

## References

- Alekh Agarwal, Akshay Krishnamurthy, John Langford, Haipeng Luo, et al. Open problem: First-order regret bounds for contextual bandits. In *Conference on Learning Theory*, pages 4–7. PMLR, 2017.
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020.
- Zeyuan Allen-Zhu, Sébastien Bubeck, and Yuanzhi Li. Make the minority great again: First-order regret bound for contextual bandits. In *International Conference on Machine Learning*, pages 186–194. PMLR, 2018.
- András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71:89–129, 2008.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- Gabriel Barth-Maron, Matthew W. Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva TB, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributional policy gradients. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=SyZipzbCb>.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International conference on machine learning*, pages 449–458. PMLR, 2017.
- Marc G Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C Machado, Subhodeep Moitra, Sameera S Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588(7836):77–82, 2020.
- Cristian Bodnar, Adrian Li, Karol Hausman, Peter Pastor, and Mrinal Kalakrishnan. Quantile qt-opt for risk-aware vision-based robotic grasping. *Robotics: Science and Systems*, 2020.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Jonathan Chang, Kaiwen Wang, Nathan Kallus, and Wen Sun. Learning bellman complete representations for offline policy evaluation. In *International Conference on Machine Learning*, pages 2938–2971. PMLR, 2022.
- Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic for offline reinforcement learning with function approximation. In *International Conference on Machine Learning*, pages 3852–3878. PMLR, 2022.
- Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for distributional reinforcement learning. In *International conference on machine learning*, pages 1096–1105. PMLR, 2018a.
- Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018b.
- Chris Dann, Yishay Mansour, Mehryar Mohri, Ayush Sekhari, and Karthik Sridharan. Guarantees for epsilon-greedy reinforcement learning with function approximation. In *International Conference on Machine Learning*, pages 4666–4689. PMLR, 2022.

418 Christoph Dann, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E  
419 Schapire. On oracle-efficient pac rl with rich observations. *Advances in neural information*  
420 *processing systems*, 31, 2018.

421 Mónica Farsang, Paul Mineiro, and Wangda Zhang. Conditionally risk-averse contextual bandits.  
422 *arXiv preprint arXiv:2210.13573*, 2022.

423 Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Moham-  
424 madamin Barekatain, Alexander Novikov, Francisco J R Ruiz, Julian Schrittwieser, Grzegorz  
425 Swirszcz, et al. Discovering faster matrix multiplication algorithms with reinforcement learning.  
426 *Nature*, 610(7930):47–53, 2022.

427 Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with  
428 regression oracles. In *International Conference on Machine Learning*, pages 3199–3210. PMLR,  
429 2020.

430 Dylan J Foster and Akshay Krishnamurthy. Efficient first-order contextual bandits: Prediction,  
431 allocation, and triangular discrimination. *Advances in Neural Information Processing Systems*, 34:  
432 18907–18919, 2021.

433 Dylan J Foster, Zhiyuan Li, Thodoris Lykouris, Karthik Sridharan, and Eva Tardos. Learning in  
434 games: Robustness of fast convergence. *Advances in Neural Information Processing Systems*, 29,  
435 2016.

436 Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of  
437 interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

438 Dylan J Foster, Akshay Krishnamurthy, David Simchi-Levi, and Yunzong Xu. Offline reinforcement  
439 learning: Fundamental barriers for value function approximation. In *Conference on Learning*  
440 *Theory*, pages 3489–3489. PMLR, 2022.

441 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
442 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
443 pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.

444 Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan  
445 Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in  
446 deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*,  
447 volume 32, 2018.

448 Audrey Huang, Jinglin Chen, and Nan Jiang. Reinforcement learning in low-rank mdps with density  
449 features. *arXiv preprint arXiv:2302.02252*, 2023.

450 Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Context-  
451 tual decision processes with low bellman rank are pac-learnable. In *International Conference on*  
452 *Machine Learning*, pages 1704–1713. PMLR, 2017.

453 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for  
454 reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879.  
455 PMLR, 2020a.

456 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement  
457 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143.  
458 PMLR, 2020b.

459 Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl  
460 problems, and sample-efficient algorithms. *Advances in neural information processing systems*,  
461 34:13406–13418, 2021a.

462 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In  
463 *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021b.

464 Sham Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information  
465 theoretic regret bounds for online nonlinear control. *Advances in Neural Information Processing*  
466 *Systems*, 33:15312–15325, 2020.

467 Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14,  
468 2001.

469 Nathan Kallus, Xiaojie Mao, Kaiwen Wang, and Zhengyuan Zhou. Doubly robust distributionally  
470 robust off-policy evaluation and learning. In *International Conference on Machine Learning*, pages  
471 10598–10632. PMLR, 2022.

472 Ramtin Keramati, Christoph Dann, Alex Tamkin, and Emma Brunskill. Being optimistic to be  
473 conservative: Quickly learning a cvar policy. In *Proceedings of the AAAI Conference on Artificial  
474 Intelligence*, volume 34, pages 4436–4443, 2020.

475 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-  
476 based offline reinforcement learning. *Advances in neural information processing systems*, 33:  
477 21810–21823, 2020.

478 J Kolter. The fixed points of off-policy td. *Advances in Neural Information Processing Systems*, 24,  
479 2011.

480 Alex Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. 2009.

481 Lucien Le Cam. *Asymptotic methods in statistical decision theory*. Springer Science & Business  
482 Media, 2012.

483 Shiao Hong Lim and Ilyas Malik. Distributional reinforcement learning for risk-sensitive policies.  
484 In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in  
485 Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=wSVEd3Ta42m>.

487 Thodoris Lykouris, Karthik Sridharan, and Éva Tardos. Small-loss bounds for online learning with  
488 partial information. *Mathematics of Operations Research*, 47(3):2186–2218, 2022.

489 Clare Lyle, Marc G Bellemare, and Pablo Samuel Castro. A comparative analysis of expected  
490 and distributional reinforcement learning. In *Proceedings of the AAAI Conference on Artificial  
491 Intelligence*, volume 33, pages 4504–4511, 2019.

492 Yecheng Ma, Dinesh Jayaraman, and Osbert Bastani. Conservative offline distributional reinforcement  
493 learning. *Advances in Neural Information Processing Systems*, 34:19235–19247, 2021.

494 Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstrac-  
495 tion and provably efficient rich-observation reinforcement learning. In *International conference on  
496 machine learning*, pages 6961–6971. PMLR, 2020.

497 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,  
498 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control  
499 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

500 Aditya Modi, Jinglin Chen, Akshay Krishnamurthy, Nan Jiang, and Alekh Agarwal. Model-free  
501 representation learning and exploration in low-rank mdps. *arXiv preprint arXiv:2102.07035*, 2021.

502 Anna Montoya, BigJek14, Bull, denisedunleavy, egrad, FleetwoodHack, Imbayoh, PadraicS,  
503 Pru\_Admin, tpitman, and Will Cukierski. Prudential life insurance assessment, 2015. URL  
504 <https://kaggle.com/competitions/prudential-life-insurance-assessment>.

505 Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine  
506 Learning Research*, 9(5), 2008.

507 Gergely Neu. First-order regret bounds for combinatorial semi-bandits. In *Conference on Learning  
508 Theory*, pages 1360–1375. PMLR, 2015.

509 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-  
510 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information  
511 Processing Systems*, 34:11702–11716, 2021.

512 Mark Rowland, Marc Bellemare, Will Dabney, Rémi Munos, and Yee Whye Teh. An analysis  
513 of categorical distributional reinforcement learning. In *International Conference on Artificial*  
514 *Intelligence and Statistics*, pages 29–37. PMLR, 2018.

515 Mark Rowland, Rémi Munos, Mohammad Gheshlaghi Azar, Yunhao Tang, Georg Ostrovski, Anna  
516 Harutyunyan, Karl Tuyls, Marc G Bellemare, and Will Dabney. An analysis of quantile temporal-  
517 difference learning. *arXiv preprint arXiv:2301.04462*, 2023.

518 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

519 Flemming Topsoe. Some inequalities for information divergence and related measures of discrimina-  
520 tion. *IEEE Transactions on information theory*, 46(4):1602–1609, 2000.

521 John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function  
522 approximation. *Advances in neural information processing systems*, 9, 1996.

523 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under  
524 partial coverage. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=tyrJsbKAe6>.  
525

526 Sara van de Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.

527 Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked science in  
528 machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013. doi: 10.1145/2641190.2641198.  
529 URL <http://doi.acm.org/10.1145/2641190.2641198>.

530 István Vincze. On the concept and measure of information contained in an observation. In *Contribu-*  
531 *tions to Probability*, pages 207–214. Elsevier, 1981.

532 Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. First-order  
533 regret in reinforcement learning with linear function approximation: A robust estimation approach.  
534 In *International Conference on Machine Learning*, pages 22384–22429. PMLR, 2022.

535 Kaiwen Wang, Nathan Kallus, and Wen Sun. Near-minimax-optimal risk-sensitive reinforcement  
536 learning with cvar. *International Conference on Machine Learning*, 2023.

537 Ruosong Wang, Dean Foster, and Sham M. Kakade. What are the statistical limits of offline RL with  
538 linear function approximation? In *International Conference on Learning Representations*, 2021a.  
539 URL <https://openreview.net/forum?id=30EvkP2aQLD>.

540 Ruosong Wang, Yifan Wu, Ruslan Salakhutdinov, and Sham Kakade. Instabilities of offline rl  
541 with pre-trained neural representation. In *International Conference on Machine Learning*, pages  
542 10948–10960. PMLR, 2021b.

543 Yuanhao Wang, Ruosong Wang, and Sham Kakade. An exponential lower bound for linearly  
544 realizable mdp with constant suboptimality gap. *Advances in Neural Information Processing*  
545 *Systems*, 34:9521–9533, 2021c.

546 Runzhe Wu, Masatoshi Uehara, and Wen Sun. Distributional offline policy evaluation with predictive  
547 error guarantees. *International Conference of Machine Learning*, 2023.

548 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent  
549 pessimism for offline reinforcement learning. *Advances in neural information processing systems*,  
550 34:6683–6694, 2021.

551 Tengyang Xie, Dylan J Foster, Yu Bai, Nan Jiang, and Sham M. Kakade. The role of coverage in online  
552 reinforcement learning. In *The Eleventh International Conference on Learning Representations*,  
553 2023. URL <https://openreview.net/forum?id=LQIjzPdDt3q>.

554 Derek Yang, Li Zhao, Zichuan Lin, Tao Qin, Jiang Bian, and Tie-Yan Liu. Fully parameterized  
555 quantile function for distributional reinforcement learning. *Advances in neural information*  
556 *processing systems*, 32, 2019.

- 557 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement  
558 learning without domain knowledge using value function bounds. In *International Conference on*  
559 *Machine Learning*, pages 7304–7312. PMLR, 2019.
- 560 Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic methods  
561 for offline reinforcement learning. *Advances in neural information processing systems*, 34:13626–  
562 13640, 2021.
- 563 Tianjun Zhang, Tongzheng Ren, Mengjiao Yang, Joseph Gonzalez, Dale Schuurmans, and Bo Dai.  
564 Making linear mdps practical via contrastive representation learning. In *International Conference*  
565 *on Machine Learning*, pages 26447–26466. PMLR, 2022.
- 566 Tong Zhang. *Mathematical Analysis of Machine Learning Algorithms*. 2023. <http://www.tongzhang-ml.org/lt-book.html>.  
567
- 568 Zhengyuan Zhou, Susan Athey, and Stefan Wager. Offline multi-action policy learning: Generalization  
569 and optimization. *Operations Research*, 71(1):148–183, 2023.

# Appendices

## A Notations

Table 2: List of Notations

$\mathcal{S}, \mathcal{A}, A$	State and action spaces, and $A =  \mathcal{A} $ .
$\Delta(S)$	The set of distributions supported by $S$ .
$\bar{d}$	The expectation of any real-valued distribution $d$ , i.e., $\bar{d} = \mathbb{E}_{y \sim d}[y]$ .
$[N]$	$\{1, 2, \dots, N\}$ for any natural number $N$ .
$Z_h^\pi(x, a)$	Distribution of $\sum_{t=h}^H c_t$ given $x_h = x, a_h = a$ rolling in from $\pi$ .
$Q_h^\pi(x, a), V_h^\pi(x)$	$Q_h^\pi(x, a) = \bar{Z}_h^\pi(x, a)$ and $V_h^\pi = \mathbb{E}_{a \sim \pi(x)}[Q_h^\pi(x, a)]$ .
$\pi^*$	Optimal policy, i.e., $\pi^* = \arg \min_{\pi} V_1^\pi(x_1)$ . Without loss of optimality, we take $\pi^* : \mathcal{X} \rightarrow \mathcal{A}$ to be Markov & deterministic.
$Z_h^*, Q_h^*, V_h^*$	$Z_h^\pi, Q_h^\pi, V_h^\pi$ with $\pi = \pi^*$ , the optimal policy.
$\mathcal{T}_h^\pi, \mathcal{T}_h^*$	The Bellman operators that act on functions.
$\mathcal{T}_h^{\pi, D}, \mathcal{T}_h^{*, D}$	The distributional Bellman operators that act on conditional distributions.
$V^\pi, Z^\pi, V^*, Z^*$	$V^\pi = V_1^\pi(x_1), Z^\pi = Z_1^\pi(x_1)$ . $V^*, Z^*$ are defined similarly with $\pi^*$ .
$d_h^\pi(x, a)$	The probability of $\pi$ visiting $(x, a)$ at time $h$ .
$C^{\tilde{\pi}}$	Coverage coefficient $\max_h \ d_h^{\tilde{\pi}} / d\nu_h\ _\infty$ .
$D_\Delta(f \parallel g)$	Triangular discrimination between $f, g$ .
$H(f \parallel g)$	Hellinger distance between $f, g$ .
$D_{KL}(f \parallel g)$	KL divergence between $f, g$ .

