

Motivation

Markov Decision Process $(S, \mathcal{A}, P, R_r, \gamma)$:

$$\max_{\pi} \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t)$$

- Single scalar reward to represent multiple contradictory aspects (e.g. complete task and avoid collisions)
- No control over the spread of the performance distribution

Constrained MDP $(S, \mathcal{A}, P, R_r, R_c, \gamma)$

- [Beutler and Ross 1985; Altman 1999]
- Introduce a cost signal R_c
- Constrained objective

$$\max_{\pi \in \mathcal{M}(\mathcal{A})^S} \mathbb{E}[G_r^{\pi} | s_0 = s] \quad \text{s.t.} \quad \mathbb{E}[G_c^{\pi} | s_0 = s] \leq \beta$$

- The cost budget β cannot be changed after training

Budgeted MDP

- [Boutillier and Lu 2016]
- We seek one general policy $\pi(s, \beta)$ that solves every CMDP for any β
- Can only be solved for finite S and known P, R_r, R_c .

Setting

Budgeted policies π

- Take a budget β as an additional input
- Output a next budget β'

$$\pi : \underbrace{(s, \beta)}_s \rightarrow \underbrace{(a, \beta')}_{\bar{a}}$$

2D signals

- Rewards $R = (R_r, R_c)$
- Returns $G^{\pi} = (G_r^{\pi}, G_c^{\pi})$
- Values $V^{\pi} = (V_r^{\pi}, V_c^{\pi})$ and $Q^{\pi} = (Q_r^{\pi}, Q_c^{\pi})$

Policy Evaluation

The Bellman Expectation equations are preserved, and the Bellman Expectation Operator \mathcal{T}^{π} is a γ -contraction.

Budgeted Optimality

Definition. In that order, we want to:

(i) Respect the budget β :

$$\Pi_a(\bar{s}) \stackrel{\text{def}}{=} \{\pi \in \Pi : V_c^{\pi}(s, \beta) \leq \beta\}$$

(ii) Maximise the rewards:

$$V_r^*(\bar{s}) \stackrel{\text{def}}{=} \max_{\pi \in \Pi_a(\bar{s})} V_r^{\pi}(\bar{s}) \quad \Pi_r(\bar{s}) \stackrel{\text{def}}{=} \arg \max_{\pi \in \Pi_a(\bar{s})} V_r^{\pi}(\bar{s})$$

(iii) Minimise the costs:

$$V_c^*(\bar{s}) \stackrel{\text{def}}{=} \min_{\pi \in \Pi_r(\bar{s})} V_c^{\pi}(\bar{s}), \quad \Pi^*(\bar{s}) \stackrel{\text{def}}{=} \arg \min_{\pi \in \Pi_r(\bar{s})} V_c^{\pi}(\bar{s})$$

We define the budgeted action-value function Q^* similarly

Acknowledgements

This work has been supported by CPER Nord-Pas de Calais/FEDER DATA Advanced data science and technologies 2015-2020, the French Ministry of Higher Education and Research, INRIA, and the French Agence Nationale de la Recherche (ANR).

References

- Eitan Altman. *Constrained Markov Decision Processes*. CRC Press, 1999.
- Frederick J. Beutler and Keith W. Ross. "Optimal policies for controlled Markov chains with a constraint". In: *Journal of Mathematical Analysis and Applications*. 1985.
- Craig Boutilier and Tyler Lu. "Budget Allocation using Weakly Coupled, Constrained Markov Decision Processes". In: *Uncertainty in Artificial Intelligence (UAI)*. 2016.

Budgeted Dynamic Programming

Theorem (Budgeted Bellman Optimality). Q^* verifies:

$$Q^*(\bar{s}, \bar{a}) = \mathcal{T}Q^*(\bar{s}, \bar{a}) \stackrel{\text{def}}{=} R(\bar{s}, \bar{a}) + \gamma \sum_{\bar{s}' \in \bar{S}} P(\bar{s}' | \bar{s}, \bar{a}) \sum_{\bar{a}' \in \bar{A}} \pi_{\text{greedy}}(\bar{a}' | \bar{s}'; Q^*) Q^*(\bar{s}', \bar{a}'), \quad (1)$$

where the greedy policy π_{greedy} is defined by:

$$\pi_{\text{greedy}}(\bar{a} | \bar{s}; Q) \in \arg \min_{\rho \in \Pi_r^Q} \mathbb{E}_{\bar{a} \sim \rho} Q_c(\bar{s}, \bar{a}), \quad (2a)$$

$$\text{where } \Pi_r^Q \stackrel{\text{def}}{=} \arg \max_{\rho \in \mathcal{M}(\bar{A})} \mathbb{E}_{\bar{a} \sim \rho} Q_r(\bar{s}, \bar{a}) \quad (2b)$$

$$\text{s.t. } \mathbb{E}_{\bar{a} \sim \rho} Q_c(\bar{s}, \bar{a}) \leq \beta \quad (2c)$$

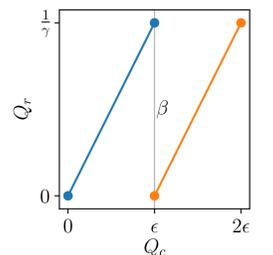
Proposition. $\pi_{\text{greedy}}(\cdot; Q^*)$ is simultaneously optimal in all states $\bar{s} \in \bar{S}$:

$$\pi_{\text{greedy}}(\cdot; Q^*) \in \Pi^*(\bar{s})$$

In particular, $V^{\pi_{\text{greedy}}(\cdot; Q^*)} = V^*$ and $Q^{\pi_{\text{greedy}}(\cdot; Q^*)} = Q^*$.

