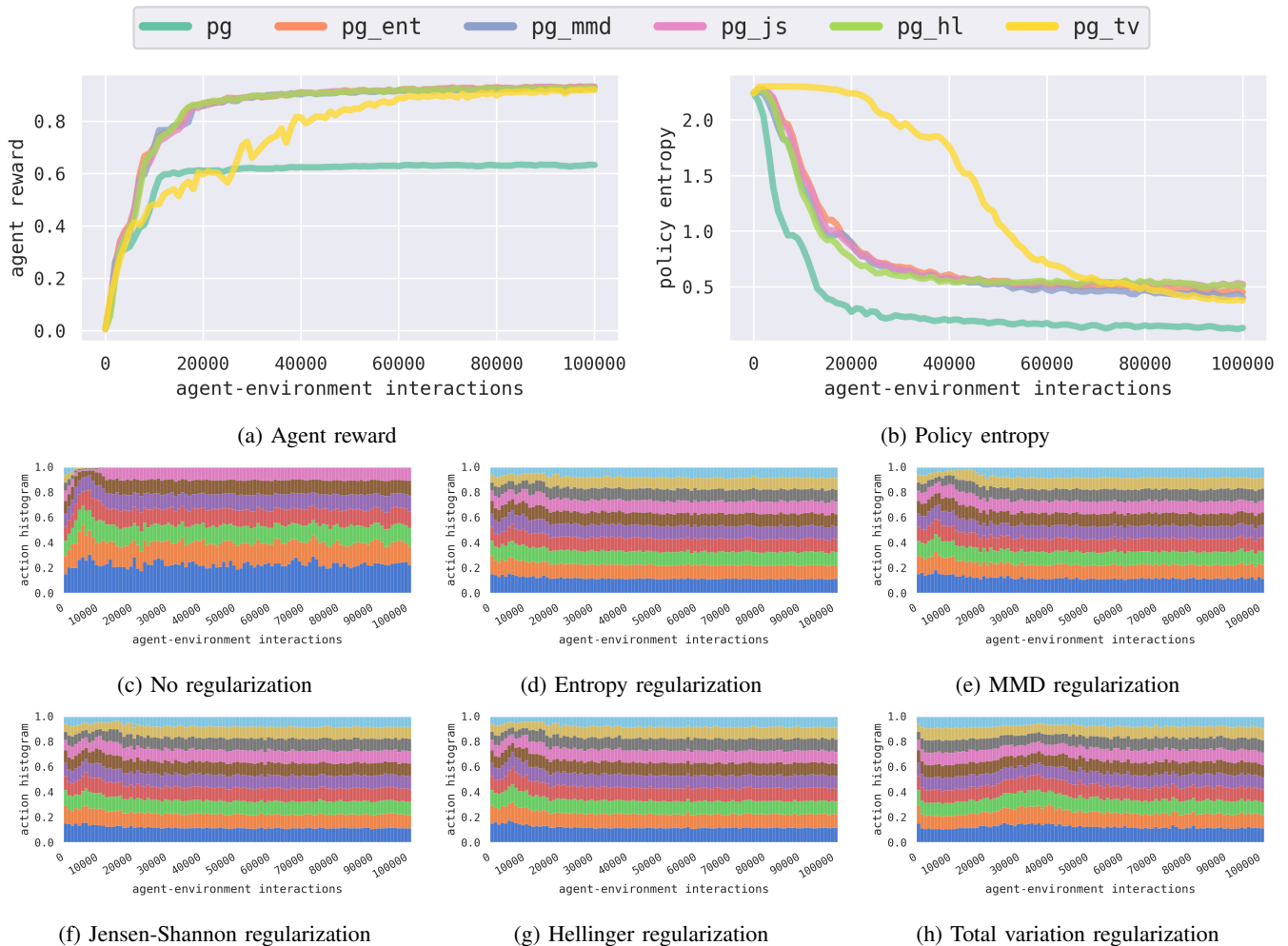


Fig. 1: Results of image classification experiment on MNIST environment.

(albeit unbalanced) action selection, while the baseline agent only selects 5 out of 10 available actions.