### 572 A.1 Statistical Distances

573 Let  $f, g$  be distributions over  $\mathcal{Y}$ . Then,

$$\begin{aligned}
 D_\Delta(f \parallel g) &= \sum_y \frac{(f(y) - g(y))^2}{f(y) + g(y)}, \\
 H(f \parallel g) &= \sqrt{\frac{1}{2} \sum_y \left( \sqrt{f(y)} - \sqrt{g(y)} \right)^2}, \\
 D_{KL}(f \parallel g) &= \sum_y f(y) \log(f(y)/g(y)), \\
 D_{TV}(f \parallel g) &= \frac{1}{2} \sum_y |f(y) - g(y)|.
 \end{aligned}$$

574 The following standard inequalities will be helpful:

$$\begin{aligned}
 H^2 &\leq D_{TV} \leq \sqrt{2}H, \\
 2H^2 &\leq D_\Delta \leq 4H^2, \\
 H &\leq \sqrt{D_{KL}}.
 \end{aligned}
 \tag{Lemma A.1}$$

575 **Lemma A.1.** For any distributions  $f, g$ , we have  $2H^2(f \parallel g) \leq D_\Delta(f \parallel g) \leq 4H^2(f \parallel g)$ .

576 *Proof.* Recall that

$$D_\Delta(f \parallel g) = \int_y \left( \frac{f(y) - g(y)}{\sqrt{f(y) + g(y)}} \right)^2.$$

577 Applying  $\frac{1}{\sqrt{f(y) + g(y)}} \leq \frac{1}{\sqrt{f(y)}} \leq \frac{\sqrt{2}}{\sqrt{f(y) + g(y)}}$  concludes the proof.  $\square$

## 578 B Omitted Algorithms

579 In this section, we present the O-DISCO algorithm with Uniform Action Exploration (UAE), as  
 580 described in [Section 5.2](#). We also present versions of O-DISCO and P-DISCO for the reward-  
 581 maximizing setting (instead of the cost-minimizing setting studied throughout the paper); if SMALL-  
 582 RETURN is turned on, we can derive small-return bounds in [Appendix I](#).

---

### Algorithm 4 O-DISCO (with UAE and small return)

---

- 1: **Input:** number of episodes  $K$ , distribution function class  $\mathcal{F}$ , threshold  $\beta$ , flag UAE, flag SMALLRETURN.
  - 2: Initialize  $\mathcal{D}_{h,0} \leftarrow \emptyset$  for all  $h \in [H]$ , and set  $\mathcal{F}_0 = \mathcal{F}$ .
  - 3: Set  $\text{op} = \max$  if SMALLRETURN else  $\text{op} = \min$ .
  - 4: **for** episode  $k = 1, 2, \dots, K$  **do**
  - 5:   Set  $f^{(k)} = \arg \text{op}_{f \in \mathcal{F}_{k-1}} \text{op}_a \bar{f}_1(x_1, a)$ .
  - 6:   Set  $\pi_h^k(x) = \arg \text{op}_a \bar{f}_h^{(k)}(x, a)$ .
  - 7:   **if** UAE **then**
  - 8:     For each  $h \in [H]$ , collect  $x_{h,k} \sim d_h^{\pi_h^k}$ ,  $a_{h,k} \sim \text{unif}(\mathcal{A})$ ,  $c_{h,k} \sim C_h(x_{h,k}, a_{h,k})$ ,  $x'_{h,k} \sim P_h(x_{h,k}, a_{h,k})$ , and augment the dataset  $\mathcal{D}_{h,k} = \mathcal{D}_{h,k-1} \cup \{(x_{h,k}, a_{h,k}, c_{h,k}, x'_{h,k})\}$ .
  - 9:   **else**
  - 10:     Roll out  $\pi^k$  and obtain a trajectory  $x_{1,k}, a_{1,k}, c_{1,k}, \dots, x_{H,k}, a_{H,k}, c_{H,k}$ .
  - 11:     For each  $h \in [H]$ , augment the dataset  $\mathcal{D}_{h,k} = \mathcal{D}_{h,k-1} \cup \{(x_{h,k}, a_{h,k}, c_{h,k}, x_{h+1,k})\}$ .
  - 12:   **end if**
  - 13:   For all  $(h, f) \in [H] \times \mathcal{F}$ , sample  $y_{h,i}^f \sim f_{h+1}(x'_{h,i}, a')$  and  $a' = \arg \text{op}_a \bar{f}_{h+1}(x'_{h,i}, a)$ , where  $(x_{h,i}, a_{h,i}, c_{h,i}, x'_{h,i})$  is the  $i$ -th datapoint of  $\mathcal{D}_{h,k}$ . Also, set  $z_{h,i}^f = c_{h,i} + y_{h,i}^f$  and define the confidence set,
  - 14:   
$$\mathcal{F}_k = \left\{ f \in \mathcal{F} : \sum_{i=1}^k \log f_h(z_{h,i}^f \mid x_{h,i}, a_{h,i}) \geq \max_{\tilde{f} \in \mathcal{F}} \sum_{i=1}^k \log \tilde{f}_h(z_{h,i}^f \mid x_{h,i}, a_{h,i}) - 7\beta, \forall h \in [H] \right\}.$$
  - 15: **end for**
  - 16: **Output:**  $\bar{\pi} = \text{unif}(\pi^{1:K})$ .
- 

---

### Algorithm 5 P-DISCO (with small return)

---

- 1: **Input:** datasets  $\mathcal{D}_1, \dots, \mathcal{D}_H$ , distribution function class  $\mathcal{F}$ , threshold  $\beta$ , policy class  $\Pi$ , flag SMALLRETURN.
- 2: For all  $(h, f, \pi) \in [H] \times \mathcal{F} \times \Pi$ , sample  $y_{h,i}^{f,\pi} \sim f_{h+1}(x'_{h,i}, \pi_{h+1}(x'_{h,i}))$ , where  $(x_{h,i}, a_{h,i}, c_{h,i}, x'_{h,i})$  is the  $i$ -th datapoint of  $\mathcal{D}_h$ . Then, set  $z_{h,i}^{f,\pi} = c_{h,i} + y_{h,i}^{f,\pi}$  and define the confidence set,

$$\mathcal{F}_\pi = \left\{ f \in \mathcal{F} : \sum_{i=1}^N \log f_h(z_{h,i}^{f,\pi} \mid x_{h,i}, a_{h,i}) \geq \max_{\tilde{f} \in \mathcal{F}} \sum_{i=1}^N \log \tilde{f}_h(z_{h,i}^{f,\pi} \mid x_{h,i}, a_{h,i}) - 7\beta, \forall h \in [H] \right\}.$$

- 3: Set  $\text{op} = \max$  if SMALLRETURN else  $\text{op} = \min$ .
  - 4: For each  $\pi \in \Pi$ , define the pessimistic estimate  $f^\pi = \arg \text{op}_{f \in \mathcal{F}_\pi} \mathbb{E}_{a \sim \pi(x_1)} [\bar{f}_1(x_1, a)]$ .
  - 5: **Output:**  $\hat{\pi} = \arg \text{op}_{\pi \in \Pi} \mathbb{E}_{a \sim \pi(x_1)} [\bar{f}_1^\pi(x_1, \pi)]$ .
-

## C Proofs for DISTCB

**Lemma C.1** (Azuma). *Let  $\{X_i\}_{i \in [N]}$  be a sequence of random variables supported on  $[0, 1]$ , adapted to filtration  $\{\mathcal{F}_i\}_{i \in [N]}$ . For any  $\delta \in (0, 1)$ , we have w.p. at least  $1 - \delta$ ,*

$$\sum_{t=1}^N \mathbb{E}[X_t \mid \mathcal{F}_{t-1}] \leq \sum_{t=1}^N X_t + \sqrt{N \log(2/\delta)}, \quad (\text{Standard Azuma})$$

$$\sum_{t=1}^N \mathbb{E}[X_t \mid \mathcal{F}_{t-1}] \leq 2 \sum_{t=1}^N X_t + 2 \log(1/\delta). \quad (\text{Multiplicative Azuma})$$

*Proof.* For standard Azuma, see [Zhang \[2023, Theorem 13.4\]](#). For multiplicative Azuma, apply [\[Zhang, 2023, Theorem 13.5\]](#) with  $\lambda = 1$ . The claim follows, since  $\frac{1}{1 - \exp(-\lambda)} \leq 2$ .  $\square$

**Theorem 4.1.** *Fix any  $\delta \in (0, 1)$  and set  $\gamma = 10A \vee \sqrt{\frac{40A(C^* + \log(1/\delta))}{112(\text{Regret}_{\log}(K) + \log(1/\delta))}}$ . Then, w.p. at least  $1 - \delta$ , DISTCB satisfies,*

$$\text{Regret}_{\text{DISTCB}}(K) \leq 232 \sqrt{AC^* \text{Regret}_{\log}(K) \log(1/\delta)} + 2300A(\text{Regret}_{\log}(K) + \log(1/\delta)),$$

where  $C^* = \sum_{k=1}^K \min_{a \in \mathcal{A}} \bar{C}(x_k, a)$  is the cumulative cost of the optimal policy.

*Proof of Theorem 4.1.* First, recall the per-step inequality of ReIGW [Foster and Krishnamurthy \[2021, Theorem 4\]](#), which states: for any  $\hat{f}$  and  $\gamma \geq 2A$ , if we set  $p = \text{ReIGW}_\gamma(\hat{f}, \gamma)$ , then, for all  $f \in [0, 1]^{\mathcal{A}}$ , we have

$$\sum_a p(a)(f(a) - f(a^*)) \leq \frac{5A}{\gamma} \sum_a p(a)f(a) + 7\gamma \sum_a p(a) \frac{(\hat{f}(a) - f(a))^2}{\hat{f}(a) + f(a)},$$

where  $a^* = \arg \min_a f(a)$ . For any  $k \in [K]$ , applying this to  $\hat{f} = \bar{f}^{(k)}(s_k, \cdot)$ ,  $p = p_k$  and  $f = \bar{C}(s_k, \cdot)$ , we have

$$\begin{aligned} \sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, a_k) - \bar{C}(s_k, \pi^*(s_k))] &\leq \sum_{k=1}^K \mathbb{E}_{a_k} \left[ \frac{5A}{\gamma} \bar{C}(s_k, a_k) + 7\gamma \frac{(\bar{f}^{(k)}(s_k, a_k) - \bar{C}(s_k, a_k))^2}{\bar{f}^{(k)}(s_k, a_k) + \bar{C}(s_k, a_k)} \right] \\ &\leq \sum_{k=1}^K \mathbb{E}_{a_k} \left[ \frac{5A}{\gamma} \bar{C}(s_k, a_k) + 7\gamma D_\Delta(f^{(k)}(s_k, a_k) \parallel C(s_k, a_k)) \right] \end{aligned} \quad (\text{Eq. } (\Delta_1))$$

Since  $D_\Delta \leq 4H^2$ , we have

$$\begin{aligned} &\sum_{k=1}^K \mathbb{E}_{a_k} [D_\Delta(f^{(k)}(s_k, a_k) \parallel C(s_k, a_k))] \\ &\leq 4 \sum_{k=1}^K \mathbb{E}_{a_k} [H^2(C(s_k, a_k) \parallel f^{(k)}(s_k, a_k))] \\ &\leq 8 \sum_{k=1}^K H^2(C(s_k, a_k) \parallel f^{(k)}(s_k, a_k)) + 8 \log(1/\delta) \quad (\text{Multiplicative Azuma, since } H^2 \in [0, 1]) \\ &\leq 8 \text{Regret}_{\log}(K) + 10 \log(1/\delta). \quad (\text{Foster et al. [2021, Lemma A.14]}) \end{aligned}$$

Hence, we have

$$\sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, a_k) - \bar{C}(s_k, \pi^*(s_k))] \leq \frac{5A}{\gamma} \sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, a_k)] + 70\gamma(\text{Regret}_{\log}(K) + \log(1/\delta)).$$



598 Finally, recalling that  $1/(1 - \varepsilon) \leq 1 + 2\varepsilon$  when  $\varepsilon \leq \frac{1}{2}$ , and the fact that  $\frac{5A}{\gamma} \leq \frac{1}{2}$ , we have

$$\sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, a_k) - \bar{C}(s_k, \pi^*(s_k))] \leq \frac{10A}{\gamma} \sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, \pi^*(s_k))] + 140\gamma(\text{Regret}_{\log}(K) + \log(1/\delta)).$$

599 By Azuma's inequality, we have

$$\begin{aligned} & \sum_{k=1}^K \bar{C}(s_k, a_k) - \bar{C}(s_k, \pi^*(s_k)) \\ & \leq 2 \sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, a_k) - \bar{C}(s_k, \pi^*(s_k))] + 2\log(1/\delta) \\ & \leq \frac{20A}{\gamma} \sum_{k=1}^K \mathbb{E}_{a_k} [\bar{C}(s_k, \pi^*(s_k))] + 140\gamma(\text{Regret}_{\log}(K) + \log(1/\delta)) + 2\log(1/\delta) \\ & \leq \frac{40A}{\gamma} (C^* + \log(1/\delta)) + 140\gamma(\text{Regret}_{\log}(K) + \log(1/\delta)) + 2\log(1/\delta). \end{aligned}$$

600 Now set  $\gamma = \sqrt{\frac{40A(C^* + \log(1/\delta))}{140(\text{Regret}_{\log}(K) + \log(1/\delta))}} \vee 10A$ .

