# Polytime Computation of Strong- and $n$-Present-Value Optimal Policies in Markov Decision Chains

by

## Michael O'Sullivan

and

## Arthur F. Veinott, Jr.

# Introduction

- **Markov Decision Chain:** $S < \infty$ **states** $s$, **finite actions** $a \in A_s$, **stationary policy** $\delta = (\delta^s)$, **reward** $r_\delta$, **Submarkov matrix** $P_\delta$, $V_\delta^\rho$ **is present value of income with interest rate** $\rho > 0$ **using** $\delta$

- **Optimality Concepts: Strengths and Weaknesses**

  - ♦ **Maximum Present Value (MPV):** $V^\rho \equiv \max_\delta V_\delta^\rho$

    **How small should** $\rho$ **be if the distant future is relevant?**

    **MPV is undefined if** $\rho = 0$.

    **What if the problem does not involve interest?**

  - ♦ **Maximum Reward Rate (MRR):** $\max_\delta \lim_{\rho \downarrow 0} \rho V_\delta^\rho$

    **If** $\delta$ **has MRR, then** $V_\delta^\rho - V^\rho = O(1)$.

    **MRR does not reflect transient effects.**

    **The rewards earned in any finite horizon are irrelevant.**

    **In transient systems, all rewards are irrelevant.**

    **MRR policies are not nearly SPV (see below) optimal.**

  - ♦ **Limiting Present Value [Bl62]:** $V_\delta^\rho - V^\rho = o(1)$

    **LPV is stronger than MRR.**

    **LPV reflects transient effects.**

    **LPV policies are nearly SPV optimal.**

LPV is the right way to define MPV when $\rho = 0$.

In transient systems, LPV is equivalent to MV.

♦ $n$-**Present Value [Ve69]:** $V_\delta^\rho - V^\rho = o(\rho^n)$, $n \quad -1$

$n$-PV polices are MRR if $n = -1$ and LPV if $n = 0$.

$n$-PV policies are nearly SPV for $n \quad 0$.

The quality of $n$-PV policies rises with $n$.

The set of $n$-PV policies diminishes with $n$.

Computation: Find $n$-PV Opt given $(n-1)$-PV Opt

♦ **Strong Present Value [Bl62]:** $V_\delta^\rho - V^\rho = 0$, $\rho^* > \rho > 0$

Such policies are robust under small changes in $\rho$.

The sets of $n$-PV & SPV policies coincide for $n \quad S$.

SPV policies are MRR and SPV for all $n \geq -1$.

But SPV policies can be more expensive to compute.

• **Extensions**

♦ **Immigration [RV92]**

♦ **Perturbation of Interest Rates [Ve69]**

♦ **Long Time Horizons: $n$-Cesàro-overtaking optimality**

♦ **Continuous Time Parameter**

♦ **Markov Population Decision Chains [RV92]**

# Background [MV69, Ve69, Ve74]

- **Laurent Expansion of** $V_\delta^\rho = \sum_{-1}^{\infty} \rho^n v_\delta^n$

- $\delta$ **is** $n$**-Optimal if** $V_\delta^n \succeq V_\gamma^n$ **all** $\gamma$, $V_\gamma^n \equiv \left( v_\gamma^{-1}, \ldots, v_\gamma^n \right)$

- $\Delta_n \equiv$ **Set of (stationary)** $n$**-Optimal Policies**

  - ♦ $\delta \in \Delta_{-1} \Leftrightarrow \delta$ **Maximum Reward Rate (MRR)**

  - ♦ $\delta \in \Delta_0 \quad \Leftrightarrow \delta$ **Present-Value Optimal** $\rho = 0$

