

---

# Score-Models for Offline Goal-Conditioned Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Motivation

2 Despite recent progress in developing methods for goal-reaching in the online setting (where  
3 environment interactions are allowed), a number of these methods are either suboptimal in the  
4 offline setting or suffer from learning difficulties. Prior GCRL algorithms can largely be classified  
5 into one of three categories: iterated behavior cloning, RL with sparse rewards, and contrastive  
6 learning. Iterated behavior cloning or goal-conditioned supervised learning approaches [16, 38] have  
7 been shown to be provably suboptimal [9] for GCRL. Modifying single-task RL methods [33, 18]  
8 for GCRL with 0-1 reward implies learning a  $Q$ -function that predicts the discounted probability  
9 of goal reaching, which makes it essentially a density model. Modeling density directly is a hard  
10 problem, an insight which has prompted the development of methods [8] that learn density-ratio  
11 instead of densities, as classification is an easier problem than density estimation. Contrastive RL  
12 approaches to GCRL [8, 10, 40] aim to do precisely this and are the main methods to enjoy success  
13 for applying GCRL in high-dimensional observation spaces. However, when dealing with offline  
14 datasets, contrastive RL approaches [10, 40] are suboptimal, as they learn a policy that is a greedy  
15 improvement over the  $Q$ -function of the data generation policy. A prior GCRL work [21] explores  
16 the insight of occupancy matching for GCRL which requires learning a discriminator. Unfortunately,  
17 errors in learned discriminators can compound and adversely affect the learned policy’s performance,  
18 especially in the offline setting where these errors cannot be corrected with further interaction with  
19 the environment. This begs the question: *How can we derive a performant GCRL method that learns  
20 optimal policies from offline datasets of suboptimal quality?*

21 Going beyond the shortcomings of the previous methods, our proposed method combines the insight  
22 of formulating GCRL as an occupancy matching problem along with an efficient, discriminator-free  
23 dual formulation that learns from offline data. The resulting algorithm SMORe forgoes learning  
24 density functions or classifiers, but instead learns unnormalized densities or *scores* that allow it to  
25 produce optimal goal-reaching policies. The scores are learned via a Bellman-regularized contrastive  
26 procedure that makes our method a desirable candidate for GCRL with high-dimensional observations,  
27 avoiding the need for density modeling.

## 2 SMORe: Score Models for Offline GCRL

29 Define a training distribution over goals  $q^{train}(g)$  and *goal-transition distribution*  $q(s, a, g)$  in  
30 a stochastic MDP as  $q(s, a, g) \propto q^{train}(g) \mathbb{E}_{s' \sim p(\cdot|s,a)} [\mathbb{I}_{\phi(s')=g}]$ . Intuitively, the distribution has  
31 probability mass on each transition that leads to a goal. Let  $\rho$  be the offline data distribution and  $d_g^\pi$   
32 denote the visitation distribution induced by goal-conditioned policy  $\pi_g$  when the goals are sampled  
33 from  $q^{train}(g)$ . To leverage offline data to learn performant goal-reaching policies, we consider a  
34 surrogate objective to the occupancy matching learning problem by matching *mixture* distributions:

$$\min_{\pi_g} \mathcal{D}_f(\text{Mix}_\beta(d^{\pi_g}, \rho)(s, a, g) \| \text{Mix}_\beta(q, \rho)(s, a, g)), \quad (1)$$

35 where for any two distributions  $\mu_1$  and  $\mu_2$ ,  $\text{Mix}_\beta(\mu_1, \mu_2)$  denotes the mixture distribution with  
36 coefficient  $\beta \in (0, 1]$  defined as  $\text{Mix}_\beta(\mu_1, \mu_2) = \beta\mu_1 + (1 - \beta)\mu_2$ . Proposition 2.1 (in appendix)  
37 shows the matching mixture distribution provably maximizes a lower bound to the Lagrangian

Submitted to 37th Conference on Neural Information Processing Systems (NeurIPS 2023). Do not distribute.

### SMORe : Mixture occupancy matching with *dual* objective

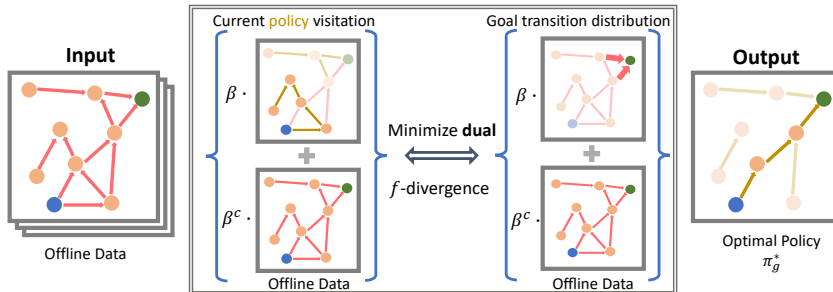


Figure 1: Illustration of the SMORe objective where  $\beta^c = 1 - \beta$ : SMORe matches a mixture distribution of current policy and offline data to a mixture of the goal-transition distribution and offline data in order to find the optimal goal reaching policy.

38 relaxation of the max-entropy GCRL objective subject to the constraint that the policy visitation is  
 39 close to the offline data visitation. Theorem 2 presents our core method SMORe, showing that we can  
 40 leverage tools from convex duality to obtain an unconstrained dual problem that does not require  
 41 computing  $d^{\pi_g}(s, a, g)$  or sampling from it, while effectively leveraging offline data.

42 **Theorem 1.** (SMORe) The dual problem to the primal occupancy matching objective (Equation 9) is  
 43 given by:

$$\max_{\pi_g} \min_S \beta(1 - \gamma) \mathbb{E}_{d_0, \pi_g} [S(s, a, g)] + \mathbb{E}_{\mathbb{H}_{\text{IxB}}(q, \rho)} [f^*(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \quad (2)$$

$$- (1 - \beta) \mathbb{E}_{\rho} [\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)],$$

44 where  $f^*$  is conjugate function of  $f$ ,  $S$  is the Lagrange dual variable defined as  $S : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow \mathbb{R}$ ,  
 45  $d_0$  is the initial state distribution and  $P^{\pi_g}$  the transition operator induced by the policy  $\pi_g$  defined  
 46 as  $P^{\pi_g} S(s, a, g) := \mathbb{E}_{s' \sim p(\cdot | s, a), a' \sim \pi_g(\cdot | s', g)} [S(s', a', g)]$ . Moreover, as strong duality holds from  
 47 Slater’s conditions the primal and dual share the same optimal solution  $\pi_g^*$  for any offline data  
 48 distribution  $\rho$ .

49 **Our contribution is a novel method for GCRL that is discriminator-free, applicable for a**  
 50 **number of  $f$ -divergences, and robust to low coverage of goals in the offline dataset.**

### 51 3 Overview of Empirical Results

Task	Occupancy Matching		Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	SMORe	GoFAR	WGCSL	GCSL		AM	IQL
Reacher (★)	28.40±0.88	19.74±1.35	17.57±0.53	15.87±1.31	16.44±0.60	23.26 ±0.14	11.70 ±1.97
SawyerReach (★)	37.67±0.12	15.34±0.64	15.15±0.44	14.25±0.7	22.32 ±0.34	23.34±0.17	35.18 ±0.29
SawyerDoor (★)	31.48±0.46	18.94±0.01	20.01±1.55	20.88±0.22	12.96±5.19	22.12 ±0.13	25.52 ±1.45
FetchReach (★)	35.08±0.54	28.2 ±0.61	21.9±2.13	20.91 ±2.78	30.07±0.07	30.1 ±0.32	34.43 ±1.00
FetchPick (★)	26.47 ±1.34	19.7 ±2.57	9.84 ±2.58	7.58±1.85	0.42±0.29	8.94 ±3.09	16.8 ±3.10
FetchPush (★)	26.83±1.21	18.2 ±3.00	14.7 ±2.65	13.4 ±3.02	2.40 ±1.28	14.0 ±2.81	22.40 ±0.74
FetchSlide	4.99±0.40	2.47 ±1.44	2.73 ±1.64	1.75 ±1.3	0.0±0.0	1.46 ±1.38	4.80 ±1.59
HandReach (★)	18.68 ±3.35	11.5 ±5.26	5.97 ±4.81	1.37 ±2.21	0.0±0.0	0.0 ±0.0	1.44 ±1.77
CheetahTgtVel-m-e (★)	136.71 ±10.59	0.0 ±0.0	0.0 ±0.0	95.98 ±15.72	0.0±0.0	0.0 ±0.0	100.38±1.22
CheetahTgtVel-r-e (★)	60.01 ±39.40	0.0 ±0.0	0.0 ±0.0	11.56 ±13.47	0.0±0.0	0.0 ±0.0	0.0±0.0
AntTgtVel-m-e	154.95±19.44	168.27±9.58	0.0 ±0.0	164.54±7.69	0.0±0.0	0.0 ±0.0	148.17 ±5.43
AntTgtVel-r-e (★)	126.22±14.40	74.36±15.97	0.0 ±0.0	104.95±6.00	0.0±0.0	0.0 ±0.0	3.06 ±2.64

Table 1: Discounted Return for the offline GCRL benchmark. Results are averaged over 10 seeds. ‘m-e’ and ‘r-e’ stands for medium-expert mixture and random-expert mixture respectively.

52 Our experiments in Table 2 show across a broad range of offline datasets and environments that  
 53 SMORe outperforms prior offline GCRL baselines. A key property of SMORe is that it learns scores  
 54 through a contrastive procedure, making it a particularly appealing choice for GCRL with large  
 55 observation spaces. Our experiments on image-observation domains in Figure 4 also demonstrate  
 56 that SMORe outperforms baselines that are designed specifically for image-based GCRL. Finally, we  
 57 show in Table 3 that the discriminator-free nature of SMORe allows to be more robust to decreasing  
 58 coverage of goal-reaching policy in the offline dataset.

## 59 References

- 60 [1] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun. Reinforcement learning: Theory and  
61 algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep*, pages 10–4, 2019. 18
- 62 [2] F. Al-Hafez, D. Tateo, O. Arenz, G. Zhao, and J. Peters. Ls-iq: Implicit reward regularization  
63 for inverse reinforcement learning. *arXiv preprint arXiv:2303.00599*, 2023. 9
- 64 [3] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,  
65 O. Pieter Abbeel, and W. Zaremba. Hindsight experience replay. *Advances in neural information  
66 processing systems*, 30, 2017. 9, 11, 13
- 67 [4] Y. Chebotar, K. Hausman, Y. Lu, T. Xiao, D. Kalashnikov, J. Varley, A. Irpan, B. Eysenbach,  
68 R. Julian, C. Finn, et al. Actionable models: Unsupervised offline reinforcement learning of  
69 robotic skills. *arXiv preprint arXiv:2104.07749*, 2021. 11, 22
- 70 [5] L. Chen, R. Paleja, and M. Gombolay. Learning from suboptimal demonstration via  
71 self-supervised reward regression. *arXiv preprint arXiv:2010.11723*, 2020. 13
- 72 [6] Y. Ding, C. Florensa, P. Abbeel, and M. Phielipp. Goal-conditioned imitation learning. *Advances  
73 in neural information processing systems*, 32, 2019. 13
- 74 [7] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and  
75 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets.  
76 *arXiv preprint arXiv:2109.13396*, 2021. 6
- 77 [8] B. Eysenbach, R. Salakhutdinov, and S. Levine. C-learning: Learning to achieve goals via  
78 recursive classification. *arXiv preprint arXiv:2011.08909*, 2020. 1, 6, 11, 12, 13, 22
- 79 [9] B. Eysenbach, S. Udatha, R. R. Salakhutdinov, and S. Levine. Imitating past successes can be  
80 very suboptimal. *Advances in Neural Information Processing Systems*, 35:6047–6059, 2022. 1,  
81 6, 13
- 82 [10] B. Eysenbach, T. Zhang, S. Levine, and R. R. Salakhutdinov. Contrastive learning as  
83 goal-conditioned reinforcement learning. *Advances in Neural Information Processing Systems*,  
84 35:35603–35620, 2022. 1, 6, 11, 12, 13
- 85 [11] M. Fang, T. Zhou, Y. Du, L. Han, and Z. Zhang. Curriculum-guided hindsight experience replay.  
86 *Advances in neural information processing systems*, 32, 2019. 13
- 87 [12] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven  
88 reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020. 11
- 89 [13] D. Garg, S. Chakraborty, C. Cundy, J. Song, and S. Ermon. Iq-learn: Inverse soft-q learning for  
90 imitation. *Advances in Neural Information Processing Systems*, 34:4028–4039, 2021. 9
- 91 [14] D. Garg, J. Hejna, M. Geist, and S. Ermon. Extreme q-learning: Maxent rl without entropy.  
92 *arXiv preprint arXiv:2301.02328*, 2023. 10
- 93 [15] S. K. S. Ghasemipour, R. Zemel, and S. Gu. A divergence minimization perspective on imitation  
94 learning methods. In *Conference on Robot Learning*, pages 1259–1277. PMLR, 2020. 6, 13
- 95 [16] D. Ghosh, A. Gupta, A. Reddy, J. Fu, C. Devin, B. Eysenbach, and S. Levine. Learning to reach  
96 goals via iterated supervised learning. *arXiv preprint arXiv:1912.06088*, 2019. 1, 6, 11, 13, 22
- 97 [17] L. P. Kaelbling. Learning to achieve goals. In *IJCAI*, volume 2, pages 1094–8. Citeseer, 1993.  
98 13
- 99 [18] I. Kostrikov, A. Nair, and S. Levine. Offline reinforcement learning with implicit q-learning.  
100 *arXiv preprint arXiv:2110.06169*, 2021. 1, 6, 10, 11, 22
- 101 [19] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement  
102 learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020. 11
- 103 [20] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning  
104 latent plans from play. In *Conference on robot learning*, pages 1113–1132. PMLR, 2020. 13

