

Appendix

VII. PROOF OF COROLLARY 1

Proof: Drawing inspiration from the proofs by Khaled and Richtarik [31] and Yuan et al. [30], we initiate our analysis by considering Assumption 1, which concerns Lipschitz smoothness. For the risk-sensitive objective J_β with a Lipschitz smoothness constant L_β , we have that:

$$\begin{aligned} J_\beta(\theta_{t+1}) &\geq J_\beta(\theta_t) + \langle \nabla J_\beta(\theta_t), \theta_{t+1} - \theta_t \rangle - \frac{L_\beta}{2} \|\theta_{t+1} - \theta_t\|^2 \\ &= J_\beta(\theta_t) + \eta \left\langle \nabla J_\beta(\theta_t), \hat{\nabla} J_\beta(\theta_t) \right\rangle - \frac{L_\beta \eta^2}{2} \left\| \hat{\nabla} J_\beta(\theta_t) \right\|^2. \end{aligned} \quad (31)$$

Take expectations on both sides conditioned on θ_t and use Assumption 2, we get:

$$\begin{aligned} \mathbb{E}_t [J_\beta(\theta_{t+1})] &\geq J_\beta(\theta_t) + \eta \langle \nabla J_\beta(\theta_t), \nabla J_\beta(\theta_t) \rangle - \frac{L_\beta \eta^2}{2} \mathbb{E}_t \left[\left\| \hat{\nabla} J_\beta(\theta_t) \right\|^2 \right] \\ &\geq J_\beta(\theta_t) + \eta \|\nabla J_\beta(\theta_t)\|^2 - \frac{L_\beta \eta^2}{2} \left(2A(J_\beta^* - J_\beta(\theta_t)) + B \|\nabla J_\beta(\theta_t)\|^2 + C \right) \\ &= J_\beta(\theta_t) + \eta \left(1 - \frac{L_\beta B \eta}{2} \right) \|\nabla J_\beta(\theta_t)\|^2 - L_\beta \eta^2 A(J_\beta^* - J_\beta(\theta_t)) - \frac{L_\beta C \eta^2}{2}. \end{aligned} \quad (32)$$

Then we subtract J_β^* from both sides,

$$\mathbb{E}_t [J_\beta(\theta_{t+1})] - J_\beta^* \geq -(1 + L_\beta \eta^2 A)(J_\beta^* - J_\beta(\theta_t)) + \eta \left(1 - \frac{L_\beta B \eta}{2} \right) \|\nabla J_\beta(\theta_t)\|^2 - \frac{L_\beta C \eta^2}{2}. \quad (33)$$

Take the expectation on both sides and rearrange the equation, we obtain:

$$\mathbb{E} [J_\beta^* - J_\beta(\theta_{t+1})] + \eta \left(1 - \frac{L_\beta B \eta}{2} \right) \mathbb{E} [\|\nabla J_\beta(\theta_t)\|^2] \leq (1 + L_\beta \eta^2 A) \mathbb{E} [J_\beta^* - J_\beta(\theta_t)] + \frac{L_\beta C \eta^2}{2}. \quad (34)$$

Define $\delta_t \stackrel{\text{def}}{=} \mathbb{E} [J_\beta^* - J_\beta(\theta_t)]$ and $r_t \stackrel{\text{def}}{=} \mathbb{E} [\|\nabla J_\beta(\theta_t)\|^2]$, we can rewrite the above inequality as

$$\eta \left(1 - \frac{L_\beta B \eta}{2} \right) r_t \leq (1 + L_\beta \eta^2 A) \delta_t - \delta_{t+1} + \frac{L_\beta C \eta^2}{2}. \quad (35)$$

Now, we introduce a sequence of weights, denoted as $w_{-1}, w_0, w_1, \dots, w_{T-1}$, based on a method used by [30, 31, 37]. We initialize w_{-1} with a positive value. We define w_t as $w_t =: \frac{w_{t-1}}{1 + L_\beta \eta^2 A}$ for all $t \geq 0$. It's important to note that when $A = 0$, all w_t are equal, i.e., $w_t = w_{t-1} = \dots = w_{-1}$. By multiplying (35) by w_t/η , we can derive:

$$\begin{aligned} \left(1 - \frac{L_\beta B \eta}{2} \right) w_t r_t &\leq \frac{w_t (1 + L_\beta \eta^2 A)}{\eta} \delta_t - \frac{w_t}{\eta} \delta_{t+1} + \frac{L_\beta C \eta}{2} w_t \\ &= \frac{w_{t-1}}{\eta} \delta_t - \frac{w_t}{\eta} \delta_{t+1} + \frac{L_\beta C \eta}{2} w_t. \end{aligned} \quad (36)$$