C. Spotify Environment

In this experiment we set up a synthetic music recommendation system proposed in [47]. We use Spotify Web API⁴ to construct a contextual bandit environment that replicates the task of track recommendation to a user. In this setting the observations (users) are synthetically generated and are represented by their preferences to various musical genres, and the actions are given by the set of tracks the agent can recommend. The reward for recommending a track to a user is either 1, -1, or 0, indicating that the user liked/disliked/did not provide feedback, respectively. See [47] for a more detailed explanation of the environment. The agent reward, policy entropy, and action selection histograms for the various regularizers on Spotify environment are shown in Figure 3.

We note that while all agents achieve satisfactory performance in terms of reward, the action selection of the baseline

agent is constrained to only 3 tracks (out of 50 possible), which is neither practical nor acceptable in real-world applications. All regularized agents provide a much more diverse action selection, while also achieving higher reward values.

A particular interest of this environment is the fact that there are infinitely many policies that achieve near-optimal performance. As an example, even the unregularized policy gradient agent almost learned the environment, while only ever taking about 6% of the available actions, with one action being taken about 40% of the time. In comparison, for regularized agents the action selection is much more diverse with fewer “favorite” actions. Most notably, the MMD-regularized agent is actively taking about 40% of the actions with the most frequent one being selected only about 8% of the time.

V. CONCLUDING REMARKS

In this effort, we consider the impact of regularization on the diversity of actions taken by policies generated from policy gradient RL agents. In the context of personalized RL there are several additional advantages that extend from this work. First, the φ -divergence and MMD-based regularization encourages