- 105 [21] Y. J. Ma, J. Yan, D. Jayaraman, and O. Bastani. How far i'll go: Offline goal-conditioned  
106 reinforcement learning via  $f$ -advantage regression. *arXiv preprint arXiv:2206.03023*, 2022. 1,  
107 6, 8, 9, 10, 11, 13, 14, 22
- 108 [22] A. S. Manne. Linear programming and sequential decisions. *Management Science*, 6(3):  
109 259–267, 1960. 9
- 110 [23] G. Molinaro and A. G. Collins. A goal-centric outlook on learning. *Trends in Cognitive Sciences*,  
111 2023. 13
- 112 [24] O. Nachum and B. Dai. Reinforcement learning via fenchel-rockafellar duality. *arXiv preprint*  
113 *arXiv:2001.01866*, 2020. 9, 13, 14
- 114 [25] T. Ni, H. Sikchi, Y. Wang, T. Gupta, L. Lee, and B. Eysenbach. f-irl: Inverse reinforcement  
115 learning via state marginal matching. In *Conference on Robot Learning*, pages 529–551. PMLR,  
116 2021. 13
- 117 [26] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding.  
118 *arXiv preprint arXiv:1807.03748*, 2018. 10
- 119 [27] K. Paster, S. A. McIlraith, and J. Ba. Planning from pixels using inverse dynamics models.  
120 *arXiv preprint arXiv:2012.02419*, 2020. 13
- 121 [28] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin,  
122 M. Chociej, P. Welinder, et al. Multi-goal reinforcement learning: Challenging robotics  
123 environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018. 10, 22
- 124 [29] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
125 Wiley & Sons, 2014. 7
- 126 [30] A. Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley*  
127 *Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory*  
128 *of Statistics*, pages 547–561. University of California Press, 1961. 16
- 129 [31] H. Sikchi, A. Saran, W. Goo, and S. Niekum. A ranking game for imitation learning. *arXiv*  
130 *preprint arXiv:2202.03481*, 2022. 13
- 131 [32] H. Sikchi, Q. Zheng, A. Zhang, and S. Niekum. Dual rl: Unification and new methods for  
132 reinforcement and imitation learning. In *Sixteenth European Workshop on Reinforcement*  
133 *Learning*, 2023. 9, 10, 13
- 134 [33] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller. Deterministic policy  
135 gradient algorithms. In *International conference on machine learning*, pages 387–395. Pmlr,  
136 2014. 1, 6
- 137 [34] R. K. Srivastava, P. Shyam, F. Mutz, W. Jaśkowski, and J. Schmidhuber. Training agents using  
138 upside-down reinforcement learning. *arXiv preprint arXiv:1912.02877*, 2019. 13
- 139 [35] G. Swamy, S. Choudhury, J. A. Bagnell, and S. Wu. Of moments and matching: A  
140 game-theoretic framework for closing the imitation gap. In *International Conference on*  
141 *Machine Learning*, pages 10022–10032. PMLR, 2021. 13
- 142 [36] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch,  
143 Q. Vuong, A. He, et al. Bridgedata v2: A dataset for robot learning at scale. *arXiv preprint*  
144 *arXiv:2308.12952*, 2023. 6
- 145 [37] Y. Wu, G. Tucker, and O. Nachum. Behavior regularized offline reinforcement learning. *arXiv*  
146 *preprint arXiv:1911.11361*, 2019. 14
- 147 [38] C. Yang, X. Ma, W. Huang, F. Sun, H. Liu, J. Huang, and C. Gan. Imitation learning from  
148 observations by minimizing inverse dynamics disagreement. *arXiv preprint arXiv:1910.04417*,  
149 2019. 1, 6, 13

- 150 [39] R. Yang, Y. Lu, W. Li, H. Sun, M. Fang, Y. Du, X. Li, L. Han, and C. Zhang. Rethinking  
151 goal-conditioned supervised learning and its connection to offline rl. *arXiv preprint*  
152 *arXiv:2202.04478*, 2022. 11, 22
- 153 [40] C. Zheng, B. Eysenbach, H. Walke, P. Yin, K. Fang, R. Salakhutdinov, and S. Levine. Stabilizing  
154 contrastive rl: Techniques for offline goal reaching. *arXiv preprint arXiv:2306.03346*, 2023. 1,  
155 6, 12, 13
- 156 [41] L. Zheng, T. Fiez, Z. Alumbaugh, B. Chasnov, and L. J. Ratliff. Stackelberg actor-critic:  
157 Game-theoretic reinforcement learning algorithms. *arXiv preprint arXiv:2109.12286*, 2021. 12

## 158 A Appendix

### 159 A.1 Introduction

160 A generalist agent must be able to leverage large amounts of offline pre-collected data to learn  
161 useful skills. Other fields of machine learning like vision and NLP have enjoyed great success by  
162 designing objectives to learn a general model from large and diverse datasets. In robot learning,  
163 offline interaction data has become more prominent in the recent past [7], with the scale of the datasets  
164 growing consistently [36]. Goal-conditioned reinforcement learning (GCRL) offers a principled way  
165 to acquire a variety of useful skills without the prohibitively difficult process of hand-engineering  
166 reward functions. In GCRL, the agent learns a policy to accomplish a variety of goals in the  
167 environment. The rewards are sparse and goal-conditioned: 1 when the agent’s state is in proximity  
168 to the goal and 0 otherwise. However, the benefit of not requiring the designer to hand-engineer  
169 dense reward functions can also be a curse, because learning from sparse rewards is difficult. Driving  
170 progress in fundamental offline GCRL algorithms thus becomes an important aspect of moving  
171 towards performant generalist agents whose skills scale with data.

172 Despite recent progress in developing methods for goal-reaching in the online setting (where  
173 environment interactions are allowed), a number of these methods are either suboptimal in the  
174 offline setting or suffer from learning difficulties. Prior GCRL algorithms can largely be classified  
175 into one of three categories: iterated behavior cloning, RL with sparse rewards, and contrastive  
176 learning. Iterated behavior cloning or goal-conditioned supervised learning approaches [16, 38] have  
177 been shown to be provably suboptimal [9] for GCRL. Modifying single-task RL methods [33, 18]  
178 for GCRL with 0-1 reward implies learning a  $Q$ -function that predicts the discounted probability  
179 of goal reaching, which makes it essentially a density model. Modeling density directly is a hard  
180 problem, an insight which has prompted the development of methods [8] that learn density-ratio  
181 instead of densities, as classification is an easier problem than density estimation. Contrastive RL  
182 approaches to GCRL [8, 10, 40] aim to do precisely this and are the main methods to enjoy success  
183 for applying GCRL in high-dimensional observation spaces. However, when dealing with offline  
184 datasets, contrastive RL approaches [10, 40] are suboptimal, as they learn a policy that is a greedy  
185 improvement over the  $Q$ -function of the data generation policy. This begs the question: *How can we*  
186 *derive a performant GCRL method that learns optimal policies from offline datasets of suboptimal*  
187 *quality?*

188 In this work, we leverage the underexplored insight of formulating GCRL as an occupancy matching  
189 problem. Occupancy matching between the joint state-action-goal visitation distribution induced  
190 by the current policy and the distribution over state-actions that transition to goals can be shown to  
191 be equivalent to optimizing a max-entropy GCRL objective. Occupancy matching has been studied  
192 extensively in imitation learning [15] and often requires learning a discriminator and using the learned  
193 discriminator for downstream policy learning through RL. Indeed, a prior GCRL work [21] explores a  
194 similar insight. Unfortunately, errors in learned discriminators can compound and adversely affect the  
195 learned policy’s performance, especially in the offline setting where these errors cannot be corrected  
196 with further interaction with the environment.

197 Going beyond the shortcomings of the previous methods, our proposed method combines the insight  
198 of formulating GCRL as an occupancy matching problem along with an efficient, discriminator-free  
199 dual formulation that learns from offline data. The resulting algorithm SMORe forgoes learning  
200 density functions or classifiers, but instead learns unnormalized densities or *scores* that allow it to  
201 produce optimal goal-reaching policies. The scores are learned via a Bellman-regularized contrastive  
202 procedure that makes our method a desirable candidate for GCRL with high-dimensional observations,  
203 avoiding the need for density modeling. Our experiments represent a wide variety of goal-reaching  
204 environments – consisting of robotic arms, anthropomorphic hands, and locomotion environments.  
205 We lay out the following contributions: 1) on the extended offline GCRL benchmark, our results  
206 demonstrate that SMORe significantly outperforms prior methods in the offline GCRL setting. 2) In  
207 line with our hypothesis, discriminator-free training makes SMORe particularly robust to decreasing  
208 goal-coverage in the offline dataset, a property we demonstrate in the experiments. 3) We test SMORe  
209 for zero-shot GCRL on a prior benchmark [40] for high dimensional vision-based GCRL where  
210 contrastive RL approaches are the only class of GCRL methods that have been successful, and show  
211 improved performance over other state-of-the-art baselines.

## 212 A.2 Problem Formulation

213 **Goal-Conditioned Reinforcement Learning:** We consider an infinite-horizon Markov decision  
 214 process (MDP) [29]  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, p, d_0, \gamma)$  with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , deterministic rewards  
 215  $r(s, a)$ , transition probabilities  $p(s' | s, a)$  from state  $s$  to  $s'$  given action  $a$ , initial state distribution  
 216  $d_0(s)$ , and discount factor  $\gamma \in (0, 1)$ . A policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  outputs a distribution over actions in a  
 217 given state. In goal-conditioned RL, the MDP additionally assumes a goal space  $\mathcal{G} := \{\phi(s) | s \in \mathcal{S}\}$ ,  
 218 where the state-to-goal mapping  $\phi : \mathcal{S} \rightarrow \mathcal{G}$  is known. The sparse reward function  $r(s, a, g)$  as well  
 219 as the policy  $\pi(a | s, g)$  depend on the commanded goal  $g \in \mathcal{G}$ . Given a distribution over desired  
 220 evaluation goals  $q^{\text{test}}(g)$ , the objective of goal-conditioned RL is to find a policy  $\pi_g^1$  that maximizes  
 221 the discounted return:

$$J(\pi_g) := \mathbb{E}_{g \sim q^{\text{test}}(g), s_0 \sim d_0, a_t \sim \pi_g(\cdot | s_t, g)} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, g) \right]. \quad (3)$$

222 We denote by  $P^{\pi_g}$  the transition operator induced by the policy  $\pi_g$  defined as  $P^{\pi_g}S(s, a, g) :=$   
 223  $\mathbb{E}_{s' \sim p(\cdot | s, a), a' \sim \pi_g(\cdot | s', g)}[S(s', a', g)]$ , for any *score* function  $S : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow \mathbb{R}$ .

224 The *goal-conditioned* discounted state-action occupancy distribution  $d^{\pi_g}(s, a | g)$  of  $\pi_g$  is given by:

$$d^{\pi_g}(s, a | g) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s, a_t = a | s_0 \sim d_0, a_t \sim \pi_g(\cdot | s_t, g), s_{t+1} \sim p(\cdot | s_t, a_t)), \quad (4)$$

225 which represents the expected discounted time spent in each state-action pair by the policy  $\pi_g$   
 226 conditioned on the goal  $g$ . It follows that  $\pi_g(a | s, g) = \frac{d^{\pi_g}(s, a | g)}{d^{\pi_g}(s | g)}$ , where  $d^{\pi_g}(s | g) :=$   
 227  $\sum_{s, a} d^{\pi_g}(s, a | g)$ . For complete generality, in GCRL, the distribution of goals the policy is  
 228 trained on often differs from the test goal distribution. To make this distinction clear we define  
 229 the training distribution  $q^{\text{train}}(g)$ , a uniform measure over goals we desire to learn to optimally  
 230 reach during training. We write  $d^{\pi_g}(s, a, g) = q^{\text{train}}(g)d^{\pi_g}(s, a | g)$  as the joint state-action-goal  
 231 visitation distribution of the policy  $\pi_g$  under the training goal distribution. A state-action-goal  
 232 occupancy distribution must satisfy the *Bellman flow constraint* in order for it to be a valid occupancy<sup>2</sup>  
 233 distribution for some stationary policy  $\pi_g, \forall s \in \mathcal{S}, a \in \mathcal{A}, g \in \mathcal{G}$ :

$$d(s, a, g) = (1 - \gamma)d_0(s, g)\pi_g(a | s, g) + \gamma \sum_{s', a'} p(s | s', a')d(s', a', g)\pi_g(a | s, g), \quad (5)$$

234 where  $d_0(s, g) = d_0(s)q^{\text{train}}(g)$ . Finally, given  $d^{\pi_g}$ , we can express the learning objective for the  
 235 GCRL agent under the training goal distribution as  $J^{\text{train}}(\pi_g) = \frac{1}{1-\gamma} \mathbb{E}_{(s, a, g) \sim d^{\pi_g}}[r(s, a, g)]$ .

