

CMPT 419/983: Theoretical Foundations of Reinforcement Learning Assignment 4

Total marks: 250

Due: December 10, 11:59 pm

Submission Instructions

- Assignments typed in Latex (the Latex source is provided) are preferred, but can be handwritten.
- You can write the code in any language, but are not allowed to make use of automatic differentiation (using Numpy + Python is preferred). The code and plots is to be submitted as a separate zip file. All code files and plots should be stored in a directory named **a4 and then zip** the directory for submission.
- Assignment (PDF + separate zip file for code and plots) is to be submitted online via Coursys.

(1) [100 marks] Proving the remaining results from class

- Prove that the **Jacobian of h** : $\mathbb{R}^A \rightarrow \mathbb{R}^A$ is given by:

$$H(\pi_\theta) \in \mathbb{R}^{A \times A} = \text{diag}(\pi_\theta) - \pi_\theta \pi_\theta^T,$$

where $\text{diag}(\pi_\theta) \in \mathbb{R}^{A \times A}$ is a diagonal matrix s.t. for any s , $[\text{diag}(\pi_\theta)]_{a,a} = \pi_\theta(a|s)$ and $\pi_\theta(\cdot|s) \in \mathbb{R}^A$ s.t. $\pi_\theta(a|s) = \frac{\exp(\theta(s,a))}{\sum_{a'} \exp(\theta(s,a'))}$. **Using the above calculation**, show that for any s', a' ,

$$\frac{\partial \log(\pi_\theta(a'|s'))}{\partial \theta(s,a)} = \mathcal{I}\{s' = s\} [\mathcal{I}\{a' = a\} - \pi_\theta(a|s)] \quad [20 \text{ marks}]$$

- For arbitrary policies $\pi, \tilde{\pi}$: prove that

$$v^\pi(\rho) - v^{\tilde{\pi}}(\rho) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^\pi} \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathbf{a}^{\tilde{\pi}}(s, a)] \quad [25 \text{ marks}]$$

- Consider the modified TRPO update: for $\theta \in \mathbb{R}^d$,

$$\theta_{t+1} = \arg \max_{\theta} \left[v^{\pi_t}(\rho) + \underbrace{\frac{\mathbb{E}_{s \sim d^{\pi_t}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \mathbf{a}^{\pi_t}(s, a)}{1 - \gamma}}_{\text{Term (i)}} - \beta \underbrace{\sum_s d_t^\pi(s) \text{KL}(\pi_t(\cdot|s) || \pi_\theta(\cdot|s))}_{\text{Term (ii)}} \right]$$

Prove that using a linear approximation of Term (i) and a quadratic approximation of Term (ii) results in the following update:

$$\theta_{t+1} = \arg \max_{\theta} \left[v^{\pi_t}(\rho) + \langle \nabla_{\theta} \text{Term (i)} |_{\theta=\theta_t}, \theta \rangle - \frac{\beta}{2} (\theta - \theta_t)^T [\nabla_{\theta}^2 \text{Term (ii)} |_{\theta=\theta_t}] (\theta - \theta_t) \right]$$

and recovers the update for natural gradient descent, i.e. $\theta_{t+1} = \theta_t + \frac{1}{\beta} F_{\theta_t}^{\dagger} \nabla J(\theta_t)$. [30 marks]

- For the finite-horizon problem in Lecture 12, prove the following result that bounds the difference in the performance of the same deterministic policy π on two different MDPs: $M = (\mathcal{S}, \mathcal{A}, \{\mathcal{P}_h\}_{h=0}^{H-1}, \{r_h\}_{h=0}^{H-1})$ and $\tilde{M} = (\mathcal{S}, \mathcal{A}, \{\tilde{\mathcal{P}}_h\}_{h=0}^{H-1}, \{\tilde{r}_h\}_{h=0}^{H-1})$. If $v_0^{\pi, M}$ is the value function of policy π on MDP M from $h=0$, assuming $v_H^{\pi, M} = v_H^{\pi, \tilde{M}} = 0$, prove that for a starting state $s_0 \in \mathcal{S}$,