⁴<https://developer.spotify.com/documentation/web-api/>

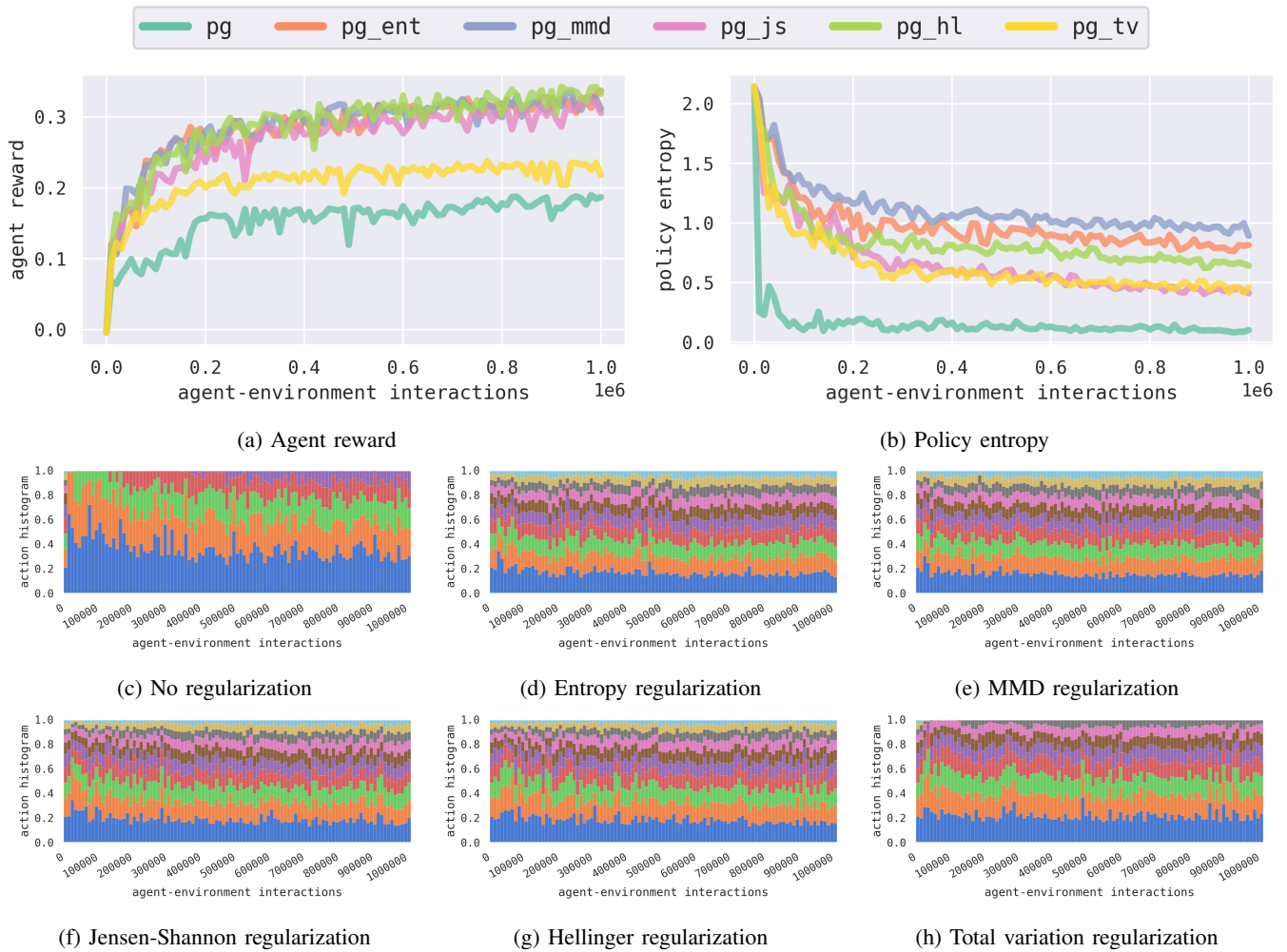


Fig. 2: Results of image classification experiment on CIFAR10 environment.

exploration and aids to prevent early convergence to sub-optimal policies. Second, the resulting policies can serve as a good (macro) initialization for a more (micro) specific behavior. Finally, the resulting policies are more robust in the face of adversarial perturbations or noise as evidenced by our various numerical examples.

However, there is much more extensive testing to be done and a supporting theory needs to be developed before any victories can be declared. As mentioned throughout, there has been extensive amounts of research by the RL community on using KL-type entropy regularization, but more advanced discrepancies such as the MMD-based approach we presented here are still in their infancy. In addition, there is a vast amount of research on optimal transport theory which, in connection with entropy-type penalization is something we also plan to investigate. These methods possess some computational challenges but have the ability to lift a ground metric from the data-space to the set of probability measures on this space and, therefore, take into account the underlying geometry of the data [48], [49], [26]. To our knowledge, this area of research has yet to be explored by the machine learning community.

REFERENCES

- [1] F. Ricci, L. Rokach, and B. Shapira, *Recommender Systems: Introduction and Challenges*. Boston, MA: Springer US, 2015, pp. 1–34. [Online]. Available: https://doi.org/10.1007/978-1-4899-7637-6_1
- [2] R. Burke, A. Felfernig, and M. H. Göker, “Recommender systems: An overview,” *AI Magazine*, vol. 32, no. 3, pp. 13–18, Jun. 2011. [Online]. Available: <https://ojs.aaai.org/index.php/aimagazine/article/view/2361>
- [3] C. A. Gomez-Uribe and N. Hunt, “The netflix recommender system: Algorithms, business value, and innovation,” *ACM Transactions on Management Information Systems (TMIS)*, vol. 6, no. 4, 2016. [Online]. Available: <https://doi.org/10.1145/2843948>
- [4] X. Amatriain and J. Basilico, *Recommender Systems in Industry: A Netflix Case Study*. Boston, MA: Springer US, 2015, pp. 385–419. [Online]. Available: https://doi.org/10.1007/978-1-4899-7637-6_11
- [5] K. Jacobson, V. Murali, E. Newett, B. Whitman, and R. Yon, “Music personalization at spotify,” in *Proceedings of the 10th ACM Conference on Recommender Systems*, ser. RecSys ’16. New York, NY, USA: Association for Computing Machinery, 2016, p. 373. [Online]. Available: <https://doi.org/10.1145/2959100.2959120>
- [6] B. Smith and G. Linden, “Two decades of recommender systems at amazon.com,” *IEEE Internet Computing*, vol. 21, no. 3, pp. 12–18, 2017.
- [7] X. Wang, Y. Wang, D. Hsu, and Y. Wang, “Exploration in interactive personalized music recommendation: A reinforcement learning approach,” *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 11, no. 1, Sep. 2014. [Online]. Available: <https://doi.org/10.1145/2623372>
- [8] T. N. T. Tran, A. Felfernig, C. Trattner, and A. Holzinger, “Recommender systems in the healthcare domain: state-of-the-art and research issues,”

