# HW LOG

CME241

Justin Lundgren, 06289145
justinlundgren@stanford.edu
Stanford University

# HW – Jan 11

- Write out the MP/MRP definitions and MRP Value Function definition (in LaTeX) in your own style/notation (so you really internalize these concepts)

- Think about the data structures/class design (in Python 3) to represent MP/MRP and implement them with clear type declarations

- Remember your data structure/code design must resemble the Mathematical/notational formalism as much as possible

- Specifically the data structure/code design of MRP should be incremental (and not independent) to that of MP

- Separately implement the $r(s, s')$ and the $R(s) = \sum_{s'} p(s, s') * r(s, s')$ definitions of MRP

- Write code to convert/cast the $r(s, s')$ definition of MRP to the $R(s)$ definition of MRP (put some thought into code design here)

- Write code to generate the stationary distribution for an MP

## MP/MRP definition

**MP:** A markov process is a chain that is memory less, i.e. it only cares about about the current state and not the past. The mathematical definition is

$$\mathbb{P}(S_{t+h} = s_{t+h} | S_0 = s_0, S_1 = s_1, \ldots, S_{t-1} = s_{t-1}, S_t = s_t) = \mathbb{P}(S_{t+h} = s_{t+h} | S_t = s_t)$$

or equivalently

$$\mathbb{E}[S_{t+h} | S_0 = s_0, S_1 = s_1, \ldots, S_{t-1} = s_{t-1}, S_t = s_t] = \mathbb{E}[S_{t+h} | S_t = s_t].$$

The Markov process is defined as $\{s, P_s\}$ where $s \in \{s_0, \ldots, s_k\}$ is the state spaces and $P_s$ is the probability distribution in each state.

**MRP:** A Markov reward process is a Markov process that has a reward $R(s)$ associated with each state and some discounting factor $\gamma \in [0, 1]$.

**Value function:** The value function is the accumulated expected reward associated with the current known state $s$. It is defines as

$$v(s) = \mathbb{E}[\sum_{i=0}^{T} R(s_{t+i}) \gamma^i | S_t = s],$$

where $T$ is the time of termination for the process.

## Data structures

- $State$ – TypeVar('State')

- $States$ – List[State]

- $\mathcal{R}(s)$ – List[float] (*)

- $r(ss')$ – List[List[float]] (*)

- $P_{MP}$ – Dict[State,Tuple[State,float]]

- $P_{MRP_A}$ – Dict[$P_{MP}$,float]

- $P_{MRP_B}$ – Dict[State,Dict[State,Tuple[float,float]]]

- $\gamma$ – float.

Thus we see that

$$
\begin{aligned}
\mathcal{R}(s) &= \mathbb{E}[R_t | S_{t-1} = s] \\
&= \sum_{s'} R_t(\{\texttt{reward after state } s'\}) \mathbb{P}(S_t = s' | S_{t-1} = s) \\
&= \sum_{s'} \mathbb{E}[R_t | S_{t-1} = s \ \cap \ S_t = s'] \mathbb{P}(S_t = s' | S_{t-1} = s) \\
&= \sum_{s'} r(s, s') p(s, s')
\end{aligned}
$$

# HW – Jan 16

> - Write the Bellman equation for MRP Value Function and code to calculate MRP Value Function (based on Matrix inversion method you learnt in this lecture)
>
> - Write out the MDP definition, Policy definition and MDP Value Function definition (in LaTeX) in your own style/notation (so you really internalize these concepts)
>
> - Think about the data structure/class design (in Python 3) to represent MDP, Policy, Value Function, and implement them with clear type definitions
>
> - The data structure/code design of MDP should be incremental (and not independent) to that of MRP
>
> - Separately implement the $r(s, s', a)$ and $R(s, a) = \sum_{s'} p(s, s', a) * r(s, s', a)$ definitions of MDP
>
> - Write code to convert/cast the $r(s, s', a)$ definition of MDP to the $R(s, a)$ definition of MDP (put some thought into code design here)
>
> - Write code to create a MRP given a MDP and a Policy
>
> - Write out all 8 MDP Bellman Equations and also the transformation from Optimal Action-Value function to Optimal Policy (in LaTeX)

## Data structures

- *Action* – TypeVar('Action')

- *Policy* – Dict[State,Tuple[Action,(float or int)]]]

- $MDP_A$ – Dict[State,Dict[Action,Dict[Tuple[State,(float or int)],(float or int)]]]]

- $MDP_B$ – Dict[State,Dict[Action,Tuple[State,Tuple[float,float]]]]

## Bellman Equations

(1) Basic Bellman for MRP (A)

$$v(s) = \mathbb{E}[\sum_{i=0}^{T} R_{t+i+1}\gamma^i \big| S_t = s]$$
$$= \mathbb{E}[R_{t+1}\big| S_t = s] + \gamma\mathbb{E}[v(S_{t+1})\big| S_t = s]$$
$$= \mathcal{R}_s + \gamma \sum_{s'} v(s')\mathbb{P}(S_{t+1} = s'\big| S_t = s).$$