601 Case 1 is when  $\sqrt{\frac{40A(C^* + \log(1/\delta))}{140(\text{Regret}_{\log}(K) + \log(1/\delta))}} \leq 10A$ , i.e.,  $(C^* + \log(1/\delta)) \leq$   
 602  $280A(\text{Regret}_{\log}(K) + \log(1/\delta))$ , we have the above is at most

$$\begin{aligned} & 4(C^* + \log(1/\delta)) + 1120A(\text{Regret}_{\log}(K) + \log(1/\delta)) + 2\log(1/\delta) \\ & \leq 2240A(\text{Regret}_{\log}(K) + \log(1/\delta)) + 2\log(1/\delta). \end{aligned}$$

603 Case 2 is when the left term dominates, then the bound is,

$$\begin{aligned} & 2\sqrt{4480A(C^* + \log(1/\delta))(\text{Regret}_{\log}(K) + \log(1/\delta)) + 2\log(1/\delta)} \\ & \leq 2\sqrt{13440AC^* \text{Regret}_{\log}(K) \log(1/\delta) + 4480A \log^2(1/\delta) + 2\log(1/\delta)} \\ & \leq 232\sqrt{AC^* \text{Regret}_{\log}(K) \log(1/\delta) + 134\sqrt{A} \log(1/\delta) + 2\log(1/\delta)}. \end{aligned}$$

604 Putting these two cases together, we have the result. □

## 605 D Placeholder

606 This section used to contain information that is no longer needed. We kept this placeholder section to  
 607 ensure the main text's references to the appendix are consistent.

## 608 E Maximum Likelihood Estimation

609 This section reviews generalization bounds for the maximum likelihood estimator (MLE). We adopt  
 610 the same sequential condition probability estimation setup as in Agarwal et al. [2020, Appendix E],  
 611 which we now recall for completeness. Let  $\mathcal{X}$  be the context/feature space and  $\mathcal{Y}$  be the label space,  
 612 and we are given a dataset  $D = \{(x_i, y_i)\}_{i \in [n]}$  from a martingale process: for  $i = 1, 2, \dots, n$ , sample  
 613  $x_i \sim \mathcal{D}_i(x_{1:i-1}, y_{1:i-1})$  and  $y_i \sim p(\cdot | x_i)$ . Let  $f^*(x, y) = p(y | x)$  and we are given a realizable,  
 614 i.e.,  $f^* \in \mathcal{F}$ , function class  $\mathcal{F} : \mathcal{X} \times \mathcal{Y} \rightarrow \Delta(\mathbb{R})$  of distributions. The MLE is an estimate for  $f^*$  that  
 615 maximizes the log-likelihood objective over our dataset:

$$\hat{f}_{\text{MLE}} = \arg \max_{f \in \mathcal{F}} \sum_{i=1}^n \log f(x_i, y_i).$$

616 For our guarantees to hold for general hypotheses classes  $\mathcal{F}$ , we use the bracketing number to quantify  
 617 the statistical complexity of  $\mathcal{F}$  [van de Geer, 2000].

618 **Definition E.1** (Bracketing Number). Let  $\mathcal{G}$  be a set of functions mapping  $\mathcal{X} \rightarrow \mathbb{R}$ . Given two  
 619 functions  $l, u$  such that  $l(x) \leq u(x)$  for all  $x \in \mathcal{X}$ , the bracket  $[l, u]$  is the set of functions  $g \in \mathcal{G}$   
 620 such that  $l(x) \leq g(x) \leq u(x)$  for all  $x \in \mathcal{X}$ . We call  $[l, u]$  an  $\varepsilon$ -bracket if  $\|u - l\| \leq \varepsilon$ . Then, the  
 621  $\varepsilon$ -bracketing number of  $\mathcal{G}$  with respect to  $\|\cdot\|$ , denoted by  $N_{[]}(\varepsilon, \mathcal{G}, \|\cdot\|)$  is the minimum number of  
 622  $\varepsilon$ -brackets needed to cover  $\mathcal{G}$ .

623 Since the triangular discrimination is equivalent to squared Hellinger up to universal constants, we  
 624 now prove MLE generalization bounds in terms of squared Hellinger.

625 **Lemma E.2.** Let  $f_1 : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$  and  $f_2 : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  satisfying  $\sup_{x \in \mathcal{X}} \int_{\mathcal{Y}} f_2(x, y) dy \leq s$ ,  
 626 then for any distribution  $\mathcal{D} \in \Delta(\mathcal{X})$ , we have

$$\mathbb{E}_{x \sim \mathcal{D}} [H^2(f_1(x) \parallel f_2(x, \cdot))] \leq (s - 1) - 2 \log \mathbb{E}_{x \sim \mathcal{D}, y \sim f_1(x)} \exp \left( -\frac{1}{2} \log(f_1(x, y)/f_2(x, y)) \right).$$

627 *Proof.* This follows from the proof of Wu et al. [2023, Lemma C.1].  $\square$

628 **Lemma E.3.** Fix  $\delta \in (0, 1)$ . Then w.p. at least  $1 - \delta$ , for any  $f \in \mathcal{F}$ , we have

$$\begin{aligned} & \sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} [H^2(f(x, \cdot) \parallel f^*(x, \cdot))] \\ & \leq 6n\epsilon|\mathcal{Y}| + 2 \sum_{i=1}^n \log(f^*(x_i, y_i)/f(x_i, y_i)) + 8 \log(N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{\infty})/\delta). \end{aligned} \quad (2)$$

629 *Rearranging, we also have*

$$\sum_{i=1}^n \log(f(x_i, y_i)/f^*(x_i, y_i)) \leq 3n\epsilon|\mathcal{Y}| + 4 \log(N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{\infty})/\delta). \quad (3)$$

630 *Proof.* We take an  $\epsilon$ -bracketing of  $\mathcal{F}$ ,  $\{[l_i, u_i] : i = 1, 2, \dots\}$ , and denote  $\tilde{\mathcal{F}} = \{u_i : i = 1, 2, \dots\}$ .  
 631 Applying Lemma 24 of Agarwal et al. [2020] to function class  $\tilde{\mathcal{F}}$  and using Chernoff method, w.p. at  
 632 least  $1 - \delta$ , for all  $\tilde{f} \in \tilde{\mathcal{F}}$ , we have

$$\underbrace{-\log \mathbb{E}_{D'} \exp(L(\tilde{f}(D), D'))}_{(i)} \leq \underbrace{-L(\tilde{f}(D), D) + 2 \log(N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{\infty})/\delta)}_{(ii)}. \quad (4)$$

633 Now, fix any  $f \in \mathcal{F}$  and pick  $\tilde{f} \in \tilde{\mathcal{F}}$  as the upper bracket, i.e.,  $f \leq \tilde{f}$ . Now set  $L(f, D) =$   
 634  $\sum_{i=1}^n -1/2 \log(f^*(x_i, y_i)/f(x_i, y_i))$ . Then the right hand side of (4) is

$$\begin{aligned} (ii) &= \frac{1}{2} \sum_{i=1}^n \log(f^*(x_i, y_i)/\tilde{f}(x_i, y_i)) + 2 \log(N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{\infty})/\delta) \\ &\leq \frac{1}{2} \sum_{i=1}^n \log(f^*(x_i, y_i)/f(x_i, y_i)) + 2 \log(N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{\infty})/\delta). \end{aligned}$$

635 On the other hand, since  $H$  is a metric, we have

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} H^2(f(x, \cdot), f^*(x, \cdot)) &\leq \sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} \left( H\left(f(x, \cdot), \tilde{f}(x, y)\right) + H\left(\tilde{f}(x, y), f^*(x, \cdot)\right) \right)^2 \\ &\leq 2 \underbrace{\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} H^2\left(f(x, \cdot), \tilde{f}(x, y)\right)}_{\text{(iii)}} + 2 \underbrace{\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} H^2\left(\tilde{f}(x, y), f^*(x, \cdot)\right)}_{\text{(iv)}}. \end{aligned}$$

636 For (iii), by the definition, we have  $\tilde{f}(x, y) - f(x, y) \in [0, \epsilon]$  for all  $x$ , so

$$\text{(iii)} = \sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} H^2\left(f(x, \cdot), \tilde{f}(x, y)\right) \leq \sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} 2 \int_y \left| f(x, y) - \tilde{f}(x, y) \right| dy \leq 2n\epsilon|\mathcal{Y}|.$$

637 For (iv), we apply [Lemma E.2](#) with  $f_1 = f^*$  and  $f_2 = \tilde{f}$  (thus  $s = 1 + \epsilon|\mathcal{Y}|$ ) and get

$$\begin{aligned} \text{(iv)} &= n\epsilon|\mathcal{Y}| - 2 \sum_{i=1}^n \log \mathbb{E}_{x, y \sim f^*(x, \cdot)} \exp\left(-\frac{1}{2} \log\left(f^*(x, y)/\tilde{f}(x, y)\right)\right) \\ &= n\epsilon|\mathcal{Y}| - 2 \sum_{i=1}^n \log \mathbb{E}_{x, y \sim \mathcal{D}_i} \exp\left(-\frac{1}{2} \log\left(f^*(x, y)/\tilde{f}(x, y)\right)\right) \\ &= n\epsilon|\mathcal{Y}| - 2 \log \mathbb{E}_{x, y \sim \mathcal{D}'} \left[ \exp\left(\sum_{i=1}^n -\frac{1}{2} \log\left(f^*(x, y)/\tilde{f}(x, y)\right)\right) \right] \\ &= n\epsilon|\mathcal{Y}| + 2 \cdot \text{(i)}. \end{aligned}$$

638 By plugging (iii) and (iv) back we get

$$\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} H^2(f(x, \cdot), f^*(x, \cdot)) \leq 6n\epsilon|\mathcal{Y}| + 4 \cdot \text{(i)}.$$

639 Notice that (i)  $\leq$  (ii), so we complete the proof by plugging (ii) into the above.  $\square$

640 We first state the MLE generalization result for finite  $\mathcal{F}$ .

641 **Theorem E.4.** Suppose  $\mathcal{F}$  is finite. Fix any  $\delta \in (0, 1)$ , set  $\beta = \log(|\mathcal{F}|/\delta)$  and define

$$\hat{\mathcal{F}} = \left\{ f \in \mathcal{F} : \sum_{i=1}^n \log f(x_i, y_i) \geq \max_{\tilde{f} \in \mathcal{F}} \sum_{i=1}^n \tilde{f}(x_i, y_i) - 4\beta \right\}.$$

642 Then w.p. at least  $1 - \delta$ , the following holds:

- 643 (1) The true distribution is in the version space, i.e.,  $f^* \in \hat{\mathcal{F}}$ .  
 644 (2) Any function in the version space is close to the ground truth data-generating distribution,  
 645 i.e., for all  $f \in \hat{\mathcal{F}}$

$$\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} [H^2(f(x, \cdot) \parallel f^*(x, \cdot))] \leq 22\beta.$$

646 *Proof.* These two claims follow from [Lemma E.3](#) with  $\epsilon = 0$ , and so  $N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_\infty) = |\mathcal{F}|$ . For

