# Markov Decision Problem Linear Programming Method

author

This Hypercard stack was prepared by: Dennis L. Bricker, Dept. of Industrial Engineering, University of Iowa, Iowa City, Iowa 52242 e-mail: dbricker@icaen.uiowa.edu





# Linear Programming Algorithm without Discounting

Optimizes the "average", i.e., expected, cost or return per period in steady state.



# Linear Programming Algorithm with Discounting

Optimizes the present value of all future expected costs

# LP model of MDP

Assume that, using the optimal policy, a steady state distribution exists.

Define "randomized" or "mixed" strategies:

 $X_i^k$  = joint probability, in steady state, of being in state i and selecting action k  $\epsilon$  K<sub>i</sub>



### LP Model

Maximize 
$$\sum_{i \in S} \sum_{k \in K_i} C_i^k X_i^k$$

$$\sum_{k \in K_{j}} X_{j}^{k} = \sum_{i \in S} \sum_{k \in K_{i}} p_{ij}^{k} X_{i}^{k} \quad \forall j \in S$$

$$\sum_{i \in S} \sum_{k \in K_{i}} X_{i}^{k} = 1$$

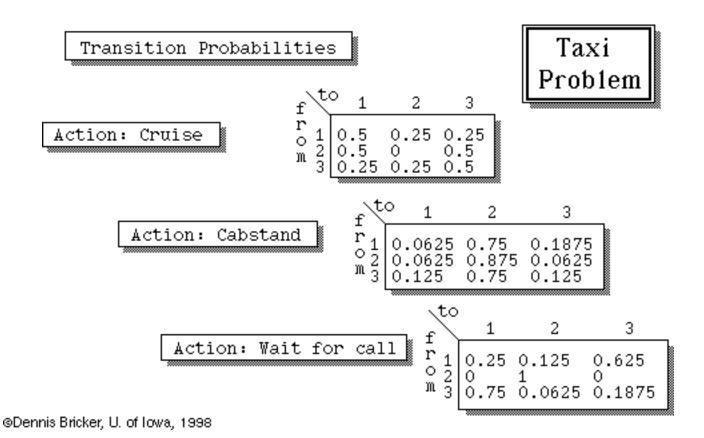
$$\sum_{i \in S} \sum_{k \in K_{i}} X_{i}^{k} = 1$$

$$X_{i}^{k} \ge 0$$

$$2 \text{ and constant } C$$

$$2 \text{ and constant } C$$

-One constraint is redundant, and can be eliminated



Taxi Problem

#### Cost Matrix

k	name	1	2	3
1	Cruise	-8	-16	-7
2	Cabstand	-2.75	-15	-4
3	Wait for call	-4.25	999	-4.5

(Rows ~ actions, Columns ~ states)

A value of 999 above signals an infeasible action in a state.

Expected returns for each 1&k

LP Tableau

Taxi Problem

k:	1	2	3	1	2	1	2	3	R
i:	1	1	1	2	2	3	3	3	н S
Min	-8	-2.75	-4.25	<sup>-</sup> 16	<sup>-</sup> 15	-7	-4	<sup>-</sup> 4.5	
	0.5 -0.25	0.9375 -0.75	0.75 -0.125	-0.5 1	-0.0625 0.125		-0.125 -0.75	-0.75 -0.0625	00
	1	1	1	1	1	1	1	1	1

Iteration 0

LP Tableau

### Initial basic feasible solution

basic	*			*		*			
k: [	1	2	3	1	2	1	2	3	R
i:	1	1	1	2	2	3	3	3	H S
Min	0 1 0 0	2.1 1.45 -0.4 -0.05	5.01667 1.36667 0.1 -0.466667	0	-4.95 0.35 0.3 0.35	0	-0.566667 0.0333333 -0.4 1.36667	3.23333 -0.616667 0.15 1.46667	9.2 0.4 0.2 0.4

Initial policy: in each city, select "cruise" ("greedy" policy)

@Dennis Bricker, U. of Iowa, 1998

Iteration 0
Policy: (Cost= -9.2)

$$\begin{array}{l}
\text{initial} \\
\text{basic} \\
\text{solution}
\end{array} \begin{cases}
X_1^1 = 0.4 \\
X_2^1 = 0.2 \\
X_3^1 = 0.4
\end{cases}$$

State	Action	P{i}
1 Town A	1 Cruise	0.4
2 Town B	1 Cruise	0.2
3 Town C	1 Cruise	0.4

