

8.2.2 확률적 LQR

확률적 LQR 문제에서는 동역학 모델이 다음과 같이 선형 확률 동적 시스템으로 주어진다.

$$\begin{aligned} x_{t+1} &= A_t x_t + B_t u_t + f_{ct} + n, t = 0, \dots, T \\ &= F_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + f_{ct} + n_t \end{aligned}$$

여기서 $F_t = [A_t, B_t]$ 이고 A_t, B_t 는 확정된 행렬.

n_t 는 프로세스 노이즈로 평균이 0이고 분산이 Σ_t 인 가우시안 분포를 갖는 화이트 시퀀스 즉,

$$p(n_t) = N(0, \Sigma_t), \quad E[n_i n_j^T] = \Sigma_i \delta_{ij}$$

여기서 δ_{ij} 는 크로네커 델타 함수

$$\delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

프로세스 노이즈가 랜덤 시퀀스이기 때문에 상태변수 x_t 는 랜덤 시퀀스가 됨

초기값 x_0 도 n_t 와 독립이며 가우시안 분포를 갖는 것으로 가정

$$p(x_0) = N(\bar{x}_0, X_0), \quad E[(x_0 - \bar{x}_0) n_i^T] = 0$$

선형 확률 동적 시스템의 상태전이 확률밀도함수는 다음과 같이 전파된다.

$$p(x_{t+1} | x_t, u_t) = N(A_t x_t + B_t u_t + f_{ct}, \Sigma_t)$$

확률적 시스템에서는 무작위적 특성 때문에 비용함수의 값도 랜덤 변수 \rightarrow 목적함수를 항상 최소화할 수 있는 제어 또는 행동 시퀀스를 구할 수 없음 \rightarrow 비용함수의 기대값을 최소로 만드는 제어 시퀀스를 구하는 것을 목적으로 한다.

$$J_0 = E_{\tau \sim p(\tau)} \left[\sum_{t=0}^T c(x_t, u_t) \right]$$

확정적 시스템에서는 최적제어의 시퀀스를 구했다면 확률적 시스템에서는 현재 상태 변수가 주어진 조건에서 최적 정책 $\pi_t(u_t | x_t)$ 의 시퀀스 $(\pi_0, \pi_1, \dots, \pi_T)$ 를 구하는 것이 목적.

정책이 시변함수이기 때문에 목적함수는 다음과 같이 시변 상태가치 함수가 된다.

$$V(x_0) = E_{\tau_{u_0} \sim p(\tau_{u_0} | x_0)} \left[\sum_{t=0}^T c(x_t, u_t) \right]$$

$$J_0 = E_{x_0 \sim p(x_0)} [V(x_0)]$$

DP 를 적용하기 위해 상태 가치 함수를 시간 스텝 t 를 기준으로 전개

$$\begin{aligned}
 V(x_t) &= E_{\tau_{0:t} \sim p(\tau_{0:t}|x_t)} \left[\sum_{k=t}^T c(x_k, u_k) \right] \\
 &= \int_{\tau_{0:t}} \left(\sum_{k=t}^T c(x_k, u_k) \right) p(\tau_{0:t}|x_t) d\tau_{0:t} \\
 &= \int_{u_t} \left[\int_{\tau_{0:t+1}} \left(\sum_{k=t}^T c(x_k, u_k) \right) p(\tau_{0:t+1}|x_t, u_t) d\tau_{0:t+1} \right] \pi(u_t|x_t) du_t \\
 &= E_{u_t \sim \pi(u_t|x_t)} [Q(x_t, u_t)]
 \end{aligned} \tag{8.41}$$

대괄호 항은 어떤 상태변수 x_t 에서 제어 u_t 를 선택하고 그로부터 정책 π 로 기대할 수 있는 미래 비용의 기대값이므로 시변 행동 가치 함수 $Q(x_t, u_t)$ 가 된다.