236 **Offline GCRL.** In offline GCRL, the agent cannot interact with the environment  $\mathcal{M}$  and is  
 237 equipped with a static dataset of logged transitions  $\mathcal{D} := \{\tau_i\}_{i=1}^N$ , where each trajectory  $\tau^{(i)} =$   
 238  $(s_0^{(i)}, a_0^{(i)}, r_0^{(i)}, s_1^{(i)}, \dots; g^{(i)})$  with  $s_0^{(i)} \sim d_0$ . The trajectories are not necessarily generated by a  
 239 goal-directed agent and are relabelled with the  $q^{\text{train}}(g)$  during learning. We denote the joint  
 240 state-action-goal distribution of the offline dataset  $\mathcal{D}$  as  $\rho(s, a, g)$ .

## 241 A.3 Score-models for Offline Goal Conditioned Reinforcement Learning

242 In this section, we introduce our method in two parts: First, we build up the equivalence of the GCRL  
 243 objective to the occupancy matching problem in Section A.3.1, and then we derive a discriminator-free  
 244 dual objective for solving the occupancy matching problem using off-policy data in Section A.3.2.  
 245 Finally, we present the algorithm for SMORe under practical considerations in Section A.3.3.

### 246 A.3.1 GCRL as an occupancy matching problem

247 Define a *goal-transition distribution*  $q(s, a, g)$  in a stochastic MDP as  
 248  $q(s, a, g) \propto q^{\text{train}}(g)\mathbb{E}_{s' \sim p(\cdot | s, a)}[\mathbb{I}_{\phi(s')=g}]$ . Intuitively, the distribution has probability mass  
 249 on each transition that leads to a goal. We formulate the GCRL problem as an occupancy matching  
 250 problem by searching for the policy  $\pi_g$  that minimizes the discrepancy between its state-action-goal  
 251 occupancy distribution and the goal-transition distribution  $q(s, a, g)$ :

$$\text{Occupancy matching problem: } \mathcal{D}_f(d^{\pi_g}(s, a, g) \| q(s, a, g)), \quad (6)$$

<sup>1</sup>We use the subscript  $g$  to make the policy's conditioning on  $g$  explicit.

<sup>2</sup>We will use ‘‘occupancy’’ and ‘‘visitation’’ interchangeably.

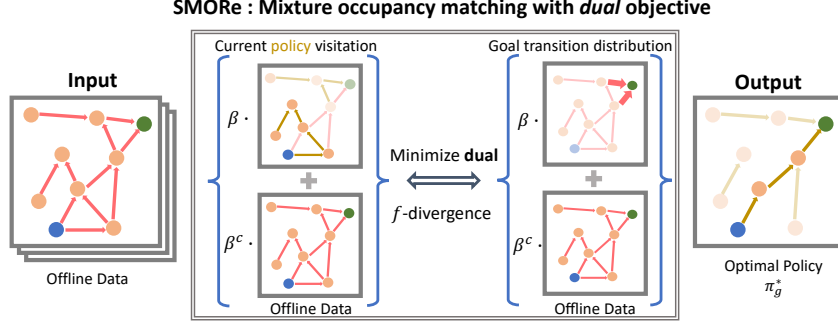


Figure 2: Illustration of the SMORe objective where  $\beta^c = 1 - \beta$ : SMORe matches a mixture distribution of current policy and offline data to a mixture of the goal-transition distribution and offline data in order to find the optimal goal reaching policy.

where  $D_f$  denotes an  $f$ -divergence with generator function  $f$ . Note that the  $q$  distribution is potentially unachievable by any goal-conditioned policy  $\pi_g$ . Firstly, it does not account for the initial transient phase that the policy must navigate to reach the desired goal. Secondly, even if we consider only the stationary regime (when  $\gamma \rightarrow 1$ ), it may not be dynamically possible for the policy to continuously remain at the goal and rather necessitate cycling around the goal. However, in Proposition 1.1, we show that the occupancy matching in Eq. 6 offers a principled objective since it forms a lower bound to the max-entropy GCRL problem.

**Proposition 1.1.** Consider a stochastic MDP, a stochastic policy  $\pi$ , and a sparse reward function  $r(s, a, g) = \mathbb{E}_{s' \sim p(\cdot|s,a)}[\mathbb{I}(\phi(s') = g, q^{\text{train}}(g) > 0)]$  where  $\mathbb{I}$  is an indicator function. Define a soft goal transition distribution to be  $q(s, a, g) \propto \exp(\alpha r(s, a, g))$ . The following bounds hold for any  $f$ -divergence that upper bounds KL-divergence (eg.  $\chi^2$ , Jensen-Shannon):

$$J^{\text{train}}(\pi_g) + \frac{1}{\alpha} \mathcal{H}(d^{\pi_g}) \geq -\frac{1}{\alpha} \mathcal{D}_f(d^{\pi_g}(s, a, g) \| q(s, a, g)) + C, \quad (7)$$

where  $\mathcal{H}$  denotes the entropy,  $\alpha$  is a temperature parameter and  $C$  is the partition function for  $e^{R(s,a,g)}$ . Furthermore, the bound is tight when  $f$  is the KL-divergence.

Proposition 1.1 extends the insights of formulating GCRL as an imitation learning problem from [21] for goal-transition distributions when matching state-action-goal visitations.

How does converting a GCRL objective to an imitation learning objective make learning easier? Estimating the  $f$ -divergence still requires estimating the joint policy visitation probabilities  $d^{\pi_g}(s, a, g)$ , which itself presents a challenging problem. We show in the following section that we can leverage convex duality to transform the imitation learning problem into an off-policy optimization problem, removing the need to sample from  $d^{\pi_g}(s, a, g)$  whilst being able to leverage offline data collected from arbitrary sources.

### 273 A.3.2 SMORe: A Dual Formulation for Occupancy Matching

274 The previous section establishes GCRL as an occupancy matching problem (Eq. 6) but provides no  
 275 way to use offline data whose joint visitation distribution is given by  $\rho(s, a, g)$ . To leverage offline  
 276 data to learn performant goal-reaching policies, we consider a surrogate objective to the occupancy  
 277 matching learning problem by matching *mixture* distributions:

$$\min_{\pi_g} \mathcal{D}_f(\text{Mix}_\beta(d^{\pi_g}, \rho)(s, a, g) \| \text{Mix}_\beta(q, \rho)(s, a, g)), \quad (8)$$

278 where for any two distributions  $\mu_1$  and  $\mu_2$ ,  $\text{Mix}_\beta(\mu_1, \mu_2)$  denotes the mixture distribution with  
 279 coefficient  $\beta \in (0, 1]$  defined as  $\text{Mix}_\beta(\mu_1, \mu_2) = \beta\mu_1 + (1 - \beta)\mu_2$ . Proposition 2.1 (in appendix)  
 280 shows the matching mixture distribution<sup>3</sup> provably maximizes a lower bound to the Lagrangian  
 281 relaxation of the max-entropy GCRL objective subject to the constraint that the policy visitation is  
 282 close to the offline data visitation. We can rewrite the mixture occupancy matching objective as a

<sup>3</sup>Note that Eq. 8 shares the same global optima as the previous occupancy matching objective at  $d_g^\pi(s, a, g) = q(s, a, g)$  when  $q$  is an achievable visitation under some policy and recovers the original objective in Eq. 6 when  $\beta = 1$ .



283 convex program with linear constraints [22, 24]:

$$\max_{\pi_g, d} -\mathcal{D}_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g))$$

$$\text{s.t. } d(s, a, g) = (1 - \gamma)d_0(s, g)\pi(a|s) + \gamma \sum_{s' \in \mathcal{S}} d(s', a', g)p(s|s', a')\pi(a'|s', g), \quad \forall s \in \mathcal{S}. \quad (9)$$

284 An illustration of this objective can be found in Figure 2. Effectively, we have simply rewritten  
 285 Eq. 8 into an equivalent problem by considering an arbitrary probability distribution  $d(s, a, g)$  in  
 286 the optimization objective, only to later constrain it to be a valid probability distribution induced by  
 287 some policy  $\pi_g$  using the *Bellman-flow constraints*. The motivation behind this construction of the  
 288 primal form is that we have made computing the Lagrangian-dual easier as this objective is convex  
 289 with linear constraints. Theorem 2 shows that we can leverage tools from convex duality to obtain an  
 290 unconstrained dual problem that does not require computing  $d^{\pi_g}(s, a, g)$  or sampling from it, while  
 291 effectively leveraging offline data.

292 **Theorem 2.** *The dual problem to the primal occupancy matching objective (Equation 9) is given by:*

$$\max_{\pi_g} \min_S \beta(1 - \gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + \mathbb{E}_{\text{Mix}_\beta(q, \rho)}[f^*(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \quad (10)$$

$$- (1 - \beta)\mathbb{E}_\rho[\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)],$$

293 where  $f^*$  is conjugate function of  $f$  and  $S$  is the Lagrange dual variable defined as  $S : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow$   
 294  $\mathbb{R}$ . Moreover, as strong duality holds from Slater’s conditions the primal and dual share the same  
 295 optimal solution  $\pi_g^*$  for any offline transition distribution  $\rho$ .

296 To our knowledge, the closest prior works to our proposed method are GoFAR [21] and Dual-RL [32].  
 297 GoFAR considers the special case of KL-divergence for the imitation formulation and derives a dual  
 298 objective that requires learning the density ratio  $\frac{\rho(s, g)}{q(s, g)}$  in the form of a discriminator and using this  
 299 as a pseudo-reward. This leads to compounding errors in the downstream RL optimization when  
 300 learning the density ratio is challenging, e.g. in the case of low coverage between  $\rho(s, a, g)$  and  
 301  $q(s, a, g)$ . We show this phenomenon experimentally in Section A.4.3. Dual-RL [32] uses convex  
 302 duality for matching visitation distribution of realizable expert demonstrations and does not deal  
 303 with the GCRL setting. *Our contribution is a novel method for GCRL that is discriminator-free,*  
 304 *applicable for a number of  $f$ -divergences, and robust to low coverage of goals in the offline dataset.*

305 **Sampling from the goal-transition distribution:** Goal relabelling is an effective technique to  
 306 address reward sparsity by widening the training goal distribution  $q^{\text{train}}(g)$ . It utilizes knowledge  
 307 about reaching other goals, possibly unrelated to test goals, to help in reaching the test distribution  
 308 of goals  $q^{\text{test}}(g)$ . In the most general case,  $q^{\text{train}}(g)$  can be set to a uniform distribution over  
 309 goals corresponding to all the states in the offline data. A common method, Hindsight Experience  
 310 Replay (HER) [3] chooses a training goal distribution that depends on the current sampled state from  
 311 the offline dataset as well as the data-collecting policies. In this setting, the sampling distribution  
 312 used for training Eq 10,  $\rho(s, a, g)$ , can no longer be factorized into  $\rho(s, a)$  and  $q^{\text{train}}(g)$ , as goals  
 313 are conditionally dependent on state-actions. However, our formulation can naturally account for  
 314 learning from such relabelled data as the SMORe objective in Eq 10 is derived considering the joint  
 315 distribution  $\rho(s, a, g)$ . In this setting, we construct our goal transition distribution  $q(s, a, g)$  as the  
 316 uniform distribution over all transitions that lead to the goals selected by the HER procedure — in  
 317 practice, this amounts to first selecting  $g$  through HER and then selecting  $\{s, a\}$  that transitions to  
 318 the selected goal from the offline dataset to get a sample  $\{s, a, g\}$  from goal transition distribution.  
 319 We emphasize that relabelling does not change the test distribution of goals, which is an immutable  
 320 property of the environment.

### 321 A.3.3 Practical Algorithm

322 To devise a stable learning algorithm we consider the Pearson  $\chi^2$  divergence. Pearson  $\chi^2$  divergence  
 323 has been found to lead to distribution matching objectives that are stable to train as a result of a  
 324 smooth quadratic generator function  $f$  [13, 2, 32]. Our dual formulation SMORe simplifies to the  
 325 following objective:

$$\max_{\pi_g} \min_S \underbrace{\beta(1 - \gamma)\mathbb{E}_{(s, g) \sim d_0, a \sim \pi_g(\cdot|s, g)}[S(s, a, g)]}_{\text{Decrease score at transitions under current policy } \pi_g} + \underbrace{\beta\gamma\mathbb{E}_{(s, a, g) \sim q, s' \sim p(\cdot|s, a), a' \sim \pi_g(\cdot|s', g)}[S(s', a', g)]}_{\text{Smoothness/Bellman regularization}}$$

$$- \underbrace{\beta\mathbb{E}_{(s, a, g) \sim q}[S(s, a, g)]}_{\text{Increase score at the proposed goal transition distribution}} + \underbrace{0.25\mathbb{E}_{(s, a, g) \sim \text{Mix}_\beta(q, \rho)}[(\gamma S(s', \pi_g(s'), g) - S(s, a, g))^2]}_{\text{Smoothness/Bellman regularization}}. \quad (11)$$

Equation 11 suggests a contrastive procedure, maximizing the score at the goal-transition distribution and minimizing the score at the offline data distribution under the current policy with Bellman regularization. The Bellman regularization has the interpretation of discouraging neighboring  $S$  values from deviating far and smoothing the score landscape. Instantiating with KL divergence results in an objective with similar intuition while resembling an InfoNCE [26] objective. Although Propositions 1.1 and 2.1 suggest that KL divergence gives an objective that is a tighter bound to the GCRL objective, prior work has found KL divergence to be unstable in practice [32, 14] for dual optimization. It is important to note that  $S$ -function is not grounded to any rewards and does not serve as a probability density of reaching goals, but is rather a score function learned via a *Bellman-regularized contrastive learning procedure*.

