

A Properties of PMD

We present lemmas relevant to the analysis of PMD. Key to the analysis is the Three-Point Descent Lemma, that relates the improvement of the proximal gradient update compared to an arbitrary point. It originally comes from [37] (Lemma 3.2) where a proof can be found, though we use a slightly modified version from [7] (Lemma 6).

Lemma A.1 (Three-Point Descent Lemma, Lemma 6 in [7]). *Suppose that $\mathcal{C} \subset \mathbb{R}^n$ is a closed convex set, $\phi : \mathcal{C} \rightarrow \mathbb{R}$ is a proper, closed convex function, $D_h(\cdot, \cdot)$ is the Bregman divergence generated by a function h of Legendre type and $\text{rint dom } h \cap \mathcal{C} \neq \emptyset$. For any $x \in \text{rint dom } h$, let*

$$x^+ = \operatorname{argmin}_{u \in \mathcal{C}} \{\phi(u) + D_h(u, x)\}. \quad (14)$$

Then $x^+ \in \text{rint dom } h \cap \mathcal{C}$ and $\forall u \in \mathcal{C}$,

$$\phi(x^+) + D_h(x^+, x) \leq \phi(u) + D_h(u, x) - D_h(u, x^+) \quad (15)$$

The update (4) of PMD is an instance of the proximal minimisation (14) with $\mathcal{C} = \Delta(\mathcal{A})$, $x = \pi_s^k$ and $\phi(x) = -\eta_k \langle Q_s^k, x \rangle$. Plugging these into (15), Lemma A.1 relates the decrease in the proximal objective of π_s^{k+1} to any other policy, i.e. $\forall p \in \Delta(\mathcal{A})$,

$$-\eta_k \langle Q_s^k, \pi_s^{k+1} \rangle + D_h(\pi_s^{k+1}, \pi_s^k) \leq -\eta_k \langle Q_s^k, p \rangle + D_h(p, \pi_s^k) - D_h(p, \pi_s^{k+1}). \quad (16)$$

This equation is key to the analysis in Section 6. In particular, it allows us to prove the following lemma regarding the monotonic improvement in action-value of PMD iterates. This is an extension of Lemma 7 in [7].

Lemma A.2. *Consider the policies produced by the iterative updates of PMD in (4). Then for any $k \geq 0$,*

$$Q^{k+1}(s, a) \geq Q^k(s, a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

A.1 Proof of Lemma A.2

We first present Lemma 7 from [7], from which Lemma A.2 almost immediately follows.

Lemma A.3 (Descent Property of PMD, Lemma 7 in [7]). *Consider the policies produced by the iterative updates of PMD in (4). Then for any $k \geq 0$*

$$\langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle \geq 0, \quad \forall s \in \mathcal{S},$$

$$V^{k+1}(\rho) \geq V^k(\rho), \quad \forall \rho \in \Delta(\mathcal{S}).$$

Proof. From [7]. Recall that the Three-Point Descent Lemma states that $\forall p \in \Delta(\mathcal{A})$,

$$-\eta_k \langle Q_s^k, \pi_s^{k+1} \rangle + D_h(\pi_s^{k+1}, \pi_s^k) \leq -\eta_k \langle Q_s^k, p \rangle + D_h(p, \pi_s^k) - D_h(p, \pi_s^{k+1}).$$

Using this with $p = \pi_s^k$,

$$D_h(\pi_s^k, \pi_s^{k+1}) + D_h(\pi_s^{k+1}, \pi_s^k) \leq \eta_k \langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle$$

and since the Bregman divergences are non-negative and $\eta_k > 0$,

$$0 \leq \langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle$$

and the result follows by an application of the performance difference lemma (Appendix B)

$$\begin{aligned} V^{k+1}(\rho) - V^k(\rho) &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\rho^{k+1}} [\langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle] \\ &\geq 0. \end{aligned}$$

Note that we use the performance difference lemma here because it gives a simple concise proof, but we do not actually need to. To maintain our claim that we avoid the use of the performance difference lemma, we can get the same result without it. We sketch how to do this as follows. From the first part of the lemma, we have

$$\langle Q_s^k, \pi_s^{k+1} \rangle \geq \langle Q_s^k, \pi_s^k \rangle = V^k(s),$$

498 in all states s . Now note that the left hand side above is

$$\begin{aligned}\langle Q_s^k, \pi_s^{k+1} \rangle &= \sum_a \pi^{k+1}(a|s) Q^k(s, a) \\ &= \sum_a \pi^{k+1}(a|s) \left(r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^k(s') \right)\end{aligned}$$

499 and we can then apply $\langle Q_{s'}^k, \pi_{s'}^{k+1} \rangle \geq V^k(s')$ at state s' :

$$\begin{aligned}V^k(s) &\leq \sum_a \pi^{k+1}(a|s) \left(r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^k(s') \right) \\ &\leq \sum_a \pi^{k+1}(a|s) \left(r(s, a) + \gamma \sum_{s'} p(s'|s, a) \sum_{a'} \pi^{k+1}(a'|s') \left(r(s', a') + \gamma \sum_{s''} p(s''|s', a') V^k(s'') \right) \right)\end{aligned}$$

500 and as proceed iteratively in the limit you get exactly $V^{k+1}(s)$. ■

501 Since Lemma [A.3](#) holds for any $\rho \in \Delta(\mathcal{S})$, it guarantees that the value in each state is non-decreasing
502 for an update of PMD, i.e for all $s \in \mathcal{S}$,

$$V^{k+1}(s) - V^k(s) \geq 0.$$

503 Using this, we get

$$Q^{k+1}(s, a) - Q^k(s, a) = \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) \left(V^{k+1}(s') - V^k(s') \right) \geq 0,$$

504 which concludes the proof. ■

505 A.2 Extension of Lemma [A.2](#) to inexact setting:

506 As in the exact case, we first present Lemma 12 from [\[7\]](#) which is the extension of Lemma [A.3](#) to the
507 inexact case. We note that in the inexact case, we lose the monotonic increase of values due to the
508 inaccuracy in our estimate \hat{Q}_s^k of Q_s^k .

