# MDPREFINED → MRP

## ASHWIN RAO (STANFORD CME 241)

We start with a MDPRefined given by transitions function $\mathcal{P}_{s,s'}^a$ denoting the probability of transition from state $s$ to state $s'$ upon taking action $a$ and given by rewards function $\mathcal{R}_{s's'}^a$ denoting the reward obtained when transitioning from state $s$ to state $s'$ upon taking action $a$.

First we construct an MDP from $(\mathcal{P}, \mathcal{R})$. The reward function of this MDP will be:

$$\sum_{s'} \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a$$

Next we construct an MRPRefined from $(\mathcal{P}, \mathcal{R})$ and a policy $\pi(a|s)$. The transitions function of this MRPRefined will be:

$$\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a$$

The reward function of this MRPRefined will be:

$$\frac{\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a}{\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a}$$

Next we construct MRP1 from the above MDP. The reward function of this MRP1 will be:

$$\sum_{a} \pi(a|s) \cdot (\sum_{s'} \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a) = \sum_{a} \sum_{s'} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a = \sum_{s'} \sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a$$

Next we construct MRP2 from the above MRPRefined. The reward function of this MRP2 will be:

$$\sum_{s'} (\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a)(\frac{\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a}{\sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a}) = \sum_{s'} \sum_{a} \pi(a|s) \cdot \mathcal{P}_{s,s'}^a \cdot \mathcal{R}_{s,s'}^a$$

Hence, MRP1 and MRP2 are the same MRP.