We now derive a practical approach for SMORe in the offline GCRL setting. We use parameterized functions:  $S_\phi(s, a, g)$ ,  $M_\psi(s, g)$ ,  $\pi_\theta(a|s, g)$ . The offline learning regime necessitates measures to constrain the learning policy to the offline data support in order to prevent overestimation due to maximizing  $\pi_g$  in Eq. 11 over potentially out-of-distribution actions. Inspired by prior work [18], we use implicit maximization to constrain the learning algorithm to learn expectiles using the observed empirical samples. More concretely, we use expectile regression:

$$\min_{\psi} \mathcal{L}(\psi) := \mathbb{E}_{(s,a,g) \sim \rho} [L_2^\tau(M_\psi(s, g) - S_\phi(s, a, g))], \quad (12)$$

where  $L_2^\tau(u) = |\tau - 1(u < 0)|u^2$ . Intuitively, this step implements the maximization w.r.t  $\pi$  by using expectile regression. With the above practical considerations, our objective for learning  $S_\phi$  reduces to:

$$\begin{aligned} \min_{\phi} \mathcal{L}(\phi) := & \beta(1 - \gamma) \mathbb{E}_{(s,g) \sim \mathcal{D}, a \sim \pi_g(\cdot|a,g)} [S_\phi(s, \pi_g(s), g)] + \beta\gamma \mathbb{E}_{(s,a,g) \sim q, s' \sim p(\cdot|s,a)} [S_\phi(s', \pi_g(s'), g)] \\ & - \beta \mathbb{E}_{(s,a,g) \sim q} [S_\phi(s, a, g)] + \mathbb{E}_{(s,a,g) \sim \text{Mix}_{\beta}(q,\rho)} [(\gamma M_\psi(s', g) - S_\phi(s, a, g))^2], \end{aligned} \quad (13)$$

where we have set the offline data distribution as our initial state distribution. Finally, the policy is extracted via advantage-weighted regression that learns in-distribution actions maximizing the score  $S(s, a, g)$ :

$$\min_{\theta} \mathcal{L}(\theta) := \mathbb{E}_{(s,a,g) \sim \rho} [\exp(\alpha(S_\phi(s, a, g) - M_\psi(s, g))) \log(\pi_\theta(a|s, g))], \quad (14)$$

where  $\alpha$  is the temperature parameter. Algorithm 1 details the practical implementation.

## A.4 Experiments

Our experiments study the effectiveness of proposed GCRL algorithm SMORe on a set of simulated benchmarks against other GCRL methods that employ behavior cloning, RL with sparse reward, and contrastive learning. We also analyze if SMORe is robust to environment stochasticity — a number of prior methods are based on an assumption of deterministic dynamics. Then, we study if the discriminator-free nature of SMORe is indeed able to prevent performance degradation in the face of low expert coverage in offline data. Finally, we analyze if SMORe’s score-modeling approach helps SMORe scale to a vision-based manipulation offline GCRL benchmark, as density modeling and discriminator learning become increasingly difficult with high-dimensional observations. Hyperparameter ablations can be found in Appendix E.

### A.4.1 Experimental Setup

Our experiments will use a suite of simulated goal-conditioned tasks extending the tasks from previous work [21, 28]. In particular we consider the following environments: Reacher, Robotic arm environments - [SawyerReach, SawyerDoor, FetchReach, FetchPick, FetchPush, FetchSlide], Anthropomorphic hand environment - HandReach and Locomotion environments

---

#### Algorithm 1: SMORe

---

- 1: Init  $S_\phi$ ,  $M_\psi$ , and  $\pi_\theta$
  - 2: Params: expectile  $\tau$ , mixture ratio  $\beta$ , temperature  $\alpha$
  - 3: Let  $\mathcal{D} = \hat{\rho} = \{(s, a, s', g)\}$  be an offline dataset and  $q$  be goal-transition distribution
  - 4: **for**  $t = 1..T$  iterations **do**
  - 5:   Train  $S_\phi$  via Eq. 13
  - 6:   Train  $M_\psi$  via Eq. 12
  - 7:   Update  $\pi_\theta$  via Eq. 14
  - 8: **end for**
-

371 -[CheetahTgtVel-me,CheetahTgtVel-re,AntTgtVel-me,AntTgtVel-re]. Tasks in all  
 372 environments are specified by a sparse reward function. Depending on whether the task involves  
 373 object manipulation, the goal distribution is defined over valid configurations in robot or object space.  
 374 The offline dataset for manipulation tasks consists of transitions collected by a random policy or  
 375 mixture of 90% random policy and 10% expert policy. For locomotion tasks, we generate our dataset  
 376 using the D4RL benchmark [12], combining a random or medium dataset with 30 episodes of expert  
 377 data. Note that the policies used to collect the expert locomotion datasets have a different objective  
 378 than the tasks here, which are to achieve and maintain a particular desired velocity.

#### 379 A.4.2 Offline Goal-conditioned RL benchmark

380 **Baselines.** We compare to state-of-art offline GCRL algorithms, consisting of both regression-based  
 381 and actor-critic methods. The occupancy-matching based methods are: (1) **GoFar** [21], which derives  
 382 a dual objective for GCRL based on a coverage assumption. The behavior cloning based methods  
 383 are: (1) **GCSL** [16], which incorporates hindsight relabeling in conjunction with behavior cloning to  
 384 clone actions that lead to a specified goal, and (2) **WGCSL** [39], which improves upon GCSL by  
 385 incorporating discount factor and advantage weighting into the supervised policy learning update.  
 386 **Contrastive RL** [10] generalizes C-learning [8] and represents contrastive GCRL approaches. The  
 387 RL with sparse reward methods are (1) **IQL** [18] where we use a state-of-the-art offline RL method  
 388 repurposed for GCRL along with HER [3] goal sampling, and (2) **ActionableModel (AM)** [4], which  
 389 incorporates conservative Q-Learning [19] as well as goal-chaining on top of an actor-critic method.

390 The results for all baselines are tuned individually, particularly the best HER ratio was searched  
 391 among  $\{0.2, 0.5, 0.8, 1.0\}$  for each task. SMORe shares the same network architecture for baselines  
 392 and uses a mixture ratio of  $\beta = 0.5$ . Each method is trained for 10 seeds. Complete architecture and  
 393 hyperparameter table as well as additional training details are provided in Appendix D.

394 Table 2 reports the **discounted return** obtained  
 395 by the learned policy with a sparse binary task  
 396 reward. (\*) denotes statistically significant  
 397 improvement over the second best method under  
 398 a two-sample t-test. This metric allows us  
 399 to compare the algorithms on a finer scale to  
 400 understand which methods reach the goal as fast  
 401 as possible and stay in the goal region thereafter  
 402 for the longest time. Additional results on  
 403 metrics like success rate and final distance to  
 404 goal can be found in the appendix. These  
 405 additional metrics do not take into consideration  
 406 how *precisely* and *consistently* a goal is being  
 407 reached. In Table 2, we see that SMORe enjoys  
 408 a high-performance gain consistently across all  
 409 tasks in the extended offline GCRL benchmark.

410 **Robustness to environment stochasticity:** We consider a noisy version of the FetchReach  
 411 environment in this experiment. Gaussian zero-mean noise is added to generate different variants  
 412 of the environment with standard deviations of  $\{0.5, 1.0, 1.5\}$ . Datasets for these environments  
 413 are obtained from prior work [21]. As we see in Figure 3, SMORe is robust to stochasticity  
 414 in the environment, outperforming baselines in terms of discounted return. Behavior cloning  
 415 based approaches assume deterministic dynamics and are therefore over-optimistic in stochastic  
 416 environments.

#### 417 A.4.3 Robustness of Occupancy-Matching Methods to Decreasing Expert Coverage

418 We posit that the discriminator-free nature of SMORe makes it more robust to decreasing goal coverage,  
 419 as it does not suffer from cascading errors stemming from a learned discriminator. In this section, we  
 420 set out to test this hypothesis by decreasing the amount of expert data in the offline goal-reaching  
 421 dataset. We compare with GoFar in Table 3 due to the similarity between methods and GoFar’s  
 422 restrictive assumption on coverage of expert data in the suboptimal dataset. Comparison against all  
 423 the baselines can be found in Appendix E.

424 Our hypothesis holds true as we see in Table 3, the performance of the discriminator-based method  
 425 GoFar rapidly decays as expert data is decreased in the offline dataset – 28.4% with 2.5% and 36.15%

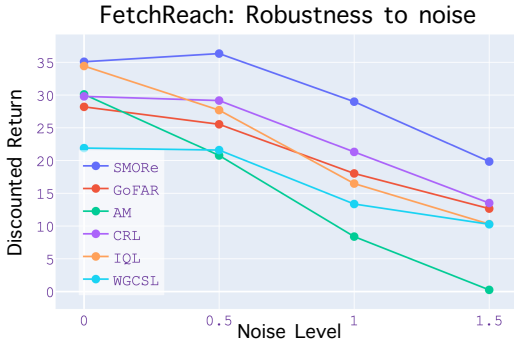


Figure 3: SMORe is robust in stochastic environments. With increasing noise, SMORe still outperforms prior methods.

Task	Occupancy Matching		Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	SMORe	GoFAR	WGCSL	GCSL		AM	IQL
Reacher (*)	28.40±0.88	19.74±1.35	17.57±0.53	15.87±1.31	16.44±0.60	23.26 ±0.14	11.70 ±1.97
SawyerReach (*)	37.67±0.12	15.34±0.64	15.15±0.44	14.25±0.7	22.32 ±0.34	23.34±0.17	35.18 ±0.29
SawyerDoor (*)	31.48±0.46	18.94±0.01	20.01±1.55	20.88±0.22	12.96±5.19	22.12 ±0.13	25.52 ±1.45
FetchReach (*)	35.08±0.54	28.2 ±0.61	21.9±2.13	20.91 ±2.78	30.07±0.07	30.1 ±0.32	34.43 ±1.00
FetchPick (*)	26.47 ±1.34	19.7 ±2.57	9.84 ±2.58	7.58±1.85	0.42±0.29	8.94 ±3.09	16.8 ±3.10
FetchPush (*)	26.83±1.21	18.2 ±3.00	14.7 ±2.65	13.4 ±3.02	2.40 ±1.28	14.0 ±2.81	22.40 ±0.74
FetchSlide	4.99±0.40	2.47 ±1.44	2.73 ±1.64	1.75 ±1.3	0.0±0.0	1.46 ±1.38	4.80 ±1.59
HandReach (*)	18.68 ±3.35	11.5 ±5.26	5.97 ±4.81	1.37 ±2.21	0.0±0.0	0.0 ±0.0	1.44 ±1.77
CheetahTgtVel-m-e (*)	136.71 ±10.59	0.0±0.0	0.0±0.0	95.98±15.72	0.0±0.0	0.0±0.0	100.38±1.22
CheetahTgtVel-r-e (*)	60.01 ±39.40	0.0±0.0	0.0±0.0	11.56 ±13.47	0.0±0.0	0.0±0.0	0.0±0.0
AntTgtVel-m-e	154.95±19.44	168.27±9.58	0.0±0.0	164.54±7.69	0.0±0.0	0.0±0.0	148.17 ±5.43
AntTgtVel-r-e (*)	126.22±14.40	74.36±15.97	0.0±0.0	104.95±6.00	0.0±0.0	0.0±0.0	3.06 ±2.64

Table 2: Discounted Return for the offline GCRL benchmark. Results are averaged over 10 seeds. 'm-e' and 'r-e' stands for medium-expert mixture and random-expert mixture respectively.

Task	5 % expert data		2.5 % expert data		1 % expert data	
	SMORe	GoFAR	SMORe	GoFAR	SMORe	GoFAR
Reacher	22.43±3.46	16.86 ±1.26	17.92 ±0.93	12.20±0.81	19.61±1.56	11.52 ±0.52
SawyerReach	36.35±0.37	13.20 ±1.36	36.74±0.62	11.57 ±1.79	35.44 ±0.27	9.34±0.17
SawyerDoor	32.82±0.88	20.07±0.01	25.69±0.21	19.54±1.32	23.78±2.88	18.04 ±1.80
FetchReach	36.00±0.01	27.66 ±0.55	35.58 ±0.47	27.84 ±0.82	35.97 ±0.25	28.01 ±0.20
FetchPick	26.43±1.95	16.21 ±1.46	26.17±3.37	3.21 ±2.22	15.38 ±1.52	0.31 ±0.31
FetchPush	23.81±0.37	18.2 ±3.00	22.75±1.08	5.17 ±2.01	19.04±2.79	4.23±3.96
FetchSlide	4.05±1.12	1.08 ±0.06	3.11 ±1.61	0.96 ±0.73	3.50±0.97	0.86 ±1.22
Average Performance	25.98	16.18	23.99	11.49	21.81	10.33
Avg. Perf. Drop	0	0	-7.6%	-28.4%	-16%	-36.15%

Table 3: Discounted Return for the offline GCRL benchmark with 5%, 2.5% and 1% expert data in offline dataset. Results are averaged over 10 seeds.

with 1% expert data (i.e. optimal policy's coverage) respectively. SMORe shows a much slower decay in performance, 7.6% with 2.5% and 16% with 1% expert data, attesting to the method's robustness under decreasing expert coverage in the offline dataset.