509 **Lemma A.4.** Consider the policies produced by the iterative updates of IPMD in [\(9\)](#). For any $k \geq 0$,
510 if $\|\hat{Q}_s^k - Q_s^k\|_\infty \leq \tau$, then

$$\begin{aligned}511 \quad \langle \hat{Q}_s^k, \pi_s^{k+1} - \pi_s^k \rangle &\geq 0, \quad \forall s \in \mathcal{S}, \\ V^{k+1}(\rho) &\geq V^k(\rho) - \frac{2\tau}{1-\gamma}, \quad \forall \rho \in \Delta(\mathcal{S}).\end{aligned}$$

512 *Proof.* From [\[7\]](#). The Three-Point Descent Lemma applied to the IPMD update [\(9\)](#) gives $\forall p \in \Delta(\mathcal{A})$,

$$-\eta_k \langle \hat{Q}_s^k, \pi_s^{k+1} \rangle + D_h(\pi_s^{k+1}, \pi_s^k) \leq -\eta_k \langle \hat{Q}_s^k, p \rangle + D_h(p, \pi_s^k) - D_h(p, \pi_s^{k+1}).$$

513 Using this with $p = \pi_s^k$,

$$D_h(\pi_s^k, \pi_s^{k+1}) + D_h(\pi_s^{k+1}, \pi_s^k) \leq \eta_k \langle \hat{Q}_s^k, \pi_s^{k+1} - \pi_s^k \rangle$$

514 and since the Bregman divergences are none-negative and $\eta_k > 0$,

$$0 \leq \langle \hat{Q}_s^k, \pi_s^{k+1} - \pi_s^k \rangle,$$

515 which proves the first inequality. Now we cannot use the above inequality directly with the perfor-
516 mance difference lemma since \hat{Q}_s^k is not the true action-value. Instead, we have

$$\begin{aligned}V^{k+1}(\rho) - V^k(\rho) &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{k+1}} \left[\langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{k+1}} \left[\langle Q_s^k - \hat{Q}_s^k, \pi_s^{k+1} - \pi_s^k \rangle + \langle \hat{Q}_s^k, \pi_s^{k+1} - \pi_s^k \rangle \right] \\ &\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{k+1}} \left[-\|Q_s^k - \hat{Q}_s^k\|_\infty \|\pi_s^{k+1} - \pi_s^k\|_1 \right] \\ &\geq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{k+1}} \left[-2\tau \right] \\ &= -\frac{2\tau}{1-\gamma}\end{aligned}$$

517 which concludes the proof. ■

518 Using the above lemma, we can state and prove the extension of Lemma A.2 to the inexact setting.

519 **Lemma A.5.** *Consider the policies produced by the iterative updates of IPMD in (9). For any $k \geq 0$,*
 520 *if $\|\hat{Q}_s^k - Q_s^k\|_\infty \leq \tau$, then*

$$\hat{Q}^{k+1}(s, a) \geq \hat{Q}^k(s, a) - \frac{2\tau\gamma}{1-\gamma}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

521 *Proof.* As in the exact case, since Lemma A.4 holds for any $\rho \in \Delta(\mathcal{S})$, it applies to each state, i.e for
 522 all $s \in \mathcal{S}$,

$$V^{k+1}(s) - V^k(s) \geq -\frac{2\tau}{1-\gamma}.$$

523 Using this, we immediately have

$$Q^{k+1}(s, a) - Q^k(s, a) = \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) (V^{k+1}(s') - V^k(s')) \geq \frac{-2\tau\gamma}{1-\gamma},$$

524 which concludes the proof. ■

525 B Performance difference lemma

526 **Lemma B.1** (Performance Difference Lemma). *For any $\pi, \pi' \in \Pi$, we have*

$$V^\pi(\rho) - V^{\pi'}(\rho) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^\pi} [\langle Q_s^{\pi'}, \pi_s - \pi'_s \rangle].$$

527 The performance difference lemma [14] is a property that relates the difference in values of policies
 528 to the policies themselves. The proof can be found in their paper under Lemma 6.1.

529 C Guarantees of Theorem 4.1 for various step-size choices

530 We give here two more choices of $\{c_k\}_{k \in \mathbb{Z}_{\geq 0}}$ for the step-size [5] of PMD and their corresponding
 531 guarantees from Theorem 4.1:

- 532 • $c_i = c_0$ for some $c_0 > 0$ yields a step-size with a constant component. The resulting bound
 533 is

$$\|V^* - V^k\|_\infty \leq \gamma^k \|V^* - V^0\|_\infty + \frac{c_0}{1-\gamma},$$

534 which converges linearly up to some accuracy controlled by c_0 .

- 535 • $c_i = \gamma^{i+1} c_0$ for some initial $c_0 > 0$ will yield a step-size with a component that is
 536 geometrically increasing as in [7], though at a slower rate than the one discussed in Section
 537 [4]. The resulting bound is

$$\|V^* - V^k\|_\infty \leq \gamma^k (\|V^* - V^0\|_\infty + k c_0),$$

538 which converges linearly with the sought-for γ -rate, though in early iterations the k factor
 539 may dominate.

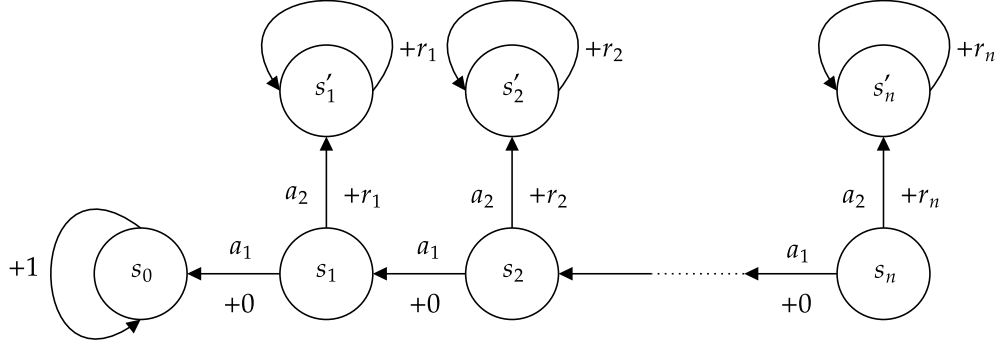


Figure 1: Example MDP used in the proof of Theorem 4.2

540 D Proof of Theorem 4.2

541 Fix $n > 0$ and $\delta \in (0, (1 - \gamma)\gamma^n)$. Consider the MDP shown in Figure 1. The state space is
 542 $\mathcal{S} = \{s_0, s_1, s'_1, \dots, s_n, s'_n\}$ and the action space is $\mathcal{A} = \{a_1, a_2\}$. There is a chain of states of length
 543 $n + 1$ with the states indexed from 0 to n . The left-most state (s_0) is absorbing with reward +1. In
 544 the other states in the chain (s_i for $i = 1, \dots, n$), the agent can take action a_1 and move left (to s_{i-1})
 545 with reward of 0, or take action a_2 and move to an additional absorbing state unique to the state it is
 546 currently in (s'_i) with reward $r_i = \gamma^{i+1} + \delta$ (that the agent also receives in that state for all future
 547 time-steps). Summarising, we have for $1 \leq i \leq n$

