

---

**Algorithm 1** The policy mirror descent (PMD) method

---

**Input:** initial points  $\pi_0$  and stepsizes  $\eta_k \geq 0$ .

**for**  $k = 0, 1, \dots$ , **do**

$$\pi_{k+1}(\cdot|s) = \arg \min_{p(\cdot|s) \in \Delta_{|\mathcal{A}|}} \left\{ \eta_k [\langle Q^{\pi_k}(s, \cdot), p(\cdot|s) \rangle + h^p(s)] + D_{\pi_k}^p(s) \right\}, \forall s \in \mathcal{S}. \quad (3.5)$$

**end for**

---