

# A Compromise Programming Approach to Multiobjective Markov Decision Processes

Włodzimierz Ogryczak<sup>1</sup>, Patrice Perny and Paul Weng

LIP6 - UPMC, Paris, France

June, 13-17 2011

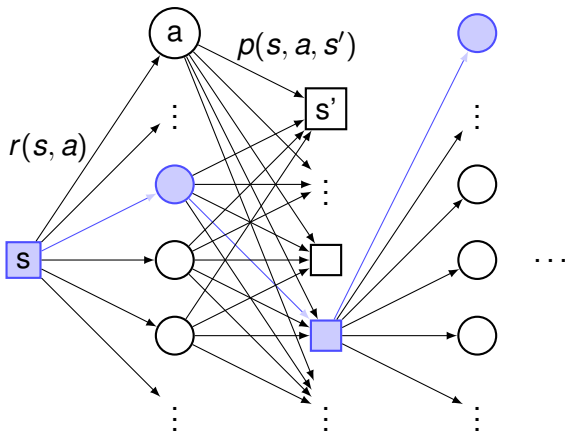
21st International Conference on MCDM  
Jyväskylä, Finland



---

<sup>1</sup> on leave from Warsaw University of Technology, Poland

# Sequential Decision Making under Uncertainty



# Markov Decision Processes (MDPs)

## Definition

- $S$  set of states
- $A$  set of actions
- $p : S \times A \times S \rightarrow [0, 1]$
- $r : S \times A \rightarrow \mathbb{R}$

## Solution

- Pure/randomized decision rule  $\delta$
- (Stationary) pure policy  $\pi$

# Value Functions and Solution Methods

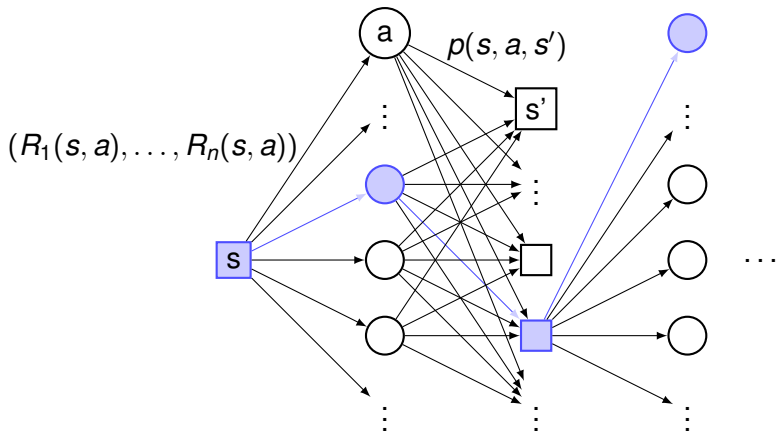
## Value functions

- $v_t^\pi(\mathbf{s}) = r(\mathbf{s}, \delta_t(\mathbf{s})) + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}, \delta_t(\mathbf{s}), \mathbf{s}') v_{t-1}^\pi(\mathbf{s}')$
- $\pi \succsim \pi' \Leftrightarrow \forall \mathbf{s}, v^\pi(\mathbf{s}) \geq v^{\pi'}(\mathbf{s})$
- $v^*(\mathbf{s}) = \max_{a \in A} r(\mathbf{s}, a) + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{s}, a, \mathbf{s}') v^*(\mathbf{s}')$

## Family of solution methods

- Value/Policy iterations
- LP

# Multiobjective MDPs (MMDPs)



# Multiobjective MDPs (MMDPs)

## Definition

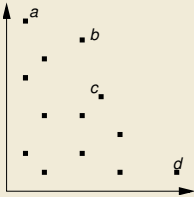
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^n$  ( $n$  criteria)
- $V^\pi(\mathbf{s}) \in \mathbb{R}^n$

## Value functions

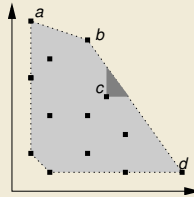
- $$V_t^\pi(\mathbf{s}) = R(\mathbf{s}, \delta_t(\mathbf{s})) + \gamma \sum_{s' \in \mathcal{S}} p(\mathbf{s}, \delta_t(\mathbf{s}), s') V_{t-1}^\pi(s')$$
- $\pi \succsim \pi' \Leftrightarrow \forall \mathbf{s}, V^\pi(\mathbf{s}) \geq_P V^{\pi'}(\mathbf{s})$