When we sum up both sides for $t = 0, 1, \dots, T-1$, we get:

$$\begin{aligned} \left(1 - \frac{L_\beta B \eta}{2} \right) \sum_{t=0}^{T-1} w_t r_t &\leq \frac{w_{-1}}{\eta} \delta_0 - \frac{w_{T-1}}{\eta} \delta_T + \frac{L_\beta C \eta}{2} \sum_{t=0}^{T-1} w_t \\ &\leq \frac{w_{-1}}{\eta} \delta_0 + \frac{L_\beta C \eta}{2} \sum_{t=0}^{T-1} w_t. \end{aligned} \quad (37)$$

We can define W_T as $W_T =: \sum_{t=0}^{T-1} w_t$. By dividing both sides of the equation by W_T , we obtain:

$$\left(1 - \frac{L_\beta B \eta}{2} \right) \min_{0 \leq t \leq T-1} r_t \leq \frac{1}{W_T} \cdot \left(1 - \frac{L_\beta B \eta}{2} \right) \sum_{t=0}^{T-1} w_t r_t \leq \frac{w_{-1}}{W_T} \frac{\delta_0}{\eta} + \frac{L_\beta C \eta}{2}. \quad (38)$$

Note that,

$$W_T = \sum_{t=0}^{T-1} w_t \geq \sum_{t=0}^{T-1} \min_{0 \leq i \leq T-1} w_i = T w_{T-1} = \frac{T w_{-1}}{(1 + L_\beta \eta^2 A)^T}. \quad (39)$$

Use this in (38),

$$\left(1 - \frac{L_\beta B \eta}{2}\right) \min_{0 \leq t \leq T-1} r_t \leq \frac{(1 + L\eta^2 A)^T}{\eta T} \delta_0 + \frac{LC\eta}{2}. \quad (40)$$

By substituting r_t in (40) with $\mathbb{E} [\|\nabla J_\beta(\theta_t)\|^2]$, we obtain:

$$\begin{aligned} \left(1 - \frac{L_\beta B \eta}{2}\right) \min_{0 \leq t \leq T-1} \mathbb{E} [\|\nabla J_\beta(\theta_t)\|^2] &\leq \frac{(1 + L\eta^2 A)^T}{\eta T} \delta_0 + \frac{L_\beta C \eta}{2}, \\ \min_{0 \leq t \leq T-1} \mathbb{E} [\|\nabla J_\beta(\theta_t)\|^2] &\leq \frac{2\delta_0(1 + L_\beta \eta^2 A)^T}{\eta T(2 - L_\beta B \eta)} + \frac{L_\beta C \eta}{2 - L_\beta B \eta}. \end{aligned} \quad (41)$$

The choice of our step size ensures that for both cases, whether $B > 0$ or $B = 0$, we have $1 - \frac{L_\beta B \eta}{2} > 0$. ■

VIII. PROOF OF THEOREM 1

Lemma 1: Subject to Assumption 4, for any non-negative integer t , and for any state-action pair $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ at time t within a trajectory τ sampled under the parameterized policy π_θ , we have the following:

$$\mathbb{E}_{\tau \sim p(\cdot|\theta)} [\|\nabla_\theta \log \pi_\theta(a_t | s_t)\|^2] \leq F_1^2, \quad (42)$$

$$\mathbb{E}_{\tau \sim p(\cdot|\theta)} [\|\nabla_\theta^2 \log \pi_\theta(a_t | s_t)\|] \leq F_2. \quad (43)$$

Proof: For $t > 0$ and $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$, we have

$$\mathbb{E}_\tau [\|\nabla_\theta \log \pi_\theta(a_t | s_t)\|^2] = \mathbb{E}_{s_t} [\mathbb{E}_{a_t \sim \pi_\theta(\cdot|s_t)} [\|\nabla_\theta \log \pi_\theta(a_t | s_t)\|^2 | s_t]] \stackrel{22}{\leq} F_1^2, \quad (44)$$

where the first equality is obtained by the Markov property. Similarly, we have

$$\mathbb{E}_\tau [\|\nabla_\theta^2 \log \pi_\theta(a_t | s_t)\|] = \mathbb{E}_{s_t} [\mathbb{E}_{a_t \sim \pi_\theta(\cdot|s_t)} [\|\nabla_\theta^2 \log \pi_\theta(a_t | s_t)\| | s_t]] \stackrel{23}{\leq} F_2. \quad (45)$$

Then Lemma 1 is then used for the derivation of L and L_β .

