

A THEORETICAL PERSPECTIVE ON SEQUENTIAL DECISION MAKING WITH PREFERENCE FEEDBACK





SIMONE DRAGO, MARCO MUSSI, AND ALBERTO MARIA METELLI

{simone.drago, marco.mussi, albertomaria.metelli}@polimi.it

SETTING

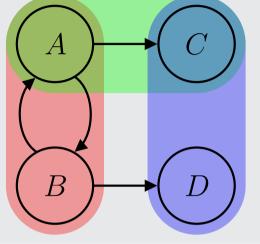
PREFERENCE

UTILITY

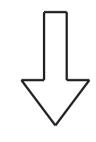
REWARD

PrefMDP = $(\underbrace{\mathcal{S}, \mathcal{A}, H, p, \mu}, \leq_{\mathcal{T}})$

 $\leq_{\mathcal{T}} \subseteq \mathcal{T} \times \mathcal{T}$ is a **partial (pre)order** over the trajectory space \mathcal{T}



- (<) C is preferred over A
- (\approx) A is equivalent to B
- (\parallel) C is incomparable to D



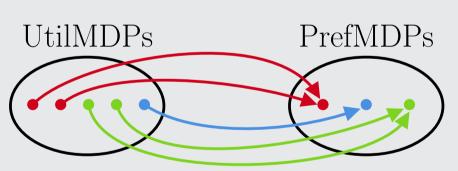
Assumption: Human expresses preferences based on an underlying (unknown) utility

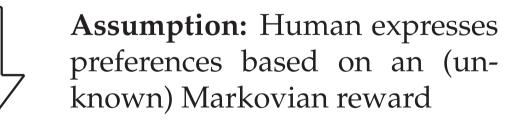
$$\text{UtilMDP} = (\underbrace{\mathcal{S}, \mathcal{A}, H, p, \mu}_{\text{MDP} \setminus R}, \boldsymbol{u})$$

 $oldsymbol{u}:\mathcal{T} o\mathbb{R}^m$ is a multi-dimensional utility

Expected utility of a policy $\pi \in \Pi$: $\boldsymbol{J}(\pi; \boldsymbol{u}) \coloneqq \sum_{\tau \in \mathcal{T}} d_{\pi}(\tau) \boldsymbol{u}(\tau) = \langle d_{\pi}, \boldsymbol{u} \rangle,$

where d_{π} is the distribution over trajectories induced by π

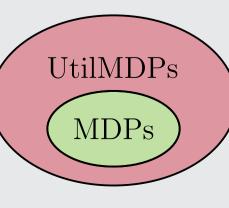




$$MDP = (\underbrace{\mathcal{S}, \mathcal{A}, H, p, \mu}_{MDP \setminus R}, \mathbf{r})$$

 $r := (r_h)_{h \in \llbracket H \rrbracket}$ is a stage-dependent multi-dimensional reward function

Trajectory return: $u_{\boldsymbol{r}}(\tau) \coloneqq \sum_{h=1}^{H} r_h(s_h, a_h)$ Expected policy return: $J(\pi; \boldsymbol{r}) \coloneqq J(\pi; \boldsymbol{u_r})$

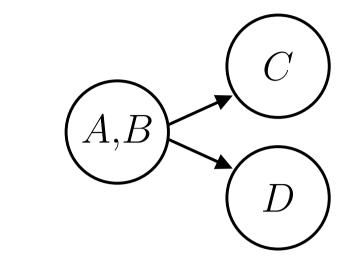


COMPATIBLE UTILITIES

COMPATIBLE UTILITY

u is compatible with $\leq_{\mathcal{T}}$ if $\forall \tau, \tau' \in \mathcal{T}$: $\tau \leq_{\mathcal{T}} \tau' \Rightarrow u(\tau) \leq u(\tau')$ (element-wise)

| Realizer Dimension | Exists? | Computational Complexity |
|--|---------|--------------------------------------|
| $<\dim(\leq_{\mathcal{T}})$ | X | |
| $=\dim(\leq_{\mathcal{T}})$ | | NP-hard |
| $> \dim(\leq_{\mathcal{T}})$ | | $\operatorname{Poly}(\mathcal{T})$ |
| where $\dim(\cdot)$ is the order dimension | | |