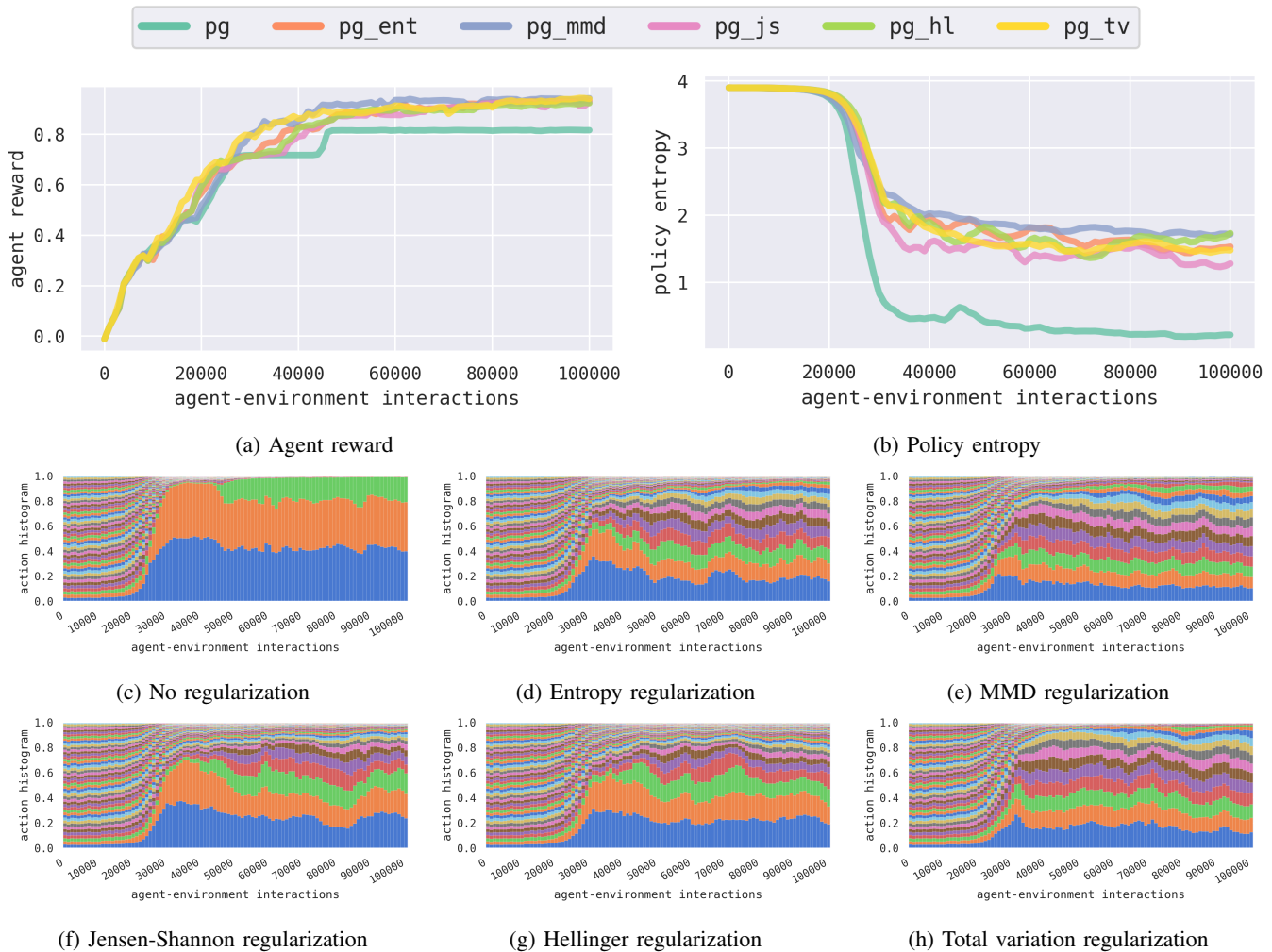


Fig. 3: Results of image classification experiment on Spotify environment.