# Scalarizing Function for Compromise Search

## Example



Pure policies



Randomized policies

## Scalarizing Function $\psi$

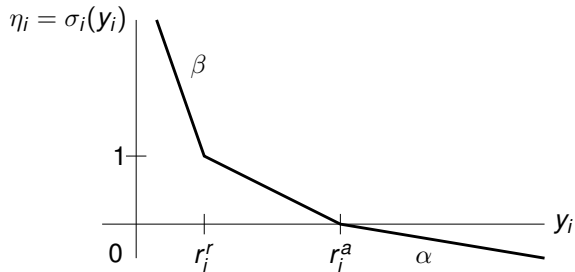
- $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$  monotonic w.r.t. Pareto dominance
- $v(\mathbf{s}) = \psi(V_1(\mathbf{s}), \dots, V_n(\mathbf{s}))$
- Weighted sum does not provide any control on tradeoffs

# Reference Point Method (RPM)

Generic Scalarizing Achievement Function (Wierzbicki, 82)

$$\psi_\varepsilon(\mathbf{y}) = (1 - \varepsilon) \max_{i=1 \dots n} \sigma_i(\mathbf{y}_i) + \frac{\varepsilon}{n} \sum_{i=1 \dots n} \sigma_i(\mathbf{y}_i)$$

$$\sigma_i(\mathbf{y}_i) = \frac{1}{r_i^f - r_i^a} \max \{ \beta y_i + (1 - \beta) r_i^f - r_i^a, y_i - r_i^a, \alpha (y_i - r_i^a) \}$$





# RPM with an OWA

## OWA

$$OWA(\eta) = \sum_{i=1}^n \omega_i \eta_{\langle i \rangle} \quad \text{where} \quad \eta_i = \sigma_i(y_i) \quad \forall i = 1 \dots n$$

where  $\omega_1 > \omega_2 > \dots > \omega_n > 0$  and  $\eta_{\langle 1 \rangle} \geq \eta_{\langle 2 \rangle} \geq \dots \geq \eta_{\langle n \rangle}$

## Example

$$r^r = (0, 0, 0) \quad r^a = (10, 10, 10) \quad w = (5/10, 3/10, 2/10)$$

$y$	$\eta$	$\eta_{\langle 1 \rangle}$	$\eta_{\langle 2 \rangle}$	$\eta_{\langle 3 \rangle}$	OWA	$\psi_0$	$\psi_\epsilon$
(4, 5, 9)	(6, 5, 1)	6	5	1	4.7	6	$6 + 4\epsilon$
(4, 8, 6)	(6, 2, 4)	6	4	2	4.6	6	$6 + 4\epsilon$
(4, 7, 7)	(6, 3, 3)	6	3	3	4.5	6	$6 + 4\epsilon$

# Main properties of OWA

- Symmetry:

$$OWA(\eta_1, \dots, \eta_n) = OWA(\eta_{\tau(1)}, \dots, \eta_{\tau(n)})$$

- Pareto-Monotonicity:

$$\eta \succ_P \eta' \Rightarrow OWA(\eta) > OWA(\eta')$$

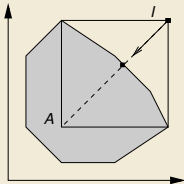
- Fairness (Monotonicity w.r.t Pigou-Dalton transfers):

$$\forall i, j \in \{1, \dots, n\} \text{ s.t. } \eta_i > \eta_j, \forall \varepsilon \in (0, \eta_i - \eta_j),$$

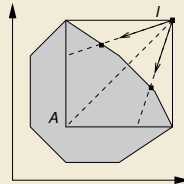
$$OWA(\eta_1, \dots, \eta_i - \varepsilon, \dots, \eta_j + \varepsilon, \dots, \eta_n) < OWA(\eta_1, \dots, \eta_n)$$

# RPM with a Weighted OWA

OWA symmetric on regrets  $\Rightarrow$  WOWA



Compromise Solution



Different Weights

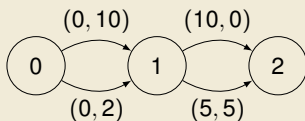
RPM WOWA

$$WOWA(\eta) = \sum_{i=1}^n w_i(\lambda, \eta) \eta_{\langle i \rangle}$$

where  $w_i(\lambda, \eta) = \varphi(\sum_{k \leq i} \lambda_{\tau(k)}) - \varphi(\sum_{k < i} \lambda_{\tau(k)})$  with  $\varphi$  a monotone increasing function that interpolates points  $(\frac{i}{m}, \sum_{k \leq i} \omega_k)$  together with the point  $(0.0)$ .