Initial policy: in each city, select "cruise" ("greedy" policy)

#### Iteration 0

# LP Tableau

basic	*			*		*			
k: [	1	2	3	1	2	1	2	3	R
1:	1	1	1	2	2	3	3	3	H S
Min	0 1 0 0	2.1 1.45 -0.4 -0.05	5.01667 1.36667 0.1 -0.466667	0 0 1 0	-4.95 0.35 0.3 0.35	0 0 0 1	-0.566667 0.0333333 -0.4 1.36667	3.23333 -0.616667 0.15 1.46667	9.2 0.4 0.2 0.4



minimum 
$$\left\{ \frac{0.4}{0.35}, \frac{0.2}{0.3}, \frac{0.4}{0.35} \right\} = \frac{0.2}{0.3}$$

X<sub>2</sub> enters the basis, replacing X<sub>2</sub><sup>1</sup>

@Dennis Bricker, U. of Iowa, 1998

Ę	Ite ★	ration 1			*	*		LP T	ab1eau
k:	1	2	3	1	2	1	2	3	
i:	1	1	1	2	2	3	3	3	rhs
Min	0 1 0 0	-4.5 1.9166 -1.3333 0.4166	1.25 0.3333			0		5.70833 -0.79166 0.5 1.29167	12.5 0.1666 0.6666 0.1666

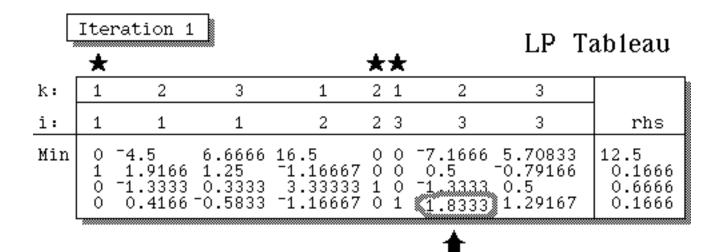
basic 
$$\begin{cases} X_1^1 = \frac{1}{6} \\ X_2^2 = \frac{2}{3} \\ X_3^1 = \frac{1}{6} \end{cases}$$

@Dennis Bricker, U. of Iowa, 1998

## Iteration 1

Policy: (Cost= ~12.5 )

State	Action	P{i}
1 Town A	1 Cruise	0.166667
2 Town B	2 Cabstand	0.666667
3 Town C	1 Cruise	0.166667



minimum 
$$\left\{ \frac{0.166}{0.5}, \frac{0.1666}{1.833} \right\} = \frac{0.1666}{1.8333}$$

 $X_3^2$  enters the basis, replacing  $X_3^1$ 

	Ιtε	eration	2					LP 7	Γableau
	*				*		*		
k:	1	2	3	1	2	1	2	3	
i:	1	1	1	2	2	3	3	3	rhs
Min	0 1 0 0	-2.8712 1.8030 -1.0303 0.2272	1.4090	11.9394 -0.8484 2.4848 -0.6363	0	-0.2727 0.7272	0	10.7576 -1.1439 1.4393 0.7045	0.1212 0.7878

Note that for every state, there is a variable in the basis for only one action!

# Iteration 2

Policy: (Cost= ~13.1515 )

State	Action	P{i}
1 Town A	1 Cruise	0.121212
2 Town B	2 Cabstand	0.787879
3 Town C	2 Cabstand	0.0909091

Iteration 2

LP Tableau

@Dennis Bricker, U. of Iowa, 1998

		*			*		*	LP Ta	ıbleau
k:	1	2	3	1	2	1	2	3	
i:	1	1	1	2	2	3	3	3	rhs
Min	1.59244 0.55462 0.57142 -0.12605	1	0.78151	-0.4705 2	0 0 1 0	-0.1512 0.5714		8.9359 -0.6344 0.7857 0.8487	13.3445 0.06722 0.85714 0.07563

Reduced costs are all nonnegative... the optimality condition is satisfied! Optimal Policy

#### Iteration 3

Policy: (Cost= -13.3445 )

State	Action	P{i}
1 Town A	2 Cabstand	0.0672269
2 Town B	2 Cabstand	0.857143
3 Town C	2 Cabstand	0.0756303

The optimal policy found by the simplex LP algorithm is deterministic, not randomized, i.e., for each state, only one action is specified.