$$\begin{aligned} p(s_{i-1}|s_i, a_1) &= 1, & r(s_i, a_1) &= 0, \\ p(s'_i|s_i, a_2) &= 1, & r(s_i, a_2) &= r_i = \gamma^{i+1} + \delta, \\ p(s'_i|s'_i, a) &= 1, & r(s'_i, a) &= r_i = \gamma^{i+1} + \delta \quad \forall a \in \mathcal{A}. \end{aligned}$$

548 The value of δ is carefully restricted so that the optimal action in all the states of the chain is a_1 . The
 549 proof will consist in showing that if the agent starts with an initial policy that places most probability
 550 mass on the sub-optimal action a_2 , then it has to learn that a_1 is the optimal action in the state directly
 551 to the left before it can start switching from action a_2 to a_1 in the current state. And this can at best
 552 happen one iteration at a time starting from the left-most state. In particular, we consider π^0
 553 s.t $\pi^0(a_1|s) = \alpha$, $\pi^0(a_2|s) = 1 - \alpha$ for all states and some α s.t $0 < \alpha \leq \delta(1 - \gamma)$. We make the
 554 following claim from which the result will follow straightforwardly.

555 **Claim:** Fix $k < n$. The policies produced by PMD satisfy $\pi^k(a_1|s_i) \leq \alpha$ for $k < i \leq n$.

556 We prove this claim by induction.

557 **Base Case:** We want to show that $\pi^1(a_1|s_i) \leq \alpha$ for $i > 1$. We do this by showing that $Q^0(s_i, a_1) \leq$
 558 $Q^0(s_i, a_2)$ for $i > 1$ so that the probability of $\pi^1(a_1|s_i)$ cannot increase w.r.t $\pi^0(a_1|s_i)$, which is α

559 (this follows from $\langle Q_s^k, \pi_s^{k+1} - \pi_s^k \rangle \geq 0$ for all iterations of PMD). We have:

$$\begin{aligned}
Q^0(s_i, a_1) &= \gamma V^0(s_{i-1}) \\
&= \gamma \left(\alpha Q^0(s_{i-1}, a_1) + (1 - \alpha) Q^0(s_{i-1}, a_2) \right) \\
&\leq \gamma \left(\alpha \frac{\gamma^{i-1}}{1 - \gamma} + \frac{r_{i-1}}{1 - \gamma} \right) \\
&\stackrel{(a)}{\leq} \gamma \left(\delta(1 - \gamma) \frac{\gamma^{i-1}}{1 - \gamma} + \frac{\gamma^i + \delta}{1 - \gamma} \right) \\
&= \frac{\gamma^{i+1}}{1 - \gamma} + \frac{\delta\gamma(1 + \gamma^{i-1} - \gamma^i)}{1 - \gamma} \\
&\stackrel{(b)}{\leq} \frac{\gamma^{i+1}}{1 - \gamma} + \frac{\delta}{1 - \gamma} \\
&= Q^0(s_i, a_2),
\end{aligned}$$

560 where we used $\alpha \leq \delta(1 - \gamma)$ in (a) and $\gamma(1 + \gamma^{i-1} - \gamma^i) < 1$ for $\gamma \in [0, 1)$ in (b). This concludes
561 the base case.

562 **Inductive Step:** Now assume that the claim is true for k and we want to show that $\pi^{k+1}(a_1|s_i) \leq \alpha$
563 for $i > k + 1$. We do this in the same way as the base case by showing that $Q^k(s_i, a_1) \leq Q^k(s_i, a_2)$
564 for $i > k + 1$ so that the probability of $\pi^{k+1}(a_1|s_i)$ cannot increase w.r.t $\pi^k(a_1|s_i)$, which is less
565 than or equal to α by the inductive hypothesis. We have:

$$\begin{aligned}
Q^k(s_i, a_1) &= \gamma V^k(s_{i-1}) \\
&= \gamma \left(\pi^k(a_1|s_{i-1}) Q^k(s_{i-1}, a_1) + \pi^k(a_2|s_{i-1}) Q^k(s_{i-1}, a_2) \right) \\
&\stackrel{(a)}{\leq} \gamma \left(\alpha Q^k(s_{i-1}, a_1) + Q^k(s_{i-1}, a_2) \right) \\
&\leq \gamma \left(\alpha \frac{\gamma^{i-1}}{1 - \gamma} + \frac{r_{i-1}}{1 - \gamma} \right) \\
&\stackrel{(b)}{\leq} \gamma \left(\delta(1 - \gamma) \frac{\gamma^{i-1}}{1 - \gamma} + \frac{\gamma^i + \delta}{1 - \gamma} \right) \\
&= \frac{\gamma^{i+1}}{1 - \gamma} + \frac{\delta\gamma(1 + \gamma^{i-1} - \gamma^i)}{1 - \gamma} \\
&\stackrel{(c)}{\leq} \frac{\gamma^{i+1}}{1 - \gamma} + \frac{\delta}{1 - \gamma} \\
&= Q^k(s_i, a_2),
\end{aligned}$$

566 where we used in (a) that $\pi^k(a_1|s_{i-1}) \leq \alpha$ for $i > k + 1$, which is true by the inductive hypothesis
567 since $i - 1 > k$, in (b) that $\alpha \leq \delta(1 - \gamma)$ and in (c) that $\gamma(1 + \gamma^{i-1} - \gamma^i) < 1$ for $\gamma \in [0, 1)$. This
568 concludes the proof of the claim.