647 (1), apply [Eq. \(3\)](#) to  $f = \hat{f}_{\text{MLE}}$  to see that  $f^* \in \hat{\mathcal{F}}$ . For (2), apply [Eq. \(2\)](#) and note that the sum term  
 648 is at most  $4\beta$ . Thus, the right hand side of [Eq. \(2\)](#) is at most  $(6 + 8 + 8)\beta = 22\beta$ .  $\square$

649 We now state the result for infinite  $\mathcal{F}$  using bracketing entropy.

650 **Theorem E.5.** Fix any  $\delta \in (0, 1)$ , set  $\beta = \log(N_{[]}((n|\mathcal{Y}|)^{-1}, \mathcal{F}, \|\cdot\|_\infty)/\delta)$  and define

$$\widehat{\mathcal{F}} = \left\{ f \in \mathcal{F} : \sum_{i=1}^n \log f(x_i, y_i) \geq \max_{\tilde{f} \in \mathcal{F}} \sum_{i=1}^n \tilde{f}(x_i, y_i) - 7\beta \right\}.$$

651 Then w.p. at least  $1 - \delta$ , the following holds:

- 652 (1) The true distribution is in the version space, i.e.,  $f^* \in \widehat{\mathcal{F}}$ .  
 653 (2) Any function in the version space is close to the ground truth data-generating distribution,  
 654 i.e., for all  $f \in \widehat{\mathcal{F}}$

$$\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} [H^2(f(x, \cdot) \| f^*(x, \cdot))] \leq 28\beta.$$

655 *Proof.* These two claims follow from Lemma E.3 with  $\epsilon = 1/n|\mathcal{Y}|$ . For (1), apply Eq. (3) to  $f = \widehat{f}_{\text{MLE}}$   
 656 to see that  $f^* \in \widehat{\mathcal{F}}$ . For (2), apply Eq. (2) and note that the sum term is at most  $7\beta$ . Thus, the right  
 657 hand side of Eq. (3) is at most  $(6 + 14 + 8)\beta = 28\beta$ .  $\square$

## 658 F Confidence set construction with general function class

659 In this section, we extend the confidence set construction of O-DISCO and P-DISCO to general  $\mathcal{F}$ ,  
 660 which can be infinite. Our procedure constructs the confidence set by performing the thresholding  
 661 scheme on an  $\varepsilon$ -net of  $\mathcal{F}$ . While constructing an  $\varepsilon$ -net for  $\mathcal{F}$  is admittedly a computationally hard  
 662 procedure, this is still information theoretically possible and our focus in O-DISCO and P-DISCO  
 663 is to show that distributional RL information-theoretically leads to small-loss bounds.

664 We first define some notations. Let  $\mathcal{F}^\downarrow$  and  $\mathcal{F}^\uparrow$  denote a lower and upper  $\varepsilon$ -bracketing of  $\mathcal{F}$ ,  
 665 i.e., for any  $f \in \mathcal{F}$ , there exists an  $\varepsilon$ -bracket  $[f^\downarrow, f^\uparrow]$  such that for all  $h$ ,  $f_h^\downarrow \leq f_h \leq f_h^\uparrow$  with  
 666  $f^\downarrow \in \mathcal{F}^\downarrow, f^\uparrow \in \mathcal{F}^\uparrow$ . Recall that a lower bracket  $g \in \mathcal{F}^\downarrow$  may not be a valid distribution, but since  
 667 elements of  $\mathcal{F}$  map to non-negative values, we can assume  $g$  has non-negative entires as well. Also,  
 668 we have  $\alpha_h^g(x, a) := \int g_h(z | x, a) \geq 1 - \varepsilon$ , so for  $\varepsilon$  small enough,  $g$  is normalizable. Hence, define  
 669  $\tilde{g}(z | x, a) = \alpha_h^g(x, a)^{-1} g(z | x, a)$  as the normalized version, which is a valid distribution that we  
 670 can sample from.

671 Now, consider any martingale  $\{x_{h,i}, a_{h,i}, c_{h,i}\}_{i \in [n], h \in [H]}$ , which could be the online data up to  
 672 episode  $k$  or the offline data (consisting of  $N$  i.i.d. samples). We define the MLE with re-  
 673 spect to a lower bracket element as follows. For any  $h \in [H], g \in \mathcal{F}^\downarrow, \pi \in \Pi$ , sample  
 674  $y_{h,i}^{g,\pi} \sim \tilde{g}_{h+1}(x'_{h,i}, \pi(x'_{h,i}))$ , and  $z_{h,i}^{g,\pi} = c_{h,i} + y_{h,i}^{g,\pi}$ , define the MLE solution for  $(g, \pi)$  at time  
 675  $h$  as,

$$\text{MLE}_h^{g,\pi} = \arg \max_{f \in \mathcal{F}} \sum_{i=1}^n \log f_h(z_{h,i}^{g,\pi} | x_{h,i}, a_{h,i}).$$

676 Also, define the version space with respect to the above MLE as,

$$\mathcal{F}_{g,\pi,h} = \left\{ f \in \mathcal{F} : \sum_{i=1}^n \log f_h(z_{h,i}^{g,\pi} | x_{h,i}, a_{h,i}) \geq \sum_{i=1}^n \log \text{MLE}_h^{g,\pi}(z_{h,i}^{g,\pi} | x_{h,i}, a_{h,i}) - \beta \right\}.$$

677 We now prove a key result that implies that  $\mathcal{T}_h^\pi f_{h+1}^\downarrow$  falls into the confidence set  $\mathcal{F}_{f^\downarrow, \pi, h}$ .

678 **Theorem F.1.** For any  $\delta \in (0, 1)$  and suppose  $n \geq 2$ . Then, w.p. at least  $1 - \delta$ , for any  $h \in [H], g \in$   
 679  $\mathcal{F}, f^\downarrow \in \mathcal{F}^\downarrow, \pi \in \Pi$ , we have

$$\sum_{i=1}^n \log g_h(z_{h,i}^{f^\downarrow, \pi} | x_{h,i}, a_{h,i}) - \log \mathcal{T}_h^\pi f_{h+1}^\downarrow(z_{h,i}^{f^\downarrow, \pi} | x_{h,i}, a_{h,i}) \leq \log(e^4 N_{[]} (n^{-1}, \mathcal{F}, \|\cdot\|_\infty)^2 |\Pi| / \delta).$$

680 where  $z_{h,i}^{f^\downarrow, \pi} = c_{h,i} + y_{h,i}^{f^\downarrow, \pi}$  and  $y_{h,i}^{f^\downarrow, \pi} \sim \tilde{f}_{h+1}^\downarrow(\cdot | x'_{h,i}, \pi_{h+1}(x'_{h,i}))$ .

681 *Proof of Theorem F.1.* Consider a  $\varepsilon$ -bracketing of  $\mathcal{F}$  where  $\varepsilon \leq 1/n \leq 1/2$ ; we will study each  
 682 element and conclude with a union bound. For any lower bracket  $l$  and upper bracket  $u$  in the  
 683 bracketing (note  $l, u$  need not correspond to the same bracket). Recall that  $\alpha_{h+1}^l(x, a) := \int l_{h+1}(z |$   
 684  $x, a)$ , so we have  $1 - \varepsilon \leq \alpha_{h+1}^l \leq 1$  since  $l$  is a lower  $\varepsilon$ -bracket of distributions. Therefore, we have

$$\mathbb{E} \left[ \exp \sum_{i=1}^n \log \left( \frac{u_h(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})}{\mathcal{T}_h^{\pi} l_{h+1}(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})} \right) \right] = \prod_{i=1}^n \mathbb{E}_{\nu_{h,i}} \left[ \frac{u_h(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})}{\mathcal{T}_h^{\pi} l_{h+1}(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})} \right],$$

685 where  $\nu_{h,i}$  is the distribution of data from  $i$ -th round and time  $h$ . Note that  $\nu_{h,i}(x, a, c, x') =$   
 686  $d_{h,i}(x, a) C_h(c | x, a) P_h(x' | x, a)$  for some distribution  $d_{h,i}(x, a)$ . Now focus on each  $i$ , so for all  $i$ ,  
 687 we have

$$\begin{aligned} & \mathbb{E}_{\nu_{h,i}} \left[ \frac{u_h(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})}{\mathcal{T}_h^{\pi} l_{h+1}(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})} \right] \\ &= \int_{x,a,c,x',y} \nu_{h,i}(x, a, c, x') \tilde{l}_{h+1}(y | x', \pi(x')) \frac{u_h(c+y | x, a)}{\int_{c,x'} \nu_{h,i}(c, x' | x, a) l_{h+1}(y | x', \pi(x'))} \\ &= \int_{x,a,z} d_{h,i}(x, a) \int_z u_h(z | x, a) \\ &\quad \times \int_{c,x'} \nu_{h,i}(c, x' | x, a) \tilde{l}_{h+1}(z-c | x', \pi(x')) \frac{1}{\int_{c,x'} \nu_{h,i}(c, x' | x, a) l_{h+1}(z-c | x', \pi(x'))} \\ &= \int_{x,a,z} d_{h,i}(x, a) \int_z u_h(z | x, a) \alpha_{h+1}^l(x, a)^{-1} \\ &\leq \frac{1+\varepsilon}{1-\varepsilon} = 1 + \frac{2\varepsilon}{1-\varepsilon} \leq 1 + \frac{4}{n}. \end{aligned}$$

688 Therefore,

$$\mathbb{E} \left[ \exp \sum_{i=1}^n \log \left( \frac{u_h(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})}{\mathcal{T}_h^{\pi} l_{h+1}(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})} \right) \right] \leq (1 + 4/n)^n \leq e^4.$$

689 Thus, by Markov's inequality, w.p. at least  $1 - \delta$ , we have

$$\sum_{i=1}^n \log \left( \frac{u_h(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})}{\mathcal{T}_h^{\pi} l_{h+1}(z_{h,i}^{l,\pi} | x_{h,i}, a_{h,i})} \right) \leq \ln(e^4/\delta).$$

690 To conclude, apply union bound to get this result for all brackets.  $\square$

691 For the remainder of this section, we assume the policy class  $\Pi$  is finite. However, it is possible  
 692 to extend our results using policy covers in the Hamming distance; in that case,  $\log|\Pi|$  would be  
 693 replaced by the log covering number or entropy integral of  $\Pi$  [as in Zhou et al., 2023, Kallus et al.,  
 694 2022]. We note that for the *online* case, we rely on the assumption that for any  $f \in \mathcal{F}$  we have  
 695  $\pi^f \in \Pi$ , where recall that  $\pi_h^f(x) = \arg \min_a \bar{f}_h(x, a)$ . This is because  $\mathcal{T}^{\star,D}$  is not a contraction  
 696 so we cannot operate with  $\mathcal{T}^{\star,D}$  directly and instead operate with  $\mathcal{T}^{\pi^f,D}$ . We highlight that this  
 697 assumption is automatically satisfied in tabular MDPs, since the whole policy space is finite, and  
 698  $\log|\Pi| = \mathcal{O}(X \log(A))$  is lower order compared to log of the bracketing entropy of  $\mathcal{F}_{tab}$ , which is  
 699  $\mathcal{O}(X^2 A^2)$ . In contrast, in non-distributional methods such as GOLF, the regular Bellman optimality  
 700 operator is a contraction so standard Lipschitz arguments for covering go through. We note that it is  
 701 also possible to construct covers of  $\mathcal{F}$  in the Hellinger distance, but the metric entropy of  $\mathcal{F}_{tab}$  seems  
 702 to be on the same order as its bracketing entropy.

703 We now describe the version space construction for general  $\mathcal{F}$ , first for the online setting. Fix any  $k$ ,  
 704 and define the set

$$\mathcal{F}_{f^\downarrow, \pi, h} = \left\{ f \in \mathcal{F} : \sum_{i=1}^k \log f_h(z_{h,i}^{f^\downarrow, \pi} | x_{h,i}, a_{h,i}) \geq \sum_{i=1}^k \log \text{MLE}_h^{f^\downarrow, \pi}(z_{h,i}^{f^\downarrow, \pi} | x_{h,i}, a_{h,i}) - \beta \right\}$$

705 Then, construct the version space as

$$\mathcal{F}_k = \{ f \in \mathcal{F} : f_h \in \mathcal{F}_{f^\downarrow, \pi^f, h}, \forall h \in [H] \}.$$

**Theorem F.2.** Fix any  $\delta \in (0, 1)$  and suppose [Assumption 5.1](#). Set  $\beta = \log(KH \cdot N_{[]} (K^{-1}, \mathcal{F}, \|\cdot\|_\infty) |\Pi| / \delta)$ . Then, w.p. at least  $1 - \delta$ , the following holds:

- (1) The optimal cost distribution is in the version space, i.e.,  $Z^* \in \mathcal{F}_k$ .  
 (2) For all  $f \in \mathcal{F}_k$  and  $h \in [H]$ ,

$$\sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{*,D} f_{h+1}(x_h, a_h)) \right] \leq 60\beta.$$

*Proof.* First, we want to verify that  $Z^* \in \mathcal{F}_k$ . Let  $f^\downarrow$  be the lower bracket of  $Z^*$  and set  $g = \text{MLE}_h^{f^\downarrow, \pi^*} \in \mathcal{F}$ ; note  $\pi^* = \pi^{Z^*}$ . By [Theorem F.1](#), we have  $\sum_{i=1}^k \log \text{MLE}_h^{f^\downarrow, \pi^*}(z_{h,i}^{f^\downarrow, \pi^*} \mid x_{h,i}, a_{h,i}) - \log \mathcal{T}_h^{\pi^*,D} f_{h+1}^\downarrow(z_{h,i}^{f^\downarrow, \pi^*} \mid x_{h,i}, a_{h,i}) \leq \mathcal{O}(\beta)$ . Therefore, noting that  $Z_h^* = \mathcal{T}_h^{\pi^*,D} Z_{h+1}^* \geq \mathcal{T}_h^{\pi^*,D} f_{h+1}^\downarrow$  shows that  $Z_h^* \in \mathcal{F}_{f^\downarrow, \pi^*, h}$  for every  $h$ , implying that  $Z^* \in \mathcal{F}_k$ .