(2) Basic Bellman for MRP (A) in matrix form is then

$$v = \mathcal{R} + \gamma\mathcal{P}v.$$

(3) For the action-value function with policy $\pi$ we have

$$q_\pi(s, a) = \mathbb{E}[\sum_{i=0}^{T} R_{t+i+1}\gamma^i \big| S_t = s \ \cap \ A_t = a]$$

which have the same solution in

$$q_\pi(s, a) = \mathbb{E}[R_{t+1}\big| S_t = s \ \cap \ A_t = a] + \gamma\mathbb{E}[q_\pi(S_{t+1}, A_{t+1})\big| S_t = s \ \cap \ A_t = a]$$
$$= \mathcal{R}_s^a + \gamma \sum_{s', a'} q_\pi(s', a')\mathbb{P}(S_{t+1} = s' \ \cap \ A_{t+1} = a'\big| S_t = s \ \cap \ A_t = a)$$

where $\mathcal{R}_s^a$ is $\mathcal{R}_s$ for action $a$.

(4) Likewise for a MDP with a policy $\pi$ we can create a value MRP with value function

$$v_\pi(s) = \sum_{a'} \pi(a'|s) q_\pi(s, a')$$

where $\pi(a'|s)$ is the probability of taking action $a'$ in state $s$.

(5) Now, we can combine (3) and (4) to express $v_\pi(s)$ as

$$v_\pi(s) = \sum_{a'} \pi(a'|s) \Big( \mathcal{R}_s^a + \gamma \sum_{s'} \mathbb{P}(S_t = s'|S_t = s \cap A_t = a') v_\pi(s') \Big)$$

(6) Combining (4) and (5) we can express $q_\pi(s, a)$ as

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s'} \mathbb{P}(S_{t+1} = s'|S_t = s \cap A_t = a) \sum_{a'} \pi(a'|s') q_\pi(s', a').$$

# HW – Jan 18

- Write code for Policy Evaluation (tabular) algorithm

- Write code for Policy Iteration (tabular) algorithm

- Write code for Value Iteration (tabular) algorithm

- Those familiar with function approximation (deep networks, or simply linear in featues) can try writing code for the above algorithms with function approximation (a.k.a. Approximate DP)

The first three are pretty much done (at least for MDP_A)

# HW – Jan 23

> - Work out (in LaTeX) the equations for Absolute/Relative Risk Premia for CARA/CRRA respectively
>
> - Write the solutions to Portfolio Applications covered in class with precise notation (in LaTeX)

## CARA

For CARA we have

$$U(x) = -\frac{1}{a}e^{-ax}, \ a \neq 0.$$

Thus we have

$$\frac{dU(x)}{dx} = e^{-ax} \text{ and } \frac{d^2U(x)}{dx^2} = -ae^{-ax}.$$

For the Arrow-Pratt risk aversion coefficient $A$ we have

$$A = -\frac{U''(x)}{U'(x)}$$
$$= a.$$

## CRRA

For CRRA we have

$$U(x) = \frac{x^{1-\gamma}}{1-\gamma} \ \gamma \neq 1.$$

Thus we have

$$\frac{dU(x)}{dx} = x^{-\gamma} \text{ and } \frac{d^2U(x)}{dx^2} = -\gamma x^{-\gamma-1}$$

For the relative Arrow-Pratt risk aversion coefficient $A$ we have

$$A = -\frac{xU''(x)}{U'(x)}$$
$$= \gamma.$$

## Portfolio Application Solution

**CARA:** We have two assets $r_a \sim N(\mu, \sigma^2)$ and $r_f \sim N(r, 0)$. We invest a fraction $\rho_a$ in $r_a$ and $\rho_f$ in $r_f$. The objective is then to

$$\max \ \mathbb{E}[U(N(\rho_a\mu + \rho_f r, \rho_r^2\sigma^2)]$$
$$\text{s.t. } \rho_a + \rho_f = 1.$$

Now, substituting $\rho_f = 1 - \rho_a$ and using the PDF of the normal distribution we can set this up as

$$\max_{\rho_a}\left\{ -\frac{1}{a}\int_{\mathbb{R}} \exp(-ax)\frac{1}{\sqrt{2\pi\rho_a\sigma^2}} \exp\big(\frac{(x-(\rho_a\mu+(1-\rho_a)r))^2}{2\rho_a^2\sigma^2}\big)\right\}. \tag{1}$$

Differentiating (1) wrt $\rho_a$ and setting to zero gives

$$\rho_a^* = \frac{\mu - r}{a\sigma^2}.$$

**CRRA:** The setup is very similar but now we assume that $\log(r_a) \sim N(\mu, \sigma^2)$ instead. This gives the solution

$$\rho_a^* = \frac{\mu - r}{\gamma\sigma^2}$$

as optimal allocation.

# HW – Jan 25

- Model Merton's Portfolio problem as an MDP (write the model in LaTeX)

- Implement this MDP model in code

- Try recovering the closed-form solution with a DP algorithm that you implemented previously