

HW 4

① So that y^π is an equilibrium point, we must prove that,

$$\dot{y}^\pi = 0$$

$$\begin{aligned}\dot{y}^\pi &= \mathbb{E}_\pi \left[c_t + (\gamma - 1) y^\pi \right] = \\ &= \mathbb{E}_\pi \left[c_t + (\gamma - 1) \frac{c_\pi}{1 - \gamma} \right] =\end{aligned}$$

$$= \mathbb{E}_\pi [c_t - c_\pi] = c_\pi - c_\pi = 0$$

↓

In average, policy π will lead to cost c_π !
(definition)

$$\textcircled{2} \quad E = \underbrace{(y^t - y^\pi)}_{> 0} (\dot{y}^t - \dot{y}^\pi)$$

Since the cost to go for an arbitrary policy must necessarily be higher than for the optimal policy.

On the other hand,

$$\dot{y}^t - \dot{y}^\pi = \mathbb{E}_\pi [C_t + (r-1)y^t]$$

$$- C_\pi - (r-1)y^\pi =$$

$$= \cancel{y^t} + (r-1)y^t - \cancel{C_\pi} - (r-1)y^\pi =$$

↳ for the same reason as before

$$= \underbrace{(r-1)}_{\leq 0} (\underbrace{y^t - y^\pi}_{> 0})$$

$$\Rightarrow \dot{y}^t - \dot{y}^\pi < 0 \Rightarrow$$

$$\Rightarrow \dot{E} < 0$$

Along all trajectories of the ode described above.

③ If the trajectories of the TD algorithm follows indeed the ode described in this exercise, we know that the energy will decrease in each iteration, therefore our cost to go will monotonically converge to the optimal one.