# Lecture 06

## * Policy Evaluation

Given a policy $\pi = \{ \mu_0, \ldots, \mu_{N-1} \}$, the value function of $\pi$ is defined as

$$J^{\pi}(x_0) = \underset{w_k \,;\, k=0,1,\ldots,N-1}{\mathbb{E}} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right]$$

So $J^{\pi} : S_0 \longrightarrow \mathbb{R}$ takes the total cost if policy $\pi$ is applied at initial state $x_0$.

Corollary 1 : For every initial state $x_0$, the value of policy $\pi$, $J^{\pi}(x_0)$ of the basic problem is equal to $J_0(x_0)$ when the function is given by the last step of the following algorithm, which proceeds backward in time from period $N-1$ to period $0$.

$$J_N(x_N) = g_N(x_N)$$

$$J_k(x_k) = \underset{w_k}{\mathbb{E}} \left[ g_k(x_k, \mu_k(x_k), w_k) + J_{k+1}\big(f_k(x_k, \mu_k(x_k), w_k)\big) \right]$$

$$k = 0, 1, \ldots, N-1$$

Proof : Just expand.

— This is policy evaluation algorithm.