

Regret based Solutions for Uncertain MDPs

Anonymous Author(s)

Affiliation

Address

email

Proposition 1. For a policy $\bar{\pi}^0 : 0 \leq \text{reg}(\bar{\pi}^0) - \text{creg}(\bar{\pi}^0) \leq \left[\max_s R^*(s) - \min_s R^*(s) \right] \cdot \frac{(1-\gamma^H)}{1-\gamma}$

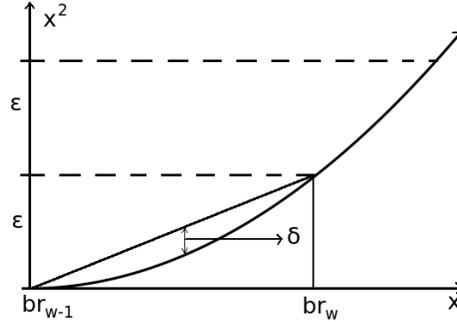
Proof. We can rewrite Equation (1) as follows:

$$\text{creg}(\bar{\pi}^0) = v^{0,\#}(\bar{\pi}^0) - v^0(\bar{\pi}^0), \text{ where } v^{0,\#}(\bar{\pi}^0) = \sum_s \alpha(s) v^{0,\#}(s, \bar{\pi}^0) \text{ and}$$

$$v^{t,\#}(s, \bar{\pi}^t) = \sum_a \pi^t(s, a) \cdot \left[R^*(s) + \gamma \sum_{s'} T(s, a, s') \cdot v^{t+1,\#}(s', \bar{\pi}^{t+1}) \right]$$

Since the value for any policy cannot exceed the value of optimal policy, we have:
 $\text{reg}(\bar{\pi}^0) - \text{creg}(\bar{\pi}^0) = v^0(\bar{\pi}^*) - v^{0,\#}(\bar{\pi}^0) \geq 0$. The difference in value of optimal policy (a deterministic one) and any other policy is because of the states visited by using the policy. In the worst case for creg , the optimal policy visits the state with highest R^* and $\bar{\pi}^0$ visits the states with the lowest R^* at every time step. Sum of a geometric progression over the time steps yields

$$\text{reg}(\bar{\pi}^0) - \text{creg}(\bar{\pi}^0) \leq \left[\max_s R^*(s) - \min_s R^*(s) \right] \cdot \frac{(1-\gamma^H)}{1-\gamma}$$



(a)

Figure 1: Error

The proposition below provides the proof for footnote 2.

Footnote 3. In the approximation of x^2 function using piecewise linear components, $\lambda(w)$, the maximum approximation in any interval $[br_{w-1}, br_w]$ occurs at the mid-point.

Proof. Without loss of generality, let us consider any point y in the interval $[br_{w-1}, br_w]$. From Equation 7, we have

$$y = \lambda_{w-1} br_{w-1} + \lambda_w br_w$$

Since, we have the sum constraint in Equation 9, the above equation can be modified as:

$$y = (1 - \lambda_w) br_{w-1} + \lambda_w br_w$$

$$\implies \lambda_w = \frac{y - br_{w-1}}{br_w - br_{w-1}}$$

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

The error is given by the difference between LHS and RHS in Equation 8:

$$\delta = y^2 - [\lambda_{w-1}(br_{w-1})^2 + \lambda_w(br_w)^2]$$

Substituting value of λ_w :

$$\delta = y^2 - \left[(br_w)^2 \cdot \frac{y - br_{w-1}}{br_w - br_{w-1}} - (br_{w-1})^2 \cdot \frac{y - br_w}{br_w - br_{w-1}} \right]$$

When δ is maximum, we have $\frac{d\delta}{dy} = 0$. Therefore:

$$2y - \frac{((br_w)^2 - (br_{w-1})^2)}{(br_w - br_{w-1})} = 0$$

$$y = \frac{br_w + br_{w-1}}{2}$$

Hence proved. ■

Proposition 3. Let $\hat{v}_{\xi_q}^t(s, \bar{\pi}^t)$ denote the approximation of $v_{\xi_q}^t(s, \bar{\pi}^t)$. Then

$$v_{\xi_q}^t(s, \bar{\pi}^t) - \frac{|\mathcal{A}| \cdot \epsilon \cdot (1 - \gamma^{H-1})}{4 \cdot (1 - \gamma)} \leq \hat{v}_{\xi_q}^t(s, \bar{\pi}^t) \leq v_{\xi_q}^t(s, \bar{\pi}^t) + \frac{|\mathcal{A}| \cdot \epsilon \cdot (1 - \gamma^{H-1})}{4 \cdot (1 - \gamma)}$$

