

Inria

Budgeted Reinforcement Learning in Continuous State Space

Nicolas Carrara¹, Edouard Leurent^{1,2},
Tanguy Urvoy³, Romain Laroche⁴,
Odalric Maillard¹, Olivier Pietquin^{1,5}

¹Inria Sequel, ²Renault Group,
³Orange Labs, ⁴Microsoft Montréal,
⁵Google Research, Brain Team

Contents

- 01.. Motivation and Setting
- 02.. Budgeted Dynamic Programming
- 03.. Budgeted Reinforcement Learning
- 04.. Experiments

01

Motivation and Setting

Optimal Decision-Making

Which action a_t should we choose in state s_t to maximise a cumulative reward R ?

$$\max_{\pi} \mathbb{E}_{a_t \sim \pi(a_t | s_t)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

Optimal Decision-Making

Which action a_t should we choose in state s_t to maximise a cumulative reward R ?

$$\max_{\pi} \mathbb{E}_{a_t \sim \pi(a_t | s_t)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

- ✓ A very general formulation
- ✓ Widely used in the industry

Optimal Decision-Making

Which action a_t should we choose in state s_t to maximise a cumulative reward R ?

$$\max_{\pi} \mathbb{E}_{a_t \sim \pi(a_t | s_t)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

- ✓ A very general formulation
- ✗ **Not** widely used in the industry

Optimal Decision-Making

Which action a_t should we choose in state s_t to maximise a cumulative reward R ?

$$\max_{\pi} \mathbb{E}_{a_t \sim \pi(a_t | s_t)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

- ✓ A very general formulation
- ✗ **Not** widely used in the industry
 - > Sample efficiency
 - > Trial and error
 - > Unpredictable behaviour

Reinforcement learning relies on a **single reward function R**

Reinforcement learning relies on a **single reward function R**

- ✓ A convenient formulation, but;

Reinforcement learning relies on a **single reward function R**

- ✓ A convenient formulation, but;
- ✗ R is not always easy to design.

Reinforcement learning relies on a **single reward function R**

- ✓ A convenient formulation, but;
- ✗ R is not always easy to design.

Conflicting Objectives

Complex tasks require multiple **contradictory** aspects. Typically:

Task completion vs **Safety**

Reinforcement learning relies on a **single reward function R**

- ✓ A convenient formulation, but;
- ✗ R is not always easy to design.

Conflicting Objectives

Complex tasks require multiple **contradictory** aspects. Typically:

Task completion vs **Safety**

For example...

Dialogue systems

A slot-filling problem: the agent fills a form by asking the user each slot. It can either:

- ask to answer using **voice** (safe/slow);
- ask to answer with a **numeric pad** (unsafe/fast).

Dialogue systems

A slot-filling problem: the agent fills a form by asking the user each slot. It can either:

- ask to answer using **voice** (safe/slow);
- ask to answer with a **numeric pad** (unsafe/fast).

Autonomous Driving

The agent is driving on a two-way road with a car in front of it,

- it can **stay behind** (safe/slow);
- it can **overtake** (unsafe/fast).



Reinforcement learning relies on a single reward function R

- ✓ A convenient formulation, but;
- ✗ R is not always easy to design.

Conflicting Objectives

Complex tasks require multiple contradictory aspects. Typically:

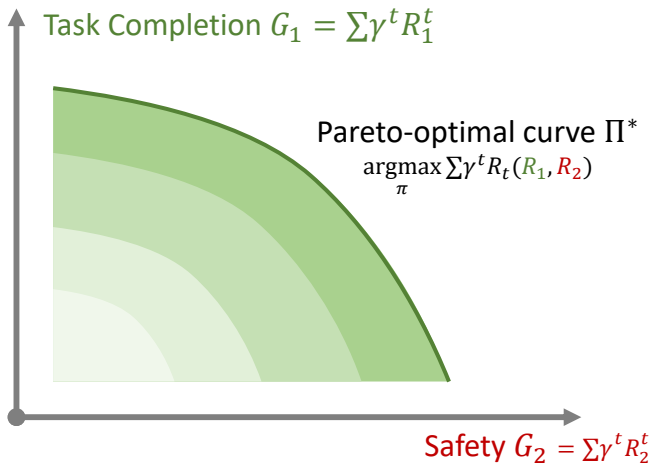
Task completion vs Safety

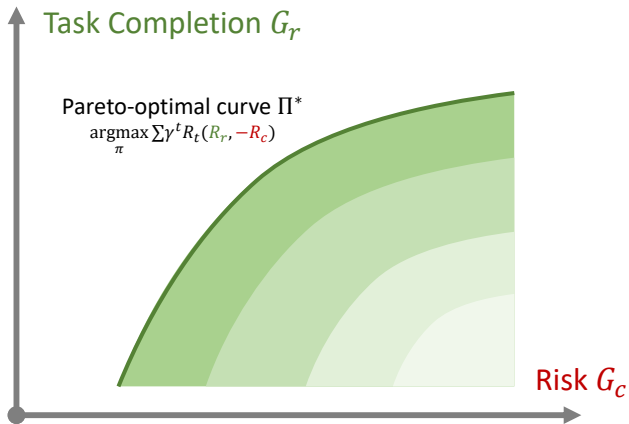
For example...

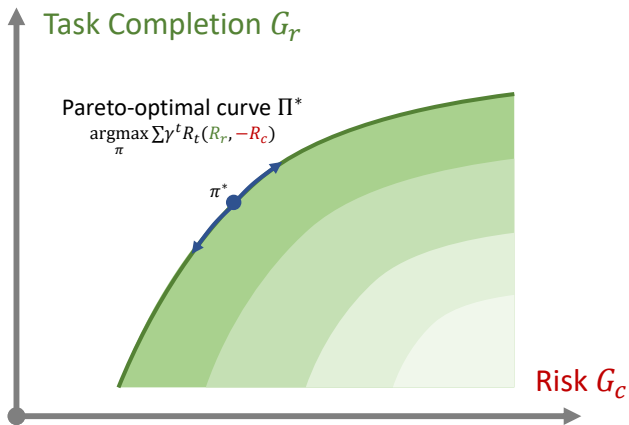
For a fixed reward function R ,

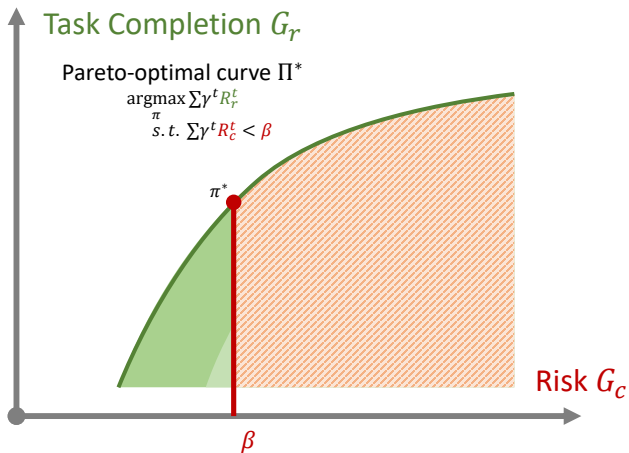
↳ no control over the $\frac{\text{Task Completion}}{\text{Safety}}$ trade-off

↳ π^* is only guaranteed to lie on a Pareto-optimal curve Π^*









Markov Decision Process

An MDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, \gamma)$ with:

- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$

Objective

Maximise rewards

$$\max_{\pi \in \mathcal{M}(\mathcal{A})^{\mathcal{S}}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s \right]$$