Algorithm 1: Budgeted Value Iteration

Data: P, R_r, R_c
Result: Q^*
1 $Q_0 \leftarrow 0$
2 repeat
3 $Q_{k+1} \leftarrow \mathcal{T}Q_k$
4 until convergence



Theorem (Contractivity). For any BMDP $(S, \mathcal{A}, P, R_r, R_c, \gamma)$ with $|\mathcal{A}| \geq 2$, \mathcal{T} is not a contraction.

$$\forall \epsilon > 0, \exists Q^1, Q^2 \in (\mathbb{R}^2)^{\bar{S}\bar{A}} : \|\mathcal{T}Q^1 - \mathcal{T}Q^2\|_{\infty} \geq \frac{1}{\epsilon} \|Q^1 - Q^2\|_{\infty}$$

Despite those theoretical limitations, we observed empirical convergence in our experiments

Theorem (Contractivity on smooth Q-functions). \mathcal{T} is a contraction when restricted to the subset \mathcal{L}_{γ} of Q-functions such that " Q_r is L -Lipschitz with respect to Q_c ", with $L < \frac{1}{\gamma} - 1$.

Budgeted Reinforcement Learning

We address several limitations of Algorithm 1.

1. The BMDP is unknown

- Work with a batch of samples $\mathcal{D} = \{(\bar{s}_i, \bar{a}_i, r_i, \bar{s}'_i)\}_{i \in [0, N]}$

2. \mathcal{T} contains an expectation $\mathbb{E}_{\bar{s}' \sim \bar{P}}$ over next states \bar{s}'

- Replace it with a sampling operator $\hat{\mathcal{T}}$:

$$\hat{\mathcal{T}}Q(\bar{s}_i, \bar{a}_i, r_i, \bar{s}'_i) \stackrel{\text{def}}{=} r_i + \gamma \sum_{\bar{a}'_i \in \bar{A}_i} \pi_{\text{greedy}}(\bar{a}'_i | \bar{s}'_i; Q) Q(\bar{s}'_i, \bar{a}'_i).$$

3. S is continuous

- Employ function approximation Q_{θ} , and minimise a regression loss

$$\mathcal{L}(Q_{\theta}, Q_{\text{target}}; \mathcal{D}) = \sum_{\mathcal{D}} \|Q_{\theta}(\bar{s}, \bar{a}) - Q_{\text{target}}(\bar{s}, \bar{a}, r, \bar{s}')\|_2^2$$

Algorithm 2: Budgeted Fitted-Q Iteration

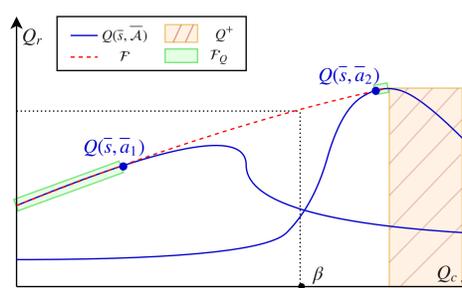
Data: \mathcal{D}
Result: Q^*
1 $Q_{\theta_0} \leftarrow 0$
2 repeat
3 $\theta_{k+1} \leftarrow \arg \min_{\theta} \mathcal{L}(Q_{\theta}, \hat{\mathcal{T}}Q_{\theta_k}; \mathcal{D})$
4 until convergence

4. How to collect the batch \mathcal{D} ?

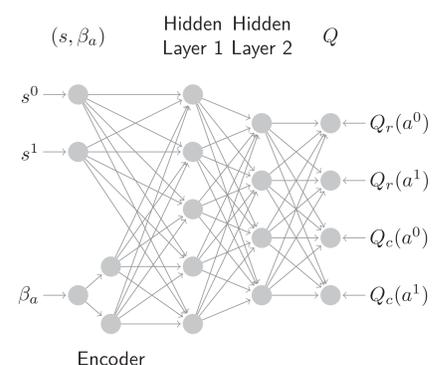
- We propose a risk-sensitive exploration procedure

Scalable Implementation

How to compute the greedy policy?



Function approximation



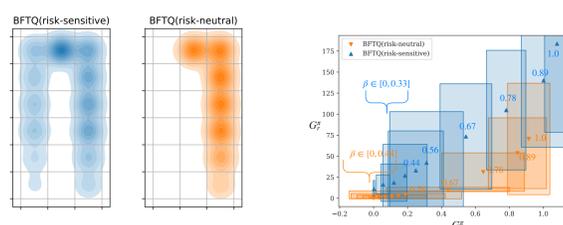
Proposition (Hull policy). π_{greedy} in (2) can be computed explicitly, as a mixture of two points that lie on the convex hull of Q .

Parallel computing

Experience collection and computation of π_{greedy} can be distributed over several cores.

Experiments

Risk-sensitive exploration



Pareto frontier