#### A.4.4 Offline GCRL with image observations

SMORe provides an effective algorithm for offline GCRL in high-dimensional observation spaces by learning unnormalized scores using a contrastive procedure as opposed to prior works that learn normalized densities [8] which are difficult to learn or density ratios [10, 40] which do not optimize for the optimal goal-conditioned policy in the offline GCRL setting. Similar to prior work [10], we consider the following structure in S-function parameterization to learn performant and generalizable policies:  $S(s, a, g) = \phi(s, a)^T \psi(g)$ . The S-function can be interpreted as the similarity between the two representations given by  $\phi$  and  $\psi$ . Our network architecture for both representations is similar to [40] and is kept the same across all baselines to ensure a fair comparison of the underlying GCRL method.

We use the offline GCRL benchmark from [41] which learns goal-reaching policies from an image-observation dataset of 250K transitions with the horizon ranging from 50-100. The benchmark adds another layer of complexity by testing on goals absent from the dataset — the dataset contains

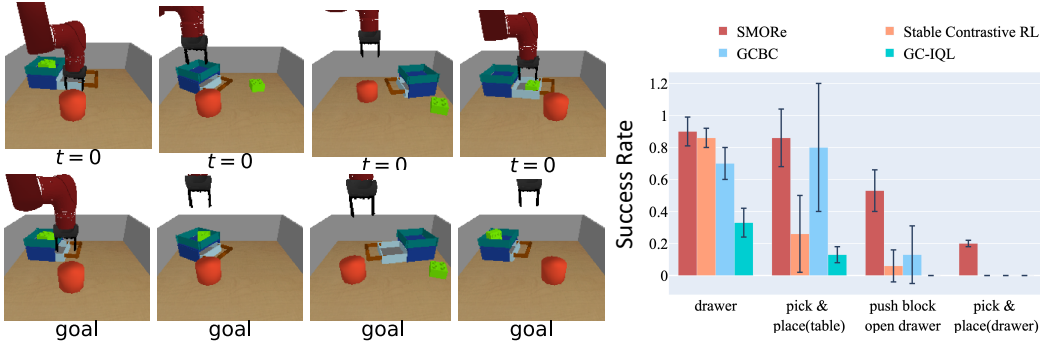


Figure 4: Evaluation on simulated manipulation tasks with image observations. The left image shows the starting state at the top and the goal at the bottom for evaluation tasks. SMORe outperforms prior methods on all the tasks we considered.

442 primitive behaviors like picking up objects and pushing drawers but no behavior that completes the  
443 compound task we consider from the initial state. The observations and goals are 48x48x3 RGB  
444 images.

445 **Baselines** We compare to the best performing GCRL algorithms from Section 1 as well as a recent  
446 state-of-the-art work, stable contrastive RL [40]. Stable contrastive RL features a number of  
447 improvements over contrastive RL by changing design decisions in neural network architecture,  
448 layer normalization, and data augmentation. Since our objective is to compare the quality of the  
449 underlying GCRL algorithm, we keep these design decisions consistent across the board.

450 **Results** Figure 4 shows the success rate on a variety of unseen tasks for all the methods. SMORe  
451 achieves the highest success rate across all the methods, even for the most challenging task of pick,  
452 place and closing the drawer. We note that our results differ from [40] for the baselines as we apply  
453 the same design decisions for all methods whereas [40] focuses on ablating design decisions.

## 454 A.5 Related Works

455 **Offline Goal Conditioned Reinforcement Learning.** Learning to achieve goals in the environment  
456 optimally forms the basis of goal-condition RL problems. Studies in cognitive science [23] underscore  
457 the importance goal-achieving plays in human development. Offline GCRL approaches are typically  
458 catered to designing learning algorithms for addressing the sparsity of reward function in the  
459 offline setting. One of the most successful techniques in this setting has been hindsight relabelling.  
460 Hindsight-experience relabelling (HER) [17, 3] suggests relabelling any experience with some  
461 commanded goal to the goal that was actually achieved in order to leverage generalization. HER  
462 has been investigated in the setting of learning from demonstrations [6] and exploration [11] to  
463 validate its effectiveness. A number of prior works [16, 38, 5, 6, 20, 27, 34] have investigated using  
464 goal-conditioned behavior cloning, a strategy that uses relabelling to learn goal-conditioned policies,  
465 as a way to learn performant policies. Eysenbach et al. [9] shows that this line of work has a limitation  
466 of learning suboptimal policies that do not consistently improve over the policy that collected the  
467 dataset. The simplest strategy of applying single-task RL to the problem of multi-task goal reaching  
468 requires learning a  $Q$ -function which represents normalized densities over the state-action space.  
469 Contrastive RL [10, 8, 40] emerged as another alternative for GCRL which relabels trajectories and,  
470 rather than use that relabelling to learn policies, learns a  $Q$ -function using a contrastive procedure.  
471 While these approaches learn optimal policies in the online setting, they fall behind in the offline  
472 setting where they only learn a policy that greedily improves over the  $Q$ -function of the data collecting  
473 policy. Our work learns optimal policies by presenting an off-policy objective that solves GCRL  
474 and furthermore learns scores (or unnormalized densities) that alleviate the learning challenges of  
475 normalized density estimation.

476 **Distribution matching.** Our approach is inspired by the distribution matching approach [15, 25, 31,  
477 35, 32] prominent in imitation learning. Ghasemipour et al. [15], Ni et al. [25] takes the problem of  
478 imitating an expert demonstrator in the environment and converts it into a problem of distribution  
479 matching between the current policy’s state-action visitation distribution and the expert policy’s  
480 visitation distribution. Indeed, prior work [21] creates one such distribution matching problem and  
481 presents a new optimization problem for GCRL in the form of an off-policy dual [24, 32]. Such an  
482 off-policy dual is very appealing for the offline RL setup, as optimizing for this dual only requires  
483 sampling from the offline data distribution. A limitation of their dual construction is the fact that  
484 they require learning a discriminator and use that discriminator as the pseudo-reward for solving the  
485 GCRL objective. Our approach presents a new construction for GCRL as a distribution matching  
486 along with a dual construction that leads to a more performant discriminator-free off-policy approach  
487 for GCRL.

## 488 A.6 Conclusion

489 Prior work in performant online goal-conditioned RL often relies on iterated behavior cloning  
490 or contrastive RL. However, these approaches are suboptimal for the offline setting. Existing  
491 methods specifically derived for offline GCRL require learning a discriminator and using it as  
492 a pseudo-reward, enabling compounding errors that make the resulting policy ineffective. We  
493 present an occupancy-matching approach to offline GCRL that provably optimizes a lower bound  
494 to the regularized GCRL objective. Our method is discriminator-free, applicable to a number of  
495  $f$ -divergences, and learns unnormalized scores over actions at a state to reach the goal. We show that

496 these positive aspects of our algorithm allow us to empirically outperform prior methods, stay robust  
 497 under decreasing goal coverage, and scale to high-dimensional observation space for GCRL.

## 498 B Supplementary Materials

### 499 B.1 Theory

500 In this section, we first show the equivalence of the GCRL problem and the distribution-matching  
 501 objective of imitation learning. Then, we show how the mixture distribution objective relates to  
 502 offline GCRL objective. Finally, we derive the dual objective for mixture distribution matching that  
 503 leads to our method SMORe.

#### 504 B.1.1 Reduction of GCRL to distribution matching

505 **Proposition 1.1.** *Consider a stochastic MDP, a stochastic policy  $\pi$ , and a sparse reward function*  
 506  *$r(s, a, g) = \mathbb{E}_{s' \sim p(\cdot|s,a)}[\mathbb{I}(\phi(s') = g, q^{\text{train}}(g) > 0)]$  where  $\mathbb{I}$  is an indicator function. Define a soft*  
 507 *goal transition distribution to be  $q(s, a, g) \propto \exp(\alpha r(s, a, g))$ . The following bounds hold for any*  
 508  *$f$ -divergence that upper bounds KL-divergence (eg.  $\chi^2$ , Jensen-Shannon):*

$$J^{\text{train}}(\pi_g) + \frac{1}{\alpha} \mathcal{H}(d^{\pi_g}) \geq -\frac{1}{\alpha} \mathcal{D}_f(d^{\pi_g}(s, a, g) \| q(s, a, g)) + C, \quad (7)$$

509 where  $\mathcal{H}$  denotes the entropy,  $\alpha$  is a temperature parameter and  $C$  is the partition function for  
 510  $e^{R(s,a,g)}$ . Furthermore, the bound is tight when  $f$  is the KL-divergence.

511 *Proof.* This proof is adapted from [21] for goal transition distributions and state-action distributions.  
 512 Let  $Z = \int e^{R(s,a,g)} ds da dg$  and  $\alpha > 0$  be the temperature parameter. Note that  $q(s, a, g) =$   
 513  $e^{r(s,a,g)}$  where  $r$  is defined in the proposition, strictly generalizes the original definition  $q(s, a, g) =$   
 514  $q^{\text{train}}(g) \mathbb{E}_{s' \sim p(\cdot|s,a)}[\mathbb{I}(\phi(s') = g)]$  and recovers it when  $\alpha \rightarrow \infty$ . Starting with the true GCRL  
 515 objective:

$$\alpha J(\pi_g) = \mathbb{E}_{d^{\pi_g}} [\alpha R(s, a, g)] \quad (15)$$

$$= \mathbb{E}_{d^{\pi_g}} \left[ \log e^{\alpha R(s,a,g)} \right] \quad (16)$$

$$= \mathbb{E}_{d^{\pi_g}} \left[ \log \left( \frac{e^{\alpha R(s,a,g)}}{Z} \frac{d^{\pi_g}(s, a, g)}{d^{\pi_g}(s, a, g)} Z \right) \right] \quad (17)$$

$$= \mathbb{E}_{d^{\pi_g}} \left[ \log \left( \frac{q(s, a, g)}{d^{\pi_g}(s, a, g)} Z \right) \right] + \mathbb{E}_{d^{\pi_g}} [\log d^{\pi_g}] \quad (18)$$

$$= -D_{KL}(d^{\pi_g}(s, a, g) \| q(s, a, g)) - \mathcal{H}(d^{\pi_g}) + \log(Z) \quad (19)$$

516 Rearranging terms we get:

$$J(\pi_g) + \frac{1}{\alpha} \mathcal{H}(d^{\pi_g}) = -\frac{1}{\alpha} D_{KL}(d^{\pi_g}(s, a, g) \| q(s, a, g)) + C \quad (20)$$

517 For any  $f$ -divergence that upper bounds the KL divergence we have:

$$J(\pi_g) + \frac{1}{\alpha} \mathcal{H}(d^{\pi_g}) = -\frac{1}{\alpha} D_{KL}(d^{\pi_g}(s, a, g) \| q(s, a, g)) + C \geq -\frac{1}{\alpha} \mathcal{D}_f(d^{\pi_g}(s, a, g) \| q(s, a, g)) + C \quad (21)$$

518  $\square$

519 **A (dataset) regularized GCRL objective:** Define a regularized objective for GCRL as follows:

$$J_{\text{offline}}(\pi) = \alpha_1 \mathbb{E}_{d^\pi} [e^{r(s,a,g)}] + \alpha_2 \mathbb{E}_{d^\pi(s,a,g)} [\rho(s, a, g)]. \quad (22)$$

520 The above offline objective mimics the classical offline RL objective [37, 24] in constraining the  
 521 visitation of the learned policy, as the second objective is minimized when  $d^\pi(s, a, g) = \rho(s, a, g)$ .  
 522 Also, a constraint of  $\mathbb{E}_{d^\pi(s,a,g)} [\rho(s, a, g)] > 1 - \delta$  implies that  $d^\pi(s, a, g)$  has atleast  $1 - \delta$  coverage  
 523 of the offline data distribution.

524 Proposition 2.1 derives the connection between the offline GCRL objective and SMORe:

525 **Proposition 2.1.** Consider a stochastic MDP, a stochastic policy  $\pi$ , and a sparse reward function  
526  $r(s, a, g) = \mathbb{E}_{s' \sim p(\cdot|s,a)} [\mathbb{I}(\phi(s') = g, q^{train}(g) > 0)]$  where  $\mathbb{I}$  is an indicator function, define a soft  
527 goal transition distribution to be  $q(s, a, g) \propto \exp(\alpha r(s, a, g))$  the following bounds hold for any  
528  $f$ -divergence that upper bounds KL-divergence (eg.  $\chi^2$ , Jensen-Shannon):

$$\log J_{offline}(\pi_g) + \mathcal{H}(\text{Mix}_\beta(d, \rho)(s, a, g)) + C \geq -\mathcal{D}_f(\text{Mix}_\beta(d, \rho)(s, a, g) \| \text{Mix}_\beta(q, \rho)(s, a, g)), \quad (23)$$

529 where  $\mathcal{H}$  denotes the entropy,  $\alpha$  is a temperature parameter,  $\alpha_1 = \beta^2$ ,  $\alpha_2 = \beta(1 - \beta)Z$  and  $C$  is a  
530 positive constant. Furthermore, the bound is tight when  $f$  is the KL-divergence.