569 Now using the claim

$$\begin{aligned}
V^k(s_{k+1}) &= \pi^k(a_1|s_{k+1}) Q^k(s_{k+1}, a_1) + \pi^k(a_2|s_{k+1}) Q^k(s_{k+1}, a_2) \\
&\leq \alpha \frac{\gamma^{k+1}}{1 - \gamma} + \frac{r_{k+1}}{1 - \gamma} \\
&= \alpha \frac{\gamma^{k+1}}{1 - \gamma} + \frac{\gamma^{k+2} + \delta}{1 - \gamma},
\end{aligned}$$

570 SO

$$\begin{aligned}
V^*(s_{k+1}) - V^k(s_{k+1}) &\geq \frac{\gamma^{k+1}}{1-\gamma} - \alpha \frac{\gamma^{k+1}}{1-\gamma} - \frac{\gamma^{k+2} + \delta}{1-\gamma} \\
&= \frac{\gamma^{k+1}(1-\alpha)}{1-\gamma} - \frac{\alpha\gamma^{k+1} + \delta}{1-\gamma} \\
&\geq \gamma^{k+1} - \frac{\alpha + \delta}{1-\gamma} \\
&\geq \gamma^{k+1} - \frac{2\delta}{1-\gamma},
\end{aligned} \tag{17}$$

571 where we used that $\alpha \leq \delta$. Now note that

$$\begin{aligned}
V^0(s_1) &= \alpha Q^0(s_1, a_1) + (1-\alpha)Q^0(s_1, a_2) \\
&= \alpha \frac{\gamma}{1-\gamma} + (1-\alpha) \frac{\gamma^2 + \delta}{1-\gamma},
\end{aligned}$$

572 SO

$$\begin{aligned}
V^*(s_1) - V^0(s_1) &= \frac{\gamma}{1-\gamma} - \alpha \frac{\gamma}{1-\gamma} - (1-\alpha) \frac{\gamma^2 + \delta}{1-\gamma} \\
&= (1-\alpha) \frac{\gamma}{1-\gamma} - (1-\alpha) \frac{\gamma^2 + \delta}{1-\gamma} \\
&= \frac{1-\alpha}{1-\gamma} (\gamma - \gamma^2 - \delta) \\
&\leq \frac{1-\alpha}{1-\gamma} (\gamma - \gamma^2) \\
&= \gamma \frac{1-\alpha}{1-\gamma} (1-\gamma) \\
&= \gamma(1-\alpha) \\
&\leq \gamma
\end{aligned}$$

573 and by induction we can show this is the case for all states (above is base case), the inductive step is
574 as follows (assuming $V^*(s_k) - V^0(s_k) \leq \gamma$),

$$\begin{aligned}
V^*(s_{k+1}) - V^0(s_{k+1}) &= \frac{\gamma^{k+1}}{1-\gamma} - (1-\alpha) \frac{\gamma^{k+2} + \delta}{1-\gamma} - \alpha \gamma V^0(s_k) \\
&= (1-\alpha) \left[\frac{\gamma^{k+1} - \gamma^{k+2} - \delta}{1-\gamma} \right] + \alpha \gamma [V^*(s_k) - V^0(s_k)] \\
&\leq (1-\alpha) \gamma^{k+1} + \alpha \gamma^2 \\
&\leq \gamma
\end{aligned}$$

575 and so

$$\|V^* - V^0\|_\infty \leq \gamma,$$

576 which combining with (17) gives,

$$\begin{aligned}
V^*(s_{k+1}) - V^k(s_{k+1}) &\geq \gamma^k \|V^* - V^0\|_\infty - \frac{2\delta}{1-\gamma} \\
\implies \|V^* - V^k\|_\infty &\geq \gamma^k \|V^* - V^0\|_\infty - \frac{2\delta}{1-\gamma},
\end{aligned}$$

577 which concludes the proof. ■

578 E Proof of Theorem 4.3

579 Consider the same MDP as in the proof of Theorem 4.2 in Appendix D (see Figure I). Denote
 580 $c = \frac{1-\gamma}{8}$ and note that $c < \frac{\sqrt{\gamma}}{1+\sqrt{\gamma}} \frac{1-\gamma}{2}$ since $\frac{1}{4} < \frac{\sqrt{\gamma}}{1+\sqrt{\gamma}}$ for $\gamma > 0.2$.

581 Suppose you consider NPG updates with initial policy $\pi^0(a_1|s_i) = \alpha$. Recall that NPG is the instance
 582 of PMD with relative entropy as the mirror map. It can be shown that NPG has the closed form update

$$\pi^{k+1}(a|s) = \frac{\pi^k(a|s)e^{\eta_k Q^k(s,a)}}{\sum_{a'} \pi^k(a'|s)e^{\eta_k Q^k(s,a')}}.$$

583 We know from the proof of Theorem D that for any step-size regime, for $i > k + 1$

$$Q^k(s_i, a_1) \leq Q^k(s_i, a_2).$$

584 Now, $\|V^* - V^0\|_\infty = V^*(s_1) - V^0(s_1) \leq \gamma - \frac{\delta}{1-\gamma}$ (see Section E.1 below). The idea of the proof
 585 is to show that satisfying the bound given in the statement of the theorem will imply that a certain
 586 condition on the step-size.

587 Fix a state s_k and let k_0 be the first iteration where $Q^{k_0}(s_k, a_1) > Q^{k_0}(s_k, a_2)$. By the above, we
 588 must have $k \leq k_0 + 1$, or $k_0 \geq k - 1$. By the proof of Theorem D, we also have $\pi^{k_0}(a_1|s_k) \leq \alpha$
 589 (before iteration k_0 , $Q(s_k, \cdot)$ favors a_2 , so $\pi^{k_0}(a_1|s_k)$ has not increased compared to $\pi^0(a_1|s_k) = \alpha$).