- **Selectivity Increases with** $n$: $\Delta_n \supseteq \Delta_{n+1}$

- $S$**-Optimal** $\Leftrightarrow$ **Strong Present-Value Optimal**

- **Comparison** $\qquad g_{\gamma\delta}^n \equiv r_\gamma^n + Q_\gamma v_\delta^n - v_\delta^{n-1}, \, n \quad -1$

$$Q_\gamma \equiv P_\gamma - I, \, r_\gamma^0 \equiv r_\gamma, \, r_\gamma^n \equiv 0 \text{ all } n \neq 0$$

$$r_\gamma + P_\gamma V_\delta^\rho - (1+\rho)V_\delta^\rho = r_\gamma + Q_\gamma V_\delta^\rho - \rho V_\delta^\rho = \sum_{-1}^{\infty} \rho^n g_{\gamma\delta}^n$$

- $r_\delta^n + Q_\delta v^n = v^{n-1}, \, n \leq m+1 \Rightarrow V^m = V_\delta^m$

- $\delta \in \Delta_n \Rightarrow G_{\gamma\delta}^n \equiv (g_{\gamma\delta}^{-1}, \ldots, g_{\gamma\delta}^n) \preceq 0 \text{ all } \gamma$

- $\delta \in \Delta_n \Leftarrow G_{\gamma\delta}^{n+1} \preceq 0 \text{ all } \gamma \qquad n$**-Optimality Condition**

- **Policy-Improvement Method**

- $\mathcal{E}_\delta^n \equiv \{\gamma : G_{\gamma\delta}^n = 0\}, \, n \geq -1$

- $\delta \in \Delta_{n-1} \Rightarrow \mathcal{E}_\delta^n \subseteq \Delta_{n-1} \subseteq \mathcal{E}_\delta^{n-1}$

# Maximum Transient Value (MTV)
## [Ve69], [EV72]

- **State $s$ is Transient or Recurrent under $\delta$ according as $s^{th}$ Column of $P_\delta^*$ is 0 or Not.**

- **$\delta$ is Transient or Recurrent according as $P_\delta^*$ is 0 or Not.**

- **If $\delta$ is Transient, its Value is $V_\delta = \sum\limits_{i=1}^{\infty} P_\delta^{i-1} r_\delta$.**

- **$\delta$ has MTV if $V_\delta \geq V_\gamma$, all Transient $\gamma$.**

- **Characterization of MTV. If $\delta$ Transient, $\delta$ has MTV $\Leftrightarrow$ $V = V_\delta$ is Least Fixed (resp., Excessive) Point of $\mathcal{R}$,**

$$\mathcal{R}V \equiv \max_\gamma \left(r_\gamma + P_\gamma V\right), \text{ all } V \in \Re^S.$$

- **Dual Linear Program to Find MTV. Choose $V$ to Minimize $1V$ Subject to $V \geq r_\gamma + P_\gamma V$, all $\gamma$. Each Optimal Basis Corresponds to MTV $\delta$.**

# Maximum Reward-Rate (MRR)

## [Ho60], [Ba61], [Bl62], [DF68], [HK78]

- **Policy Improvement for $n = -1$**

- **Dual Linear Programming**

  **Given $w \gg 0$, find $v^{-1}$ and $v^0$ that**

  **minimize** $\qquad wv^{-1}$

  **subject to** $\qquad v^{-1} - Q_\gamma v^0 \qquad r_\gamma$

  $\qquad\qquad\qquad -Q_\gamma v^{-1} \qquad\qquad\quad 0$

  **for all $\gamma \in \Delta$.**

  **This linear program has optimal solution $(v^{-1} \quad v^0)$; $v^{-1} = \max_\gamma v_\gamma^{-1}$. There is a procedure to find a $\delta$ such that $v^{-1} = v_\delta^{-1}$.**