# LP Algorithm for MDP with discounting

Determining a policy which minimizes the *present value* of all future costs over an infinitely long planning horizon.

Note: existence of a steady state distribution is *not* assumed!



The present value of future costs (i.e., the discounted future costs) will depend upon the initial state of the system.

### Define

 $\alpha_j$  = probability that system in initially in state j

Note: If the initial state is known, then  $\alpha = [0, 0, ..., 0, 1, 0, ..., 0]$ 

MDP LP Algorithm 8/20/00 page 21

### Decision variables

$$\lambda_i^k(n) = \text{Joint probability that} \\ \text{and} \\ \text{artion } k \in K_j \text{ is selected}$$

Note that this definition of the decision variables does not assume that the same policy is optimal for every stage!

#### Define

$$\beta$$
 = discount factor =  $\frac{1}{1+r}$  where  $r$  = rate of return per stage

Then the present value of a cost Y which is incurred 1 period hence is βY

2 periods hence is β<sup>2</sup> Y

i n periods hence is β<sup>n</sup> Y

If  $C_j^k = cost of action k in state j$ 

then

$$\sum_{j} \sum_{k \in K_{j}} C_{j}^{k} \lambda_{j}^{k}(n) = \text{expected cost during} \\ \text{stage (period) } n$$

and

$$\sum_{n=0}^{\infty} \beta^n \sum_{j} \sum_{k \in K_j} C_j^k \lambda_j^k(n) = \text{present value of} \\ \text{all costs in periods} \\ \text{n=0, 1, 2, ....}$$

Our objective is therefore to minimize the discounted future expected costs:

$$\sum_{j} \sum_{k \in K_{j}} \left[ \sum_{n=0}^{\infty} \beta^{n} C_{j}^{k} \lambda_{j}^{k}(n) \right]$$

# Constraints

For each state j at stage n=0:  $\sum_{k \in K_i} \lambda_j^k(0) = \alpha_j$ 

$$\sum_{k \in K_j} \lambda_j^k(0) = \alpha_j$$

For each state j at stage n, n=1,2,...

$$\sum_{\mathbf{k}\in\mathbf{K}_{\mathbf{j}}}\lambda_{\mathbf{j}}^{\mathbf{k}}(\mathbf{n})=$$

Probability that system is in state j at stage n

$$\sum_{k \in K_j} \lambda_j^k(\mathbf{n}) = \sum_i \sum_{k \in K_i} \mathbf{p}_{ij}^k \lambda_i^k(\mathbf{n}-1)$$

Probability that system makes transition from state i in stage n-1 to state j in stage n

Note that there is an infinite number of as infinitely many constraints, as well

In order to reduce the size of the LP to finite proportions, we will utilize the z - transform.

The z-transform of the sequence  $\{a_n\}_{n=0}^{\infty}$  is the *function* 

$$F(z) = \sum_{n=0}^{\infty} z^n a_n$$

[See Queueing Systems, Vol. 1, Appendix 1 by L. Kleinrock]

Note that, given F, we can reconstruct the sequence:

$$\mathbf{a}_{n} = \frac{1}{n!} \frac{d^{n} F(0)}{dz^{n}}$$

For each pair of state i and action k, consider the sequence of probabilities

$$\left\{\lambda_{j}^{k}(n)\right\}_{n=0}^{\infty}$$

 $\left\{\,\lambda_{\,j}^{\,k}(n)\right\}_{n=0}^{\infty}$  Its z-transform is  $F\left(z\right)=\sum_{n=0}^{\infty}\,z^n\,\lambda_{\,j}^{\,k}(n)$ 

Define a new set of decision variables

$$x_j^k = \sum_{n=0}^\infty \beta^n \, \lambda_j^k(n)$$
 i.e., the z-transform of  $\left\{\lambda_j^k(n)\right\}_{n=0}^\infty$ 

evaluated at B

We are then able to rewrite our objective function

$$\sum_{j} \sum_{k \in K_{j}} \left[ \sum_{n=0}^{\infty} \beta^{n} C_{j}^{k} \lambda_{j}^{k}(n) \right]$$

with a finite number of terms:

$$\sum_{j} \sum_{k \in K_{j}} C_{j}^{k} x_{j}^{k}$$

where

$$\mathbf{x}_{j}^{k} = \sum_{n=0}^{\infty} \beta^{n} \lambda_{j}^{k}(\mathbf{n})$$