For the second claim, fix any  $f \in \mathcal{F}_k$  and  $h \in [H]$ . Then,

$$\begin{aligned} & \sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{*,D} f_{h+1}(x_h, a_h)) \right] \\ &= \sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{\pi^f, D} f_{h+1}(x_h, a_h)) \right] \\ &\leq 2 \sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{\pi^f, D} \tilde{f}_{h+1}^\downarrow(x_h, a_h)) + H^2(\mathcal{T}_h^{\pi^f, D} \tilde{f}_{h+1}^\downarrow(x_h, a_h) \parallel \mathcal{T}_h^{\pi^f, D} f_{h+1}(x_h, a_h)) \right] \\ &\leq 2(28\beta + 3k\varepsilon). \end{aligned}$$

The  $\beta$  comes from [Theorem E.5](#), and for  $\varepsilon$ , we used the fact that  $H^2 \leq H \leq TV$ , and

$$\begin{aligned} & \sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ TV(\mathcal{T}_h^{\pi^f, D} \tilde{f}_{h+1}^\downarrow(x_h, a_h) \parallel \mathcal{T}_h^{\pi^f, D} f_{h+1}(x_h, a_h)) \right] \\ &= \sum_{i=1}^k \mathbb{E}_{\pi^i} \int_z \left| \mathcal{T}_h^{\pi^f, D} \tilde{f}_{h+1}^\downarrow(z \mid x_h, a_h) - \mathcal{T}_h^{\pi^f, D} f_{h+1}(z \mid x_h, a_h) \right| \\ &= \sum_{i=1}^k \mathbb{E}_{\pi^i} \int_z \sum_{c, x'} \nu(c, x' \mid x_h, a_h) \left| \tilde{f}_{h+1}^\downarrow(z - c \mid x', \pi^f(x')) - f_{h+1}(z - c \mid x', \pi^f(x')) \right| \\ &\leq \sum_{i=1}^k 3\varepsilon = 3k\varepsilon, \end{aligned}$$

since for any  $x, a$ , we have  $\int_z \left| \tilde{f}_{h+1}^\downarrow(z \mid x, a) - f_{h+1}(z \mid x, a) \right| \leq 3\varepsilon$ . There are two cases. If  $\tilde{f}_{h+1}^\downarrow(z \mid x, a) \geq f_{h+1}(z \mid x, a)$ , then  $\tilde{f}_{h+1}^\downarrow(z \mid x, a) - f_{h+1}(z \mid x, a) \leq (1 - \varepsilon)^{-1} f_{h+1}^\downarrow(z \mid x, a) - f_{h+1}(z \mid x, a) \leq 2\varepsilon f_{h+1}(z \mid x, a)$  since  $(1 - \varepsilon)^{-1} \leq 1 + 2\varepsilon$ . If  $\tilde{f}_{h+1}^\downarrow(z \mid x, a) < f_{h+1}(z \mid x, a)$ , then  $f_{h+1}(z \mid x, a) - \tilde{f}_{h+1}^\downarrow(z \mid x, a) \leq f_{h+1}(z \mid x, a) - f_{h+1}^\downarrow(z \mid x, a) \leq \varepsilon$ . Thus,  $\int_z \max(2\varepsilon f_{h+1}(z \mid x, a), \varepsilon) \leq \int_z 2\varepsilon f_{h+1}(z \mid x, a) + \varepsilon = 3\varepsilon$ . Thus, setting  $\varepsilon = 1/K$  gives

$$\sum_{i=1}^k \mathbb{E}_{\pi^i} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{*,D} f_{h+1}(x_h, a_h)) \right] \leq 59\beta.$$

□

For the offline setting, fix any  $\pi$  and define its general version space as,

$$\mathcal{F}_\pi = \{f \in \mathcal{F} : f_h \in \mathcal{F}_{f^\downarrow, \pi, h}, \forall h \in [H]\}.$$

723

724 **Theorem F.3.** Fix any  $\delta \in (0, 1)$  and suppose [Assumption 5.1](#). Set  $\beta = \log(H|\Pi| \cdot$   
725  $N_{\Pi}((n|\mathcal{Y}|)^{-1}, \mathcal{F}, \|\cdot\|_{\infty})/\delta)$ . Then, w.p. at least  $1 - \delta$ , the following holds for all policies  $\pi \in \Pi$ :

726 (1) The policy cost distribution is in the version space, i.e.,  $Z^{\pi} \in \mathcal{F}_{\pi}$ .

727 (2) Any function in the version space has bounded triangular discrimination with the ground  
728 truth data-generating distribution, i.e., for all  $f \in \mathcal{F}_{\pi}$  and  $h \in [H]$ ,

$$\mathbb{E}_{\nu_h} \left[ H^2(f_h(x_h, a_h) \parallel \mathcal{T}_h^{\pi, D} f_{h+1}(x_h, a_h)) \right] \leq 60\beta N^{-1}.$$

729 *Proof.* The proof is the same as in [Theorem F.2](#), but instead of  $\pi^f$ , we fix any  $\pi$ . □

## 730 G Proofs for Online RL

### 731 G.1 Preliminary Lemmas

732 **Lemma G.1.** *For any policy  $\pi$ , conditional distribution  $d$  and  $h \in [H]$ , we have*

$$\begin{aligned}\overline{\mathcal{T}_h^{\pi,D}d(x,a)} &= \mathcal{T}_h^{\pi}\bar{d}(x,a), \\ \overline{\mathcal{T}_h^{\star,D}d(x,a)} &= \mathcal{T}_h^{\star}\bar{d}(x,a).\end{aligned}$$

*Proof.*

$$\begin{aligned}\overline{\mathcal{T}_h^{\pi,D}d(x,a)} &= \mathbb{E}_{y \sim \mathcal{T}_h^{\pi,D}d(x,a)}[y] \\ &= \mathbb{E}_{c \sim C_h(x,a), x' \sim P_h(x,a), a' \sim \pi_{h+1}(x'), y' \sim d(x',a')}[c + y'] \\ &= \bar{C}_h(x,a) + \mathbb{E}_{x' \sim P_h(x,a), a' \sim \pi_{h+1}(x'), y' \sim d(x',a')}[y'] \\ &= \bar{C}_h(x,a) + \mathbb{E}_{x' \sim P_h(x,a), a' \sim \pi_{h+1}(x')}[\bar{d}(x',a')] \\ &= \mathcal{T}_h^{\pi}\bar{d}(x,a).\end{aligned}$$

733

$$\begin{aligned}\overline{\mathcal{T}_h^{\star,D}d(x,a)} &= \mathbb{E}_{y \sim \mathcal{T}_h^{\star,D}d(x,a)}[y] \\ &= \mathbb{E}_{c \sim C_h(x,a), x' \sim P_h(x,a), a' = \arg \min_{\bar{a}} \bar{d}(x', \bar{a}), y' \sim d(x', a')}[c + y'] \\ &= \bar{C}_h(x,a) + \mathbb{E}_{x' \sim P_h(x,a), a' = \arg \min_{\bar{a}} \bar{d}(x', \bar{a}), y' \sim d(x', a')}[y'] \\ &= \bar{C}_h(x,a) + \mathbb{E}_{x' \sim P_h(x,a), a' = \arg \min_{\bar{a}} \bar{d}(x', \bar{a})}[\bar{d}(x', a')] \\ &= \bar{C}_h(x,a) + \mathbb{E}_{x' \sim P_h(x,a)}\left[\min_{a'} \bar{d}(x', a')\right] \\ &= \mathcal{T}_h^{\star}\bar{d}(x,a).\end{aligned}$$

734

□

735 **Lemma G.2** (Performance Difference Lemma (PDL)). *For any  $f : (\mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R})^H$  and policies*  
736  *$\pi, \pi'$ , we have*

$$V^{\pi} - \mathbb{E}_{a \sim \pi'(x_1)}[f_1(x_1, a)] = \sum_{h=1}^H \mathbb{E}_{\pi} \left[ \mathcal{T}_h^{\pi'} f_{h+1}(x_h, a_h) - f_h(x_h, \pi') \right]. \quad (5)$$

737 *Proof.* We proceed by inducting on the following claim: for all  $h = H+1, H, \dots, 1$ ,

$$V_h^{\pi}(x_h) - f_h(x_h, \pi') = \sum_{t=h}^H \mathbb{E}_{\pi, x_h} \left[ \mathcal{T}_t^{\pi'} f_{t+1}(x_t, a_t) - f_t(x_t, \pi') \right].$$

738 The base case of  $H+1$  is trivially true as everything is 0. Now fix any  $h$  and suppose the IH at  $h+1$   
739 is true. Then

$$\begin{aligned}V_h^{\pi}(x_h) - f_h(x_h, \pi') &= \mathbb{E}_{\pi, x_h} [c_h + V_{h+1}^{\pi}(x_{h+1}) - f_{h+1}(x_{h+1}, \pi') + f_{h+1}(x_{h+1}, \pi') - f_h(x_h, \pi')] \\ &= \mathbb{E}_{\pi, x_h} [V_{h+1}^{\pi}(x_{h+1}) - f_{h+1}(x_{h+1}, \pi')] + \mathbb{E}_{\pi, x_h} [c_h + f_{h+1}(x_{h+1}, \pi') - f_h(x_h, \pi')].\end{aligned}$$

740 By the IH, the first term is equal to  $\sum_{t=h+1}^H \mathbb{E}_{\pi, x_h} [\mathcal{T}_t^{\pi'} f_{t+1}(x_t, a_t) - f_t(x_t, \pi')]$ . The second term  
741 is exactly  $\mathbb{E}_{\pi, x_h} [\mathcal{T}_h^{\pi'} f_{h+1}(x_h, a_h) - f_h(x_h, \pi')]$ , which concludes the proof. □



## 742 G.2 General Regret and PAC Bounds

743 For our analysis, we define a complexity measure inspired by the Sequential Extrapolation Coefficient  
 744 (SEC) of [Xie et al. \[2023\]](#). The SEC measures how well a function can be extrapolated on the  $k$ -th  
 745 episode, using data from the first  $k - 1$  episodes, and has interesting connections to the coverability  
 746 of the MDP. Recall the definition of SEC for function class  $\Psi$ , distribution class  $\mathcal{D}$ , both indexed by  
 747  $h$ , and number of episodes  $K$ :

$$\text{SEC}(\Psi, \mathcal{D}, K) = \max_{\forall k: f^{(k)} \in \Psi, d^{(k)} \in \mathcal{D}} \sum_{k=1}^K \frac{(\mathbb{E}_{d^{(k)}}[f^{(k)}(z)])^2}{1 \vee \sum_{i < k} \mathbb{E}_{d^{(i)}}[f^{(k)}(z)^2]}.$$

748 [Xie et al. \[2023\]](#) showed that the regret of standard (non-distributional) GOLF can be captured by the  
 749 SEC. However, for our distributional algorithm, we need to define a slightly different term, which we  
 750 call the *Linear SEC* (LSEC):

$$\text{LSEC}(\Psi, \mathcal{D}, K) := \max_{\forall k: f^{(k)} \in \Psi, d^{(k)} \in \mathcal{D}} \sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}}[f^{(k)}(z)]}{1 \vee \sum_{i < k} \mathbb{E}_{d^{(i)}}[f^{(k)}(z)]}. \quad (6)$$

751 The difference with the SEC is that our quantity does not have squares, hence we call it “linear”. By  
 752 Jensen’s inequality, we have  $\text{SEC}(\{f^2 : f \in \Psi\}, \mathcal{D}, K) \leq \text{LSEC}(\Psi, \mathcal{D}, K)$ , which shows that our  
 753 LSEC is in general a larger quantity. Nonetheless, we will show that it is controlled for tabular MDPs.  
 754 For our regret bound, the function class and distribution class are instantiated as, for each  $h$ ,

$$\begin{aligned} \mathcal{D}_h(\Pi) &= \{z \mapsto d^\pi(z) : \pi \in \Pi\} \\ \Psi_h &= \{z \mapsto D_\Delta(f(z) \parallel \mathcal{T}^{\star, D} f(z)) : f \in \mathcal{F}\}, \end{aligned} \quad (7)$$

755 where  $z = (s, a)$ . So let us denote  $\text{LSEC}(K) = \max_h \text{LSEC}(\Psi_h, \mathcal{D}_h(\Pi), K)$ . This quantity will  
 756 appear in our small-loss regret bounds.