- Journal of Intelligent Information Systems*, vol. 57, no. 1, pp. 171–201, 2021.
- [9] S. Ferretti, S. Mirri, C. Prandi, and P. Salomoni, “Automatic web content personalization through reinforcement learning,” *Journal of Systems and Software*, vol. 121, pp. 157–169, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121216000443>
- [10] F. Ricci, L. Rokach, and B. Shapira, “Introduction to recommender systems handbook,” in *Recommender Systems Handbook*, 2011.
- [11] L. Lasalvia, “Personalization and standardization: Can we have it all?” *Journal of Precision Medicine— Volume*, vol. 6, no. 1, 2020.
- [12] A. Vatian, S. Dudorov, A. Ivchenko, K. Smirnov, E. Chikhshova, A. Lobantsev, V. Parfenov, A. Shalyto, and N. Gusarova, “Design patterns for personalization of healthcare process,” in *Proceedings of the 2019 2nd International Conference on Geoinformatics and Data Analysis*, ser. ICGDA 2019. New York, NY, USA: Association for Computing Machinery, 2019, p. 83–88. [Online]. Available: <https://doi.org/10.1145/3318236.3318249>
- [13] H. Lei, A. Tewari, and S. A. Murphy, “An actor-critic contextual bandit algorithm for personalized mobile health interventions,” *arXiv preprint arXiv:1706.09090*, 2017.
- [14] F. Zhu, J. Guo, R. Li, and J. Huang, “Robust actor-critic contextual bandit for mobile health (mhealth) interventions,” in *Proceedings of the 2018 acm international conference on bioinformatics, computational biology, and health informatics*, 2018, pp. 492–501.
- [15] A. e. Hassouni, M. Hoogendoorn, M. v. Otterlo, and E. Barbaro, “Personalization of health interventions using cluster-based reinforcement learning,” in *International Conference on Principles and Practice of Multi-Agent Systems*. Springer, 2018, pp. 467–475.
- [16] C. Tan, R. Han, R. Ye, and K. Chen, “Adaptive learning recommendation strategy based on deep q-learning,” *Applied psychological measurement*, vol. 44, no. 4, pp. 251–266, 2020.
- [17] M. Aspinall and R. Hamermesh, “Realizing the promise of personalized medicine,” *Harvard business review*, vol. 85, pp. 108–117, 165, 11 2007.
- [18] G. S. Ginsburg and J. J. McCarthy, “Personalized medicine: revolutionizing drug discovery and patient care,” *Trends in Biotechnology*, vol. 19, no. 12, pp. 491–496, 2001. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167779901018145>
- [19] F. den Hengst, E. Grua, A. el Hassouni, and M. Hoogendoorn, “Reinforcement learning for personalization: A systematic literature review,” *Data Science*, vol. 3, pp. 1–41, 04 2020.
- [20] A. Dereventsov, C. G. Webster, and J. Daws, “An adaptive stochastic gradient-free approach for high-dimensional blackbox optimization,” Montreal, Canada, 2022.
- [21] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, “Soft actor-critic algorithms and applications,” *CoRR*, vol. abs/1812.05905, 2018. [Online]. Available: <http://arxiv.org/abs/1812.05905>
- [22] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, 10–15 Jul 2018, pp. 1861–1870. [Online]. Available: <https://proceedings.mlr.press/v80/haarnoja18b.html>