590 We want a γ -contraction at every iteration, i.e. we assume the following is satisfied:

$$V^*(s_k) - V^{k_0+1}(s_k) \leq \gamma^{k_0+1}(\|V^* - V^0\|_\infty + c) \leq \gamma^{k_0+1}(\gamma - \frac{\delta}{1-\gamma} + c).$$

591 Now, by direct computation,

$$\begin{aligned} V^*(s_k) - V^{k_0+1}(s_k) &= \pi^{k_0+1}(a_1|s_k)\gamma(V^*(s_{k-1}) - V^{k_0+1}(s_{k-1})) + \pi^{k_0+1}(a_2|s_2)\frac{\gamma^k - r_k}{1-\gamma} \\ &\geq \pi^{k_0+1}(a_2|s_2)\frac{\gamma^k - r_k}{1-\gamma} = \pi^{k_0+1}(a_2|s_2)(\gamma^k - \frac{\delta}{1-\gamma}). \end{aligned}$$

592 Putting this together with the above (this is an implication as this is about the necessity rather than
 593 sufficiency), we must have:

$$\begin{aligned} \pi^{k_0+1}(a_2|s_2)(\gamma^k - \frac{\delta}{1-\gamma}) &\leq \gamma^{k_0+1}(\gamma - \frac{\delta}{1-\gamma} + c) \\ \implies \pi^{k_0+1}(a_2|s_2) &\leq \frac{\gamma^{k_0+1}(\gamma - \frac{\delta}{1-\gamma} + c)}{(\gamma^k - \frac{\delta}{1-\gamma})} = \beta \end{aligned}$$

594 If we choose $\delta < \frac{1}{2}(1-\gamma)(1-\sqrt{\gamma})\gamma^k$ then $\beta < \sqrt{\gamma}$ and require

$$\pi^{k_0+1}(a_2|s_2) \leq \sqrt{\gamma}.$$

595 To see this, start from $\beta \leq \sqrt{\gamma}$, this is equivalent to

$$\begin{aligned}
& \frac{\gamma^{k_0+1}(\gamma - \frac{\delta}{1-\gamma} + c)}{(\gamma^k - \frac{\delta}{1-\gamma})} \leq \sqrt{\gamma} \\
& \iff \frac{\gamma^k(\gamma - \frac{\delta}{1-\gamma} + c)}{(\gamma^k - \frac{\delta}{1-\gamma})} \leq \sqrt{\gamma} \quad \text{since } k_0 + 1 \geq k \\
& \iff \gamma^k(\gamma - \frac{\delta}{1-\gamma} + c) \leq \sqrt{\gamma}(\gamma^k - \frac{\delta}{1-\gamma}) \\
& \iff \gamma^k(\gamma + c) \leq \sqrt{\gamma}(\gamma^k - \frac{\delta}{1-\gamma}) \\
& \iff \gamma^{k-\frac{1}{2}}(\gamma + c) \leq \gamma^k - \frac{\delta}{1-\gamma} \\
& \iff \frac{\delta}{1-\gamma} \leq \gamma^{k-\frac{1}{2}}(\sqrt{\gamma} - \gamma - c) \\
& \iff \frac{\delta}{1-\gamma} \leq \gamma^{k-\frac{1}{2}}(\sqrt{\gamma} - \gamma - \frac{\sqrt{\gamma}}{1+\sqrt{\gamma}} \frac{1-\gamma}{2}) \quad \text{since } -c > -\frac{\sqrt{\gamma}}{1+\sqrt{\gamma}} \frac{1-\gamma}{2} \\
& \iff \frac{\delta}{1-\gamma} \leq \gamma^{k-\frac{1}{2}}(\sqrt{\gamma} - \gamma - \frac{\sqrt{\gamma}}{2}(1-\sqrt{\gamma})) \\
& \iff \frac{\delta}{1-\gamma} \leq \gamma^{k-\frac{1}{2}}(\sqrt{\gamma} - \gamma - \frac{\sqrt{\gamma}}{2} + \frac{\gamma}{2}) \\
& \iff \frac{\delta}{1-\gamma} \leq \gamma^{k-\frac{1}{2}}(\frac{\sqrt{\gamma}}{2} - \frac{\gamma}{2}) \\
& \iff \delta \leq \frac{1}{2}\gamma^k(1-\sqrt{\gamma})(1-\gamma),
\end{aligned}$$

596 which is the condition for δ we imposed initially.

597 To achieve the above condition $\pi^{k_0+1}(a_2|s_2) \leq \sqrt{\gamma}$, recalling that $\pi^{k_0}(a_2|s_2) \geq 1-\alpha$, η_{k_0} has to
598 satisfy

$$\eta_{k_0} \geq \frac{1}{Q^{k_0}(s_k, a_1) - Q^{k_0}(s_k, a_2)} \left[\log((1-\alpha)(1-\sqrt{\gamma})) + KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}) \right]$$

599 To see this, again start from $\pi^{k_0+1}(a_2|s_2) \leq \sqrt{\gamma}$, this is equivalent to (use $k_0 = m$ for simplicity of
600 notation) using the closed-form update of NPG:

$$\begin{aligned}
& \pi^m(a_2|s_2) \exp(\eta_m Q^m(s_k, a_2)) \leq \\
& \sqrt{\gamma}(\pi^m(a_2|s_2) \exp(\eta_m Q^m(s_k, a_2)) + \pi^m(a_1|s_2) \exp(\eta_m Q^m(s_k, a_1))) \\
& \iff \frac{1}{\sqrt{\gamma}} \leq 1 + \frac{\pi^m(a_1|s_2)}{\pi^m(a_2|s_2)} \exp(\eta_m(Q^m(s_k, a_1) - Q^m(s_k, a_2))) \\
& \iff \frac{1-\sqrt{\gamma}}{\sqrt{\gamma}} \frac{\pi^m(a_2|s_2)}{\pi^m(a_1|s_2)} \leq \exp(\eta_m(Q^m(s_k, a_1) - Q^m(s_k, a_2))) \\
& \iff \eta_m(Q^m(s_k, a_1) - Q^m(s_k, a_2)) \geq \log\left(\frac{1-\sqrt{\gamma}}{\sqrt{\gamma}} \frac{\pi^m(a_2|s_2)}{\pi^m(a_1|s_2)}\right) \\
& \iff \eta_m \geq \frac{1}{Q^m(s_k, a_1) - Q^m(s_k, a_2)} \left[\log\left(\frac{1-\sqrt{\gamma}}{\sqrt{\gamma}} \frac{\pi^m(a_2|s_2)}{\pi^m(a_1|s_2)}\right) + \log\left(\frac{1}{\pi^m(a_1|s_2)}\right) \right] \\
& \implies \eta_m \geq \frac{1}{Q^m(s_k, a_1) - Q^m(s_k, a_2)} \left[\log\left((1-\alpha)\frac{1-\sqrt{\gamma}}{\sqrt{\gamma}}\right) + KL(\tilde{\pi}_{s_k}^{m+1}, \pi_{s_k}^m) \right] \\
& \implies \eta_m \geq \frac{1}{Q^m(s_k, a_1) - Q^m(s_k, a_2)} \left[\log\left((1-\alpha)(1-\sqrt{\gamma})\right) + KL(\tilde{\pi}_{s_k}^{m+1}, \pi_{s_k}^m) \right].
\end{aligned}$$