757 We can also define V-type analogs of LSEC, which we will use for obtaining small-loss PAC bounds  
 758 for latent variable models. The key difference in the V-type LSEC is that the distributions in  $\mathcal{D}_h(\Pi)$   
 759 are in the form  $d^\pi(s) \cdot \text{unif}(a)$ , i.e.,

$$\begin{aligned} \mathcal{D}_{h,v}(\Pi) &= \{(s, a) \mapsto d^\pi(s)/A : \pi \in \Pi\} \\ \text{LSEC}_v(\Psi, \mathcal{D}, K) &= \max_h \text{LSEC}(\Psi_h, \mathcal{D}_{h,v}(\Pi), K). \end{aligned} \quad (8)$$

760 We now prove the our main regret bound.

761 **Theorem G.3.** Assume [Assumption 5.1](#). Fix any  $\delta \in (0, 1)$  and set  $\beta = \log(HK|\mathcal{F}|/\delta)$  and  
 762  $\beta' = 60\beta$ . Then, w.p. at least  $1 - \delta$ , running O-DISCO ([Algorithm 4](#)) with  $\text{UAE} = \text{FALSE}$  yields  
 763 the following small-loss regret bound,

$$\text{Regret}_{\text{O-DISCO}}(K) \leq 5H\sqrt{KV^\star \text{LSEC}(K)\beta'} + 18H^2 \text{LSEC}(K)\beta'.$$

764 If instead  $\text{UAE} = \text{TRUE}$ , the outputted policy  $\bar{\pi}$  enjoys the following small-loss PAC bound,

$$V^{\bar{\pi}} - V^\star \leq 5H\sqrt{\frac{AV^\star \text{LSEC}_v(K)\beta'}{K}} + 18H^2 \frac{A \text{LSEC}_v(K)\beta'}{K}.$$

765 *Proof.* We first prove the regret bound ( $\text{UAE} = \text{FALSE}$ ); the PAC bound follows from the  
 766 same argument. For shorthand, let  $\delta_{h,k}(x, a) := D_\Delta(f_h^{(k)}(x, a) \parallel \mathcal{T}_h^{\star, D} f_{h+1}^{(k)}(x, a))$  and  
 767  $\Delta_k := \sum_{h=1}^H \mathbb{E}_{\pi^k}[\delta_{h,k}(x_h, a_h)]$ . Notice that since  $\pi_{h+1}^k(x) = \arg \min_a \bar{f}_{h+1}^{(k)}(x, a)$ , we have  
 768  $\mathcal{T}_h^{\pi^k, D} f_{h+1}^{(k)}(x, a) = \mathcal{T}_h^{\star, D} f_{h+1}^{(k)}(x, a)$ , so  $\delta_{h,k}(x, a) = D_\Delta(f_h^{(k)}(x, a) \parallel \mathcal{T}_h^{\pi^k, D} f_{h+1}^{(k)}(x, a))$  as well.

769 By [Theorem F.2](#), we have the following two facts for all  $k \in [K]$ ,

- 770 (i) Optimism:  $\min_a \bar{f}_1^{(k)}(x_1, a) \leq V^\star$  (since  $Z^\star \in \mathcal{F}_k$ ) and
- 771 (ii)  $\sum_{i < k} \mathbb{E}_{\pi^i}[\delta_{h,k}(s_h, a_h)] \leq \beta'$  for all  $h$ . If  $\text{UAE} = \text{TRUE}$ , then  $a_h$  is sampled from  $\text{unif}(\mathcal{A})$  rather  
 772 than  $\pi^i$ , i.e., we have  $\sum_{i < k} \mathbb{E}_{s_h \sim \pi^i, a_h \sim \text{unif}(\mathcal{A})}[\delta_{h,k}(s_h, a_h)] \leq \beta'$ , where  $\beta' \lesssim \beta$ . [Theorem F.2](#) and  
 773 the fact that  $D_\Delta \leq 4H^2$  certifies that  $\beta' = 240\beta$  is sufficient.

774 Now, fix any episode  $k \in [K]$ .

$$\begin{aligned}
& V^{\pi^k} - V^* \\
& \leq V^{\pi^k} - \min_a \bar{f}_1^{(k)}(x_1, a) && \text{(Fact (i))} \\
& = \sum_{h=1}^H \mathbb{E}_{\pi^k} \left[ \mathcal{T}_h^{\pi^k} \bar{f}_{h+1}^{(k)}(x_h, a_h) - \bar{f}_h^{(k)}(x_h, \pi_h^k(x_h)) \right] && \text{(PDL Lemma G.2)} \\
& = \sum_{h=1}^H \mathbb{E}_{\pi^k} \left[ \overline{\mathcal{T}_h^{\pi^k, D} f_{h+1}^{(k)}}(x_h, a_h) - \bar{f}_h^{(k)}(x_h, a_h) \right] && \text{(Lemma G.1)} \\
& \leq \sum_{h=1}^H \sqrt{\mathbb{E}_{\pi^k} \left[ 4\bar{f}_h^{(k)}(x_h, a_h) + \delta_{h,k}(x_h, a_h) \right]} \cdot \sqrt{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]} && \text{(Eq. (\Delta_2))} \\
& \leq \sum_{h=1}^H \sqrt{4eV^{\pi^k} + 17H \sum_{t=h}^H \mathbb{E}_{\pi^k} [\delta_{t,k}(x_t, a_t)]} \cdot \sqrt{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]} \\
& && \text{(Lemma G.4 and } \mathbb{E}_{\pi} [Q_h^{\pi}(s_h, a_h)] \leq V^{\pi}) \\
& \leq \sqrt{4eV^{\pi^k} + 17H\Delta_k} \cdot \sqrt{H\Delta_k} && (\star) \\
& \leq \sqrt{4eHV^{\pi^k}\Delta_k} + 5H\Delta_k \\
& \leq 2\sqrt{H}\eta^{-1}V^{\pi^k} + 2\sqrt{H}\eta\Delta_k + 5H\Delta_k.
\end{aligned}$$

775 In  $\star$ , we used Cauchy Schwartz. Setting  $\eta = 4\sqrt{H}$  and rearranging, we have

$$V^{\pi^k} \leq 2V^* + 16H\Delta_k + 10H\Delta_k \leq 2V^* + 26H\Delta_k.$$

776 Plugging this into  $\star$ , and noting  $104e + 17 \leq 300$ , we have

$$V^{\pi^k} - V^* \leq \sqrt{8eV^* + 300H\Delta_k} \sqrt{H\Delta_k}.$$

777 Thus, summing the instantaneous regrets over all episodes, we get

$$\begin{aligned}
\sum_{k=1}^K V^{\pi^k} - V^* & \leq \sum_{k=1}^K \sqrt{8eV^* + 300H\Delta_k} \sqrt{H\Delta_k} \\
& \leq \sqrt{8eKV^* + 300H \sum_k \Delta_k} \sqrt{H \sum_k \Delta_k} && \text{(Cauchy-Schwartz)} \\
& \leq 5\sqrt{HKV^* \sum_k \Delta_k} + 18H \sum_k \Delta_k.
\end{aligned}$$

778 Finally it remains the bound the sum of  $\Delta_k$ ,

$$\begin{aligned}
\sum_{k=1}^K \Delta_k & = \sum_{h=1}^H \sum_{k=1}^K \frac{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]}{1 \vee \sum_{i=1}^{k-1} \mathbb{E}_{\pi^i} [\delta_{h,k}(s_h, a_h)]} \cdot \left( 1 \vee \sum_{i=1}^{k-1} \mathbb{E}_{\pi^i} [\delta_{h,k}(s_h, a_h)] \right) \\
& \leq H \text{LSEC}(K) \cdot \beta'. && \text{(Fact (ii))}
\end{aligned}$$

779 If UAE=TRUE, we instead bound the sum of  $\Delta_k$  using the V-type LSEC:

$$\begin{aligned}
\sum_{k=1}^K \Delta_k & \leq \sum_{h=1}^H \sum_{k=1}^K \frac{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]}{1 \vee \sum_{i=1}^{k-1} \mathbb{E}_{s_h \sim \pi^i, a_h \sim \text{unif}(\mathcal{A})} [\delta_{h,k}(s_h, a_h)]} \cdot \left( 1 \vee \sum_{i=1}^{k-1} \mathbb{E}_{s_h \sim \pi^i, a_h \sim \text{unif}(\mathcal{A})} [\delta_{h,k}(s_h, a_h)] \right) \\
& \leq AH \text{LSEC}_v(K) \cdot \beta'. && \text{(Fact (ii))}
\end{aligned}$$

780 This concludes the proof for both the regret and PAC bounds.  $\square$

781 **Lemma G.4** (Self-bounding lemma). *Let  $f \in \mathcal{F}$  and let  $\pi$  be any policy. Let us denote  $\delta_h(x, a) :=$   
782  $D_\Delta(f_h(x, a) \parallel \mathcal{T}_h^{\pi, D} f_{h+1}(x, a))$ . Then, for all  $h \in [H]$ , for all  $x_h, a_h$ , we have*

$$\bar{f}_h(x_h, a_h) \leq eQ_h^\pi(x_h, a_h) + 4H \sum_{t=h}^H \mathbb{E}_{\pi, x_h, a_h} [\delta_t(x_t, a_t)].$$

783 *Proof.* We prove the following refined subclaim inductively: for all  $h \in [H]$ , for all  $x_h, a_h$ , we have

$$\bar{f}_h(x_h, a_h) \leq \sum_{t=h}^H \left(1 + \frac{1}{H}\right)^{t-h} \mathbb{E}_{\pi, x_h, a_h} [\bar{c}_t(x_t, a_t) + 2H\delta_t(x_t, a_t)]. \quad (\text{IH})$$

784 For  $H + 1$  this is trivially true. Now fix any  $h$  and suppose IH is true for  $h + 1$ . By Eq. ( $\Delta_2$ ), for any  
785  $h, x_h, a_h$ , we have,

$$\begin{aligned} \bar{f}_h(x_h, a_h) - \mathcal{T}_h^\pi \bar{f}_{h+1}(x_h, a_h) &\leq \sqrt{4\mathcal{T}_h^\pi \bar{f}_{h+1}(x_h, a_h) + \delta_h(x_h, a_h)} \sqrt{\delta_h(x_h, a_h)} \\ &\leq \sqrt{4\mathcal{T}_h^\pi \bar{f}_{h+1}(x_h, a_h) \delta_h(x_h, a_h) + \delta_h(x_h, a_h)} \\ &\leq \frac{1}{H} \mathcal{T}_h^\pi \bar{f}_{h+1}(x_h, a_h) + (H + 1) \delta_h(x_h, a_h). \quad (\text{AM-GM}) \end{aligned}$$

786 In particular, we have that

$$\begin{aligned} &\bar{f}_h(x_h, a_h) \\ &\leq \left(1 + \frac{1}{H}\right) \mathcal{T}_h^\pi \bar{f}_{h+1}(x_h, a_h) + 2H\delta_h(x_h, a_h) \\ &= \left(1 + \frac{1}{H}\right) \left(\bar{c}_h(x_h, a_h) + \mathbb{E}_{x_{h+1} \sim P_h^*(x_h, a_h)} [\bar{f}_{h+1}(x_{h+1}, \pi)]\right) + 2H\delta_h(x_h, a_h) \\ &\leq \left(1 + \frac{1}{H}\right) \left(\bar{c}_h(x_h, a_h) + \mathbb{E}_{x_{h+1} \sim P_h^*(x_h, a_h)} \left[ \sum_{t=h+1}^H \left(1 + \frac{1}{H}\right)^{t-h-1} \mathbb{E}_{\pi, x_{h+1}} [\bar{c}_t(x_t, a_t) + 2H\delta_t(x_t, a_t)] \right] \right) \\ &\quad \quad \quad (\text{IH}) \\ &\quad + 2H\delta_h(x_h, a_h), \end{aligned}$$

787 which proves the inductive claim. Noting that  $\sum_{t=1}^H (1 + 1/H)^t \leq e$ , we have proven the lemma.  $\square$

### 788 G.3 Bounding the LSEC

789 In this section, we show that the LSEC quantity is bounded for tabular MDPs and latent variable  
790 models. First, recall the notion of Coverability from Xie et al. [2023],

$$C_{\text{Cov}} := \inf_{\mu} \max_{\pi} \max_{h, x, a} \frac{d_h^\pi(x, a)}{\mu_h(x, a)}.$$

791 Let  $\mu^*$  be the measure that realizes this infimum.  $C_{\text{Cov}}$  was shown to be equivalent to  
792  $\max_h \sum_{x, a} \sup_{\pi} d_h^\pi(x, a)$  by Xie et al. [2023, Lemma 3]. For example, in tabular MDPs with  
793  $X$  states and  $A$  actions, we have  $C_{\text{Cov}} \leq XA$ , and in low-rank MDPs (and hence latent variable  
794 models) with rank  $d$ , we have  $C_{\text{Cov}} \leq d$  [Huang et al., 2023, Proposition 3].