행동 가치 함수 $Q(x_t, u_t)$ 를 한 스텝 더 전개 한다.

$$\begin{aligned}
 Q(x_t, u_t) &= \int_{\tau_{0:t+1}} c(x_t, u_t) p(\tau_{0:t+1}|x_t, u_t) d\tau_{0:t+1} \\
 &\quad + \int_{\tau_{0:t+1}} \left(\sum_{k=t+1}^T c(x_k, u_k) \right) p(\tau_{0:t+1}|x_t, u_t) d\tau_{0:t+1} \\
 &= c(x_t, u_t) + Q_1
 \end{aligned} \tag{8.42}$$

Q_1 을 정리해 보면

$$\begin{aligned}
 Q_1 &= \int_{x_{t+1}} \int_{\tau_{0:t+1}} \left(\sum_{k=t+1}^T c(x_k, u_k) \right) p(\tau_{0:t+1}|x_t, u_t, x_{t+1}) p(x_{t+1}|x_t, u_t) d\tau_{0:t+1} dx_{t+1} \\
 &= \int_{x_{t+1}} \left[\int_{\tau_{0:t+1}} \left(\sum_{k=t+1}^T c(x_k, u_k) \right) p(\tau_{0:t+1}|x_{t+1}) d\tau_{0:t+1} \right] p(x_{t+1}|x_t, u_t) dx_{t+1}
 \end{aligned} \tag{8.43}$$

위의 식에서 대괄호는 $V(x_{t+1})$ 이므로

$$Q_1 = \int_{x_{t+1}} V(x_{t+1}) p(x_{t+1}|x_t, u_t) dx_{t+1}$$

이다. 따라서 행동가치 함수 $Q(x_t, u_t)$ 는 다음과 같이 된다.

$$\begin{aligned}
 Q(x_t, u_t) &= c(x_t, u_t) + \int_{x_{t+1}} V(x_{t+1}) p(x_{t+1}|x_t, u_t) dx_{t+1} \\
 &= c(x_t, u_t) + E_{x_{t+1} \sim p(x_{t+1}|x_t, u_t)} [V(x_{t+1})]
 \end{aligned} \tag{8.44}$$

위의 식을 상태 가치 식 8.41에 대입하면

$$V(x_t) = E_{u_t \sim \pi_t(u_t|x_t)} [c(x_t, u_t) + E_{x_{t+1} \sim p(x_{t+1}|x_t, u_t)} [V(x_{t+1})]]$$

확정적 정책 $u_t = \pi(x_t)$ 를 가정하면

$$V(x_t) = c(x_t, u_t) + E_{x_{t+1} \sim p(x_{t+1}|x_t, u_t)} [V(x_{t+1})]$$

벨만의 최적성 원리에 의하면 현재 시간스텝 t 의 최적 제어는 다음을 만족해야 한다.

$$V^*(x_t) = \min_{u_t} ((x_t, u_t) + E_{x_{t+1} \sim p(x_{t+1}|x_t, u_t)} [V(x_{t+1})])$$

확률적 LQR문제에서는 비용함수가 다음과 같이 2차 함수로 주어진다.

$$c(x_t, u_t) = \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T C_T \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T c_t$$

최종 시간 $t = T$ 에서의 목적함수는

$$\begin{aligned} V(x_T) &= \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T C_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T c_T \\ &= \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T Q_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T q_T \\ &= Q(x_T, u_T) \end{aligned} \tag{8.50}$$

여기서 $Q_T = C_T$, $q_T = c_T$ 로 놓았다. 위의 목적함수를 최소로 만드는 최적제어를 구하기 위해 다음과 같이 미분한다.

$$\begin{aligned} 0 &= \frac{\partial V(x_T)}{\partial u_T} = \frac{\partial Q(x_T, u_T)}{\partial u_T} \\ &= Q_{uxT} x_T + Q_{uuT} u_T + Q_{uT} \end{aligned} \tag{8.51}$$

여기서

$$\begin{aligned} Q_T &= \begin{bmatrix} Q_{xxT} & Q_{xuT} \\ Q_{uxT} & Q_{uuT} \end{bmatrix} \\ q_T &= \begin{bmatrix} Q_{xT} \\ Q_{uT} \end{bmatrix} \end{aligned}$$

이다. 그러면 u'_T 는

$$u'_T = -Q_{uuT}^{-1} Q_{uxT} x_T - Q_{uuT}^{-1} Q_{uT}$$

칼만 게인을 다음과 같이 정의를 하면

$$\begin{aligned} K_T &= -Q_{uuT}^{-1} Q_{uxT} \\ k_T &= -Q_{uuT}^{-1} Q_{uT} \end{aligned} \tag{8.54}$$

$t = T$ 에서의 최적제어는 다음과 같다.

$$u'_T = K_T x_T + k_T$$

위의 식을 이용해 $t = T$ 에서의 최소 목적함수의 값을 구하면 다음과 같다.

