

Week-6

L2 Stochastic Inventory

- When the demand of the prod. is not known with certainty - only probab. distribu" is known. At each time period, the decision maker obs. the current inventory level & decides how much to order. C_o & C_c
- Demand is discrete ($D_t = 1, 2 \dots$), delivery is instantaneous. Demand is IID random var., with a stationary PMF of p_j ($P_{D_t}(D_t = j) = p_j$). Warehouse capacity is M units. Inventory is non-perishable. Demand not met on the day is lost.
- $t = 1, 2 \dots N$, s_t , a_t

reward func"

$$r_t(s_t, a_t) = E[\min(D_t, s_t + a_t) * \text{price}] - C_o(a_t) - C_c(s_t + a_t)$$

- terminal

$$r_N(s_N, a_N) = g(s_N) = \text{salvage value.}$$

- MDP

$$s_{t+1} = \max(s_t + a_t - D_t, 0)$$

- Obj. \rightarrow Max. of reward over time horizon.
- State transition possibilities -

$$P_{s'}(s' | s, a) = \begin{cases} 0 & \text{if } s' \notin (s+a, M) \\ p_{s+a-s}, & \text{if } s' \in [0, s+a] \text{ & } s+a \leq M \\ \sum_{k>s+a} p_k & \text{if } s=0 \text{ & } s+a \leq M \end{cases}$$

Comp. of the obj. func" -

$$C_0(a_t) = \begin{cases} K + c(a_t), & \text{if } a_t > 0 \\ 0, & \text{if } a_t = 0 \end{cases}$$

Revenue,

$$F_t(s_t + a_t) = E[\min(D_t, s_t + a_t) * \text{price}]$$

- If $u = s_t + a_t > D_t = j$, let revenue $f(j) = \text{price}_j$
This happens with probab. p_j .

If demand exceeds stock, then revenue is $f(u)$
with probability $q_u = \sum_{j=u+1}^{\infty} p_j$.

$$F_t(u) = F_t(s_t + a_t) = \sum_{j=0}^{u-1} f(j) + p(j) + f(u) \times q_u$$

- Policies \rightarrow Markov decisions rule an order quantity to be ordered in each time period (a_t) to each possible starting inventory for that time period. A policy is a seq. of such decision rules ($d_n = a_n(1), a_n(2)$)

- Value Iteration

\rightarrow Finite horizon optimality (Bellman eq") -

$$v_n(s) = \max_{a \in A_s} \left\{ u(s, a) + \sum_{j \in S} \lambda^* p(j|s, a) * v_{n-1}(j) \right\}$$

Passing to limit, the ∞ horizon eq" becomes -

$$v(i) = \max_{a \in A_i} \left\{ u(i, a) + \sum_{j \in S} \lambda^* p(j|i, a) * v(j) \right\}$$

- Q-learning

- model-based MDP - when probab. transi" matrix is known ($\text{Par}(s'|s, a)$), model-free - when these probab. are not known.
- def. seq. of vectors $\{Q_t \in R^{S \times A} : t = 0, 1, 2 \dots\}$
- this is Q learning.