Constrained Markov Decision Process

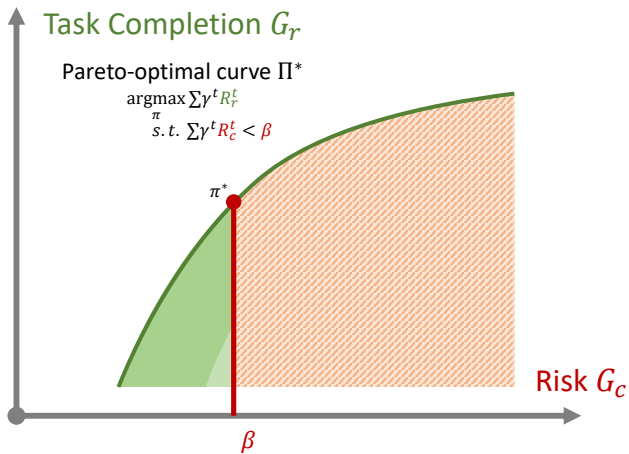
A CMDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, R_c, \gamma, \beta)$ with:

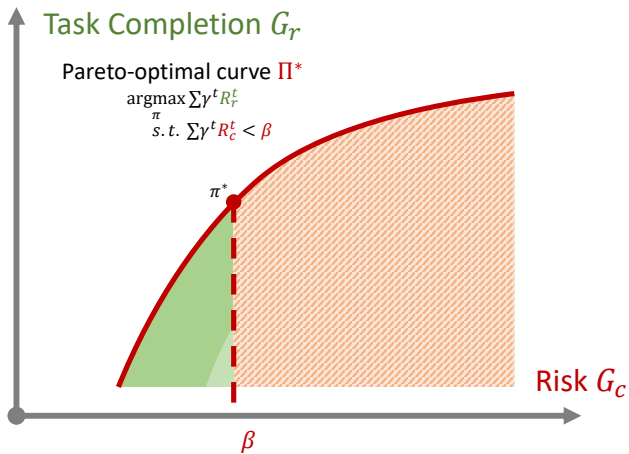
- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Costs $R_c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Budget β

Objective

Maximise rewards while keeping costs under a fixed budget

$$\begin{aligned} \max_{\pi \in \mathcal{M}(\mathcal{A})^{\mathcal{S}}} \quad & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s \right] \\ \text{s.t.} \quad & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_c(s_t, a_t) \mid s_0 = s \right] \leq \beta \end{aligned}$$





Budgeted Markov Decision Process

A BMDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, R_c, \gamma, \mathcal{B})$ with:

- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Costs $R_c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Budget space \mathcal{B}

Objective

Maximise rewards while keeping costs under an **adjustable** budget.

$\forall \beta \in \mathcal{B},$

$$\begin{aligned} \max_{\pi \in \mathcal{M}(\mathcal{A} \times \mathcal{B})^{\mathcal{S} \times \mathcal{B}}} & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s, \beta_0 = \beta \right] \\ \text{s.t.} & \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_c(s_t, a_t) \mid s_0 = s, \beta_0 = \beta \right] \leq \beta \end{aligned}$$

Budgeted policies π

- Take a budget β as an additional input
- Output a next budget β'
- $\pi : \underbrace{(s, \beta)}_{\bar{s}} \rightarrow \underbrace{(a, \beta')}_{\bar{a}}$

↳ Augment the spaces with the budget β

Definition (Augmented spaces)

- States $\bar{\mathcal{S}} = \mathcal{S} \times \mathcal{B}$.
- Actions $\bar{\mathcal{A}} = \mathcal{A} \times \mathcal{B}$.
- Dynamics \bar{P}

state (s, β) , action $(a, \beta_a) \rightarrow$ next state $\begin{cases} s' \sim P(s'|s, a) \\ \beta' = \beta_a \end{cases}$

Definition (Augmented signals)

