

Q. 2.1:

We denote \bar{V}^π the old value function. Then we have by definition:

$$\bar{V}^\pi = R^\pi + \gamma P^\pi \bar{V}^\pi \quad (1)$$

$$\bar{V}^\pi = (I - \gamma P^\pi)^{-1} R^\pi \quad (2)$$

We now apply affine transformation to reward function:

$$R_{\text{new}}^\pi = \alpha R^\pi + \beta \mathbb{1} \quad (3)$$

We plug (3) into (2):

$$V^\pi = \alpha R^\pi + \beta \mathbb{1} + \gamma P^\pi V^\pi$$

$$V^\pi - \gamma P^\pi V^\pi = \alpha R^\pi + \beta \mathbb{1}$$

$$\Rightarrow V^\pi = (I - \gamma P^\pi)^{-1} (\alpha R^\pi + \beta \mathbb{1})$$

$$\Rightarrow V^\pi = (I - \gamma P^\pi)^{-1} \alpha R^\pi + (I - \gamma P^\pi)^{-1} \beta \mathbb{1}$$

$$V^\pi = \alpha \bar{V}^\pi + (I - \gamma P^\pi)^{-1} \beta \mathbb{1}$$

So essentially our new value function $V^\pi = \alpha \bar{V}^\pi + c$, where c is a constant

Optimal policy is not preserved. We give a counter-example with affine transformation where $\alpha = -1$, $\beta = 0$. π^* is not optimal for V^π obviously.