$$\begin{aligned} V(x_T) &= \frac{1}{2} \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix}^T Q_T \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix} + \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix}^T q_T \\ &= \frac{1}{2} x_T^T V_T x_T + x_T^T v_T + const \end{aligned} \quad (8.56)$$

여기서 const는

$$\begin{aligned} V_T &= Q_{xxT} + Q_{xuT} K_T + K_T^T Q_{uxT} + K_T^T Q_{uuT} K_T \\ v_T &= Q_{xT} + K_T^T Q_{uT} + Q_{xuT} k_T + K_T^T Q_{uuT} k_T \end{aligned} \quad (8.57)$$

다음으로 역방향 시간시스템 $t = T - 1$ 에서의 목적함수는 다음과 같다.

$$\begin{aligned} V(x_{T-1}) &= \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T C_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T c_{T-1} \\ &\quad + E_{xT \sim p(x_T | x_{T-1}, u_{T-1})} \left[-\frac{1}{2} x_T^T V_T x_T + x_T^T v_T + const \right] \\ &= \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T (C_{T-1} + F_{T-1}^T V_T F_{T-1}) \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} \\ &\quad + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T (c_{T-1} + F_{T-1}^T V_T f_{eT-1} + F_{T-1}^T v_T) + const \\ &\quad + E_{xT \sim p(x_T | x_{T-1}, u_{T-1})} \left[-\frac{1}{2} n_{T-1}^T V_T n_{T-1} \right] \end{aligned} \quad (8.58)$$

$$\begin{aligned} E_{xT \sim p(x_T | x_{T-1}, u_{T-1})} \left[-\frac{1}{2} n_{T-1}^T V_T n_{T-1} \right] &= E_{xT \sim p(x_T | x_{T-1}, u_{T-1})} \left[-\frac{1}{2} tr(V_T n_{T-1} n_{T-1}^T) \right] \\ &= -\frac{1}{2} tr(V_T E_{xT \sim p(x_T | x_{T-1}, u_{T-1})} [n_{T-1} n_{T-1}^T]) \\ &= -\frac{1}{2} tr(V_T \Sigma_{T-1}) \end{aligned} \quad (8.59)$$

여기서 식 8.59에 의하면 또 다른 상수항을 추가한 것에 불가하므로 목적함수는 다음과 같이 된다.

$$\begin{aligned} V(x_{T-1}) &= \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T (C_{T-1} + F_{T-1}^T V_T F_{T-1}) \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} \\ &\quad + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T (c_{T-1} + F_{T-1}^T V_T f_{eT-1} + F_{T-1}^T v_T) + const \\ &= \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T Q_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T q_{T-1} + const \\ &= Q(x_{T-1}, u_{T-1}) \end{aligned} \quad (8.60)$$

$$Q_{T-1} = C_{T-1} + F_{T-1}^T V_T F_{T-1} \quad [8.61]$$

$$q_{T-1} = c_{T-1} + F_{T-1}^T V_T f_{cT-1} + F_{T-1}^T v_T$$

목적함수를 최소로 만드는 최적제어를 구하기 위해 다음과 같이 미분을 수행한다.

$$\begin{aligned} 0 &= \frac{\partial V(x_{T-1})}{\partial u_{T-1}} \\ &= Q_{uxT-1} x_{T-1} + Q_{uuT-1} u_{T-1} + Q_{uT-1} \end{aligned} \quad [8.62]$$

$$\begin{aligned} Q_{T-1} &= \begin{bmatrix} Q_{xxT-1} & Q_{xuT-1} \\ Q_{uxT-1} & Q_{uuT-1} \end{bmatrix} \\ q_{T-1} &= \begin{bmatrix} Q_{xT-1} \\ Q_{uT-1} \end{bmatrix} \end{aligned}$$

$$u'_{T-1} = -Q_{uuT-1}^{-1} Q_{uxT-1} x_{T-1} = Q_{uuT-1}^{-1} Q_{uT-1}$$

칼만 게임을 다음과 같이 정의 하면!

$$K_{T-1} = -Q_{uuT-1}^{-1} Q_{uxT-1}$$

$$k_{T-1} = -Q_{uuT-1}^{-1} Q_{uT-1}$$

$t = T - 1$ 에서의 최적제어는

$$u'_{T-1} = K_{T-1} x_{T-1} + k_{T-1}$$

위 식을 이용해 $t = T - 1$ 에서의 최소 목적함수 값을 구하면 다음과 같다.