531 *Proof.* We first consider the following two objectives for GCRL and show that they are equivalent.  
532 This reduction will later help in proving a connection to mixture occupancy matching. We consider  
533  $\alpha = 1$  w.l.o.g. Here are two objectives we consider:

$$J(\pi) = \mathbb{E}_{d^\pi} [r(s, a, g)] \quad (24)$$

534

$$J'(\pi) = \mathbb{E}_{d^\pi} [e^{r(s,a,g)}] \quad (25)$$

535 In GCRL reward functions are sparse and binary. We show the equivalence of first two objectives in  
536 find the optimal goal conditioned policy via two arguments. First, notice that the rewards for goal  
537 transition states for objective  $J'(\pi)$  is  $e$  and 1 for all other transitions. This is in contrast to  $J(\pi)$   
538 which considers a reward function 1 at goal transitions states and 0 otherwise. Under our assumption  
539 of infinite horizon discounted MDP, we can translate the rewards while keeping the optimal policy  
540 same in MDP considered by  $J'(\pi)$  to  $e - 1$  at goal transitions states and 0 otherwise. Further we can  
541 scale the rewards by  $1/(e - 1)$  and recover an MDP with same optimal policy that has reward of 1  
542 at goal-transition states and 0 otherwise. This concludes the equivalence of maximizing  $J'(\pi)$  as an  
543 alternative to  $J(\pi)$  while recovering the same optimal policy.

544 We now consider a regularized (pessimistic/offline) GCRL problem with the shifted reward functions  
545  $e^{r(s,a,g)}$  that maximizes the reward while ensuring the policy visitation stays close to offline data  
546 visitation in  $\chi^2$  divergence.

$$J_{offline}(\pi) = \alpha_1 \mathbb{E}_{d^\pi} [e^{r(s,a,g)}] + \alpha_2 \mathbb{E}_{d^\pi(s,a,g)} [\rho(s, a, g)]. \quad (26)$$

547 With a particular instantiation of hyperparameters we show that the  $J_{offline}(\pi)$  objective can be  
548 simplified to an equivalent objective  $J'_{offline}(\pi)$  by setting  $\alpha_1 = \beta^2$  and  $\alpha_2 = \beta(1 - \beta)Z$  where  $Z$   
549 is the partition function for  $e^{r(s,a,g)}$  over entire  $\mathcal{S} \times \mathcal{A} \times \mathcal{G}$ .

$$J'_{offline}(\pi) = \mathbb{E}_{\text{Mix}_\beta(d,\rho)(s,a,g)} [\beta e^{r(s,a,g)} + (1 - \beta)\rho(s, a, g) \cdot Z] \quad (27)$$

$$J'_{offline}(\pi) = \mathbb{E}_{\text{Mix}_\beta(d,\rho)(s,a,g)} [\beta e^{r(s,a,g)} + (1 - \beta)\rho(s, a, g) \cdot Z] \quad (28)$$

$$= \beta^2 \mathbb{E}_{d^\pi} [e^{r(s,a,g)}] + \beta(1 - \beta)Z \mathbb{E}_{d^\pi} [\rho(s, a, g)] \quad (29)$$

$$+ (1 - \beta) \mathbb{E}_{d^\pi} [\beta e^{r(s,a,g)} + (1 - \beta)\rho(s, a, g) \cdot Z] \beta \quad (30)$$

550

$$= \beta^2 \mathbb{E}_{d^\pi} [e^{r(s,a,g)}] + \beta(1 - \beta)Z \mathbb{E}_{d^\pi} [\rho(s, a, g)] + C' \quad (32)$$

$$= J_{offline}(\pi) + C' \quad (33)$$

551 Now that we have shown  $J'_{offline}(\pi) \equiv J_{offline}(\pi)$  and hence solving the same optimization  
552 problem, we proceed to derive connections with mixture occupancy matching which follows through  
553 an application of Jensen's inequality:

$$\log J'_{offline}(\pi) = \log \mathbb{E}_{\text{Mix}_\beta(d, \rho)(s, a, g)} \left[ \beta e^{r(s, a, g)} + (1 - \beta) \rho(s, a, g) \cdot Z \right] \quad (34)$$

$$\geq \mathbb{E}_{\text{Mix}_\beta(d, \rho)(s, a, g)} \left[ \log(\beta e^{r(s, a, g)} + (1 - \beta) \rho(s, a, g) \cdot Z) \right] \quad (35)$$

$$= \mathbb{E}_{\text{Mix}_\beta(d, \rho)(s, a, g)} [\log(\beta q(s, a, g) + (1 - \beta) \rho(s, a, g))] + \log Z \quad (36)$$

$$= -D_{KL}[\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)] - \mathcal{H}(\text{Mix}_\beta(d, \rho)(s, a, g)) + \log Z \quad (37)$$

$$= -D_{KL}[\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)] - \mathcal{H}(\text{Mix}_\beta(d, \rho)(s, a, g)) + \log Z \quad (38)$$

For any  $f$ -divergence that upperbounds the KL divergence since  $Z \geq 1$  we have:

$$\log J'_{offline}(\pi) + \frac{1}{\alpha} \mathcal{H}(\text{Mix}_\beta(d, \rho)(s, a, g)) \geq -\frac{1}{\alpha} D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (39)$$

Further simplifying using Eq 33:

$$\log J_{offline}(\pi) + \frac{1}{\alpha} \mathcal{H}(\text{Mix}_\beta(d, \rho)(s, a, g)) + C \geq -\frac{1}{\alpha} D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (40)$$

□

Optimizing the mixture distribution matching objective of SMORe maximizes a variant of *offline* GCRL objective where the entropy for distribution  $\text{Mix}_\beta(d, \rho)(s, a, g)$  is jointly maximized. Therefore we have shown that the minimizing discrepancy of mixture distribution occupancy maximizes a lower bounds to an offline variant of maxent GCRL objective.

## B.2 Convex Conjugates and $f$ -divergences

We first review the basics of duality in reinforcement learning. Let  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  be a convex function. The convex conjugate  $f^* : \mathbb{R}_+ \rightarrow \mathbb{R}$  of  $f$  is defined by:

$$f^*(y) = \sup_{x \in \mathbb{R}_+} [xy - f(x)]. \quad (41)$$

The convex conjugates have the important property that  $f^*$  is also convex and the convex conjugate of  $f^*$  retrieves back the original function  $f$ . We also note an important relation regarding  $f$  and  $f^*$ :  $(f^*)' = (f')^{-1}$ , where the  $'$  notation denotes first derivative.

Going forward, we would be dealing extensively with  $f$ -divergences. Informally,  $f$ -divergences [30] are a measure of distance between two probability distributions. Here's a more formal definition:

Let  $P$  and  $Q$  be two probability distributions over a space  $\mathcal{Z}$  such that  $P$  is absolutely continuous with respect to  $Q$ <sup>4</sup>. For a function  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  that is a convex lower semi-continuous and  $f(1) = 0$ , the  $f$ -divergence of  $P$  from  $Q$  is

$$D_f(P \parallel Q) = \mathbb{E}_{z \sim Q} \left[ f \left( \frac{P(z)}{Q(z)} \right) \right]. \quad (42)$$

Table 4 lists some common  $f$ -divergences with their generator functions  $f$  and the conjugate functions  $f^*$ .

## B.3 SMORe: Dual objective for Offline Goal conditioned reinforcement learning

In this section, we derive the dual objective for solving the multi-task occupancy problem formulation for GCRL. First, we derive the original variant of SMORe for the GCRL problem and later derive the action-free SMORe variant for the interested readers.

**Theorem 2.** *The dual problem to the primal occupancy matching objective (Equation 9) is given by:*

$$\max_{\pi_g} \min_S \beta(1 - \gamma) \mathbb{E}_{d_0, \pi_g} [S(s, a, g)] + \mathbb{E}_{\text{Mix}_\beta(d, \rho)(s, a, g)} [f^*(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \quad (10)$$

$$- (1 - \beta) \mathbb{E}_\rho [\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)],$$

where  $f^*$  is conjugate function of  $f$  and  $S$  is the Lagrange dual variable defined as  $S : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow \mathbb{R}$ . Moreover, as strong duality holds from Slater's conditions the primal and dual share the same optimal solution  $\pi_g^*$  for any offline transition distribution  $\rho$ .

<sup>4</sup>Let  $z$  denote the random variable. For any measurable set  $Z \subseteq \mathcal{Z}$ ,  $Q(z \in Z) = 0$  implies  $P(z \in Z) = 0$ .



Divergence Name	Generator $f(x)$	Conjugate $f^*(y)$
KL (Reverse)	$x \log x$	$e^{(y-1)}$
Squared Hellinger	$(\sqrt{x} - 1)^2$	$\frac{y}{1-y}$
Pearson $\chi^2$	$(x - 1)^2$	$y + \frac{y^2}{4}$
Total Variation	$\frac{1}{2} x - 1 $	$y$ if $y \in [-\frac{1}{2}, \frac{1}{2}]$ otherwise $\infty$
Jensen-Shannon	$-(x + 1) \log(\frac{x+1}{2}) + x \log x$	$-\log(2 - e^y)$

Table 4: List of common  $f$ -divergences.

583 *Proof.* Recall that:  $\text{Mix}_\beta(d, \rho)(s, a, g) := \beta d(s, a, g) + (1 - \beta)\rho(s, a, g)$  and  $\text{Mix}_\beta(q, \rho)(s, a, g) :=$   
584  $\beta q(s, a, g) + (1 - \beta)\rho(s, a, g)$ .  $\text{Mix}_\beta(d, \rho)(s, a, g)$  denotes the mixture between the current  
585 agent’s joint-goal visitation distribution with an offline transition dataset potentially suboptimal  
586 and  $\text{Mix}_\beta(q, \rho)(s, a, g)$  is the mixture between the expert’s visitation distribution with arbitrary  
587 experience from the offline transition dataset. Minimizing the divergence between these visitation  
588 distributions still solves the occupancy problem, i.e  $d^{\pi_g} = q$  when  $q$  is achievable. We start with the  
589 primal formulation from Eq 9 for mixture divergence regularization:

$$\begin{aligned} & \max_{d(s,a,g) \geq 0, \pi(a|s)} -D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \\ \text{s.t. } & d(s, a, g) = (1 - \gamma)\rho_0(s, g) \cdot \pi(a|s, g) + \gamma \pi(a|s, g) \sum_{s', a'} d(s', a', g) p(s|s', a'). \end{aligned}$$

590 Applying Lagrangian duality and convex conjugate (41) to this problem, we can convert it to an  
591 unconstrained problem with dual variables  $S(s, a, g)$  defined for all  $s, a \in \mathcal{S} \times \mathcal{A} \times \mathcal{G}$ :

$$\begin{aligned} & \max_{\pi, d \geq 0} \min_{S(s,a,g)} -D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \\ & + \sum_{s,a,g} S(s, a, g) \left( (1 - \gamma)d_0(s, g) \cdot \pi(a|s, g) + \gamma \sum_{s', a'} d(s', a', g) p(s|s', a') \pi(a|s, g) - d(s, a, g) \right) \end{aligned} \quad (43)$$

$$\begin{aligned} & = \max_{\pi, d \geq 0} \min_{S(s,a,g)} (1 - \gamma) \mathbb{E}_{d_0(s,g), \pi(a|s,g)} [S(s, a, g)] \\ & + \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s', a'} p(s'|s, a) \pi(a'|s') S(s', a', g) - S(s, a, g) \right] \end{aligned} \quad (44)$$

$$- D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (45)$$

592

593

$$\begin{aligned} & = \max_{\pi, d \geq 0} \min_{S(s,a,g)} \beta (1 - \gamma) \mathbb{E}_{d_0(s,g), \pi(a|s,g)} [S(s, a, g)] \\ & + \beta \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s', a'} p(s'|s, a) \pi(a'|s') S(s', a', g) - S(s, a, g) \right] \\ & + (1 - \beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s', a'} p(s'|s, a) \pi(a'|s') S(s', a', g) - S(s, a, g) \right] \\ & - (1 - \beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s', a'} p(s'|s, a) \pi(a'|s', g) S(s', a', g) - S(s, a, g) \right] \end{aligned} \quad (46)$$

$$- D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (47)$$

594 Now using the fact that strong duality holds in this problem we can swap the inner max and min  
 595 resulting in:

$$\begin{aligned}
 &= \max_{\pi} \min_{S(s,a,g)} \max_{\text{Mix}_{\beta}(d,\rho)(s,a,g) \geq 0} \beta(1-\gamma) \mathbb{E}_{d_0(s,g), \pi(a|s,g)} [S(s,a,g)] \\
 &+ \beta \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s',a'} p(s'|s,a) \pi(a'|s') S(s',a',g) - S(s,a,g) \right] \\
 &+ (1-\beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s',a'} p(s'|s,a) \pi(a'|s') S(s',a',g) - S(s,a,g) \right] \\
 &- (1-\beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s',a'} p(s'|s,a) \pi(a'|s',g) S(s',a',g) - S(s,a,g) \right] \tag{48}
 \end{aligned}$$

$$- D_f(\text{Mix}_{\beta}(d,\rho)(s,a,g) \parallel \text{Mix}_{\beta}(q,\rho)(s,a,g)) \tag{49}$$

$$\tag{50}$$

596 We can now apply the convex conjugate (Eq. (41)) definition to obtain a closed form for the inner  
 597 maximization problem simplifying to:

$$\begin{aligned}
 &\max_{\pi(a|s,g)} \min_{S(s,a,g)} \beta(1-\gamma) \mathbb{E}_{d_0(s,g), \pi(a|s,g)} [S(s,a,g)] \\
 &+ \mathbb{E}_{s,a,g \sim \text{Mix}_{\beta}(q,\rho)(s,a,g)} \left[ f^* \left( \gamma \sum_{s',a'} p(s'|s,a,g) \pi(a'|s') S(s',a',g) - S(s,a,g) \right) \right] \\
 &- (1-\beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s',a'} p(s'|s,a,g) \pi(a'|s') S(s',a',g) - S(s,a,g) \right] \tag{51}
 \end{aligned}$$