# Difficulty of Compromise Search in MDPs

## Example: Not dynamically consistent



- Best Compromise in 1 :  
Down  $\succ$  Up
- Best compromise in 0 :  
**(Up, Down)**  $\prec$  **(Up, Up)**

- Dependent on the initial state

- $\pi \succ \pi' \not\Rightarrow (\delta, \pi) \succ (\delta, \pi')$

$\Rightarrow$  Optimal policies cannot be obtained from optimal sub-policies (true for Tchebycheff, OWA, WOWA. . .)

# Solving Standard MDPs by LP

## Primal LP

$$\begin{aligned} \min \quad & \sum_{s \in S} \mu(s) v(s) \\ \text{s.t.} \quad & v(s) \geq r(s, a) + \gamma \sum_{s' \in S} p(s, a, s') v(s') \quad \forall s \in S, \forall a \in A \end{aligned}$$

## Dual LP

$$\begin{aligned} \max \quad & \sum_{s \in S} \sum_{a \in A} r(s, a) x(s, a) \\ \text{s.t.} \quad & \begin{cases} \sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S} \sum_{a \in A} x(s', a) p(s', a, s) = \mu(s) \quad \forall s \in S \\ x(s, a) \geq 0 \quad \forall s \in S, \forall a \in A \end{cases} \end{aligned}$$

# MMDP as MOLP

(Viswanathan et al., 77)

$$\begin{array}{ll} \max & y_i = \sum_{s \in S} \sum_{a \in A} R_i(s, a) x(s, a) \quad \forall i = 1, \dots, n \\ \text{s.t.} & \left\{ \begin{array}{l} \sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S} \sum_{a \in A} x(s', a) p(s', a, s) = \mu(s) \quad \forall s \in S \\ x(s, a) \geq 0 \quad \forall s \in S, \forall a \in A \end{array} \right. \end{array}$$

# Solving Method (1/2)

## NLP formulation

$$\begin{array}{ll}
 \min & \text{WOWA}(\eta) \\
 \text{s.t.} & \left\{ \begin{array}{l}
 \eta_i = \max \left\{ \beta \frac{y_i - r_i^r}{r_i^r - r_i^a} + 1, \frac{y_i - r_i^a}{r_i^r - r_i^a}, \alpha \frac{y_i - r_i^a}{r_i^r - r_i^a} \right\} \quad \forall i = 1 \dots n \\
 y_i = \sum_{s \in S} \sum_{a \in A} R_i(s, a) x(s, a) \quad \forall i = 1 \dots n \\
 \sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S} \sum_{a \in A} x(s', a) p(s', a, s) = \mu(s) \quad \forall s \in S \\
 x(s, a) \geq 0 \quad \forall s \in S, \forall a \in A
 \end{array} \right.
 \end{array}$$

# Solving Method (2/2)

## LP reformulation

$$\begin{aligned}
 \min \quad & \sum_{k=1}^n \omega'_k z_k \\
 \text{s.t.} \quad & \left\{ \begin{array}{ll}
 z_k = kt_k + n \sum_{i=1 \dots n} \lambda_i d_{ik} & \forall k = 1 \dots n \\
 \eta_i \leq t_k + d_{ik}, d_{ik} \geq 0 & \forall i = 1 \dots n, \forall k = 1 \dots n \\
 \eta_i \geq \beta(y_i - r_i^r)/(r_i^r - r_i^a) + 1 & \forall i = 1 \dots n \\
 \eta_i \geq (y_i - r_i^a)/(r_i^r - r_i^a) & \forall i = 1 \dots n \\
 \eta_i \geq \alpha(y_i - r_i^a)/(r_i^r - r_i^a) & \forall i = 1 \dots n \\
 y_i = \sum_{s \in S} \sum_{a \in A} R_i(s, a)x(s, a) & \forall i = 1 \dots n
 \end{array} \right. \\
 & \sum_{a \in A} x(s, a) - \gamma \sum_{s' \in S} \sum_{a \in A} x(s', a)p(s', a, s) = \mu(s) \quad \forall s \in S \\
 & x(s, a) \geq 0 \quad \forall s \in S, \forall a \in A
 \end{aligned}$$



