

Problem 1 Political campaign entering final stage. Candidate can find ads for 5 commercials to get possible extra votes.

Commercials	A1	A2	A3	A4
0	0	0	0	0
1	4	6	5	3
2	7	8	9	7
3	9	10	11	12
4	12	11	10	14
5	15	12	9	16

stage - # of commercials purchased overall

state - # of commercials per station

controls - choose any of {A1, A2, A3, A4}

dynamic system - follow the chart

cost function - extra votes by choosing new area

$$\text{maximize } \sum_{n=1}^4 A_n(x_n) \text{ such that } \sum_{n=1}^4 x_n = 5 = w$$

The goal of the system is to maximize votes

This is a general resource allocation problem

work backwards

$$\text{Define } V_j(w) = \sum_{n=j}^4 A_n(x_n)$$

$$w = \sum_{n=1}^4 x_n$$

$$V_n(w) = \max [r_n(x) + V_{n+1}(w-x)] \text{ try } V_3(w)$$

$$V_5(0) = 0$$

$$V_4(0) = 0, V_4(1) = 3, V_4(2) = 7, V_4(3) = 12, V_4(4) = 14, V_4(5) = 17$$

$$V_3(0) = 0$$

$$V_3(1) = \max [5, 3] = 5$$

$$V_3(2) = \max [9, 5+3, 7] = 9$$

$$V_3(3) = \max [11, 9+3, 5+7, 12] = 12$$

$$V_4(4) = \max [10, 11+3, 9+7, 5+12, 14] = 17$$

$$V_4(5) = \max [9, 10+3, 11+7, 9+12, 5+14, 16] = 23$$

lastly, try  $V_1(w)$

$$V_1(0) = 0$$

$$V_1(1) = \max [4, 6] = 6$$

$$V_1(2) = \max [7, 4+6, 11] = 11$$

$$V_1(3) = \max [9, 7+6, 4+11] = 15$$

$$V_1(4) = \max [12, 9+6, 7+11, 4+15] = 19$$

$$V_1(5) = \max [15, 12+6, 9+11, 7+15, 4+19, 23] = 23$$

$$= 23$$

try  $V_2(w)$

$$V_2(0) = 0$$

$$V_2(1) = \dots$$

$$V_2(1) = \max [6, 5] = 6$$

$$V_2(2) = \max [8, 6+5, 9] = 11$$

$$V_3(0) \quad V_3(1)$$

$$V_2(3) = \max [10, 8+5, 6+9, 12] = 15$$

$$V_2(4) = \max [11, 10+5, 8+9, 6+12, 16] = 18$$

$$V_2(5) = \max [12, 11+5, 10+9, 8+12, 6+17, 21] = 23$$

Upon finding Area 2 - 2 commercials  
Area 3 - 1 commercial we get a total of  
2 Area 4 - 3 commercials  
23k additional votes

Problem 2) Farmer J has \$5000 & 10 tons of wheat

stage - j

state - # of  $C_j$  cash &  $W_j$  wheat &  $P_j$  price of wheat per month

controls - sell upto  $M_j$ ,  $S_j$  & Buy upto  $10 - M_j$ ,  $B_j$

dynamical system -  $C_{j+1} = C_j + S_j P_j - B_j P_j$

cost function -  $g(x) = \dots$

Rules

- 1) can't borrow debt
- 2) can't lie to buyer
- 3)  $M \leq 10$

dynamic programming to maximize  $N$

Using Dynamic Programming we get the following descriptions from above

- Rules equate to
- ①  $B_j P_j \leq C_j$
  - ②  $S_j \leq W_j$
  - ③  $W_j + B_j \leq 10$