  **Each such $\delta$ has MRR.**

# Sequential Decomposition Method (OV)

# For finding an $n$-optimal policy

- **Given $\delta$ is $(n-1)$-Optimal**

    **(a) Find $\zeta$ Satisfying $(n-1)$-Optimality Conditions.**

    **(b) Find $\eta$ that is $n$-Optimal on States Recurrent under some $n$-Optimal Policy.**

    **(c) Find $\theta$ that is $n$-Optimal.**

**Theorem. Sequential Decomposition.**

**(a) Given $(n-1)$-Optimal $\delta$,**

- **Solve MTV Problem: Find Least Solution $v^n$ of**

$$v^n = \max_{\gamma \in \mathcal{E}^{n-1}_\delta} (r^n_\gamma - v^{n-1}_\delta + P_\gamma v^n) \vee v^n_\delta,$$

**so $v^n = V_\kappa$ has MTV in this System $(\mathcal{E}^{n-1}_\delta \supseteq \Delta_{n-1})$.**

- **Set $\zeta^s = \kappa^s$ if $V_{\kappa s} > v^n_{\delta s}$ & $\zeta^s = \delta^s$ if $V_{\kappa s} = v^n_{\delta s}$.**

- **$v^n_\zeta = V_\kappa$ & $\zeta$ Satisfies $(n-1)$-Optimality Conditions.**

**(b) Given $\zeta$ Satisfying $(n{-}1)$-Optimality Conditions,**

- $$\max_{\gamma\in\mathcal{E}^n_\zeta} v^n_\gamma = v^n_\zeta + \max_{\gamma\in\mathcal{E}^n_\zeta} P^*_\gamma(r^{n+1}_\gamma - v^n_\zeta).$$

- **Each $\eta$ solving above MRR Problem is $(\mathcal{E}^n_\zeta \subseteq \Delta_{n-1})$**

  - **$(n{-}1)$-Optimal and**

  - **$n$-Optimal on States Recurrent under some $n$-Optimal Policy.**

**(c) Given $\eta$ Satisfying Last Two Conditions,**

- **Solve MTV Prob: Find Least Solution $v^n$ of**

$$v^n = \max_{\gamma\in\mathcal{E}^{n-1}_\eta}\left(r^n_\gamma - v^{n-1}_\eta + P_\gamma v^n\right) \vee v^n_\eta,$$

  **so $v^n = V_\kappa$ is MTV in this System $(\mathcal{E}^{n-1}_\eta \supseteq \Delta_{n-1})$.**

- **Set $\theta^s = \kappa^s$ if $V_{\kappa s} > v^n_{\eta s}$ & $\theta^s = \eta^s$ if $V_{\kappa s} = v^n_{\eta s}$.**

- **$v^n_\theta = V_\kappa$ and $\theta$ is $n$-Optimal.**

# Remarks: $n$-Optimality-Problem

- $2n + 3$ MTV and $n + 2$ MRR Subproblems

  - ◆ Each Solvable in Polytime by LP in Original Data

    To see why, consider Theorem 1a. The inputs

    to the MTV problem there are $v_\delta^{n-1}$ and $v_\delta^n$.

    They can be found in polytime by solving

    $$r_\delta^m + Q_\delta v^m = v^{m-1}, \ m \leq n+1 \ \Rightarrow V^n = V_\delta^n.$$

  - ◆ Each Solvable by Policy Improvement

  - ◆ Subproblem Solutions Distinct

  - ◆ Subproblems Independent of Solution Method

  - ◆ State-Classification Refinements

  - ◆ More Efficient Algorithms for Special Structures

- $n$-Optimality-Problem Solvable in Polytime

- $n$-Optimality-Conditions Problem

# Unique Transition Systems [OV]

- **Unique Transition System: One in which Action is Next (Non-Stopped) State Visited**

- **Includes Standard Deterministic Systems**

- **Finding $(n-2)$-Optimal Policy is Solvable in Strongly Polynomial Time, viz.,**

  - ♦ $O(nS[A \wedge S^2] + A)$ **Time**                 **Undiscounted**

  - ♦ $O(nS^2[A \wedge S^2] + A)$ **Time**                 **Discounted**

  - ♦ $O(nS^2A \log S)$ **Time**                 **General**

    $(A = \#$ **state-action pairs)**