# Idea of the Linearization: the Example of OWA

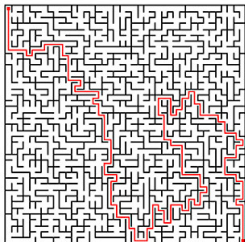
- $$OWA(\eta) = \sum_{k=1}^n w'_k L_k(\eta) \quad \left\{ \begin{array}{l} L_k(\eta) = \sum_{i=1}^k \eta_{(i)} \\ w'_i = w_i - w_{i+1} \quad \forall i = 1 \dots n-1 \\ w'_n = w_n \end{array} \right.$$

- $$L_k(\eta) = \max_{u_{ik}} \left\{ \sum_{i=1}^n \eta_i u_{ik} : \sum_{i=1}^n u_{ik} = k, 0 \leq u_{ik} \leq 1, \forall i \right\}$$

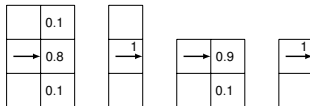
- $$L_k(\eta) = \min_{t_k, d_{ik}} \left\{ kt_k + \sum_{i=1}^n d_{ik} : \eta_i \leq t_k + d_{ik}, d_{ik} \geq 0, \forall i \right\}$$

- $$\min OWA(\eta) = \left\{ \begin{array}{l} \min \sum_{k=1}^n w'_k (kt_k + \sum_{i=1}^n d_{ik}) \\ \text{s.t.} \quad \left\{ \begin{array}{l} \eta_i \leq t_k + d_{ik} \quad \forall i, k \\ d_{ik} \geq 0 \quad \forall i, k \end{array} \right. \end{array} \right.$$

# Experiments: Navigation in a Grid

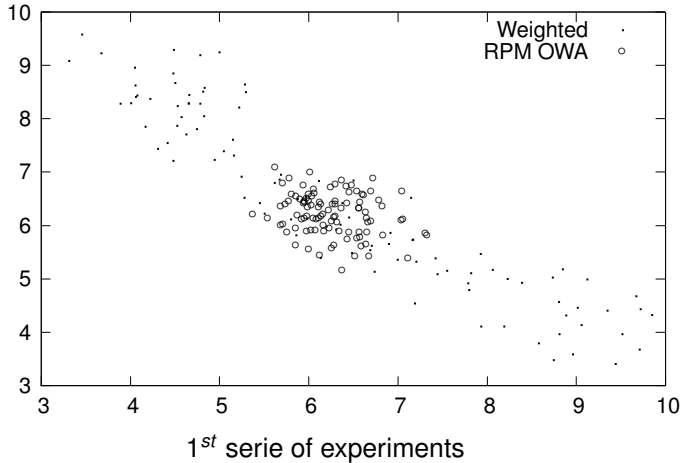


- $S = N \times N$
- $A = \{U, D, L, R\}$
- $p$
- $R$  randomly drawn in  $[0, 1]^2$

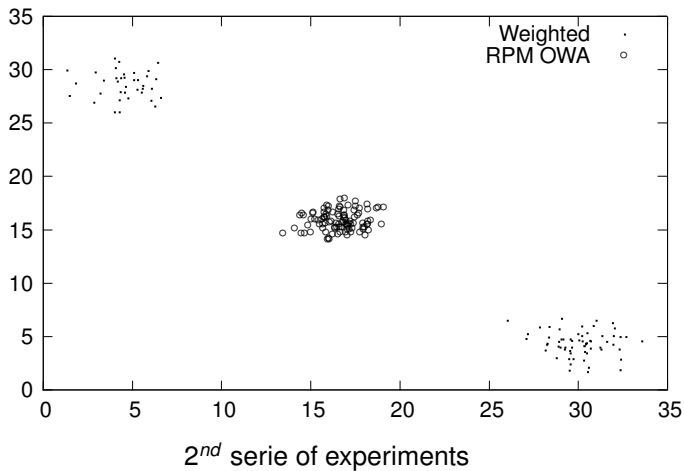


Transitions

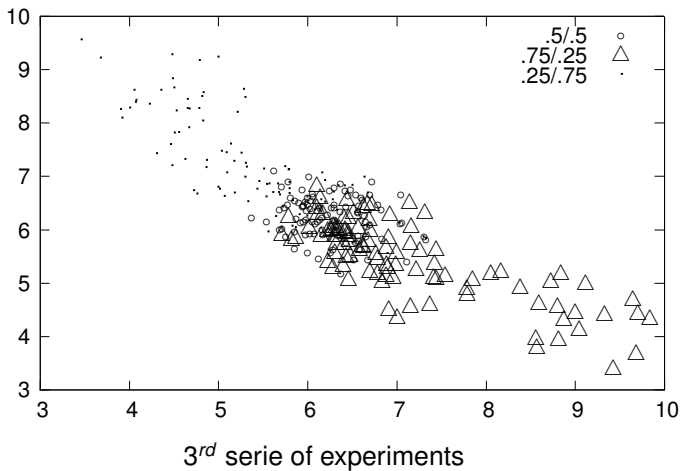
# Results on the Navigation Problem (1/3)