$$\begin{aligned} & v_0^{\pi, M}(s_0) - v_0^{\pi, \tilde{M}}(s_0) \\ &= \mathbb{E}_{\substack{s_{h+1} \sim \tilde{\mathcal{P}}_h(\cdot | s_h, a_h) \\ a_h = \pi_h(s_h)}} \sum_{h=0}^{H-1} \left[r_h(s_h, a_h) - \tilde{r}_h(s_h, a_h) \right] + \langle \mathcal{P}_h(\cdot | s_h, a_h) - \tilde{\mathcal{P}}_h(\cdot | s_h, a_h), v_{h+1}^{\pi, M} \rangle \end{aligned} \quad [25 \text{ marks}]$$

- (2) [70 marks] **A new PG method** Consider the tabular softmax PG update, i.e. for $\theta \in \mathbb{R}^{\mathcal{S}\mathcal{A}}$, $\pi_{\theta}(a|s) = \frac{\exp(\theta(s,a))}{\sum_{a'} \exp(\theta(s,a'))}$. The update $\theta_{t+1} = \theta_t + \eta \nabla_{\theta} J(\theta_t)$ can be alternatively written as:

$$\theta_{t+1} = \arg \max_{\theta} \left[\langle \theta - \theta_t, \nabla_{\theta} J(\theta_t) \rangle - \frac{1}{\eta} \|\theta - \theta_t\|_2^2 \right]$$

As in Lecture 8, the update can be generalized to measure the distance between θ and θ_t using a Bregman divergence, resulting in the following update:

$$\theta_{t+1} = \arg \max_{\theta} \left[\langle \theta - \theta_t, \nabla_{\theta} J(\theta_t) \rangle - \frac{1}{\eta} D_{\Psi}(\theta, \theta_t) \right], \quad (1)$$

where ψ is the mirror map.

- For $\theta(s, \cdot) \in \mathbb{R}^{\mathcal{A}}$, prove that choosing $\psi(\theta(s, \cdot)) = \log(\sum_a \exp(\theta(s, a)))$ results in $D_{\psi}(\theta(s, \cdot), \theta'(s, \cdot)) = \text{KL}(\pi'(\cdot|s) || \pi(\cdot|s))$, where $\pi_{\theta}(a|s) = \frac{\exp(\theta(s,a))}{\sum_{a'} \exp(\theta(s,a'))}$ [25 marks]
- For the update in Eq. (1), we will use the following Bregman divergence: $D_{\Psi}(\theta, \theta') = \sum_s d_{\pi'}^{\theta}(s) D_{\psi}(\theta(s, \cdot), \theta'(s, \cdot))$. If $\pi_t := \pi_{\theta_t}$, prove that this results in the following update rule:

$$\theta_{t+1} = \arg \max_{\theta} \left[\mathbb{E}_{s \sim d^{\pi_t}} \mathbb{E}_{a \sim \pi_t(\cdot|s)} \left[\left(\frac{\alpha^{\pi_t}(s, a)}{1 - \gamma} + \frac{1}{\eta} \right) \log \left(\frac{\pi_{\theta}(a|s)}{\pi_t(a|s)} \right) \right] \right] \quad [30 \text{ marks}]$$

- For the tabular policy parameterization, prove that the above problem is equivalent to the following problem:

$$\pi_{t+1}(\cdot|s) = \arg \max_{\pi(\cdot|s) \in \Delta_{\mathcal{A}}} \left[\mathbb{E}_{s \sim d^{\pi_t}} \left[\mathbb{E}_{a \sim \pi_t(\cdot|s)} \left(\frac{\alpha^{\pi_t}(s, a)}{1 - \gamma} \log \left(\frac{\pi(a|s)}{\pi_t(a|s)} \right) \right) - \frac{1}{\eta} \text{KL}(\pi_t(\cdot|s) || \pi(\cdot|s)) \right] \right] \quad [10 \text{ marks}]$$

- Compare the above expression to the regularized TRPO update (v1) for the tabular parameterization where Π_{θ} consists of distributions $\pi(\cdot|s) \in \mathbb{R}^{\mathcal{A}}$ for each state s . What are the advantages/disadvantages of using one over the other? [5 marks]

