

### Problem 1:

#### • Problem details

- Borrow from another bank  $y \geq 0$  at constant daily interest rate  $R$ 
  - ∴ Amount owed next day = Amount Borrowed  $\cdot (1 + R)$
- Invest a portion of bank's cash in risky asset
- Both above have no transaction cost
- If cash amount  $c < C$ , pay penalty  $K \cdot \min\left(\frac{C-c}{2C}, 1\right)$ 
  - ↳ constraint  $C \geq 0, K \geq 0$

#### • Assumptions:

- ①  $c \leq K \min\left(\frac{C-c}{2C}, 1\right)$  when  $c \leq C$
  - ② Neglect deposits
  - ③ First half of the day sees withdrawals
  - ④ If insufficient cash for withdrawal, withdrawal attempted again next day
  - ⑤ All quantities continuous
- Today planning horizon
  - Goal: Frame the NLP
  - Define the state space  $S$  by:

$$S = \{(c, j, i, w) : c \in \mathbb{R}_{\geq 0}, j \in \mathbb{R}_{\geq 0}, i \in \mathbb{R}_{\geq 0}, w \in \mathbb{R}_{\geq 0}\}$$

where  $c$  is the amount of money the bank has in cash that day,  $j$  is the amount of debt the bank is carrying that day,  $i$  is the amount of money the bank has in risky assets that day, and  $w$  is the amount of unsatisfied withdrawal requests which are carried over to the next day.

- Define the action space by :

$$A = \{(d^*, i^*) : d^* \in [-\epsilon, \alpha], i^* \in [-i, c]\}$$

where  $d^*$  and  $i^*$  are changes in debt and investment amounts to be applied the following day. This means  $i_{t+1} = i_t + i^*$  and  $d_{t+1} = d_t + d^*$ .

- The reward function accounts for the amount of interest accumulated, fees paid, and investment returns from one day to the next. Assume that daily returns from the risky asset are normally distributed with mean  $\mu = f(i) \in \mathbb{R}$  and standard deviation  $\sigma = f(i) \in \mathbb{R}_{>0}$ . Assume there is a penalty  $p \cdot w$  for asking a customer to return the next day to make withdrawal  $w$ , where  $p \in \mathbb{R}_+$ . The reward as we've defined it is independent of the action. However, the action taken affects the state at the following day, which in turn affects reward. The reward function  $R(s, a)$  for specific state  $s = (c, d, i, w)$  in  $S$  and any action  $a \in A$  is given by:

$$R(s, a) = \begin{cases} N(\mu, \sigma^2) - pw - K \cot\left(\frac{\pi c}{2C}\right) & c \leq C \\ N(\mu, \sigma^2) - pw & c > C \end{cases}$$

- The money made on the risky investment  $x \in \mathbb{R}$  during the day represents one of the stochastic components of the transition probability function. Deposits and withdrawals made to the bank represent other stochastic components of the transition probability function. Assume that deposits for the day  $a \sim N(\mu_a, \sigma_a^2)$  and withdrawals for the day are given by  $b \sim N(\mu_b, \sigma_b^2)$ . Also let  $f(x, a, b)$  be the joint probabilities of the money made from investments  $x$ , the deposit amount  $a$ , and the withdrawal amount  $b$  for the day (i.e.,  $f(x, a, b) = \text{pdf}(x, N(\mu_x, \sigma_x^2)) \cdot \text{pdf}(a, N(\mu_a, \sigma_a^2)) \cdot \text{pdf}(b, N(\mu_b, \sigma_b^2))$ ). Therefore, transitions  $P(s, a, s')$  from state  $s_0 = (c_0, d_0, i_0, w_0) \in S$  and  $s_1 = (c_1, d_1, i_1, w_1) \in S$  given action  $a = (d^*, i^*) \in A$  is given by:

$$P(s_{t+1}, a_t, s_{t+1}) = \begin{cases} f(x, a_t, b) & c_m = \max(0, c_m + a - b + x - u_{(t)} - p u_{(t)} + d' - i') \\ \text{where and} & d_{(t)} = d_{(t)} + d'; \quad i_{(t)} = i_{(t)} + i'; \quad u_{(t)} = u_{(t)} + \\ & \rightarrow \min(0, c_m + a - b + x - u_{(t)} - p u_{(t)} + d' - i') + \max(0, u_i \\ & - (c_m + a - b + x - p u_{(t)} + d' - i')); \quad c_{(t)} > C \quad \text{What could not be paid off} \\ & \text{OR} \\ & c_m = \max(0, c_m + a - b + x - u_{(t)} - p u_{(t)} + d' - i' - K \text{cat}[\frac{\pi_e}{2C}]); \\ & d_{(t)} = d_{(t)} + d'; \quad i_{(t)} = i_{(t)} + i'; \quad u_{(t)} = u_{(t)} + \\ & \min(0, c_m + a - b + x - u_{(t)} - p u_{(t)} + d' - i' - K \text{cat}[\frac{\pi_e}{2C}]) + \max(0, u_i \\ & - (c_m + a - b + x - p u_{(t)} + d' - i' - K \text{cat}[\frac{\pi_e}{2C}])); \quad c_{(t)} \leq C \\ 0 & \text{otherwise} \end{cases}$$

- How Would I solve this MDP Control Problem?

• Our state and action space are continuous (non-finite, intractably large). Therefore, I would solve this MDP control problem with Approximate DP. I would represent the policy and value function using function approximation (e.g. linear, neural network), which accept continuous states as inputs. I would then solve for the optimal policy with Approximate Policy Evaluation. In approximate policy iteration, the policy evaluation step of the algorithm is performed by iteratively applying the Bellman operator on the function approximator of the value function. In an iterative process, "source" and "next" states are sampled, reward is calculated, and the value function at the source state is approximated by the avg. return + discounted VF at next state (also obtained by the function approximator). During the policy improvement step of the algorithm, changes to the policy are made by updating parameters of the policy function approximator in the direction that selects actions corresponding to the path of greatest ascent of the value function approximation.