

486 **A Proof of Lemma 1**

488 We show that $\zeta^{\mathcal{L}}(s, a; D_{s,a}, k)$ defined in Eq. (7) decays to 0 at a rate of $1/n_{s,a}^2$. We present the
 489 proof for $k = 1$. The extension to $k > 1$ is straightforward.

490 In finite MDPs, we have to learn separate transition probability tables $\mathcal{T}(s' \mid s, a)$ for each s, a .
 491 For simplicity, we focus on one fixed (s, a) and investigate how to estimate the distribution $\mathcal{T}(s')$.
 492 We consider a Dirichlet learner where α denotes the Dirichlet posterior based on the data $D_{s,a} =$
 493 $\{s'_i\}_{i=1}^{n_{s,a}}$ and α' be the posterior based on the reduced data set $D_{s,a}^{k=1}$, that is, $\alpha'_c + 1 = \alpha_c$ where
 494 c is the outcome of the missing experience in $D_{s,a}^{k=1}$, and $\alpha'_j = \alpha_j$ for all $j \neq c$. Given Dirichlet
 495 parameters α , the MAP model $\hat{\mathcal{T}}_\alpha(s')$ is given by the vector $\alpha/\bar{\alpha}$, $\bar{\alpha} = \sum_i \alpha_i$, and we estimate ζ
 496 using Eq. (7). The log-likelihood of the data $D_{s,a}$ under $\hat{\mathcal{T}}_\alpha$ is
 497

$$498 \quad L^+ := \log P(D_{s,a} \mid \hat{\mathcal{T}}_\alpha) \\ 499 \quad = \log \prod_{i=1}^{n_{s,a}} \frac{\alpha_{s'_i}}{\bar{\alpha}} = \sum_{i=1}^{n_{s,a}} \log \alpha_{s'_i} - n_{s,a} \log \bar{\alpha} \quad (14)$$

503 The likelihood of the data $D_{s,a}$ under $\hat{\mathcal{T}}_{\alpha'}$ is

$$505 \quad L^- := \log P(D_{s,a} \mid \hat{\mathcal{T}}_{\alpha'}) \\ 506 \quad = \log \left(\prod_{i=1, s'_i \neq c}^{n_{s,a}} \frac{\alpha_{s'_i}}{\bar{\alpha} - 1} \right) \left(\prod_{i=1, s'_i = c}^{n_{s,a}} \frac{\alpha_c - 1}{\bar{\alpha} - 1} \right) \\ 507 \quad = \sum_{i=1, s'_i \neq c}^{n_{s,a}} \log \alpha_{s'_i} + \sum_{i=1, s'_i = c}^{n_{s,a}} \log(\alpha_c - 1) - n_{s,a} \log(\bar{\alpha} - 1) \quad (15)$$

513 The average difference is

$$515 \quad \zeta^{\mathcal{L}}(s, a; k=1) = \frac{1}{n_{s,a}} |L^+ - L^-| \\ 516 \quad = \frac{1}{n_{s,a}} \sum_{i=1, s'_i = c}^{n_{s,a}} (\log \alpha_c - \log(\alpha_c - 1)) - \log \bar{\alpha} + \log(\bar{\alpha} - 1) \\ 517 \quad = \frac{n_{s,a,c}}{n_{s,a}} \log \frac{1}{1 - \frac{1}{\alpha_c}} + \log \left(1 - \frac{1}{\bar{\alpha}} \right). \quad (16)$$

523 Since $\bar{\alpha} \propto n_{s,a}$, by taking the derivative of the expected value of $\zeta^{\mathcal{L}}(s, a; k=1)$ we can verify that

$$524 \quad E_{D_{s,a}}(\zeta(s, a; D_{s,a}, k=1)) = O\left(\frac{1}{n_{s,a}^2}\right). \square$$

526
527
528
529
530
531
532
533
534
535
536
537
538
539