$$C_0 = 5000$$

$$W_0 = 10$$

The dynamical system is

$$C_{j+1} = C_j - B_j P_j + S_j P_j = C_j + (S_j - B_j) P_j$$

$$\& \quad W_{j+1} = W_j + B_j - S_j = W_j - (S_j - B_j)$$

Want to maximize  $C_j$  at  $j=3$  where  $j=0,1,2,3$

from month 2 to 3

$$C_3 = C_2 + (S_2 - B_2) P_2 \quad \& \quad W_3 = W_2 - (S_2 - B_2)$$

to maximize  $C_3 \Rightarrow (S_2 - B_2)_{\max} \Rightarrow S_2 = W_2 \& \underline{B_2 = 0}$  or essentially sell all

from month 1 to 3

$$C_2 = C_1 + (S_1 - B_1) P_1 \quad \& \quad W_2 = W_1 - (S_1 - B_1) \quad \text{let } P_1 = S_1 - B_1$$

$$C_3 = C_2 + S_2 P_2 \quad \& \quad W_3 = W_2 - (S_2) S_2$$

to maximize  $C_3$

$$C_3 = C_1 + (S_1 - B_1) P_1 + S_2 P_2 \quad \& \quad W_3 = W_1 - (S_1 - B_1) - S_2 = W_1 + B_1 - (S_1 + S_2) = 0$$

$$W_1 + B_1 = S_1 + S_2$$

if  $P_1 = P_2$  then best do nothing

$$\text{or } P_1 > P_2 \quad C_3 = C_1 + W_1 P_1 + S_2 (P_2 - P_1) \rightarrow C_3 = C_1 + R P_1 + S_2 P_2 \quad \text{if } P_1 > P_2$$

then maximize  $W_1 - S_2$

if  $P_2 > P_1$  then minimize  $W_1 - S_2$

from month 0 to 3

$$C_1 = C_0 + (S_0 - B_0) P_0 \quad \& \quad W_1 = W_0 - (S_0 - B_0)$$

$$C_3 = 5000 + (S_0 - B_0) P_0 + (S_1 - B_1) P_1 + S_2 P_2 \quad W_3 = 10 - (S_0 - B_0) - (S_1 - B_1) - S_2 = 0$$

$$(10 + B_0 + B_1) = S_0 + S_1 + S_2$$

using DP.

maximize  $\{C_3\}$   
 $P_j \in$

$$\{C_0 = 5000, W_0 = 10\}$$

for delete  
to buyer

3) stage - Dec.

Problem 2) Attempt #2  
Redoing #2

let stage: months finished;  $j = 0, 1, 2, 3$

state: cash, wheat, price:  $C_j, W_j, P_j$

controls: # of sold - # bought:  $D_j$  s.t.  $W_j - D_j \leq 10$   $D_j \leq W_j$   
 $-D_j P_j \leq C_j$

Dynamical system:  $C_{j+1} = C_j + D_j P_j$   $W_{j+1} = W_j - D_j$

Cost function  $g(x) = D_j P_j$

Dynamic programming equation

$$V_j(x) = \max_{D_j} [C_j + V_{j+1}(C_j + D_j P_j)]$$

at  $j=3$   $V_{j+1}=0$

$$V_3(x) = \max_{D_3} [C_3 + \dots]$$

3) stage - period,  $n$  up to  $N$   
 state variables - stock level;  $x_n$

control - quantity ordered;  $u_n \geq 0$

$$\text{minimize } E \left[ \sum_{n=0}^{N-1} c(x_n, u_n; Y_{n+1}) \right]$$

dynamics -  $x_{n+1} = x_n + u_n - Y_{n+1}$  - demand at  $n$

cost function - operation cost  $C(x_n, u_n; Y_{n+1}) = \underbrace{p u_n}_{\text{production cost}} + \underbrace{h(x_n + u_n - Y_{n+1})^+}_{\text{unit holding cost}} + \underbrace{b(x_n + u_n - Y_{n+1})^-}_{\text{backlogging cost}}$