598 This completes our derivation of the SMORe objective. Since strong duality holds (objective convex,  
 599 constraints linear and feasible), SMORe and the primal mixture occupancy matching share the same  
 600 global optima  $\pi_g^*$ .  $\square$

#### 601 B.4 Action-free SMORe: Dual-V objective for offline goal conditioned reinforcement learning

602 The primal problem in Equation 9 is over-constrained. The objective determines the visitation  
 603 distribution  $d$  uniquely under a fixed policy. It turns out we can further relax this constraint to get an  
 604 objective that results in the same optimal solution [1]  $\pi_g^*$  by rewriting our primal formulation as:

$$\begin{aligned}
 &\max_{d(s,a,g) \geq 0} -D_f(\text{Mix}_{\beta}(d,\rho)(s,a,g) \parallel \text{Mix}_{\beta}(q,\rho)(s,a,g)) \\
 &\text{s.t. } \sum_a d(s,a,g) = (1-\gamma)\rho_0(s,g) + \gamma \sum_{s',a'} d(s',a',g)p(s|s',a'). \tag{52}
 \end{aligned}$$

605 **Theorem 3.** Let  $y(s,a,g) = \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} [S(s',g)] - S(s,g)$ . The action-free dual problem to the  
 606 multi-task mixture occupancy matching objective (Equation 52) is given by:

$$\begin{aligned}
 &\min_{S(s,g)} \beta(1-\gamma) \mathbb{E}_{d_0(s,g)} [S(s,g)] \\
 &+ \mathbb{E}_{s,a,g \sim \text{Mix}_{\beta}(q,\rho)(s,a,g)} \left[ \max(0, (f')^{-1}(y(s,a,g))) y(s,a,g) - f(\max(0, (f')^{-1}(y(s,a,g)))) \right] \\
 &- (1-\beta) \mathbb{E}_{s,a,g \sim \rho} \left[ \gamma \sum_{s'} p(s'|s,a) S(s',g) - S(s,g) \right]
 \end{aligned}$$

607 where  $S$  is the lagrange dual variable defined as  $S : \mathcal{S} \times \mathcal{G} \rightarrow \mathbb{R}$ . Moreover, strong duality holds  
 608 from Slater's conditions the primal and dual share the same optimal solution  $\pi_g^*$  for any offline  
 609 transition distribution  $d^O$ .

610 *Proof.* Proceeding as before and applying Lagrangian duality and convex conjugate (41) to this  
 611 problem, we can convert it to an unconstrained problem with dual variables  $S(s,g)$  defined for all  
 612  $s, g \in \mathcal{S} \times \mathcal{G}$ :

$$\begin{aligned} & \max_{d \geq 0} \min_{S(s,g)} -D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \\ & + \sum_{s,g} S(s, g) \left( (1 - \gamma)d_0(s, g) + \gamma \sum_{s', a', g} d(s', a', g)p(s|s', a', g) - \sum_a d(s, a, g) \right) \end{aligned} \quad (53)$$

$$\begin{aligned} & = \max_{d \geq 0} \min_{S(s,g)} (1 - \gamma)\mathbb{E}_{d_0(s,g)}[S(s, g)] \\ & + \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s'} p(s'|s, a)\pi(a'|s')S(s', g) - S(s, g) \right] \end{aligned} \quad (54)$$

$$- D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (55)$$

613

614

$$\begin{aligned} & = \max_{d \geq 0} \min_{S(s,g)} \beta(1 - \gamma)\mathbb{E}_{d_0(s,g)}[S(s, g)] \\ & + \beta \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \\ & + (1 - \beta)\mathbb{E}_{s,a,g \sim d^0} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \\ & - (1 - \beta)\mathbb{E}_{s,a,g \sim d^0} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \end{aligned} \quad (56)$$

$$- D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (57)$$

615 Now using the fact that strong duality holds in this problem we can swap the inner max and min  
616 resulting in:

$$\begin{aligned} & = \min_{S(s,g)} \max_{\text{Mix}_\beta(d, \rho)(s, a, g) \geq 0} \beta(1 - \gamma)\mathbb{E}_{d_0(s,g)}[S(s, g)] \\ & + \beta \mathbb{E}_{s,a,g \sim d} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \\ & + (1 - \beta)\mathbb{E}_{s,a,g \sim d^0} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \\ & - (1 - \beta)\mathbb{E}_{s,a,g \sim d^0} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \end{aligned} \quad (58)$$

$$- D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \quad (59)$$

617 Unlike previous case where constraints uniquely define a valid  $d$  for any given  $\pi$ , in this case we  
618 need to take into account the hidden constraint  $d \geq 0$  or equivalently  $\text{Mix}_\beta(d, \rho)(s, a, g) \geq 0$ .  
619 To incorporate the non-negativity constraints we consider the inner maximization separately and  
620 derive a closed-form solution that adheres to the non-negativity constraints. Recall  $y(s, a, g) =$   
621  $\mathbb{E}_{s' \sim p(s,a)}[S(s', g)] - S(s, g)$ .

$$\begin{aligned} & \max_{\text{Mix}_\beta(d, \rho)(s, a, g) \geq 0} \mathbb{E}_{s,a,g \sim \text{Mix}_\beta(d, \rho)(s, a, g)} \left[ \gamma \sum_{s'} p(s'|s, a)S(s', g) - S(s, g) \right] \\ & - D_f(\text{Mix}_\beta(d, \rho)(s, a, g) \parallel \text{Mix}_\beta(q, \rho)(s, a, g)) \end{aligned}$$

622 We can now construct the Lagrangian dual to incorporate the constraint  $\text{Mix}_\beta(d, \rho)(s, a, g) \geq 0$  in its  
623 equivalent form  $w(s, a, g) \geq 0$  and obtain the following where  $w \triangleq \frac{\text{Mix}_\beta(d, \rho)(s, a, g)}{\text{Mix}_\beta(q, \rho)(s, a, g)}$ :

$$\max_{w(s,a,g)} \max_{\lambda \geq 0} \mathbb{E}_{s,a \sim \text{Mix}_\beta(q,\rho)(s,a,g)} [w(s,a,g)y(s,a,g)] - \mathbb{E}_{\text{Mix}_\beta(q,\rho)(s,a,g)} [f(w(s,a,g))] + \sum_{s,a,g} \lambda(w(s,a,g) - 0) \quad (60)$$

624 Since strong duality holds, we can use the KKT constraints to find the solutions  $w^*(s, a, g)$  and  
625  $\lambda^*(s, a, g)$ .

- 626 1. **Primal feasibility:**  $w^*(s, a, g) \geq 0 \quad \forall s, a$
- 627 2. **Dual feasibility:**  $\lambda^* \geq 0 \quad \forall s, a$
- 628 3. **Stationarity:**  $\text{Mix}_\beta(q, \rho)(s, a, g)(-f'(w^*(s, a, g)) + y(s, a, g) + \lambda^*(s, a, g)) = 0 \quad \forall s, a$
- 629 4. **Complementary Slackness:**  $(w^*(s, a, g) - 0)\lambda^*(s, a, g) = 0 \quad \forall s, a$

630 Using stationarity we have the following:

$$f'(w^*(s, a, g)) = y(s, a, g) + \lambda^*(s, a, g) \quad \forall s, a, g \quad (61)$$

631 Now using complementary slackness, only two cases are possible  $w^*(s, a, g) \geq 0$  or  $\lambda^*(s, a, g) \geq 0$ .  
632 Combining both cases we arrive at the following solution for this constrained optimization:

$$w^*(s, a) = \max\left(0, f'^{-1}(y(s, a, g))\right) \quad (62)$$

633 Using the optimal closed-form solution ( $w^*$ ) for  $\text{Mix}_\beta(d, \rho)(s, a, g)$  of the inner optimization in  
634 Eq. (58) we obtain

$$\begin{aligned} & \min_{S(s,a)} \beta(1 - \gamma)\mathbb{E}_{d_0(s)}[S(s, g)] \\ & + \mathbb{E}_{s,a,g \sim \text{Mix}_\beta(q,\rho)(s,a,g)} \left[ \max\left(0, (f')^{-1}(y(s, a, g))\right) y(s, a, g) - \alpha f\left(\max\left(0, (f')^{-1}(y(s, a, g))\right)\right) \right] \\ & - (1 - \alpha)\mathbb{E}_{s,a \sim \rho} \left[ \gamma \sum_{s'} p(s'|s, a)\pi(a'|s')S(s', g) - S(s, g) \right] \end{aligned} \quad (63)$$

635 For deterministic dynamics, this reduces to the action-free SMORe objective:

$$\begin{aligned} & \min_{S(s,a)} \beta(1 - \gamma)\mathbb{E}_{d_0(s)}[S(s, g)] \\ & + \mathbb{E}_{s,a \sim \text{Mix}_\beta(q,\rho)(s,a,g)} \left[ \max\left(0, (f')^{-1}(y(s, a, g))\right) y(s, a, g) - f\left(\max\left(0, (f')^{-1}(y(s, a, g))\right)\right) \right] \\ & - (1 - \beta)\mathbb{E}_{s,a \sim \rho} [\gamma S(s', g) - S(s, g)] \end{aligned} \quad (64)$$

636 where  $y(s, a, g) = \gamma S(s', g) - S(s, g)$ .

637 Note that we no longer need actions in the offline dataset to learn an optimal goal conditioned score  
638 function. This score function can be used to learn presentation in action-free datasets as well as for  
639 transfer of value function across differing action-modalities where agents share the same observation  
640 space (eg. images as observations).

641

□

## 642 C SMORe algorithmic details

### 643 C.1 SMORe with common $f$ -divergences

#### 644 a. KL divergence

645 We consider the reverse KL divergence and start with the general SMORe objective:

$$\begin{aligned} & \max_{\pi_g} \min_S \beta(1 - \gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + \mathbb{E}_{s,a,g \sim \text{Mix}_\beta(q,\rho)(s,a,g)} [f^*(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \\ & - (1 - \beta)\mathbb{E}_{s,a,g \sim \rho} [\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)] \end{aligned} \quad (65)$$

646 Plugging in the conjugate  $f^*$  for reverse KL divergence we get:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} \left[ e^{(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))} \right] \\ - (1-\beta)\mathbb{E}_{s, a, g \sim \rho}[\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)] \end{aligned} \quad (66)$$

647 Using the telescoping sum for the last term in the objective above, we can simplify it as follows:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} \left[ e^{(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))} \right] \\ + (1-\beta)\mathbb{E}_{s, g \sim d_0, a \sim \rho(\cdot|s, g)}[S(s, a, g)] \end{aligned} \quad (67)$$

648 With the initial state distribution  $d_0$  set to the offline dataset distribution  $\rho$ , and Since our initial state  
649 distribution is the same as offline data distribution, we get:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{\rho, \pi_g}[S(s, a, g)] + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} \left[ e^{(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))} \right] \\ + (1-\beta)\mathbb{E}_{\rho}[S(s, a, g)] \end{aligned} \quad (68)$$

650 Collecting terms together we get:

$$\begin{aligned} \max_{\pi_g} \min_Q \mathbb{E}_{\rho}[\mathbb{E}_{a \sim \pi}[\beta(1-\gamma)S(s, a, g)] + \mathbb{E}_{a \sim \rho}[(1-\beta)S(s, a, g)]] \\ + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} \left[ e^{(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))} \right] \end{aligned} \quad (69)$$

651 The objective for SMORe with reverse KL divergence pushes down the "score" of offline dataset  
652 transitions selectively (without pushing down score of the goal-transition distribution) while  
653 minimizing the term resembling bellman regularization that also encourages increasing score at  
654 the mixture dataset jointly over the offline dataset as well as the goal transition distribution.

## 655 b. Pearson chi-squared divergence

656 We consider the Pearson  $\chi^2$  and start with the general SMORe objective:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} [f^*(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \\ - (1-\beta)\mathbb{E}_{s, a, g \sim \rho}[\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)] \end{aligned} \quad (70)$$

657 With the initial state distribution  $d_0$  set to the offline dataset distribution  $\rho$ , and plugging in the  
658 conjugate  $f^*$  for Pearson  $\chi^2$  divergence we get:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + 0.25\mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} [(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))^2] \\ + \mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} [(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] - (1-\beta)\mathbb{E}_{s, a, g \sim \rho}[\gamma P^{\pi_g} S(s, a, g) - S(s, a, g)] \end{aligned} \quad (71)$$

659 Using the fact that  $\text{Mix}_{\beta}(q, \rho)(s, a, g) = \beta q(s, a, g) + (1-\beta)\rho(s, a, g)$ , we can further simplify the  
660 above equation to:

$$\begin{aligned} \max_{\pi_g} \min_S \beta(1-\gamma)\mathbb{E}_{d_0, \pi_g}[S(s, a, g)] + 0.25\mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)} [(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))^2] \\ + \beta\mathbb{E}_{s, a, g \sim q}[(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))] \end{aligned} \quad (72)$$

661 Collecting terms together we get:

$$\begin{aligned} \max_{\pi_g} \min_S & \beta(1 - \gamma)\mathbb{E}_{\rho, \pi_g}[S(s, a, g)] + \beta\mathbb{E}_{s, g \sim q, a \sim \pi_g}[\gamma P^{\pi_g} S(s, a, g)] \\ & - \beta\mathbb{E}_{s, a, g \sim q}[S(s, a, g)] + 0.25\mathbb{E}_{s, a, g \sim \text{Mix}_{\beta}(q, \rho)(s, a, g)}[(\gamma P^{\pi_g} S(s, a, g) - S(s, a, g))^2] \end{aligned} \quad (73)$$

662 Observing the equation above, we note that the first two terms decrease score at offline data distribution  
 663 as well as the goal transition distribution when actions are sampled according to the policy  $\pi_g$ .  
 664 Simultaneously the third term pushes score up for the  $\{s, a, g\}$  tuples that are sampled from goal  
 665 transition distribution. Finally the last term encouraged enforces a bellman regularization enforcing  
 666 smoothness is the scores of neighbouring states.