1. Rewards $R = (R_r, R_c)$
2. Returns $G^\pi = (G_r^\pi, G_c^\pi) \stackrel{\text{def}}{=} \sum_{t=0}^{\infty} \gamma^t R(\bar{s}_t, \bar{a}_t)$
3. Value $V^\pi(\bar{s}) = (V_r^\pi, V_c^\pi) \stackrel{\text{def}}{=} \mathbb{E}[G^\pi \mid \bar{s}_0 = \bar{s}]$
4. Q-Value $Q^\pi(\bar{s}, \bar{a}) = (Q_r^\pi, Q_c^\pi) \stackrel{\text{def}}{=} \mathbb{E}[G^\pi \mid \bar{s}_0 = \bar{s}, \bar{a}_0 = \bar{a}]$

02

Budgeted Dynamic Programming

Proposition (Budgeted Bellman Expectation)

The Bellman Expectation equations are *preserved*

$$V^\pi(\bar{s}) = \sum_{\bar{a} \in \bar{\mathcal{A}}} \pi(\bar{a} | \bar{s}) Q^\pi(\bar{s}, \bar{a})$$

$$Q^\pi(\bar{s}, \bar{a}) = R(\bar{s}, \bar{a}) + \gamma \sum_{\bar{s}' \in \bar{\mathcal{S}}} \bar{P}(\bar{s}' | \bar{s}, \bar{a}) V^\pi(\bar{s}')$$

Proposition (Budgeted Bellman Expectation)

The Bellman Expectation equations are *preserved*

$$V^\pi(\bar{s}) = \sum_{\bar{a} \in \bar{\mathcal{A}}} \pi(\bar{a} | \bar{s}) Q^\pi(\bar{s}, \bar{a})$$

$$Q^\pi(\bar{s}, \bar{a}) = R(\bar{s}, \bar{a}) + \gamma \sum_{\bar{s}' \in \bar{\mathcal{S}}} \bar{P}(\bar{s}' | \bar{s}, \bar{a}) V^\pi(\bar{s}')$$

Proposition (Contraction)

The Bellman Expectation Operator \mathcal{T}^π is a γ -contraction.

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

✓ We can **evaluate** a budgeted policy π

Definition (Budgeted Optimality)

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

Theorem (Budgeted Bellman Optimality Equation)

Q^* verifies the following equation:

$$Q^*(\bar{s}, \bar{a}) = \mathcal{T}Q^*(\bar{s}, \bar{a})$$

$$\stackrel{\text{def}}{=} R(\bar{s}, \bar{a}) + \gamma \sum_{\bar{s}' \in \bar{\mathcal{S}}} \bar{P}(\bar{s}' | \bar{s}, \bar{a}) \sum_{\bar{a}' \in \bar{\mathcal{A}}} \pi_{\text{greedy}}(\bar{a}' | \bar{s}'; Q^*) Q^*(\bar{s}', \bar{a}')$$

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

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

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

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

Proposition (Optimality of the policy)

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

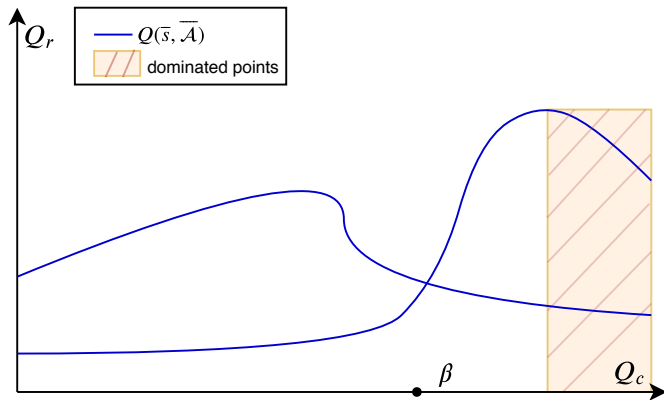
$$\pi_{\text{greedy}}(\cdot ; Q^*) \in \Pi^*(\bar{s})$$