Proof: At time step $t + 1$, the approximation error in $v_{\xi_q}^{t+1}(s, \bar{\pi}^{t+1})$ is given by $|\mathcal{A}| \cdot \delta$, ($|\mathcal{A}|$ as maximum number of actions across all states, all time steps). The maximum approximation error at time step t in $v_{\xi_q}^t(s, a, \bar{\pi}^t)$ is $\gamma \cdot |\mathcal{A}| \cdot \delta$ (due to error in value function at time step $t + 1$). We can combine equation 3 and 4 as:

$$v_{\xi_q}^t(s, \bar{\pi}^t) = \sum_a \pi^t(s, a) \cdot \left[v_{\xi_q}^t(s, a, \bar{\pi}^t) \pm \gamma \cdot |\mathcal{A}| \cdot \delta \right]$$

$$= \sum_a \pi^t(s, a) \cdot v_{\xi_q}^t(s, a, \bar{\pi}^t) \pm \gamma \cdot |\mathcal{A}| \cdot \delta$$

Now at time step t the error will be $|\mathcal{A}| \cdot \delta$ plus future error from time step $t + 1$ given by $\gamma \cdot |\mathcal{A}| \cdot \delta$. Extending to $t = 0$ we will have sum of two geometric progressions, i.e.

$$\pm \left[|\mathcal{A}| \cdot \delta + \gamma \cdot |\mathcal{A}| \cdot \delta + \gamma^2 \cdot |\mathcal{A}| \cdot \delta \dots \right]$$

Substituting $\delta = \frac{\epsilon}{4}$, we will have a positive and negative error of $\frac{|\mathcal{A}| \cdot \epsilon \cdot (1 - \gamma^{H-1})}{4 \cdot (1 - \gamma)}$. ■

Proposition 4. At time step $t - 1$, the CER corresponding to any policy, π^{t-1} will have least regret if it includes CER minimizing policy from t . Formally, if $\bar{\pi}^{*,t}$ represents the CER minimizing policy from t and $\bar{\pi}^t$ represents any arbitrary policy, then:

$$\forall s : \max_{\bar{\xi}_p^{t-1} \in \bar{\xi}^{t-1}} \text{creg}_{\bar{\xi}_p^{t-1}}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^{*,t} \rangle) \leq \max_{\bar{\xi}_p^{t-1} \in \bar{\xi}^{t-1}} \text{creg}_{\bar{\xi}_p^{t-1}}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^t \rangle)$$

$$\text{if, } \forall s : \max_{\bar{\xi}_q^t \in \bar{\xi}^t} \text{creg}_{\bar{\xi}_q^t}^t(s, \bar{\pi}^{*,t}) \leq \max_{\bar{\xi}_q^t \in \bar{\xi}^t} \text{creg}_{\bar{\xi}_q^t}^t(s, \bar{\pi}^t)$$

Proof. From Equation 12, we have:

$$\text{creg}_{\bar{\xi}_p^{t-1}}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^{*,t} \rangle) = \sum_{a \in \mathcal{A}} \pi^{t-1}(s, a) \left[\Delta \mathcal{R}_p^{t-1}(s, a) + \gamma \sum_{s'} \mathcal{T}_p^{t-1}(s, a, s') \cdot \max_{\bar{\xi}_q^t \in \bar{\xi}^t} \text{creg}_{\bar{\xi}_q^t}^t(s', \bar{\pi}^{*,t}) \right]$$

From Equation 14, we have:

$$\text{creg}_{\bar{\xi}_p^{t-1}}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^{*,t} \rangle) \leq \sum_{a \in \mathcal{A}} \pi^{t-1}(s, a) \left[\Delta \mathcal{R}_p^{t-1}(s, a) + \gamma \sum_{s'} \mathcal{T}_p^{t-1}(s, a, s') \cdot \max_{\bar{\xi}_q^t \in \bar{\xi}^t} \text{creg}_{\bar{\xi}_q^t}^t(s', \bar{\pi}^t) \right]$$

$$\leq \text{creg}_{\bar{\xi}_p^{t-1}}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^t \rangle)$$

Thus, $\max_{\bar{\xi}_q \in \bar{\xi}} \text{creg}_{\bar{\xi}_q}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^{*,t} \rangle) \leq \max_{\bar{\xi}_q \in \bar{\xi}} \text{creg}_{\bar{\xi}_q}^{t-1}(s, \langle \pi^{t-1}, \bar{\pi}^t \rangle)$. ■

Pruning dominated actions

Algorithm 1 provides the pseudo-code for pruning step discussed earlier. At each time step, for each state we maintain an upper and lower bound for the value function. Apart from pruning, this gives us tight bounds on value function that decrease the number of break points required for linearization.