## 667 D SMORe experimental details

668 **Environments:** For the offline GCRL experiments we consider the benchmark used in prior work  
 669 GoFar and extend it with locomotion tasks. For the manipulations tasks we consider the Fetch  
 670 environment and a dextrous shadow hand environment. Fetch environments [28] consists of a  
 671 manipulator with seven degrees of freedom along with a parallel gripper. The set of environments  
 672 get a sparse reward of 1 when the goal is within 5 cm and 0 otherwise. The action space is 4  
 673 dimensional (3 dimension cartesian control + 1 dimension gripper control). The shadow hand is  
 674 24 DOF manipulator with 20-dimensional action space. The goal is 15-dimension specifying the  
 675 position for each of the five fingers. The tolerance for goal reaching is 1 cm. For the locomotion  
 676 environments, the task is to achieve a particular velocity in the x direction and stay at the velocity. For  
 677 HalfCheetah, the target velocity is set to 11.0 and for Ant the target velocity is 5.0. For locomotion  
 678 environments, the tolerance for goal reaching if 0.5. The MuJoCo environments used in this work are  
 679 [licensed under CC BY 4.0](#).

680 **Offline Datasets:** We use existing datasets from the offline GCRL benchmark used in [21] for all  
 681 manipulation tasks except Reacher, SawyerReach, and SawyerDoor. For Reacher, SawyerReach, and  
 682 SawyerDoor we use existing datasets from [39]. These datasets are comprised on  $x\%$  random data  
 683 and  $(100-x)\%$  expert data depending on the coverage over goals reached in individual datasets. We  
 684 create our own datasets for locomotion by using 'random/medium/medium-replay' data as our offline  
 685 (suboptimal) data combined with 30 trajectories from corresponding 'expert' datasets. The datasets  
 686 used from D4RL are [licensed under Apache 2.0](#).

687 **Baselines:** To benchmark and analyze the performance of our proposed methods for offline imitation  
 688 learning with suboptimal data, we consider the following representative baselines in this work: GoFAR  
 689 [21], WGCSL [39], GCSL [16], and Actionable Models [4], Contrastive RL [8] and GC-IQL [18].  
 690 GoFAR is a dual occupancy matching approach to GCRL that formulates it as a weighted regression  
 691 problem. WGCSL and GSCL use goal-conditioned behavior cloning with goal relabelling as the  
 692 base algorithms and WGCL uses weights to learn improved policy over GCSL. Actionable models  
 693 uses conservative learning with goal chaining to learn goal-reaching behaviours using offline datasets.  
 694 Contrastive RL treats GCRL as a classification problem - contrastive goals that are achieved in  
 695 trajectory from random goals. Finally, GC-IQL extends the single task offline RL algorithm IQL to  
 696 GCRL.

697 The open-source implementations of the baselines GoFAR, WGCSL, GCSL, Actionable models,  
 698 Contrastive RL and IQL are provided by the authors [21] and employed in our experiments. We use  
 699 the hyperparameters provided by the authors, which are consistent with those used in the original  
 700 GoFAR paper, for all the MuJoCo locomotion and manipulation environments. We implement  
 701 contrastive learning using the code from [Contrastive RL repository](#). GC-IQL is implemented using  
 702 code from author's implementation [found here](#).

703 **Architecture and Hyperparameters** For the baselines, we use tuned hyperparameters from previous  
 704 works that were tuned on the same set of tasks and datasets. Implementation for SMORe shares the  
 705 same network architecture as baselines. GoFAR additionally requires training a discriminator. For all  
 706 experiments, all methods are trained for 10 seeds with each training run. Fetch manipulation tasks  
 707 are trained for 400k minibatch updates of size 512 whereas all other environments training is done  
 708 for 1M minibatch updates. The architectures and hyperparameters for all methods are reported in  
 709 Table 5.

Hyperparameter	Value
Policy updates $n_{pol}$	1
Policy learning rate	3e-4
Value learning rate	3e-4
MLP layers	(256,256)
LR decay schedule	cosine
Discount factor	0.99
LR decay schedule	cosine
Batch Size	512
Mixture ratio $\beta$	0.5
Expectile $\tau$	[0.65,0.7,0.8,0.85]

Table 5: Hyperparameters for SMORe.

Task	Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	WGCSL	GCSL		AM	IQL
Reacher	15.30 $\pm$ 0.58	14.01 $\pm$ 0.36	16.62 $\pm$ 2.09	23.68 $\pm$ 0.58	8.86 $\pm$ 0.61
SawyerReach	14.06 $\pm$ 0.08	12.05 $\pm$ 1.23	23.03 $\pm$ 1.17	23.37 $\pm$ 2.29	36.19 $\pm$ 0.01
SawyerDoor	16.79 $\pm$ 0.75	18.29 $\pm$ 0.94	12.26 $\pm$ 3.94	16.63 $\pm$ 0.76	29.31 $\pm$ 0.88
FetchPick	6.87 $\pm$ 0.77	6.54 $\pm$ 1.85	0.21 $\pm$ 0.29	0.45 $\pm$ 0.32	15.24 $\pm$ 1.27
FetchPush	10.62 $\pm$ 0.98	12.38 $\pm$ 1.10	3.60 $\pm$ 0.59	2.74 $\pm$ 0.70	19.95 $\pm$ 1.94
FetchSlide	2.62 $\pm$ 1.15	2.03 $\pm$ 0.01	0.41 $\pm$ 0.03	0.31 $\pm$ 0.31	3.25 $\pm$ 1.02

Table 6: Discounted Return for the offline GCRL benchmark with 5% expert data. Results are averaged over 10 seeds.

## 710 E Additional experiments

### 711 E.1 Results on offline GCRL benchmark with varying expert coverage in offline dataset

712 We ablate the effect of dataset quality on the performance of an offline GCRL method in this sections.  
713 Table 6, 7, 8 show performance of all methods with 5%, 2.5% and 1% expert data in the offline  
714 dataset respectively.

### 715 E.2 Success Rate and Final distance to goal on Manipulation tasks

716 Table 10 and Table 11 reports the success rate and final distance to goal metrics on manipulation  
717 tasks.

### 718 E.3 Robustness of mixture distribution parameter $\beta$

719 We find that SMORe is quite robust to the mixture distribution parameter  $\beta$  except in the environment  
720 FetchPush where  $\beta = 0.5$  is the most peformant. Table 9 shows this result empirically.

721

Task	Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	WGCSL	GCSL		AM	IQL
Reacher	13.03 $\pm$ 0.56	12.17 $\pm$ 0.8	19.63 $\pm$ 3.09	24.78 $\pm$ 0.23	4.44 $\pm$ 0.70
SawyerReach	11.455 $\pm$ 1.37	11.34 $\pm$ 1.18	25.35 $\pm$ 0.8	25.19 $\pm$ 0.61	35.73 $\pm$ 0.22
SawyerDoor	16.79 $\pm$ 0.29	13.20 $\pm$ 0.53	14.78 $\pm$ 5.29	16.59 $\pm$ 1.39	16.87 $\pm$ 4.21
FetchPick	4.39 $\pm$ 1.35	4.99 $\pm$ 0.11	0.21 $\pm$ 0.29	0.24 $\pm$ 0.27	11.79 $\pm$ 1.78
FetchPush	8.01 $\pm$ 1.96	8.04 $\pm$ 0.34	3.60 $\pm$ 0.59	2.02 $\pm$ 0.48	19.66 $\pm$ 1.69
FetchSlide	2.33 $\pm$ 0.23	2.37 $\pm$ 0.83	0.44 $\pm$ 0.016	0.45 $\pm$ 0.44	1.83 $\pm$ 1.31

Table 7: Discounted Return for the offline GCRL benchmark with 2.5% expert data. Results are averaged over 10 seeds.

Task	Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	WGCSL	GCSL		AM	IQL
Reacher	13.56±0.69	12.27 ±1.45	17.94±3.71	24.89±0.34	4.28 ± 0.92
SawyerReach	10.71 ±0.69	11.79±1.46	25.61±0.39	25.54±0.95	31.31 ± 2.08
SawyerDoor	15.18 ±0.81	11.89±1.51	10.26±4.61	18.04±1.8	17.11 ± 4.45
FetchPick	1.89 ± 1.22	3.30 ± 0.66	0.42 ± 0.29	0.41 ± 0.22	7.90 ± 1.22
FetchPush	6.44 ± 3.64	6.43 ± 0.56	1.69 ± 1.56	2.63± 3.04	7.11 ± 2.60
FetchSlide	1.77 ± 0.24	1.11± 0.26	0.0 ± 0.0	0.10 ± 0.11	0.80 ± 0.48

Table 8: Discounted Return for the offline GCRL benchmark with 1% expert data. Results are averaged over 10 seeds.

Task	$\beta = 0.5$	$\beta = 0.7$	$\beta = 0.8$	$\beta = 0.9$
FetchReach	35.08 ± 0.54	36.57 ± 0.20	36.59± 0.30	36.30± 0.30
FetchPick	26.47± 0.34	27.04± 0.81	27.43± 0.97	27.89 ± 1.19
FetchPush	26.83 ± 1.21	16.20± 1.11	11.50± 1.19	13.85± 5.53
FetchSlide	4.99± 0.40	3.76± 0.75	3.43± 2.4	4.10± 1.20

Table 9: Discounted Return for the offline GCRL benchmark with varying mixture coefficients in offline dataset. Results are averaged over 10 seeds.

Task	Occupancy Matching		Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	SMORe	GoFAR	WGCSL	GCSL		AM	IQL
Reacher	0.875±0.07	0.90±0.01	0.97±0.014	0.92 ± 0.08	0.76±0.74	1.0±0.1	0.26 ± 0.06
SawyerReach	0.98±0.014	0.75±0.04	1.0±0.0	0.98±0.02	0.98±0.018	1.0±0.1	0.81 ± 0.01
SawyerDoor	0.875±0.038	0.5±0.12	0.78 ± 0.10	0.5±0.12	0.22±0.11	0.3±0.11	0.84 ± 0.06
FetchReach	1.0± 0.0	1.0 ± 0.0	1.0± 0.0	0.98 ± 0.05	1.0± 0.0	1.0± 1.0	1.0 ± 0.0
FetchPick	0.925 ± 0.045	0.84 ± 0.09	0.54± 0.16	0.54 ± 0.20	0.42 ± 0.29	0.78 ± 0.15	0.86 ± 0.11
FetchPush	0.90± 0.07	0.88± 0.09	0.76± 0.12	0.72 ± 0.15	0.06± 0.03	0.67± 0.14	0.65 ± 0.052
FetchSlide	0.315± 0.07	0.18 ± 0.12	0.18± 0.14	0.17± 0.13	0.0 ± 0.0	0.11± 0.09	0.26± 0.057
HandReach	0.47± 0.11	0.40 ± 0.20	0.25± 0.23	0.047± 0.10	0.0± 0.0	0.0 ± 0.0	0.0 ± 0.0

Table 10: Success Rate for the offline GCRL benchmark with 10% expert data. Results are averaged over 10 seeds.

Task	Occupancy Matching		Behavior cloning		Contrastive RL CRL	RL+sparse reward	
	SMORe	GoFAR	WGCSL	GCSL		AM	IQL
Reacher	0.02±0.01	0.03±0.01	0.011±0.01	0.016 ±0.00	0.05±0.03	0.013±0.00	0.12 ± 0.005
SawyerReach	0.008±0.004	0.04±0.00	0.004±0.00	0.00±0.00	0.01±0.01	0.01 ± 0.00	0.053 ± 0.004
SawyerDoor	0.02±0.029	0.18±0.00	0.011±0.00	0.017±0.01	0.14±0.07	0.06 ± 0.01	0.019 ± 0.01
FetchReach	0.004± 0.0012	0.018± 0.003	0.007± 0.0043	0.008 ± 0.008	0.007 ± 0.001	0.007 ± 0.001	0.002± 0.001
FetchPick	0.04± 0.018	0.036 ± 0.013	0.094± 0.043	0.108± 0.06	0.25 ± 0.025	0.04± 0.02	0.04± 0.012
FetchPush	0.03± 0.003	0.033± 0.008	0.041± 0.02	0.042± 0.018	0.15± 0.036	0.07±0.039	0.05± 0.006
FetchSlide	0.09± 0.012	0.12 ± 0.02	0.173± 0.04	0.204± 0.051	0.42± 0.05	0.198± 0.059	0.09± 0.013
HandReach	0.039± 0.0108	0.024± 0.009	0.035 ± 0.012	0.038± 0.013	0.04 ± 0.005	0.037 ± 0.004	0.08 ± 0.005

Table 11: Final distance to goal for the offline GCRL benchmark with 10% expert data. Results are averaged over 10 seeds.