(3) **[50 marks] PG with Log-linear policies** When \mathcal{S} and \mathcal{A} is large, function approximation may be needed to reduce the dimension. For simplicity, let us consider **log-linear policy in the bandit setting (with deterministic rewards)**. A log-linear policy is such that for all $a \in [K] := \{1, 2, \dots, K\}$,

log-linear policy 中 θ 的维度是 $d=2$

$$\pi_{\theta}(a) = \frac{\exp([\Phi\theta](a))}{\sum_{a' \in [K]} \exp([\Phi\theta](a'))}$$

where $\Phi \in \mathbb{R}^{K \times d}$ is the feature matrix with full column rank $d \leq K$. When $d = K$ and the features are one-hot vectors, log-linear policies reduce to the tabular setting. The objective is:

$$J(\theta) := \mathbb{E}_{a \sim \pi_{\theta}}[r(a)] = \langle \pi_{\theta}, r \rangle$$

- Derive the **softmax policy gradient $\nabla_{\theta} J(\theta)$** for the above setting. **[10 marks]**

Consider the bandit setting with **$K = 4$** with the following **linear features** and **reward vector**:

$$\Phi = \begin{bmatrix} 0 & 2 \\ 0.4 & 0 \\ -2 & 0 \\ 0 & -0.4 \end{bmatrix} \quad r = \begin{bmatrix} 1 \\ 0.9 \\ 0.8 \\ 0.7 \end{bmatrix} \quad \text{bandit with softmax tabular setting: L9}$$

Compare **log-linear policies** with the above **linear features** to **tabular features** in **both the deterministic and stochastic settings**. For the stochastic setting, use the **importance-weighted reward estimate to construct the gradient**. In each setting initialize $\theta_0(a) = 0$ for all $a \in [K]$, grid-search over the range of $\eta \in \{10^{-3}, 10^{-2}, 10^{-1}, 1\}$ to select the step-size. L11

- In the deterministic setting, with **$T = 1000$** , plot the **suboptimality gap: $\langle \pi^*, r \rangle - \langle \pi_{\theta_t}, r \rangle$** for both **linear** and **tabular features** (with the corresponding best step-size) on the same plot. **Use a log-scale (if required) to better visualize the results**. Save the final plot as **pg.png**. **[20 marks]**
- In the stochastic setting, with **$T = 10^4$** , run the experiment for both linear and tabular features (and different step-sizes) **20 times**, and plot the average suboptimality gap and standard deviation on the same plot. Use a log-scale (if required) to better visualize the results. Save the final plot as **spg.png**. **[20 marks]**

(4) **[30 marks] Empirical Evaluation** In `mdp.py`, the following functions will generate an instance of an MDP and the corresponding v^* .

- `generate_cliffworld()` will return the following:
 - $\mathcal{S} := \{0, 1, \dots, 20\}$
 - $\mathcal{A} := \{0, 1, \dots, 3\}$
 - $r \in \mathbb{R}^{21 \times 4}$
 - $\mathcal{P} \in \mathbb{R}^{21 \times 4 \times 21}$ where $P[s][a][\text{next_s}] = \mathcal{P}[s'|s, a]$
 - $\rho \in \mathbb{R}^{21}$

tabular 的 θ : tabular softmax PG 和 NPG
 1. 更新 θ 得到新的 π ;
 2. 通过 π 计算 $v^{\pi_{\theta}}$;
 迭代

For the given MDP, let **$\gamma := 0.9$** and implement **PG and NPG**. For both PG and NPG, use **$\eta = \frac{1}{L} = \frac{(1-\gamma)^3}{8}$** and let **$\theta_0(\cdot|s) = 0$** for all $s \in \mathcal{S}$. Run each algorithm for **$T = 10^4$** iterations and plot the **log suboptimality gap: $\log(v^{\pi^*}(\rho) - v^{\pi_{\theta_t}}(\rho))$** . Save the final plot as **pg_vs_npg.png**.