$$b > p$$

$$x^+ = \max\{x, 0\} \quad \& \quad x^- = -\min\{x, 0\}$$

a) dynamic programming equation

b) value function & optimal policy at  $n=N-1$

a) The dynamic programming equation  
 the recursive equation to solve this system is

$$V_n(x) = \text{minimize}_{0 \leq u_n \leq 1} \left\{ C_n(x, u_n; y) + E[V_{n+1}(x_n + u_n - y_{n+1})] \right\}$$

$$V_n(x) = \text{minimize}_{-1 \leq u_n \leq 1} \left\{ p u_n + h(x_n + u_n - y_{n+1})^+ + b(x_n + u_n - y_{n+1})^- + E[V_{n+1}(x_n + u_n - y_{n+1})] \right\}$$

b) at  $n=N-1$

$$V_{N-1}(x_{N-1}) = \text{minimize}_{0 \leq u_{N-1} \leq 1} \left\{ C_{N-1}(x_{N-1}, u_{N-1}; y_N) + E[0] \right\}$$

no need to think about future stage

To get the optimal policy need to minimize  $C_{N-1}$  w/ respect to  $u_{N-1}$

$$\text{so take } \frac{dC_{N-1}}{du_{N-1}} = 0 \Rightarrow C(x_{N-1}, u_{N-1}; y_N) = p u_{N-1} + h(x_{N-1} + u_{N-1} - y_N)^+ + b(x_{N-1} + u_{N-1} - y_N)^-$$

$$\frac{dC_{N-1}}{du_{N-1}} = p + h^+ + b^+ = 0 \Rightarrow \boxed{p + \max(h, 0) - \min(b, 0) = 0} \quad ?$$

4) initial wealth  $x_0$  - state each timepoint  $n \in W_{n-1}$  - stage  
 get proportion  $u_n \in [0,1]$  to consume  $x_n$  & put the rest in the bank  
 interest  $r > 0$   $R = 1+r$

$$x_{n+1} = R x_n (1 - u_n) \quad n \in W_{n-1} \quad \text{get utility } g(u_n x_n) \text{ at } N \quad g_N(x_N)$$

choose controls  $\{u_n\}$  to maximize

$$v(x_0) = \max_{\{u_n\}} \left\{ g_N(x_N) + \sum_{n=0}^{N-1} g(u_n x_n) \right\}$$

a) write DPE

b) assume  $f(x) = x^\alpha$  &  $g(x) = x^\alpha$   
 $\alpha \in (0,1)$

Show  $V_n(x) = C_n x^\alpha$   
 $C_n$  defined recursively & independent of current wealth

a) The dynamic programming equation here is

$$V_n(x) = \max_{u_n} \left[ g_N(x_N) + \sum_{n=0}^{N-1} g(u_n x_n) + V_{n+1}(R x_n (1 - u_n)) \right]$$

b) assume  $g(x) = x^\alpha$  simplify  $g(x)$  by  $x^\alpha$

$$V_n(x) = \max_{u_n} \left[ x^\alpha + \sum_{n=0}^{N-1} (u_n x_n)^\alpha + V_{n+1}(R x_n (1 - u_n)) \right]$$

to get max need to take the derivative of  $V_n(x)$  w.r.t.  $u_n$  &  $= 0$

$$\frac{dV_n(x)}{du_n} = \sum_{n=0}^{N-1} \alpha u_n^{\alpha-1} x_n^\alpha + V_{n+1}'(R x_n (1 - u_n)) - R x_n u_n V_{n+1}' = 0$$

let  $C_n^\alpha = \alpha (u_n x_n)^\alpha$  & not sure but

this can be defined recursively &

$C_n$  shouldn't be a  $f(x^n)$

if  $V_n(x) = C_n x^\alpha$

$$\text{then } \frac{dV_n(x)}{du_n} = \sum_{n=0}^{N-1} \alpha u_n^{\alpha-1} x_n^\alpha + C_{n+1} x^\alpha (R x_n (1 - u_n)) - R x_n u_n C_{n+1} x^\alpha -$$

becomes 0

$$\text{then } \boxed{V_n(x) = C_n x^\alpha}$$

Sutton & Barto Tracking a Nonstationary Problem  $\leftarrow$  Multi-armed Bandits