- [23] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement learning with deep energy-based policies," in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. PMLR, 06–11 Aug 2017, pp. 1352–1361. [Online]. Available: <https://proceedings.mlr.press/v70/haarnoja17a.html>
- [24] I. Csiszar, "I-Divergence Geometry of Probability Distributions and Minimization Problems," *The Annals of Probability*, vol. 3, no. 1, pp. 146–158, 1975. [Online]. Available: <https://doi.org/10.1214/aop/1176996454>
- [25] A. Genevay, "Entropy-Regularized Optimal Transport for Machine Learning," Theses, PSL University, Mar. 2019. [Online]. Available: <https://theses.hal.science/tel-02319318>
- [26] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola, "A kernel method for the two-sample-problem," in *Advances in Neural Information Processing Systems*, B. Schölkopf, J. Platt, and T. Hoffman, Eds., vol. 19. MIT Press, 2006. [Online]. Available: <https://proceedings.neurips.cc/paper/2006/file/e9fb2eda3d9c55a0d89c98d6c54b5b3e-Paper.pdf>
- [27] A. Müller, "Integral probability metrics and their generating classes of functions," *Advances in Applied Probability*, vol. 29, no. 2, p. 429–443, 1997.
- [28] C.-L. Li, W.-C. Chang, Y. Cheng, Y. Yang, and B. Póczos, "Mmd gan: Towards deeper understanding of moment matching network," in *NIPS*, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:4685015>
- [29] M. Binkowski, D. J. Sutherland, M. Arbel, and A. Gretton, "Demystifying mmd gans," *ArXiv*, vol. abs/1801.01401, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:3531856>
- [30] R. J. Williams and J. Peng, "Function optimization using connectionist reinforcement learning algorithms," *Connection Science*, vol. 3, no. 3, pp. 241–268, 1991.
- [31] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*. PMLR, 2016, pp. 1928–1937.
- [32] Z. Ahmed, N. Le Roux, M. Norouzi, and D. Schuurmans, "Understanding the impact of entropy on policy optimization," in *International conference on machine learning*. PMLR, 2019, pp. 151–160.
- [33] A. Galashov, S. M. Jayakumar, L. Hasenclever, D. Tirumala, J. Schwarz, G. Desjardins, W. M. Czarnecki, Y. W. Teh, R. Pascanu, and N. Heess, "Information asymmetry in kl-regularized rl," *arXiv preprint arXiv:1905.01240*, 2019.
- [34] Q. Wang, Y. Li, J. Xiong, and T. Zhang, "Divergence-augmented policy optimization," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [35] Z.-W. Hong, T.-Y. Shann, S.-Y. Su, Y.-H. Chang, T.-J. Fu, and C.-Y. Lee, "Diversity-driven exploration strategy for deep reinforcement learning," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/a2802cade04644083dcde1c8c483ed9a-Paper.pdf>
- [36] M. A. Masood and F. Doshi-Velez, "Diversity-inducing policy gradient: Using maximum mean discrepancy to find a set of diverse policies," in *IJCAI*, 2019, pp. 5923–5929. [Online]. Available: <https://doi.org/10.24963/ijcai.2019/821>
- [37] J. Langford and T. Zhang, "The epoch-greedy algorithm for multi-armed bandits with side information," *Advances in neural information processing systems*, vol. 20, 2007.
- [38] L. Li, W. Chu, J. Langford, and R. E. Schapire, "A contextual-bandit approach to personalized news article recommendation," in *Proceedings of the 19th international conference on World wide web*, 2010, pp. 661–670.
- [39] L. Tang, Y. Jiang, L. Li, C. Zeng, and T. Li, "Personalized recommendation via parameter-free contextual bandits," in *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, 2015, pp. 323–332.
- [40] Z. Dou, R. Song, J.-R. Wen, and X. Yuan, "Evaluating the effectiveness of personalized web search," *IEEE Transactions on Knowledge and Data Engineering*, vol. 21, no. 8, pp. 1178–1190, 2008.
- [41] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [42] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.
- [43] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [44] M. Dudík, J. Langford, and L. Li, "Doubly robust policy evaluation and learning," in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, 2011, pp. 1097–1104.
- [45] A. Swaminathan and T. Joachims, "Counterfactual risk minimization: Learning from logged bandit feedback," in *International Conference on Machine Learning*. PMLR, 2015, pp. 814–823.
- [46] M. Chen, R. Gummadi, C. Harris, and D. Schuurmans, "Surrogate objectives for batch policy optimization in one-step decision making," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [47] A. Dereventsov, A. Starnes, and C. G. Webster, "Examining policy entropy of reinforcement learning agents for personalization tasks," *arXiv*, 2022, submitted.
- [48] L. Kantorovich, "On the transfer of masses (in russian)," *Doklady Akademii Nauk*, vol. 2, pp. 227–229, 1942.
- [49] —, "On the translocation of masses," *Journal of Mathematical Sciences*, vol. 133, pp. 1381–1382, 2006.