601 As we take $\alpha \rightarrow 0$, the KL term will dominate. In particular, note $\alpha < 1-\gamma$ so $1-\alpha > \gamma$ and
 $(1-\alpha)(1-\sqrt{\gamma}) > \gamma(1-\sqrt{\gamma})$

602 and if we further impose the condition $\alpha < \gamma^2(1 - \sqrt{\gamma})^2$ then

$$(1 - \alpha)(1 - \sqrt{\gamma}) > \sqrt{\alpha} > \sqrt{\pi^{k_0}(a_1|s_2)}$$

603 and the step-size needs to satisfy the following condition:

$$\begin{aligned} \eta_{k_0} &\geq \frac{1}{Q^{k_0}(s_k, a_1) - Q^{k_0}(s_k, a_2)} \left[\log(\sqrt{\pi^{k_0}(a_1|s_2)}) + KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}) \right] \\ &= \frac{1}{Q^{k_0}(s_k, a_1) - Q^{k_0}(s_k, a_2)} \left[-\frac{1}{2}KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}) + KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}) \right] \\ &= \frac{1}{2(Q^{k_0}(s_k, a_1) - Q^{k_0}(s_k, a_2))} KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}) \end{aligned} \quad (18)$$

604 **Distinct Iterations:** Note that the iteration $k_0(s_k)$ where $Q(\cdot, s_k)$ starts becoming bigger at a_1 that
605 a_2 is distinct for each s_k . Fix any s_k and $k_0 = k_0(s_k)$. We have

$$\begin{aligned} Q^{k_0}(s_k, a_1) &< Q^{k_0}(s_k, a_2) \\ Q^{k_0+1}(s_k, a_2) &\leq Q^{k_0+1}(s_k, a_1) \end{aligned}$$

606 then $\pi^{k_0+1}(a_1|s_k) \leq \pi^{k_0}(a_1|s_k) \leq \alpha$ (since Q^t points towards a_2 in s_k for all $t \leq k$). Then
607 applying exactly the same steps as in the proof of Theorem 4.2 we have

$$Q^{k_0+1}(s_{k+1}, a_1) < Q^{k_0+1}(s_{k+1}, a_2),$$

608 meaning that $k_0(s_k)$ is distinct to $k_0(s_{k+1})$.

609 **Upper Bounding Q-value difference:** We want to upper-bound the Q-value difference appearing in
610 the step-size condition above. We have,

$$\begin{aligned} Q^{k_0}(s_k, a_2) &= \frac{r_k}{1 - \gamma} = \frac{\gamma^{k+1} + \delta}{1 - \gamma} \\ Q^{k_0}(s_k, a_1) &= \gamma V^{k_0}(s_{k-1}) \leq \frac{\gamma^k}{1 - \gamma}. \end{aligned}$$

611 So,

$$\begin{aligned} Q^{k_0}(s_k, a_1) - Q^{k_0}(s_k, a_2) &\leq \frac{\gamma^k}{1 - \gamma} - \frac{\gamma^{k+1} + \delta}{1 - \gamma} \\ &= \gamma^k - \frac{\delta}{1 - \gamma} \\ &\leq \gamma^k. \end{aligned}$$

612 Plugging this into the above bound (18), if the iterates of NPG are to satisfy the bound with the γ -rate
613 in the statement of the theorem, the step-size must at least satisfy the following condition:

$$\eta_{k_0} \geq \frac{1}{2\gamma^k} KL(\tilde{\pi}_{s_k}^{k_0+1}, \pi_{s_k}^{k_0}),$$

614 which concludes the proof. ■

615 E.1 Largest sub-optimality gap at iteration 0

616 In this section, we prove the claim that

$$\|V^* - V^0\|_\infty = V^*(s_1) - V^0(s_1) \leq \gamma - \frac{\delta}{1 - \gamma}$$

617 **Proof:** First of all, $V^*(s_1) - V^0(s_1) = \pi^0(a_2|s_1) \frac{\gamma - r_1}{1 - \gamma} = (1 - \alpha)(\gamma - \frac{\delta}{1 - \gamma}) \leq \gamma - \frac{\delta}{1 - \gamma}$. For the
618 first part, we proceed by induction. We will use throughout that

$$\frac{\gamma^k - r_k}{1 - \gamma} = \gamma^k - \frac{\delta}{1 - \gamma} \leq V^*(s_1) - V^0(s_1) = (1 - \alpha)(\gamma - \frac{\delta}{1 - \gamma}).$$

619 This is true if (when LHS is the largest)

$$\gamma^2 - \frac{\delta}{1-\gamma} \leq (1-\alpha)(\gamma - \frac{\delta}{1-\gamma})$$

620 which holds when

$$\begin{aligned} \alpha &\leq \frac{\gamma(1-\gamma)^2}{\gamma(1-\gamma) - \delta} \\ \iff \alpha &\leq 1-\gamma \end{aligned}$$

621 **Base Case:**

$$\begin{aligned} V^*(s_2) - V^0(s_2) &= \alpha\gamma(V^*(s_1) - V^0(s_1)) + (1-\alpha)\frac{\gamma^2 - r_2}{1-\gamma} \\ &\leq \alpha\gamma(V^*(s_1) - V^0(s_1)) + (1-\alpha)(V^*(s_1) - V^0(s_1)) \\ &\leq V^*(s_1) - V^0(s_1) \end{aligned}$$

622 **Inductive Step:** Assume true for k . Then,

$$\begin{aligned} V^*(s_{k+1}) - V^0(s_{k+1}) &= \alpha\gamma(V^*(s_k) - V^0(s_k)) + (1-\alpha)\frac{\gamma^{k+1} - r_{k+1}}{1-\gamma} \\ &\leq \alpha\gamma(V^*(s_1) - V^0(s_1)) + (1-\alpha)(V^*(s_1) - V^0(s_1)) \\ &\leq V^*(s_1) - V^0(s_1), \end{aligned}$$