## Results on the Navigation Problem (2/3)



# Results on the Navigation Problem (3/3)



# Computation Times

Table: Average execution time in seconds

$n$	Size	W	RPM OWA	RPM WOVA
2	400	0.17	0.48	0.46
2	2500	5.13	15.06	15.12
2	10000	151.51	417.02	422.06
4	400	0.12	0.75	0.76
4	2500	5.2	28.21	28.27
4	10000	154	821.27	829.83
8	400	0.12	1.3	1.3
8	2500	4.96	50.62	50.72
8	10000	158.26	1514.21	1538.19

# Conclusion and Future Work

- Best compromise solution in MMDPs
  - LP approach based on a WOWA scalarizing function
- Other non linear scalarizing functions
  - Choquet integral,...
- Factorized MDP and MMDP
  - to overcome the curse of dimensionality





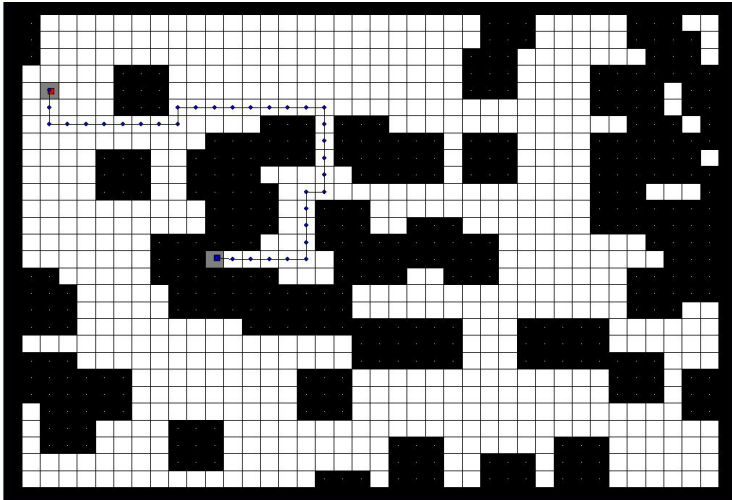
## Proposition

*For any positive and strictly decreasing normalized preferential weights  $\omega_1 > \omega_2 > \dots > \omega_n > 0$  and any positive importance weights  $\lambda_i$ , if  $\bar{y}$  is a properly nondominated with tradeoffs bounded by  $\Delta = n\bar{\lambda}\beta\omega_1/(1 - n\bar{\lambda}\omega_1)$  where  $\bar{\lambda} = \min_{i \in I} \lambda_i$ , i.e. if for any attainable outcome vector  $y$  the implication*

$$y_i > \bar{y}_i \text{ and } y_k < \bar{y}_k \Rightarrow \frac{y_i - \bar{y}_i}{\bar{y}_k - y_k} \leq \Delta = \frac{n\bar{\lambda}\beta\omega_1}{(1 - n\bar{\lambda}\omega_1)} \quad (1)$$

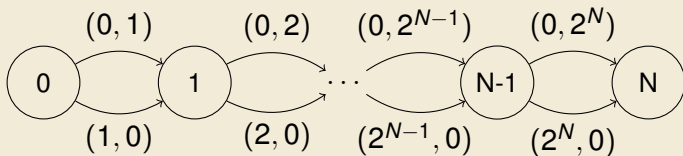
*is valid for any  $i, k \in I$ , then there exist aspirations levels  $r_i^a$  and reservation levels  $r_i^r$  such that  $\bar{y}$  is an optimal solution of the corresponding problem.*

# Example: Path-planning problems

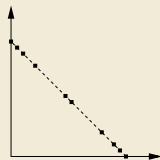


# Do we really need all the Pareto optimal solutions?

Example adapted from (Hansen, 80)



$$V_1^\pi(0) + V_2^\pi(0) = \sum_{i=0}^N 2^i = 2^{N+1} - 1$$



- The number of Pareto optimal pure policies grows exponentially with the number of states