## 2.2.4 Reduction between control policies classes

We first show a reduction from a general history dependent policies to Randomized Markovian policies. The main observation is that the only influence on the cumulative cost is the expected instantaneous cost $\mathbb{E}[c_t(s_t, a_t)]$. Namely, let

$$\rho_t^\pi(s, a) = \Pr_{h'_{t-1}}[a_t = a, s_t = s] = \mathbb{E}_{h'_{t-1}}[\mathbb{I}[s_t = s, a_t = a]|h'_{t-1}],$$

where $h'_{t-1} = (s_0, a_0, \ldots, s_{t-1}, a_{t-1})$ is the history of the first $t-1$ time steps generated using $\pi$, and the probability and expectation are taken with respect to the randomness of the policy $\pi$. Now we can rewrite the expected cost to go as,

$$\mathbb{E}[\mathcal{C}^\pi(s_0)] = \mathbb{E}[\sum_{t=1}^{T-1} \sum_{a \in \mathcal{A}_t, s \in \mathcal{S}_t} c_t(s, a)\rho_t^\pi(s, a)],$$

where $\mathcal{C}^\pi(s_0)$ is the random variable of the cost when starting at state $s_0$ and following policy $\pi$.

This implies that any two policies $\pi$ and $\pi'$ for which $\rho_t^\pi(s, a) = \rho_t^{\pi'}(s, a)$, for any time $t$, state $s$ and action $a$, would have the same expected cumulative cost for any cost function, i.e., $\mathbb{E}[\mathcal{C}^\pi(s_0)] = \mathbb{E}[\mathcal{C}^{\pi'}(s_0)]$