Assumption 1 is equivalent to $\|\nabla^2 J(\theta)\| \leq L$ for the risk-neutral REINFORCE and $\|\nabla^2 J_\beta(\theta)\| \leq L_\beta$ for the risk-sensitive REINFORCE. We first take the second order derivative of the risk-neutral objective w.r.t. θ , in order to derive the L -Lipschitz smooth constant.

$$\begin{aligned} \nabla^2 J(\theta) &\stackrel{??}{=} \nabla_\theta \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right] \\ &= \nabla_\theta \left[\int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) d\tau \right] \\ &= \int \nabla_\theta p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right)^\top d\tau + \int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) R(t) d\tau \\ &= \int p(\tau | \theta) \nabla_\theta \log p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right)^\top d\tau + \int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) R(t) d\tau \\ &= \mathbb{E}_\tau \left[\nabla_\theta \log p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right)^\top \right] + \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) R(t) \right] \\ &\stackrel{9}{=} \underbrace{\mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right)^\top \right]}_{\textcircled{1}} + \underbrace{\mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) R(t) \right]}_{\textcircled{2}}. \end{aligned} \quad (46)$$

We individually bound the aforementioned two terms for the risk-neutral REINFORCE.

For the term ①,

$$\begin{aligned}
\|\textcircled{1}\| &= \left\| \mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) R(t) \right)^\top \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right)^\top \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \left(\sum_{t'=0}^t \nabla_\theta \log \pi_\theta(a_{t'} | \theta_{t'}) \right) \left(\sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right)^\top \right] \right\| \\
&\leq \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)| \left\| \sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right\|^2 \right] \\
&\leq r_{\max} \sum_{t=0}^{\infty} \gamma^t \sum_{k=0}^t \mathbb{E}_\tau \left[\|\nabla_\theta \log \pi_\theta(a_k | \theta_k)\|^2 \right] \\
&\stackrel{42}{\leq} r_{\max} F_1^2 \sum_{t=0}^{\infty} \gamma^t (t+1) \\
&= \frac{r_{\max} F_1^2}{(1-\gamma)^2},
\end{aligned} \tag{47}$$

where the third line is due to the fact that the future actions do not depend on the past rewards.

For the term ②,

$$\begin{aligned}
\|\textcircled{2}\| &= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) R(t) \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \left(\sum_{k=0}^t \nabla_\theta^2 \log \pi_\theta(a_k | s_k) \right) \right] \right\| \\
&\leq \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)| \left(\sum_{k=0}^t \|\nabla_\theta^2 \log \pi_\theta(a_k | s_k)\| \right) \right] \\
&\leq r_{\max} \sum_{t=0}^{\infty} \gamma^t \left(\sum_{k=0}^t \mathbb{E}_\tau [\|\nabla_\theta^2 \log \pi_\theta(a_k | s_k)\|] \right) \\
&\stackrel{43}{\leq} r_{\max} F_2 \sum_{t=0}^{\infty} \gamma^t (t+1) \\
&= \frac{r_{\max} F_2}{(1-\gamma)^2},
\end{aligned} \tag{48}$$

where the third line is also due to the fact that the future actions do not depend on the past rewards.

Finally,

$$\|\nabla^2 J(\theta)\| \leq \frac{r_{\max}}{(1-\gamma)^2} (F_1^2 + F_2), \tag{49}$$

so $L = \frac{r_{\max}}{(1-\gamma)^2} (F_1^2 + F_2)$.

