1.3 a) assuming $\ell(s, \pi)$ is 0-1 loss and

cost function $C$ bounded in $[0,1]$

Let $\mathbb{E}_{s \sim d_{\pi^*}}[\ell(s, \pi)] = \epsilon$ $\quad \to \Delta C(s) \quad \nearrow C(s, \pi^*) \quad \nearrow C(s, \pi)$

\* think of $\ell(s, \pi)$ as <u>difference in cost</u> suffered by optimal policy vs. learned policy given a state

\* Then $\Delta C$ can be at most 1 b/c $C$ is bounded in $[0,1]$, so the absolute difference b/w two values that are a result of $C(\cdot)$ is $|0-1|$ or $|1-0|$ at the most

\* if $\ell(s, \pi)$ is the loss in cost suffered from not following $\pi^*$ then $\ell(s, \pi)$ is also $= \Delta C(s)$

\* since $s \sim d_{\pi^*}$, $\Delta C(s)$ can be expected to equal $\mathbb{E}_{s \sim d_{\pi^*}}[\ell(s, \pi)] = \epsilon$

First we have at a state by state basis:

$$\Delta C(s) = C(s, \pi) - C(s, \pi^*)$$

We expect these values to be

$$\mathbb{E}_{s \sim d_\pi^*}[\Delta C(s)] = \mathbb{E}_{s \sim d_\pi^*}[C(s, \pi) - C(s, \pi^*)] = \mathbb{E}_{s \sim d_\pi^*}[C(s, \pi)] - \mathbb{E}_{s \sim d_\pi^*}[C(s, \pi^*)]$$

$\downarrow \epsilon$

we rearrange to get

$$\mathbb{E}_{s \sim d_\pi^*}[C(s, \pi)] = \mathbb{E}_{s \sim d_\pi^*}[C(s, \pi^*)] + \epsilon$$

Over $T$ time steps, the expected cost from following a policy accumulates to $J(\cdot)$

↳ $J(\pi) = J(\pi^*) + T\epsilon$ $\quad \swarrow$ loss is accumulated at each step.

Since $\epsilon$ is at most $= 1$, at every step, we can accumulate a total loss of $T \cdot 1$ at most so:

$$J(\pi) \leq J(\pi^*) + T^2 \epsilon$$