795 **Bounding the LSEC in Tabular MDPs** First, consider any function class  $\Psi$  and distribution  
796 class  $\mathcal{D}$ . For all  $k$ , let  $f^{(k)} \in \Psi$  and  $d^{(k)} \in \mathcal{D}$ . Define  $\tilde{d}^{(k)} = \sum_{i < k} d^{(i)}$  and  $\tau(z) :=$   
797  $\min\{k \mid \tilde{d}^{(k)}(z) \geq C_{\text{Cov}} \mu^*(z)\}$ . Then, for any  $f \in \Psi$  and  $d \in \mathcal{D}$ , we have

$$\begin{aligned} &\sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}} [f^{(k)}]}{1 \vee \sum_{i < k} \mathbb{E}_{d^{(i)}} [f^{(k)}]} \\ &= \underbrace{\sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}} [f^{(k)}(z) \mathbb{I}[k < \tau(z)]]}{1 \vee \sum_{i < k} \mathbb{E}_{d^{(i)}} [f^{(k)}]}}_{\text{Term 1}} + \underbrace{\sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}} [f^{(k)}(z) \mathbb{I}[k \geq \tau(z)]]}{1 \vee \sum_{i < k} \mathbb{E}_{d^{(i)}} [f^{(k)}]}}_{\text{Term 2}}. \end{aligned}$$

798 Focusing on Term 1, we have it is at most,

$$\sum_{k=1}^K \mathbb{E}_{d^{(k)}} \left[ f^{(k)}(z) \mathbb{I}[k < \tau(z)] \right] \leq \sum_{k=1}^K \mathbb{E}_{d^{(k)}} [\mathbb{I}[k < \tau(z)]] \leq 2C_{\text{Cov}},$$

799 by the proof of Proposition 13 of [Xie et al. \[2023\]](#).

800 For Term 2, we need to specialize  $\mathcal{D}$ . If the MDP is tabular, we can set  $\mathcal{D}$  as defined in [Eq. \(7\)](#). Then,  
801 for  $z = (x, a)$ ,

$$\begin{aligned} & \sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}} [f^{(k)}(z) \mathbb{I}[k \geq \tau(z)]]}{\sum_{i < k} \mathbb{E}_{d^{(i)}} [f^{(k)}]} \\ &= \sum_{k=1}^K \sum_z \frac{d^{(k)}(z) f^{(k)}(z) \mathbb{I}[k \geq \tau(z)]}{\sum_z \tilde{d}^{(k)}(z) f^{(k)}(z)} \\ &\leq \sum_{k=1}^K \sum_z \frac{d^{(k)}(z) f^{(k)}(z) \mathbb{I}[k \geq \tau(z)]}{\tilde{d}^{(k)}(z) f^{(k)}(z)} \quad (\text{terms are non-negative}) \\ &= \sum_{k=1}^K \sum_z \frac{d^{(k)}(z) \mathbb{I}[k \geq \tau(z)]}{\tilde{d}^{(k)}(z)} \\ &\leq 2 \sum_z \sum_{k=1}^K \frac{d^{(k)}(z) \mathbb{I}[k \geq \tau(z)]}{\tilde{d}^{(k)}(z) + C_{\text{Cov}} \mu^*(z)} \\ &\leq 2 \sum_z 2 \log(K+1) \quad (\text{Xie et al. [2023, Lemma 4]}) \\ &= 4Z \log(K+1). \end{aligned}$$

802 Since the MDP is tabular we have  $Z = XA$ . We have proven the following lemma,

803 **Lemma G.5.** *Suppose the MDP is tabular. Then, for any  $\Psi, K$ , we have*

$$\text{LSEC}(\Psi, \mathcal{D}_h(\Pi), K) \in \mathcal{O}(XA \log(K)).$$

804 Combining this with [Theorem G.3](#) directly implies [Theorem 5.2](#).

805 **Bounding V-type LSEC in Latent Variable Models** Now suppose the MDP is a latent variable  
806 model (LVM), *i.e.*, an MDP with small non-negative rank  $d$  [Modi et al. \[2021\]](#). The sampling  
807 procedure for latent variable model is, start with a distribution over  $d$  latent states  $p_1$ , sample an  
808 unobserved latent state  $s_1 \sim p_1$ , observe  $x_1 \sim o(s_1)$ , take action  $a_1 \sim \pi_1(s_1)$  and transition to the  
809 next distribution of latent states  $p_2$ . This process repeats  $H$  times. Note that the observation set  $\mathcal{X}$   
810 can be very large or infinite, so instead of having a bound that depends on  $X$ , we'd like to depend on  
811 the number of latent states  $S$ . To do so, we make a simple modification to our previous argument.

812 Set  $\tau(s, a) = \min \left\{ k \mid \tilde{d}^{(k)}(s, a) \geq C_{\text{Cov}} \mu^*(s, a) \right\}$ , where we've abused notation to use  $s$  as input  
813 instead of  $x$ , denoting that we are considering distributions over latent states rather than observa-  
814 tions. For any distribution, we have  $d(x, a) = o(x \mid s) d(s, a)$  where  $s$  is the encoded latent state  
815 corresponding to  $x$ . Crucially,  $\tau$  depends on  $s$  rather than  $x$ .

816 In this case, we can take  $\mathcal{D}$  as the V-type distributions  $\mathcal{D}_{h,v}(\Pi)$ . So  $d^{(k)}(s, a) = d^{\pi^k}(s)/A$  and we  
817 can bound Term 2 as follows,

$$\begin{aligned}
& \sum_{k=1}^K \frac{\mathbb{E}_{d^{(k)}}[f^{(k)}(x, a) \mathbb{I}[k \geq \tau(s, a)]]}{\sum_{i < k} \mathbb{E}_{d^{(i)}}[f^{(k)}]} \\
&= \sum_{k=1}^K \sum_{s, a} \frac{d^{(k)}(s, a) \mathbb{E}_{x \sim o(s)}[f^{(k)}(x, a)] \mathbb{I}[k \geq \tau(s, a)]}{\sum_{s, a} \tilde{d}^{(k)}(s, a) \mathbb{E}_{x \sim o(s)}[f^{(k)}(x, a)]} \\
&\leq \sum_{k=1}^K \sum_{s, a} \frac{d^{(k)}(s, a) \mathbb{E}_{x \sim o(s)}[f^{(k)}(x, a)] \mathbb{I}[k \geq \tau(s, a)]}{\tilde{d}^{(k)}(s, a) \mathbb{E}_{x \sim o(s)}[f^{(k)}(x, a)]} \quad (\text{terms are non-negative}) \\
&= \sum_{k=1}^K \sum_{s, a} \frac{d^{(k)}(s, a) \mathbb{I}[k \geq \tau(s, a)]}{\tilde{d}^{(k)}(s, a)} \\
&\leq 2 \sum_{s, a} \sum_{k=1}^K \frac{d^{(k)}(s, a) \mathbb{I}[k \geq \tau(s, a)]}{\tilde{d}^{(k)}(s, a) + C_{\text{Cov}} \mu^*(s, a)} \\
&\leq 2 \sum_{s, a} 2 \log(K+1) \quad (\text{Xie et al. [2023, Lemma 4]}) \\
&= 4SA \log(K+1).
\end{aligned}$$

818 We highlight that this argument only works for the V-type LSEC, since the uniform action  $a$  does  
819 not depend on the observation generating process,  $x \sim o(s)$ , while the action from the Q-type LSEC  
820 does. This dependence in the Q-type LSEC is what prevents us from doing the decomposition in the  
821 first step. This is why uniform action exploration is needed for our theory to extend to latent variable  
822 models. Thus, we've shown the following lemma,

823 **Lemma G.6.** *Suppose the MDP is a latent variable model. Then, for any  $\Psi, K$ , we have*

$$\text{LSEC}_v(\Psi, \mathcal{D}_{h,v}(\Pi), K) \in \mathcal{O}(SA \log(K)).$$

824 Combining this with [Theorem G.3](#) directly implies [Theorem 5.4](#).

## H Proofs for Offline RL

**Theorem 6.1** (Small-Loss PAC bound for P-DISCO). Assume [Assumption 5.1](#). Fix any  $\delta \in (0, 1)$  and set  $\beta = \log(H|\Pi||\mathcal{F}|/\delta)$ . Then, w.p. at least  $1 - \delta$ , P-DISCO learns a policy  $\hat{\pi}$  such that for any comparator policy  $\tilde{\pi} \in \Pi$ , we have

$$V^{\hat{\pi}} - V^{\tilde{\pi}} \leq 9H \sqrt{\frac{C^{\tilde{\pi}} V^{\tilde{\pi}} \beta}{N}} + \frac{30H^2 C^{\tilde{\pi}} \beta}{N}.$$

*Proof of Theorem 6.1.* For shorthand, let  $\delta_h^\pi(x, a) = D_\Delta(f_h^\pi(x, a) \parallel \mathcal{T}_h^{\pi, D} f_{h+1}^\pi(x, a))$  and  $\Delta^\pi = \sum_{h=1}^H \mathbb{E}_\pi[\delta_h^\pi(x_h, a_h)]$ . Also, let  $f(x, \pi) = \mathbb{E}_{a \sim \pi(x)}[f(x, a)]$ .

By [Theorem F.3](#), we have the following two facts, for all  $\pi \in \Pi$ ,

- (i) Pessimism:  $V^\pi \leq \bar{f}_1^\pi(x_1, \pi)$  (since  $Z^\pi \in \mathcal{F}_\pi$ ) for all  $\pi \in \Pi$ , and
- (ii)  $\mathbb{E}_{\nu_h}[\delta_h^\pi(x_h, a_h)] \leq \beta' N^{-1}$  for all  $h$  where [Theorem F.3](#) and the fact that  $D_\Delta \leq 4H^2$  certifies that  $\beta' = 240\beta$  is sufficient.

With these two facts, we can bound the suboptimality of  $\hat{\pi}$  as follows:

$$\begin{aligned} V^{\hat{\pi}} - V^{\tilde{\pi}} &\leq \bar{f}_1^{\hat{\pi}}(x_1, \hat{\pi}) - V^{\tilde{\pi}} && \text{(Fact (i))} \\ &\leq \bar{f}_1^{\tilde{\pi}}(x_1, \tilde{\pi}) - V^{\tilde{\pi}} && \text{(Policy selection scheme in Algorithm 3 (Line 4))} \\ &= \sum_{h=1}^H \mathbb{E}_{\tilde{\pi}} \left[ \bar{f}_h^{\tilde{\pi}}(x_h, \tilde{\pi}) - \mathcal{T}_h^{\tilde{\pi}} \bar{f}_{h+1}^{\tilde{\pi}}(x_h, a_h) \right] && \text{(PDL Lemma G.2)} \\ &\leq \sum_{h=1}^H \sqrt{\mathbb{E}_{\tilde{\pi}} [4\bar{f}_h^{\tilde{\pi}}(x_h, a_h) + \delta_h^{\tilde{\pi}}(x_h, a_h)]} \sqrt{\mathbb{E}_{\tilde{\pi}} [\delta_h^{\tilde{\pi}}(x_h, a_h)]} && \text{(Eq. (\Delta_2))} \\ &\leq \sum_{h=1}^H \sqrt{4eV^{\tilde{\pi}} + 17H \sum_{t=h}^H \mathbb{E}_{\tilde{\pi}} [\delta_t^{\tilde{\pi}}(x_t, a_t)]} \sqrt{\mathbb{E}_{\tilde{\pi}} [\delta_h^{\tilde{\pi}}(x_h, a_h)]} && \text{(Lemma G.4)} \\ &\leq \sqrt{4eV^{\tilde{\pi}} + 17H\Delta^{\tilde{\pi}}} \sqrt{H\Delta^{\tilde{\pi}}} \\ &\leq 4\sqrt{HV^{\tilde{\pi}}\Delta^{\tilde{\pi}}} + 5H\Delta^{\tilde{\pi}}. \end{aligned}$$

Finally, we can bound  $\Delta^{\tilde{\pi}}$  by a change of measure,

$$\begin{aligned} \Delta^{\tilde{\pi}} &= \sum_{h=1}^H \mathbb{E}_{\tilde{\pi}} [\delta_h^{\tilde{\pi}}(x_h, a_h)] \\ &\leq C^{\tilde{\pi}} \sum_{h=1}^H \mathbb{E}_{\nu_h} [\delta_h^{\tilde{\pi}}(x_h, a_h)] \\ &\leq C^{\tilde{\pi}} H \cdot \beta' N^{-1}. \end{aligned} \tag{Fact (ii)}$$

Therefore,

$$V^{\hat{\pi}} - V^{\tilde{\pi}} \leq 4H \sqrt{\frac{C^{\tilde{\pi}} V^{\tilde{\pi}} \beta'}{N}} + \frac{5H^2 C^{\tilde{\pi}} \beta'}{N}.$$

838

□

## I Extension: Small-Return Bounds

In this section, we show that O-DISCO and P-DISCO can also be used to obtain small-return bounds. Compared to the algorithms presented in the main text for minimizing cost, we simply have to replace min with max (and vice versa) for maximizing reward, *i.e.*, see [Appendix B](#) and enable the SMALLRETURN flag. The proofs are also largely the same, with slight changes to the first few steps.