$$\begin{aligned} V(x_{T-1}) &= \frac{1}{2} \begin{bmatrix} x_{T-1} \\ K_{T-1} x_{T-1} + k_{T-1} \end{bmatrix}^T Q_{T-1} \begin{bmatrix} x_{T-1} \\ K_{T-1} x_{T-1} + k_{T-1} \end{bmatrix} \\ &\quad + \begin{bmatrix} x_{T-1} \\ K_{T-1} x_{T-1} + k_{T-1} \end{bmatrix}^T q_{T-1} + const \\ &= \frac{1}{2} x_{T-1}^T V_{T-1} x_{T-1} + x_{T-1}^T v_{T-1} + const \end{aligned}$$

$$V_{T-1} = Q_{xxT-1} + Q_{xuT-1} K_{T-1} + K_{T-1}^T Q_{uxT-1} + K_{T-1}^T Q_{uuT-1} K_{T-1}$$

$$v_{T-1} = Q_{xT-1} + K_{T-1}^T Q_{uT-1} + Q_{xuT-1} k_{T-1} + K_{T-1}^T Q_{uuT-1} k_{T-1}$$

확정적 LQR과 비교해 보면 상후상ㅇ에 또다른 상수항 $\frac{1}{2}tr(V, \Sigma_{t-1})$ 만 추가 됐고 나머지는 동일하다는 것을 알 수 있다.

결론

확률적 LQR 업데이트 식은 프로세스 노이즈의 공분산 Σ_t 와 초기 상태 변수의 공분산 $\backslash \mathbf{X}_0$ 와는 무관하며 확정적 LQR 식과 동일하다는 것을 알 수 있다.

8.2.3 가우시안 LQR

- 확률적 LQR의 확정적 정책을 $\pi(u_t|x_t) = N(x_t|\mu_t, S_t)$ 와 같은 확률밀도 함수를 갖는 확률적 정책으로 확장 시킨 것
- 동일한 x_t 에 대해 다른 제어 μ_t 가 선택될 수 있다.
- 목적함수

$$J_0 = E_{\tau \sim p(\tau)} \left[\sum_{t=0}^T (c(x_t, u_t) + \log \pi(u_t|x_t)) \right]$$

위의 식을 전개 하면

$$\begin{aligned} J_0 &= E_{\tau \sim p(\tau)} \left[\sum_{t=0}^T c(x_t, u_t) \right] + \sum_{t=0}^T E_{u_t \sim p(x_t), u_t \sim \pi(u_t|x_t)} [\log \pi(u_t|x_t)] \\ &= E_{\tau \sim p(\tau)} \left[\sum_{t=0}^T c(x_t, u_t) \right] - \sum_{t=0}^T E_{u_t \sim p(x_t)} [H(\pi(u_t|x_t))] \end{aligned} \quad (8.71)$$

$$H(\pi(u_t|x_t)) = - \int \pi(u_t|x_t) \log \pi(u_t|x_t) du_t$$

확률적 LQR의 목적함수에 **한계확률밀도함수 $p(x_t)$ 에 대한 정책 엔트로피 기대값** 도입됨

기존 목적함수에 확률적 정책의 엔트로피가 추가함으로써 가우시안 LQR 제어기는 **비용을 최소화함과 동시에 정책의 무작위성을 최대화** 하는 특성이 있다.

확률적 시스템에서는 상태 변수가 랜덤 변수. → 시간스텝 t 에서의 제어는 상태변수 x_t 를 기반으로 구축해야 한다는 추가 조건 필요 → 목적함수를 다음과 같이 초기 상태변수 x_0 에 대한 조건부 함수로 바꾼다.

$$V(x_0) = E_{\tau_{u_0} \sim p(\tau_{u_0}|x_0)} \left[\sum_{t=0}^T (c(x_t, u_t) + \log \pi(u_t|x_t)) \right]$$

새로운 목적함수는 상태가치 함수에 엔트로피를 추가한 것이므로 **소프트 상태가치 함수(soft state-value function)**이라 부른다.