623 which concludes the proof. ■

624 F Inexact policy mirror descent and the generative model

625 The following Lemma from [7] controls the accuracy of the estimator \hat{Q}_s^k specified in (10) of Section
 626 5 with respect to H and M_k :

627 **Lemma F.1** (Lemma 15 in [7]). *Consider using (10) to estimate Q_s^k for all state-action pairs for K*
 628 *iterations of IPMD. Then for $\delta \in (0, 1)$, if for all $k \leq K$,*

$$M_k \geq \frac{\gamma^{-2H}}{2} \log\left(\frac{2K|\mathcal{S}||\mathcal{A}|}{\delta}\right).$$

629 *Then with probability at least $1 - \delta$, we have for all $k \leq K$,*

$$\|\hat{Q}_s^k - Q_s^k\|_\infty \leq \frac{2\gamma^H}{1-\gamma}.$$

630 The proof of this result can be found in Lemma 15 of [7].

631 F.1 Proof of Theorem 5.1

632 This proof is similar to that of [7] (Theorem 14). It is also similar in structure to the proof of Theorem
 633 4.1 in Section 6.

634 Fix a state $s \in \mathcal{S}$ and an integer $k \geq 0$. For now let's assume that our Q-estimates are τ -accurate
 635 ($\tau > 0$), i.e.

$$\|Q^k - \hat{Q}^k\|_\infty \leq \tau$$

636 for all $k \geq 0$. With this assumption, we have from Lemma A.5 in Appendix A.1

$$Q^{k+1}(s, a) \geq Q^k(s, a) - \frac{2\gamma\tau}{1-\gamma}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

637 Now proceeding in a similar way to Section 6,

$$\begin{aligned} \langle \hat{Q}_s^k, \pi_s^* - \pi_s^{k+1} \rangle &= \langle Q_s^k, \pi_s^* - \pi_s^{k+1} \rangle + \langle \hat{Q}_s^k - Q_s^k, \pi_s^* - \pi_s^{k+1} \rangle \\ &\geq \langle Q_s^k, \pi_s^* \rangle - \langle Q_s^k, \pi_s^{k+1} \rangle - \|\hat{Q}_s^k - Q_s^k\|_\infty \|\pi_s^* - \pi_s^{k+1}\|_1 \\ &\geq \langle Q_s^k, \pi_s^* \rangle - \langle Q_s^{k+1}, \pi_s^{k+1} \rangle - \frac{2\gamma\tau}{1-\gamma} - 2\tau \\ &\geq \langle Q_s^k, \pi_s^* \rangle - V^{k+1}(s) - \frac{4\gamma\tau}{1-\gamma} \\ &= \langle Q_s^k - Q_s^*, \pi_s^* \rangle + V^*(s) - V^{k+1}(s) - \frac{4\gamma\tau}{1-\gamma} \\ &\geq -\|Q_s^* - Q_s^k\|_\infty + V^*(s) - V^{k+1}(s) - \frac{4\gamma\tau}{1-\gamma} \\ &\geq -\gamma\|V^* - V^k\|_\infty + V^*(s) - V^{k+1}(s) - \frac{4\gamma\tau}{1-\gamma}. \end{aligned}$$

638 Now again proceeding exactly as in Section 6 with this extra τ -term using the step-size condition
 639 ($c_k = \gamma^{2k+1}$), we end up with

$$\|V^* - V^{k+1}\|_\infty \leq \gamma\|V^* - V^k\|_\infty + \gamma^{2k+1} + \frac{4\gamma\tau}{1-\gamma}.$$

640 Unravelling this recursion yields

$$\begin{aligned} \|V^* - V^k\|_\infty &\leq \gamma^k \left(\|V^* - V^0\|_\infty + \sum_{i=1}^k \gamma^{-i} \gamma^{2(i-1)+1} \right) + \frac{4\gamma\tau}{1-\gamma} \sum_{i=0}^{k-1} \gamma^i \\ &\leq \gamma^k \left(\|V^* - V^0\|_\infty + \frac{1}{1-\gamma} \right) + \frac{4\gamma\tau}{(1-\gamma)^2}. \end{aligned}$$

641 Now using the properties of the estimator (10) in Lemma F.1, we have with probability $1 - \delta$ for all
 642 $0 \leq k \leq K$,

$$\tau = \frac{2\gamma^H}{1-\gamma},$$

643 giving

$$\begin{aligned} \|V^\star - V^k\|_\infty &\leq \gamma^k \left(\|V^\star - V^0\|_\infty + \frac{1}{1-\gamma} \right) + \frac{8\gamma^H}{(1-\gamma)^3} \\ &\leq \frac{2}{1-\gamma} \gamma^k + \frac{8\gamma^H}{(1-\gamma)^3}. \end{aligned}$$

644 This establishes the first bound. Now

$$\begin{aligned} K &> \frac{1}{1-\gamma} \log \frac{4}{(1-\gamma)\varepsilon} \implies \frac{2}{1-\gamma} \gamma^k \leq \varepsilon/2, \\ H &\geq \frac{1}{1-\gamma} \log \frac{16}{(1-\gamma)^3\varepsilon} \implies \frac{8\gamma^H}{(1-\gamma)^3} \leq \varepsilon/2 \end{aligned}$$

645 giving

$$\|V^\star - V^k\|_\infty \leq \varepsilon/2 + \varepsilon/2 = \varepsilon$$

646 as required. In terms of M , we have

$$\begin{aligned} M &\geq \frac{\gamma^{-2H}}{2} \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta} \\ &\geq \frac{1}{2} \left(\frac{16}{(1-\gamma)^3\varepsilon} \right)^2 \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta} \\ &= \frac{16^2}{2(1-\gamma)^6\varepsilon^2} \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta} \end{aligned}$$

647 and the corresponding number of calls to the sampling model, i.e. the sample complexity is (what we
 648 have shown above is actually a lower bound but can choose K, H, M so that it is of the following
 649 order),

$$|\mathcal{S}| \cdot |\mathcal{A}| \cdot K \cdot H \cdot M = \tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^8\varepsilon^2}\right),$$

650 where the notation $\tilde{O}()$ hides poly-logarithmic factors. This completes the proof. ■

651 G MDP examples

652 G.1 MDP on which distribution-mismatch coefficient scales with size of state space

653 We construct an MDP on which

$$\theta_\rho = \frac{1}{1-\gamma} \left\| \frac{d_\rho^*}{\rho} \right\|_\infty,$$

654 scales with $|\mathcal{S}|$, and hence so does the iteration complexity of the bound of [7] for exact PMD.