**Theorem I.1.** Assume [Assumption 5.1](#) and suppose we want to maximize returns (instead of minimize cost), so enable the SMALLRETURN flag. Fix any  $\delta \in (0, 1)$  and set  $\beta = \log(HK|\mathcal{F}|/\delta)$  and  $\beta' = 60\beta$ . Then, w.p. at least  $1 - \delta$ , running O-DISCO ([Algorithm 4](#)) with  $\text{UAE} = \text{FALSE}$  yields the following small-loss regret bound,

$$\text{Regret}_{\text{O-DISCO}}(K) \leq 5H\sqrt{KV^*\text{LSEC}(K)\beta'} + 18H^2\text{LSEC}(K)\beta'. \quad (9)$$

If instead  $\text{UAE} = \text{TRUE}$ , the outputted policy  $\bar{\pi}$  enjoys the following small-loss PAC bound,

$$V^* - V^{\bar{\pi}} \leq 5H\sqrt{\frac{AV^*\text{LSEC}_v(K)\beta'}{K}} + 18H^2\frac{A\text{LSEC}_v(K)\beta'}{K}.$$

*Proof.* Adopt the same notation as in the proof of [Theorem G.3](#). By [Theorem F.2](#), we have the following two facts for all  $k \in [K]$ ,

- (i) Optimism:  $V^* \leq \max_a \bar{f}_1^{(k)}(x_1, a)$  (since  $Z^* \in \mathcal{F}_k$ ) and
- (ii)  $\sum_{i < k} \mathbb{E}_{\pi^i}[\delta_{h,k}(s_h, a_h)] \leq \beta'$  for all  $h$ . If  $\text{UAE} = \text{TRUE}$ , then  $a_h$  is sampled from  $\text{unif}(\mathcal{A})$  rather than  $\pi^i$ , *i.e.*, we have  $\sum_{i < k} \mathbb{E}_{s_h \sim \pi^i, a_h \sim \text{unif}(\mathcal{A})}[\delta_{h,k}(s_h, a_h)] \leq \beta'$ , where  $\beta' \lesssim \beta$ . [Theorem F.2](#) certifies that  $\beta' = 60\beta$  is sufficient.

Fix any episode  $k \in [K]$ . Then,

$$\begin{aligned} V^* - V^{\pi^k} &\leq \max_a \bar{f}_1^{(k)}(x_1, a) - V^{\pi^k} && \text{(Fact (i))} \\ &= \sum_{h=1}^H \mathbb{E}_{\pi^k} \left[ \bar{f}_h^{(k)}(x_h, \pi_h^k(x_h)) - \mathcal{T}_h^{\pi^k} \bar{f}_{h+1}^{(k)}(x_h, a_h) \right] && \text{(PDL Lemma G.2)} \\ &= \sum_{h=1}^H \mathbb{E}_{\pi^k} \left[ \bar{f}_h^{(k)}(x_h, a_h) - \overline{\mathcal{T}_h^{\pi^k, D} \bar{f}_{h+1}^{(k)}}(x_h, a_h) \right] && \text{(Lemma G.1)} \\ &\leq \sum_{h=1}^H \sqrt{\mathbb{E}_{\pi^k} \left[ 4\bar{f}_h^{(k)}(x_h, a_h) + \delta_{h,k}(x_h, a_h) \right]} \cdot \sqrt{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]} && \text{(Eq. } (\Delta_2)) \\ &\leq \sum_{h=1}^H \sqrt{4eV^{\pi^k} + 17H \sum_{t=h}^H \mathbb{E}_{\pi^k} [\delta_{t,k}(x_t, a_t)]} \cdot \sqrt{\mathbb{E}_{\pi^k} [\delta_{h,k}(x_h, a_h)]} \\ &\hspace{15em} \text{(Lemma G.4 and } \mathbb{E}_{\pi} [Q_h^{\pi}(s_h, a_h)] \leq V^{\pi}) \\ &\leq \sqrt{4eV^{\pi^k} + 17H\Delta_k} \cdot \sqrt{H\Delta_k} && (\clubsuit) \\ &\leq \sqrt{4eV^* + 17H\Delta_k} \cdot \sqrt{H\Delta_k} \end{aligned}$$

Thus, summing the instantaneous regrets over all episodes, we get

$$\begin{aligned} \sum_{k=1}^K V^{\pi^k} - V^* &\leq \sum_{k=1}^K \sqrt{4eV^* + 17H\Delta_k} \sqrt{H\Delta_k} \\ &\leq \sqrt{4eKV^* + 17H \sum_k \Delta_k} \sqrt{H \sum_k \Delta_k} && \text{(Cauchy-Schwartz)} \\ &\leq 5\sqrt{HKV^* \sum_k \Delta_k} + 18H \sum_k \Delta_k. \end{aligned}$$

The bounds for  $\Delta_k$  are the same as in [Theorem G.3](#).  $\square$

858 In some sense, the proof for the small-returns bound is actually easier than the small-loss bound.  
 859 Recall that in the cost-minimizing setting, we needed to perform a crucial Cauchy-Schwartz step to  
 860 rearrange terms at the step labelled ♣. However, in the reward-maximizing setting, we simply bound  
 861  $V^{\pi^k} \leq V^*$ , without needing to rearrange terms.

862 **Theorem I.2.** Assume [Assumption 5.1](#) and suppose we want to maximize returns (instead of minimize  
 863 cost), so enable the SMALLRETURN flag. Fix any  $\delta \in (0, 1)$  and set  $\beta = \log(H|\Pi||\mathcal{F}|/\delta)$ . Then,  
 864 w.p. at least  $1 - \delta$ , P-DISCO ([Algorithm 4](#)) learns a policy  $\hat{\pi}$  such that for any comparator policy  
 865  $\tilde{\pi} \in \Pi$ , we have

$$V^{\tilde{\pi}} - V^{\hat{\pi}} \leq 9H\sqrt{\frac{C^{\tilde{\pi}}V^{\tilde{\pi}}\beta}{N}} + \frac{30H^2C^{\tilde{\pi}}\beta}{N}.$$

866 *Proof of Theorem I.2.* Adopt the same notation as in the proof of [Theorem 6.1](#). By [Theorem F.3](#), we  
 867 have the following two facts, for all  $\pi \in \Pi$ ,  
 868 (i) Pessimism:  $\bar{f}_1^\pi(x_1, \pi) \leq V^\pi$  (since  $Z^\pi \in \mathcal{F}_\pi$ ) for all  $\pi \in \Pi$ , and  
 869 (ii)  $\mathbb{E}_{\nu_h}[\delta_h^\pi(x_h, a_h)] \leq \beta' N^{-1}$  for all  $h$  where  $\beta' \leq 60\beta$ .  
 870 With these two facts, we can bound the suboptimality of  $\hat{\pi}$  as follows:

$$\begin{aligned} & V^{\tilde{\pi}} - V^{\hat{\pi}} \\ & \leq V^{\tilde{\pi}} - \bar{f}_1^{\tilde{\pi}}(x_1, \hat{\pi}) && \text{(Fact (i))} \\ & \leq V^{\tilde{\pi}} - \bar{f}_1^{\tilde{\pi}}(x_1, \tilde{\pi}) && \text{(Policy selection rule in Line 5)} \\ & = \sum_{h=1}^H \mathbb{E}_{\tilde{\pi}} \left[ \mathcal{T}_h^{\tilde{\pi}} \bar{f}_{h+1}^{\tilde{\pi}}(x_h, a_h) - \bar{f}_h^{\tilde{\pi}}(x_h, \tilde{\pi}) \right] && \text{(PDL Lemma G.2)} \\ & \leq \sum_{h=1}^H \sqrt{\mathbb{E}_{\tilde{\pi}} [4\bar{f}_h^{\tilde{\pi}}(x_h, a_h) + \delta_h^{\tilde{\pi}}(x_h, a_h)]} \sqrt{\mathbb{E}_{\tilde{\pi}} [\delta_h^{\tilde{\pi}}(x_h, a_h)]}. && \text{(Eq. (\Delta_2))} \end{aligned}$$

871 From here, the same argument in the proof of [Theorem 6.1](#) finishes the proof.  $\square$



## J Experiment Details

### Experiment Settings

In our experiments, as outlined in Foster and Krishnamurthy [2021], our  $\gamma$  learning rate at each time step  $t$  is set to  $\gamma_t = \gamma_0 t^p$  where  $\gamma_0$  and  $p$  are hyperparameters. We use batch sizes of 32 samples per episode, and the King County and Prudential experiments run for 5,000 episodes while the CIFAR-100 experiment runs for 15,000.

For each dataset, we select the hyperparameter configuration with the best performance for each algorithm. As we report two metrics, performance over the last 100 episodes and over all episodes, we choose the best hyperparameters for each metric as well. While it is often the same hyperparameters that give the best last 100 episodes and all episodes results for a model, that is not always the case. We use the WandB (Weights and Biases) library to run sweeps over hyperparameters.

### Oracles

For our regression oracles, we use ResNet18 [He et al., 2016], with a modified output layer (so that the output is suited for 100 prediction classes) for CIFAR-100, and a simple 2 hidden-layer neural network for the Prudential Life Insurance and King’s County Housing datasets. For DistCB, the oracle’s output layer has size  $AC$  where  $A$  is the number of actions and  $C$  is the number of potential costs. This is reshaped so that for each action, there are predictions associated with each potential cost, which then have a softmax function applied to them to represent cost probabilities. For SquareCB and FastCB, the output size is  $A$  because there is just a single prediction associated with each action. As per Foster and Krishnamurthy [2021], a sigmoid function is applied to this output layer. All experiments were implemented using PyTorch.

### Datasets

We now provide an overview table as well as additional details and context to our setups for each dataset. Note that the number of items in each dataset in the table is the count after preprocessing.

Datasets			
Dataset	Items	Number of Actions	Number of Costs
CIFAR-100	50,000	100	3
Prudential Life Insurance	59,381	8	9
King County Housing	20,148	100	101

Table 3: Overview of the three datasets and their experimental setups

**Prudential Life Insurance** This dataset is from the Prudential Life Insurance Kaggle competition [Montoya et al., 2015]. It is featured in Farsang et al. [2022], which inspires our experimental setup. The risk level in [8] directly determines the price charged to the customer. Thus, we can consider the chosen risk level as the action taken. If the model overpredicts the risk level, we get a cost of 1.0 because this is considered over charging the customer and not getting a sale. Otherwise, the model’s prediction is charging too little for the customer. To reiterate, the cost in this case is  $.1 * (y - \hat{y})$  where  $y$  is the actual risk level, and  $\hat{y}$  is the predicted risk level.

**King County Housing** The King County housing dataset is also used in Farsang et al. [2022]. An interesting part of the setup is that the cost construction in the case of not overpredicting differs from the Prudential experiment, even though they’re both effectively about predicting a price point. Here, the model’s chosen price is considered the gain, which is why the cost is 1.0 minus the chosen price. On the other hand, in the Prudential experiment, the cost is a linear function of the difference between the chosen value and the actual value.

**CIFAR-100** For the CIFAR-100 experiment, we use the training dataset of 50,000 images as our dataset. The inclusion of the superclass is critical, as it lets us delineate 3 possible costs that DISTCB can learn. Without the super class, the cost construction would be a pure binary of correct vs. incorrect. If this were the case, the ability to test the effectiveness of learning the distribution would be nullified. The distribution would just be whether an action is correct or not, which means our algorithm would essentially be predicting the mean directly.

## 915 **Results**

916 The largest advantages DISTCB had over the next best algorithm were in the Prudential experiment,  
917 with DISTCB having a .086 advantage over the last 100 episodes and a .045 advantage over all  
918 episodes. While the gaps were not as large for the other two datasets, they are still statistically  
919 significant and further showcase the benefit of distribution learning.