- (i) Represent the quotient set w.r.t. equivalence of $\leq_{\mathcal{T}}$ as a DAG
- (A,B)
- $(A,B) \longrightarrow (C)$
- (ii) Solve a minimum path cover problem to obtain a set of width($\leq_{\mathcal{T}}$) chains
- (iii) Extend each chain to obtain a linear extension of $\leq_{\mathcal{T}}$, accounting for incomparabilities

Procedure to derive a realizer of size width($\leq_{\mathcal{T}}$) $\geqslant \dim(\leq_{\mathcal{T}})$ in $\mathcal{O}(|\mathcal{T}|^3)$

POLICY DOMINANCE

POLICY DOMINANCE

Policy $\pi \leq_{\mathcal{T}}$ -strictly dominates policy π' ($\pi' <_{\Pi} \pi$) if it yields a strictly better expected utility

 $\boldsymbol{J}(\pi; \boldsymbol{u}) - \boldsymbol{J}(\pi'; \boldsymbol{u}) > \boldsymbol{0}$ (element-wise)

for every compatible utility function $oldsymbol{u}$

$\leq_{\mathcal{T}}$ -Pareto Optimality

Set of $\leq_{\mathcal{T}}$ -Pareto optimal policies:

 $\Pi^*(\leq_{\mathcal{T}}) := \{ \pi \in \Pi : \nexists \pi' \in \Pi \text{ s.t. } \pi \prec_{\Pi} \pi' \}$

How to Evaluate?

 $\pi' \leq_{\Pi} \pi$ can be verified by evaluating if:

$$\forall n \in [\![|\mathcal{T}|]\!] : \sum_{i=1}^{n} (d_{\pi}(i) - d_{\pi'}(i)) \ge 0$$

holds for every linear extension of $\leq_{\mathcal{T}}$, where index i represents the i-th trajectory sorted w.r.t. the linear extension

Tractability

OPEN QUESTION

Is there an efficient evaluation method?

Trivial solution requires the evaluation of $\mathcal{O}(|\mathcal{T}|!)$ linear extensions of $\leq_{\mathcal{T}}$

APPROXIMATION VIA MARKOVIAN REWARDS

WHY THE NEED FOR APPROXIMATION?

Preferences

Most general type of feedback

Intractable without introducing a structure

Utilities

- Assign numerical signals to each trajectory
- Complexity of learning is exponential in H

Rewards

- Enable efficient learning
- Less representational power

(CONVEX) QUADRATIC PROGRAM

Goal: Find u and r that best represent $\leq_{\mathcal{T}}$

Input: A realizer $\{\leqslant_{\mathcal{T},i}\}_{i\in[m]}$ of $\leq_{\mathcal{T}}$ of size m

Idea: Jointly choose u compatible with the realizer and r as Markovian approximation of u

Define: $B \in \{0,1\}^{|\mathcal{T}| \times |\mathcal{S}||\mathcal{A}|H}$ (binary encoding of \mathcal{T}) and $A := I_{|\mathcal{T}|} - B(B^{\top}B)^{-1}B^{\top}$ (OLS)

$$\eta^* := \min \|\mathbf{A}\mathbf{u}\|_F^2$$
s.t. $u_i(i+1) \leq u_i(i) - \varepsilon \ \forall i \in [\![|\mathcal{T}|-1]\!], \ j \in [\![m]\!]$

If $\eta^* = 0 \rightarrow u_r = u \rightarrow$ Preferences derive from a Markovian reward function

If $\eta^* > 0 \rightarrow$ Preferences cannot be expressed via Markovian rewards

- ightarrow Using r, we learn how to solve a simpler surrogate problem
- \rightarrow Such an approximation introduces a suboptimality in terms of the performance of the optimal policy bounded by $2\sqrt{m\eta^*}$