655 Consider an MDP with state-space $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ of size n and arbitrary action space \mathcal{A} . s_1 is
656 an absorbing state giving out rewards of 1 at each time-step, regardless of the action taken, i.e

$$p(s_1|s_1, a) = 1, \quad r(s_1, a) = 1 \quad \forall a \in \mathcal{A}.$$

657 All others states have an action, say a_1 , that gives out a reward of 1 and with probability $1 - \delta$ brings
658 the agent to state s_1 for some $\delta > 0$ and spreads the remaining δ probability arbitrarily amongst the
659 other states. The other actions have arbitrary rewards strictly less than 1 associated to them, and
660 arbitrary transition probabilities that place 0 mass on state s_1 , i.e

$$\begin{aligned} p(s_1|s, a_1) &= 1 - \delta, \quad r(s, a_1) = 1 \quad \forall s \neq s_1, \\ p(s_1|s, a) &= 0, \quad r(s, a) < 1 \quad \forall s \neq s_1, \forall a \neq a_1. \end{aligned}$$

661 Denote $r_{\max} = \max_{s \neq s_1, a \neq a_1} r(s, a) < 1$. The following condition ensures that a_1 is the optimal
662 action in all states,

$$\delta \leq \frac{1-\gamma}{\gamma} (1 - r_{\max})$$

663 so that $\pi^*(s) = a_1$ for all states s . To see this, consider $s_i \neq s_1, a_m \neq a_1$ and an arbitrary policy π ,

$$\begin{aligned} Q^\pi(s_i, a_1) &= 1 + \gamma \left(\frac{1-\delta}{1-\gamma} + \sum_{j=2}^n p(s_j|s_i, a_1) V^\pi(s_j) \right) \\ &\geq 1 + \gamma \frac{1-\delta}{1-\gamma} \\ Q^\pi(s_i, a_m) &= r(s_i, a_m) + \gamma \sum_{j=2}^n p(s_j|s_i, a_1) V^\pi(s_j) \\ &\leq r_{\max} + \gamma \frac{1}{1-\gamma} \end{aligned}$$

664 and solving

$$r_{\max} + \gamma \frac{1}{1-\gamma} \leq 1 + \gamma \frac{1-\delta}{1-\gamma}$$

665 will yield the condition above.

666 Then for $t \geq 1$ (abusing notation, s_t denotes the state at time t),

$$\begin{aligned} \mathbb{P}^{\pi^*}(s_t = s_1 | s_0 = s) &= \sum_{s'} \mathbb{P}^{\pi^*}(s_t = s_1, s_{t-1} = s' | s_0 = s) \\ &= \sum_{s'} p(s_1 | s', a_1) \mathbb{P}^{\pi^*}(s_{t-1} = s' | s_0 = s) \\ &\geq \sum_{s'} (1 - \delta) \mathbb{P}^{\pi^*}(s_{t-1} = s' | s_0 = s) \\ &= 1 - \delta \end{aligned}$$

667 and

$$\begin{aligned}
d_\rho^*(s_1) &= (1 - \gamma) \sum_s \rho(s) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi^*}(s_t = s_1 | s_0 = s) \\
&\geq (1 - \gamma) \sum_s \rho(s) \sum_{t=1}^{\infty} \gamma^t (1 - \delta) \\
&\geq (1 - \gamma) \sum_s \rho(s) \frac{\gamma}{1 - \gamma} (1 - \delta) \\
&= \gamma(1 - \delta).
\end{aligned}$$

668 Now

$$\left\| \frac{d_\rho^*}{\rho} \right\|_\infty \geq \frac{d_\rho^*(s_1)}{\rho(s_1)} \geq \frac{\gamma(1 - \delta)}{\rho(s_1)}$$

669 and depending on what ρ you consider, θ_ρ can be arbitrarily large. In particular, the natural choice of
670 the uniform starting-state distribution $\rho(s) = 1/n$ leads to

$$\theta_\rho \geq n \frac{\gamma(1 - \delta)}{(1 - \gamma)}$$

671 which gives an iteration complexity under the result of [7] for an ε -optimal policy that is

$$n \frac{\gamma(1 - \delta)}{(1 - \gamma)} \log \frac{2}{(1 - \gamma)\varepsilon}.$$

672 Recall that $n = |\mathcal{S}|$, so this iteration complexity scales linearly with the size of the state space.

673 G.2 Family of MDPs on which sub-optimality gaps can be made arbitrarily small

674 We present how to construct a family of MDPs on which $\Delta^k(s)$ defined in Section 4 can be made
675 arbitrarily small.

676 Consider an arbitrary MDP \mathcal{M} with state space \mathcal{S} and action space \mathcal{A} . For each state-action pair
677 $(s, a) \in \mathcal{S} \times \mathcal{A}$, create a duplicate action a' s.t the transitions from that action in that state are the
678 same as for the original pair, i.e

$$p(s' | s, a) = p(s' | s, a') \quad \forall s' \in \mathcal{S}$$

679 and the reward is shifted down by $\delta > 0$ from the original reward, i.e

$$r(s, a') = r(s, a) - \delta.$$

680 This results in a new MDP \mathcal{M}' with an augmented action space \mathcal{A}' , that is twice the size of the action
681 space of the original MDP \mathcal{M} . In terms of action-value of an arbitrary policy π , this results in

$$Q_{\mathcal{M}'}^\pi(s, a) - Q_{\mathcal{M}'}^\pi(s, a') = \delta,$$

682 where the notation $Q_{\mathcal{M}'}^\pi$ refers to action-values in the MDP \mathcal{M}' . In terms of sub-optimality gaps, this
683 gives

$$\Delta^\pi(s) \leq \delta.$$

684 Choosing δ small enough, we can make the step-size of [9] arbitrarily large, at least in early iterations.
685 The step-size condition (5) of Theorem 4.1 will be less affected by this issue as it does not depend
686 directly on $\Delta^k(s)$, and not at all in the first iteration. Beyond its generality to PMD, this illustrates
687 the benefit of our result restricted to NPG over the result of [9].