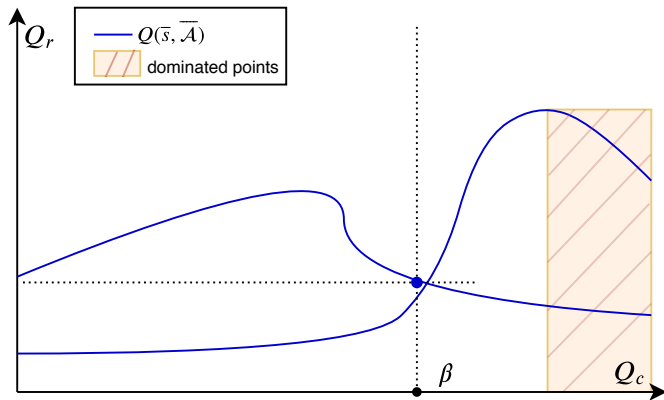
In particular, $V^{\pi_{\text{greedy}}(\cdot ; Q^*)} = V^*$ and $Q^{\pi_{\text{greedy}}(\cdot ; Q^*)} = Q^*$.

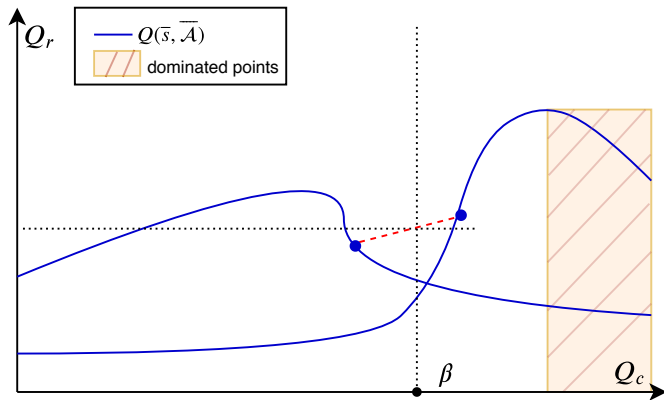
Proposition (Solving the non-linear program)

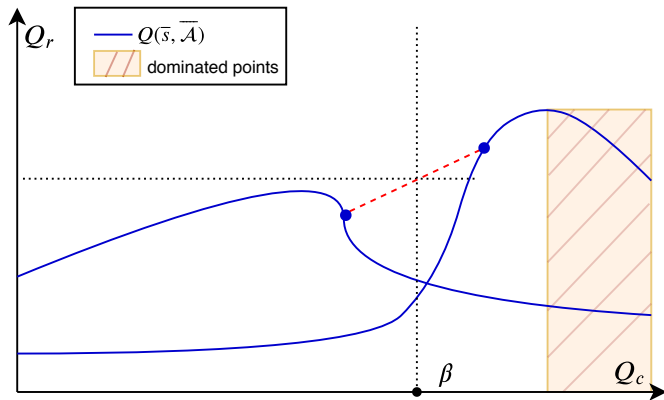
π_{greedy} can be computed *efficiently*, as a mixture π_{hull} of two points that lie on the convex hull of Q .

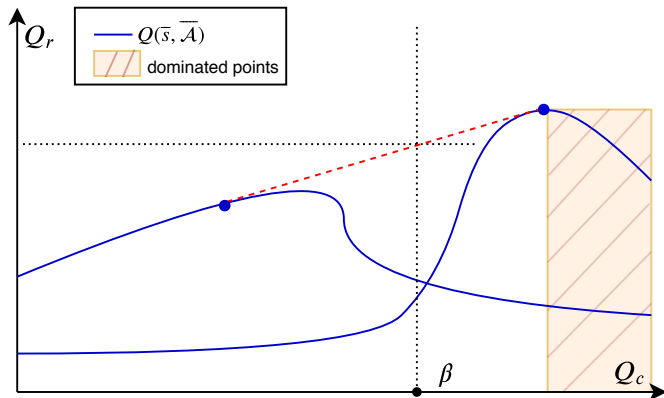
$$\pi_{\text{greedy}} = \pi_{\text{hull}}$$