2.4) if  $x_n \neq C$  then estimate  $Q_n$  is a weighted average of previously received rewards with a weighting different from (2.6)

$$Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

What is the weighting on each prior reward for the general case, analogous to (2.6) in terms of the sequence of step-size parameters

from 2.6  $Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$

from  $Q_{n+1} = Q_n + \alpha [R_n - Q_n]$

now  $Q_{n+1} = Q_n + \alpha_n [R_n - Q_n] = \alpha_n R_n + Q_n (1-\alpha_n)$

$$Q_{n+1} = \alpha_n R_n + (1-\alpha_n) [Q_{n-1} + \alpha_{n-1} [R_{n-1} - Q_{n-1}]]$$

$$Q_{n+1} = \alpha_n R_n + \alpha_{n-1} R_{n-1} (1-\alpha_n) + (1-\alpha_n)(1-\alpha_{n-1}) Q_{n-1}$$

$$Q_{n+1} = \alpha_n R_n + \alpha_{n-1} (1-\alpha_n) R_{n-1} + (1-\alpha_n)(1-\alpha_{n-1})(1-\alpha_{n-2}) Q_{n-2} + (1-\alpha_n)(1-\alpha_{n-1}) \alpha_{n-2} R_{n-2}$$

$$Q_{n+1} = Q_1 \prod_{i=1}^n (1-\alpha_i) + \sum_{i=1}^n \alpha_i R_i \prod_{j=i+1}^n (1-\alpha_j)$$

pattern from

- q' 2.5) Design & conduct an experiment to demonstrate difficulties that  
ga sample-average methods have for nonstationary problems. Modified version of  
i. 10-armed testbed  $q_*(a)$  all = . independent random walks ( $\mu=0, \sigma=0.01$ )  
 $\alpha=0.1, \epsilon=0.1, k=10,000$  steps
-

## 2.7) Optimistic Initial Values

on nonstationary problems

use step size of  $\beta_n \equiv \alpha / \bar{O}_n$   $\alpha > 0$  conventional constant  
 $\bar{O}_n$  trace of  $o_n$  that starts at 0:

$$\bar{O}_n \equiv \bar{O}_{n-1} + \alpha(1 - \bar{O}_{n-1}) \quad \text{for } n \geq 0 \text{ w/ } \bar{O}_0 \equiv 0$$

show that  $\bar{O}_n$  is an exponential recency-weighted average w/o initial bias

$$\bar{O}_{n+1} = \bar{O}_n + \alpha(1 - \bar{O}_n) = \alpha + (1 - \alpha)\bar{O}_n, \quad \bar{O}_n = \alpha + (1 - \alpha)\bar{O}_{n-1}$$

following this

$$\bar{O}_{n+1} = \alpha + (1 - \alpha)\alpha + (1 - \alpha)^2 \bar{O}_{n-1}$$

$$= \alpha + (1 - \alpha)\alpha + (1 - \alpha)^2 \alpha + (1 - \alpha)^3 \bar{O}_{n-2}$$

$$= \alpha + \alpha + \alpha + \dots + (1 - \alpha)^{n+1} \bar{O}_0 + (1 - \alpha)^n \bar{O}_1$$

$$= \alpha \sum_{i=0}^n (1 - \alpha)^i$$

$$\beta_n = \frac{\alpha}{\bar{O}_n} = \frac{\alpha}{\alpha \sum_{i=0}^{n-1} (1 - \alpha)^i} = \frac{1}{\sum_{i=0}^{n-1} (1 - \alpha)^i}$$

$$\beta_n = \frac{1}{\sum_{i=0}^{n-1} (1 - \alpha)^i}$$