Then we take the second order derivative of the risk-sensitive objective w.r.t. θ , in order to derive the L_β -Lipschitz

smoothness constant.

$$\begin{aligned}
\nabla^2 J_\beta(\theta) &\stackrel{15}{=} \nabla_\theta \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right] \\
&= \nabla_\theta \left[\int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} d\tau \right] \\
&= \int \nabla_\theta p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right)^\top d\tau + \int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} d\tau \\
&= \int p(\tau | \theta) \nabla_\theta \log p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right)^\top d\tau \\
&\quad + \int p(\tau | \theta) \sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} d\tau \\
&= \mathbb{E}_\tau \left[\nabla_\theta \log p(\tau | \theta) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right)^\top \right] + \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right] \\
&\stackrel{9}{=} \underbrace{\mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right)^\top \right]}_{\textcircled{3}} \\
&\quad + \underbrace{\mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right]}_{\textcircled{4}}. \tag{50}
\end{aligned}$$

We also bound the above two terms separately for the risk-sensitive REINFORCE algorithm.

For the term $\textcircled{3}$,

$$\begin{aligned}
\|\textcircled{3}\| &= \left\| \mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right)^\top \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{k=0}^{\infty} \nabla_\theta \log \pi_\theta(a_k | \theta_k) \left(\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'})} \right)^\top \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\beta e^{\beta \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)} \left(\sum_{t'=0}^t \nabla_\theta \log \pi_\theta(a_{t'} | \theta_{t'}) \right) \left(\sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right)^\top \right] \right\| \\
&\leq \mathbb{E}_\tau \left[|\beta| e^{|\beta| \sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)|} \left\| \sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right\|^2 \right] \\
&\stackrel{21}{=} \mathbb{E}_\tau \left[\alpha \sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)| \left\| \sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right\|^2 \right] \\
&\leq \alpha \cdot r_{\max} \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_\tau \left[\left\| \sum_{k=0}^t \nabla_\theta \log \pi_\theta(a_k | \theta_k) \right\|^2 \right] \\
&= \alpha \cdot r_{\max} \sum_{t=0}^{\infty} \gamma^t \sum_{k=0}^t \mathbb{E}_\tau \left[\|\nabla_\theta \log \pi_\theta(a_k | \theta_k)\|^2 \right] \\
&\leq \alpha \cdot r_{\max} F_1^2 \sum_{t=0}^{\infty} \gamma^t (t+1) \\
&= \alpha \cdot \frac{r_{\max} F_1^2}{(1-\gamma)^2}, \tag{51}
\end{aligned}$$

where in the third line, we use the fact that the future actions do not depend on the past rewards. In the fifth line, we use Assumption 3.

For the term ④,

$$\begin{aligned}
\|\textcircled{4}\| &= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta R(t)} \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \nabla_\theta^2 \log \pi_\theta(a_t | s_t) \cdot \beta e^{\beta \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'})} \right] \right\| \\
&= \left\| \mathbb{E}_\tau \left[\beta e^{\beta \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)} \left(\sum_{k=0}^t \nabla_\theta^2 \log \pi_\theta(a_k | s_k) \right) \right] \right\| \\
&\leq \mathbb{E}_\tau \left[|\beta| e^{|\beta| \sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)|} \left(\sum_{k=0}^t \|\nabla_\theta^2 \log \pi_\theta(a_k | s_k)\| \right) \right] \\
&\stackrel{21}{=} \mathbb{E}_\tau \left[\alpha \cdot \sum_{t=0}^{\infty} \gamma^t |r(s_t, a_t)| \left(\sum_{k=0}^t \|\nabla_\theta^2 \log \pi_\theta(a_k | s_k)\| \right) \right] \\
&\leq \alpha \cdot r_{\max} \sum_{t=0}^{\infty} \gamma^t \left(\sum_{k=0}^t \mathbb{E}_\tau [\|\nabla_\theta^2 \log \pi_\theta(a_k | s_k)\|] \right) \\
&\stackrel{43}{\leq} \alpha \cdot r_{\max} F_2 \sum_{t=0}^{\infty} \gamma^t (t+1) \\
&= \alpha \cdot \frac{r_{\max} F_2}{(1-\gamma)^2} \tag{52}
\end{aligned}$$

Finally,

$$\|\nabla^2 J_\beta(\theta)\| \leq \alpha \cdot \frac{r_{\max}}{(1-\gamma)^2} (F_1^2 + F_2), \tag{53}$$

so $L_\beta = \alpha \cdot \frac{r_{\max}}{(1-\gamma)^2} (F_1^2 + F_2)$, where $0 < \alpha < 1$.