Recall what we've shown so far:

$$\mathcal{T} \xrightarrow{\text{fixed-point}} Q^* \xrightarrow{\text{tractable}} \pi_{\text{hull}}(Q^*) \xrightarrow{\text{equal}} \pi_{\text{greedy}}(Q^*) \xrightarrow{\text{optimal}}$$

Recall what we've shown so far:

$$\mathcal{T} \xrightarrow{\text{fixed-point}} Q^* \xrightarrow{\text{tractable}} \pi_{\text{hull}}(Q^*) \xrightarrow{\text{equal}} \pi_{\text{greedy}}(Q^*) \xrightarrow{\text{optimal}}$$

We're **almost there!**

All that is left is to perform **Fixed-Point Iteration** to compute Q^* .

Recall what we've shown so far:

$$\mathcal{T} \xrightarrow{\text{fixed-point}} Q^* \xrightarrow{\text{tractable}} \pi_{\text{hull}}(Q^*) \xrightarrow{\text{equal}} \pi_{\text{greedy}}(Q^*) \xrightarrow{\text{optimal}}$$

We're **almost there!**

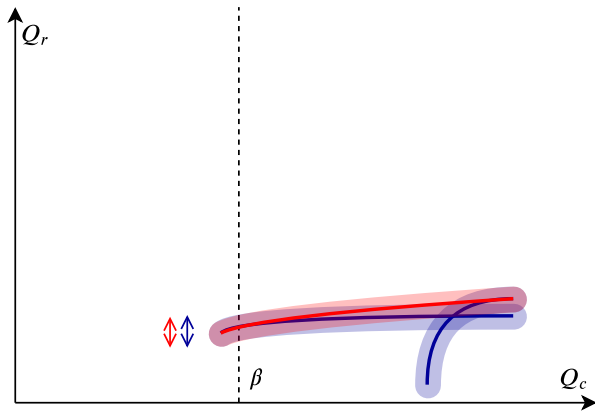
All that is left is to perform **Fixed-Point Iteration** to compute Q^* .

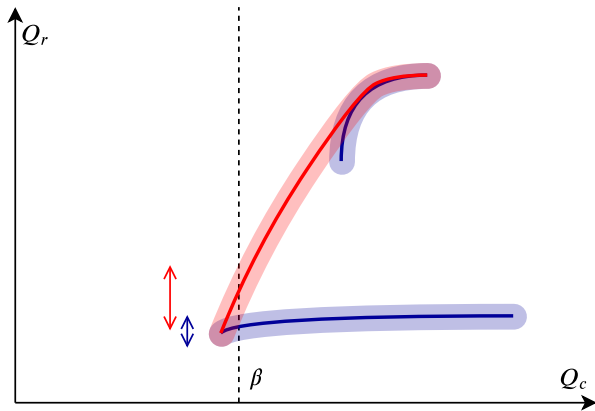
Theorem (Non-Contractivity)

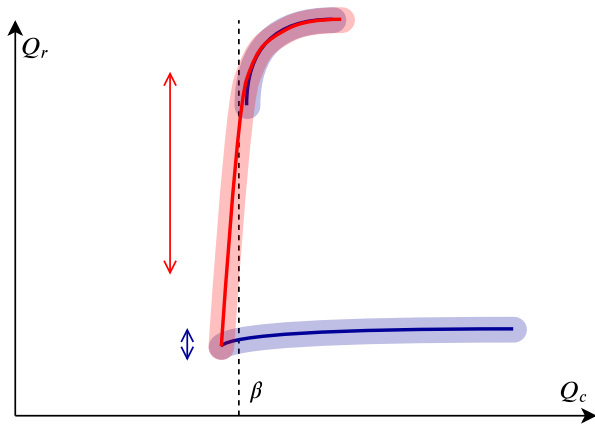
For any BMDP $(S, \mathcal{A}, P, R_r, R_c, \gamma)$ with $|\mathcal{A}| \geq 2$, \mathcal{T} is **not** a contraction.