# Constraints

In order to reduce the set of constraints to a finite number

(with finitely many variables), perform the following operations:

 $\bullet$  For each pair j & n, multiply the corresponding constraint by  $\beta^n$ 

$$\begin{cases} \beta^{o} \sum_{k \in K_{j}} \lambda^{k}_{j}(0) = \beta^{o} \alpha_{j} & \text{for each state } j \\ \\ \sum_{k \in K_{j}} \beta^{n} \lambda^{k}_{j}(n) = \beta \sum_{i} \sum_{k \in K_{i}} p^{k}_{ij} \beta^{n-1} \lambda^{k}_{i}(n-1) & \text{for each state } j \\ \\ & \& n \ge 1 \end{cases}$$

• For each state j, sum the equations over n:

$$\begin{cases} \beta^{\circ} \sum_{k \in K_{j}} \lambda_{j}^{k}(0) &= \beta^{\circ} \alpha_{j} & \text{for each state } j \\ \\ \sum_{k \in K_{j}} \beta^{n} \lambda_{j}^{k}(n) &= \beta \sum_{i} \sum_{k \in K_{i}} p_{ij}^{k} \beta^{n-1} \ \lambda_{i}^{k}(n-1) & \text{for each state } j \\ \\ &\& n \geq 1 \end{cases}$$

$$\implies \sum_{n=0}^{\infty} \sum_{k \in K_i} \beta^n \lambda_j^k(n) = \alpha_j + \beta \sum_{n=1}^{\infty} \sum_{i} \sum_{k \in K_i} p_{ij}^k \beta^{n-1} \lambda_i^k(n-1)$$

 Rearrange the order of summation in this new constraint:

$$\begin{split} &\sum_{n=0}^{\infty} \sum_{k \in K_{j}} \beta^{n} \lambda_{j}^{k}(n) = \alpha_{j} + \beta \sum_{n=1}^{\infty} \sum_{i} \sum_{k \in K_{i}} p_{ij}^{k} \beta^{n-1} \lambda_{i}^{k}(n-1) \\ \Longrightarrow &\sum_{k \in K_{j}} \sum_{n=0}^{\infty} \beta^{n} \lambda_{j}^{k}(n) = \alpha_{j} + \beta \sum_{i} \sum_{k \in K_{i}} \sum_{n=1}^{\infty} p_{ij}^{k} \beta^{n-1} \lambda_{i}^{k}(n-1) \\ \Longrightarrow &\sum_{k \in K_{j}} x_{j}^{k} = \alpha_{j} + \beta \sum_{i} \sum_{k \in K_{i}} x_{i}^{k} \qquad \text{for all } j \\ \\ \text{since} &\sum_{k \in K_{j}} p_{ij}^{k} \beta^{n-1} \lambda_{i}^{k}(n-1) = \sum_{k \in K_{i}} \beta^{n} \lambda_{i}^{k}(n) \end{split}$$

@Dennis Bricker, U. of Iowa, 1998

# LP Model

$$\begin{array}{ll} \text{Minimize} & \sum_{j} \sum_{k \in K_{j}} C_{j}^{k} x_{j}^{k} \end{array}$$

subject to

$$\begin{split} \sum_{k \in K_j} x_j^k &= \alpha_j + \beta \sum_i \sum_{k \in K_i} p_{ij}^k x_i^k &\quad \text{for all } j \\ x_j^k &\ge 0 \end{split}$$

- Note that

   sum of x is not specified to be 1
   no redundant constraint was
   eliminated from state equations

MDP LP Algorithm 8/20/00 page 33

Using the "Kronecker delta", i.e.,

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

this LP model may be rewritten:

$$\begin{aligned} & \text{Minimize } \sum_{j} \sum_{k \in K_{j}} C_{j}^{k} x_{j}^{k} \\ & \text{subject to} \\ & \sum_{i} \sum_{k \in K_{i}} \left( \delta_{ij} - \beta \, p_{i \, j}^{\, k} \, x_{i}^{k} \right) = \alpha_{j} \quad \text{ for all } j \\ & x_{j}^{\, k} \geq 0 \end{aligned}$$

If  $x^*$  is the optimal basic solution, then

$$x_j^{*k} > 0$$
 (i.e., basic)

implies that

the optimal policy is to select action k when in state j for every stage n=0,1,2,...

i.e., the optimal policy is stationary, same for every time period!