Algorithm 1: PRUNEDOMINATEDACTIONS()

```

114  $t \leftarrow H - 1$ 
115 for all  $\xi_q \in \xi, s \in \mathcal{S}$  do
116    $v_{\xi_q}^{H,min}(s) \leftarrow 0$ 
117    $v_{\xi_q}^{H,max}(s) \leftarrow 0$ 
118 while  $t \geq 0$  do
119   for all  $s \in \mathcal{S}$  do
120     for all  $\xi_q \in \xi, a \in \mathcal{A}$  do
121        $v_{\xi_q}^{t,min}(s, a) \leftarrow R_q^t(s, a) + \gamma \sum_{s'} \mathcal{T}_q^t(s, a, s') \cdot v_{\xi_q}^{t+1,min}(s')$ 
122        $v_{\xi_q}^{t,max}(s, a) \leftarrow R_q^t(s, a) + \gamma \sum_{s'} \mathcal{T}_q^t(s, a, s') \cdot v_{\xi_q}^{t+1,max}(s')$ 
123       if  $\exists a' s.t. v_{\xi_q}^{t,min}(s, a') \geq v_{\xi_q}^{t,max}(s, a) \forall \xi_q$  then
124         PRUNE  $a$ 
125          $v_{\xi_q}^{t+1,min}(s) = \min_a v_{\xi_q}^{t,min}(s, a)$ 
126          $v_{\xi_q}^{t+1,max}(s) = \max_a v_{\xi_q}^{t,max}(s, a)$ 
127        $t \leftarrow t - 1$ 

```

SAA Analysis

Each sample (scenario) is described by $i = \{i_1, i_2, i_3, \dots, i_{|T|}\}$ and belong to the set I (in the case where we consider independent transition probabilities/rewards in each stage, I is the set of samples which are cross products of independent samples in each stage). Followed from the sample average approximation (SAA) method described by [2], the steps to calculate the approximate optimality gap are as follows:

1. Generate the set of sample sets, $M = \{I_1, I_2, \dots, I_{|M|}\}$, where each sample set is of size $|I|$. Also generate a larger sample set of size $|I'| \gg |I|$.
 - For $m = 1, \dots, |M|$, solve the problem with sample set I_m to obtain the solution value $r\bar{e}g_m^*$ and policy $\bar{\pi}_m$
2. Compute the average of the objective values obtained which is a statistical lower bound of the problem and their corresponding variance as follows:

$$r\hat{e}g^* = \frac{1}{|M|} \sum_{m \in M} r\bar{e}g_m^* \text{ and } \sigma_{r\hat{e}g^*}^2 = \frac{1}{|M|(|M| - 1)} \sum_{m \in M} (r\bar{e}g_m^* - r\hat{e}g^*)^2.$$

3. Let $\bar{\pi}$ be the selected solution from the set of solutions obtained in Step 1. Denote by $reg_{I'}^*(\bar{\pi})$ the regret value of the policy $\bar{\pi}$ on the large sample set I' . This value is the sample average estimate of the true objective function of the policy $\bar{\pi}$. Also, its variance can be computed as follows:

$$\sigma_{I'}^2(\bar{\pi}) = \frac{1}{|I'|(|I'| - 1)} \sum_{i \in I'} (r\bar{e}g_i^*(\bar{\pi}) - reg_{I'}^*(\bar{\pi}))^2$$

where $r\bar{e}g_i^*(\bar{\pi})$ is the regret of the policy $\bar{\pi}$ corresponding to each sample $i \in I'$.

4. The absolute optimality gap of the solution $\bar{\pi}$ and its variance can be estimated as follows:

$$gap(\bar{\pi}) = |reg_{I'}^*(\bar{\pi}) - r\hat{e}g^*| \text{ and } \sigma_{gap}^2(\bar{\pi}) = \sigma_{I'}^2(\bar{\pi}) + \sigma_{r\hat{e}g^*}^2.$$

We can similarly perform SAA analysis for MILP-CER.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

Single Product Stochastic Inventory Control Problem

In the single product finite horizon stochastic inventory control problem [1], at the beginning of each time period and before observing the demand, the manager determines the current inventory size x^t and decides whether or not to order additional stock y^t from a supplier. We assume the cost of ordering u units is given by $k_1 \cdot u$, the cost of maintaining an inventory of u units is given by $k_2 \cdot u$ and the revenue obtained when the demand is j units is given by $k_3 \cdot j$.

Denote $D^t = \{d_0^t, d_1^t, \dots, d_q^t\}$ as the set of demand values at time step t (independent of demand in other time steps). The inventory at time step $t + 1$ for demand d_q^t is given by $x^{t+1}(d_q^t) = \max\{x^t + y^t - d_q^t, 0\} \equiv [x^t + y^t - d_q^t]^+$. Note that the reward at time step t depends on the current and subsequent inventory size and is given by $r^t(x^t, y^t, x^{t+1}(d_q^t)) = -k_1 \cdot y^t - k_2 \cdot (x^t + y^t) + k_3 \cdot ([x^t + y^t - x^{t+1}(d_q^t)]^+)$.

The discrete demand uncertainty values translate to uncertainty over reward and transition functions, which require robust solution concepts. A standard approach is to maximize the minimum expected values or *maximin solution*. In this paper, we compare DP-CER against maximin across different cost-to-revenue ratio defined as $\frac{k_1+k_2}{k_3}$.

References

- [1] Martin Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, 1994.
- [2] A. Shapiro. Monte carlo sampling methods. In A. Ruszczyński and A. Shapiro, editors, *Stochastic Programming*, volume 10 of *Handbooks in Operations Research and Management Science*, pages 353–425. Elsevier, 2003.