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

X We **cannot guarantee** the convergence of $\mathcal{T}^n(Q_0)$ to Q^*







Thankfully,

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 $L < \frac{1}{\gamma} - 1$.

$$\mathcal{L}_\gamma = \left\{ Q \in (\mathbb{R}^2)^{\overline{\mathcal{S}\mathcal{A}}} \text{ s.t. } \exists L < \frac{1}{\gamma} - 1 : \forall \bar{s} \in \overline{\mathcal{S}}, \bar{a}_1, \bar{a}_2 \in \overline{\mathcal{A}}, \right. \\ \left. |Q_r(\bar{s}, \bar{a}_1) - Q_r(\bar{s}, \bar{a}_2)| \leq L |Q_c(\bar{s}, \bar{a}_1) - Q_c(\bar{s}, \bar{a}_2)| \right\}$$

- ✓ We guarantee convergence under some (strong) assumptions
- ✓ We observe empirical convergence

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*
-

03

Budgeted Reinforcement Learning

We address several limitations of Budgeted Value-Iteration

1. If the P , R_r and R_c are unknown:

We address several limitations of Budgeted Value-Iteration

1. If the P , R_r and R_c are unknown:

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

We address several limitations of Budgeted Value-Iteration

1. If the P , R_r and R_c are unknown:

- > Work with a batch of samples $\mathcal{D} = \{(\bar{s}_i, \bar{a}_i, r_i, \bar{s}'_i)\}_{i \in [0, M]}$
- > Replace \mathcal{T} with a sampling operator $\hat{\mathcal{T}}$:

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

We address several limitations of Budgeted Value-Iteration

1. If the P , R_r and R_c are unknown:

- > Work with a batch of samples $\mathcal{D} = \{(\bar{s}_i, \bar{a}_i, r_i, \bar{s}'_i)\}_{i \in [0, M]}$
- > Replace \mathcal{T} with a sampling operator $\hat{\mathcal{T}}$:

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

2. If \mathcal{S} is continuous:

We address several limitations of Budgeted Value-Iteration

1. If the P , R_r and R_c are unknown:

- > Work with a batch of samples $\mathcal{D} = \{(\bar{s}_i, \bar{a}_i, r_i, \bar{s}'_i)\}_{i \in [0, M]}$
- > Replace \mathcal{T} with a sampling operator $\hat{\mathcal{T}}$:

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

2. If \mathcal{S} is continuous:

- > Employ function approximation Q_θ , and minimise a regression loss

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

- CPU parallel computing of the targets

$$\sum_{\bar{a}'_i \in \mathcal{A}_i} \pi_{\text{greedy}}(\bar{a}'_i | \bar{s}'_i; Q) Q(\bar{s}'_i, \bar{a}'_i)$$

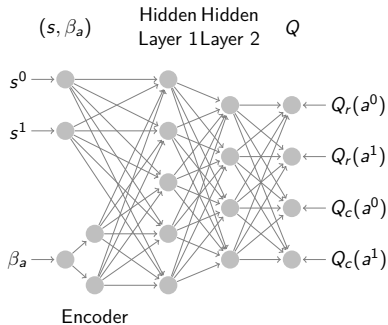
- CPU parallel computing of the targets

$$\sum_{\bar{a}'_i \in \mathcal{A}_i} \pi_{\text{greedy}}(\bar{a}'_i | \bar{s}'_i; Q) Q(\bar{s}'_i, \bar{a}'_i)$$

- Same for interactions with the environment.

- CPU parallel computing of the targets

$$\sum_{\bar{a}'_i \in \mathcal{A}_i} \pi_{\text{greedy}}(\bar{a}'_i | \bar{s}'_i; Q) Q(\bar{s}'_i, \bar{a}'_i)$$
- Same for interactions with the environment.
- Neural Network for function approximation:



04

Experiments

Lagrangian Relaxation

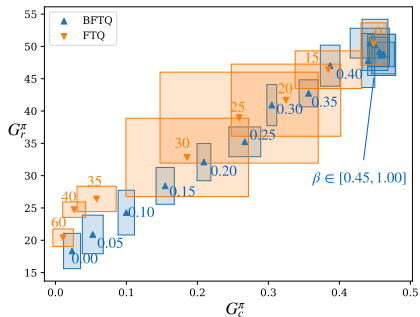
Consider the dual problem so as to replace the hard constraint by a soft constraint penalised by a **Lagrangian multiplier** λ :

$$\max_{\pi} \mathbb{E} \sum_t \gamma^t R_r(s, a) - \lambda \gamma^t R_c(s, a)$$

- Train many policies π_k with penalties λ_k and recover the cost budgets β_k
- Very data/memory-heavy

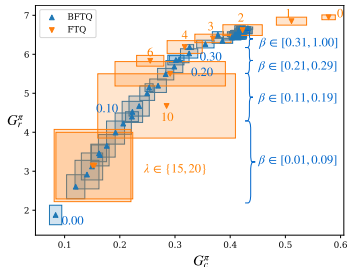
A slot-filling problem: the agent (the dialogue system) fills a form by asking the user each slot. It can either:

- ask to answer using **voice** (safe/slow);
- ask to answer with a **numeric pad** (unsafe/fast).



The agent (the car) is on a two-way road with a car in front of it,

- it can stay behind (safe/slow);
- it can overtake (unsafe/fast).



How to collect the batch \mathcal{D} ?

We propose an ε -greedy exploration procedure:

How to collect the batch \mathcal{D} ?

We propose an ε -greedy exploration procedure:

- Sample an initial budget $\beta_0 \sim \mathcal{U}(\mathcal{B})$

How to collect the batch \mathcal{D} ?

We propose an ε -greedy exploration procedure:

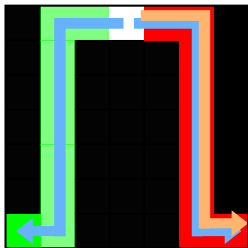
- Sample an initial budget $\beta_0 \sim \mathcal{U}(\mathcal{B})$
- At each step, where $\bar{s} = (s, \beta)$ only explore feasible budgets:

$$\bar{a} = (a, \beta_a) \sim \mathcal{U}(\Delta_{\mathcal{A}\mathcal{B}})$$

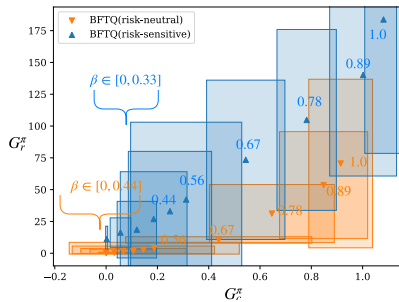
where Δ is such that $\mathbb{P}(a, \beta_a | s, \beta)$ verifies $\mathbb{E}[\beta_a] \leq \beta$

Two corridors:

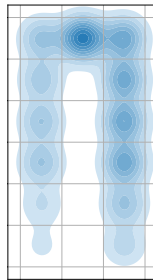
1. one with **high costs / high rewards**
2. the other with **no costs / low rewards**



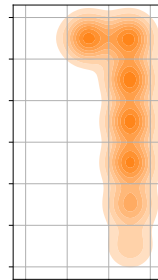
→ Validate the **risk-sensitive exploration** procedure



BFTQ(risk-sensitive)



BFTQ(risk-neutral)



Thank You!