원래 목적함수와 소프트 상태가치 함수와의 관계

$$J_0 = E_{x_0 \sim p(x_0)} [V(x_0)]$$

DP를 적용하기 위해 소프트 상태가치 함수를 시간스텝 t 를 기준으로 전개 해 보자. (식 8.75)

$$\begin{aligned}
V(x_t) &= E_{\tau_{u_t} \sim p(\tau_{u_t}|x_t)} \left[\sum_{k=t}^T (c(x_k, u_k) + \log \pi(u_k | x_k)) \right] \quad [8.75] \\
&= E_{\tau_{u_t} \sim p(\tau_{u_t}|x_t)} \left[\log \pi(u_t | x_t) + \sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_{k+1} | x_{k+1})) \right] \\
&= \int_{\tau_{u_t}} \log \pi(u_t | x_t) p(\tau_{u_t} | x_t) d\tau_{u_t} \\
&\quad + \int_{\tau_{u_t}} \left(\sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_{k+1} | x_{k+1})) \right) p(\tau_{u_t} | x_t) d\tau_{u_t} \\
&= \int_{u_t} \int_{\tau_{u_t}} \log \pi(u_t | x_t) p(\tau_{u_{t+1}} | x_t, u_t) \pi(u_t | x_t) d\tau_{u_{t+1}} du_t \\
&\quad + \int_{u_t} \left[\int_{\tau_{u_t}} \left(\sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_{k+1} | x_{k+1})) \right) p(\tau_{u_{t+1}} | x_t, u_t) d\tau_{u_{t+1}} \right] \pi(u_t | x_t) du_t \\
&= E_{u_t \sim \pi(u_t | x_t)} [\log \pi(u_t | x_t) + Q(x_t, u_t)]
\end{aligned}$$

여기서 $\tau_{u_t} = (u_t, x_{t+1}, u_{t+1}, \dots, x_T, u_T)$, $\tau_{x_{t+1}} = (x_{t+1}, u_{t+1}, \dots, x_T, u_T)$ 다. 정리를 하면,

$$V(x_t) = E_{u_t \sim \pi(u_t | x_t)} [\log \pi(u_t | x_t) + Q(x_t, u_t)]$$

$Q(x_t, u_t)$: 소프트 행동가치 함수

$$\begin{aligned}
Q(x_t, u_t) &= \int_{\tau_{u_t}} \left(c(x_t, u_t) + \sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_k | x_k)) \right) p(\tau_{u_{t+1}} | x_t, u_t) d\tau_{u_{t+1}} \quad [8.77] \\
&= c(x_t, u_t) + \int_{\tau_{u_t}} \left(\sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_k | x_k)) \right) p(\tau_{u_{t+1}} | x_t, u_t) d\tau_{u_{t+1}} \\
&= c(x_t, u_t) + Q_t
\end{aligned}$$

Q 를 정리해 보면,

$$\begin{aligned}
Q_t &= \int_{x_{t+1}} \int_{\tau_{u_t}} \left(\sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_k | x_k)) \right) p(\tau_{u_{t+1}} | x_t, u_t, x_{t+1}) p(x_{t+1} | x_t, u_t) d\tau_{u_{t+1}} dx_{t+1} \quad [8.78] \\
&= \int_{x_{t+1}} \left[\int_{\tau_{u_t}} \left(\sum_{k=t+1}^T (c(x_k, u_k) + \log \pi(u_k | x_k)) \right) p(\tau_{u_{t+1}} | x_{t+1}) d\tau_{u_{t+1}} \right] p(x_{t+1} | x_t, u_t) dx_{t+1}
\end{aligned}$$

대괄호는 $V(x_{t+1})$ 이므로

$$Q_1 = \int_{x_{t+1}} V(x_{t+1}) p(x_{t+1} | x_t, u_t) dx_{t+1}$$

따라서 소프트 행동가치 함수는 다음과 같이 된다.

$$\begin{aligned}
Q(x_t, u_t) &= c(x_t, u_t) + \int_{x_{t+1}} V(x_{t+1}) p(x_{t+1} | x_t, u_t) dx_{t+1} \quad [8.80] \\
&= c(x_t, u_t) + E_{x_{t+1} \sim p(x_{t+1} | x_t, u_t)} [V(x_{t+1})]
\end{aligned}$$

위의 식을 식(8.75)에 대입하면 다음과 같다.

$$V(x_t) = E_{u_t \sim \pi(u_t | x_t)} [\log \pi(u_t | x_t) + c(x_t, u_t) + E_{u_{t+1} \sim p(x_{t+1} | x_t, u_t)} [V(x_{t+1})]]$$

벨만의 최적성 원리에 의하면 현재 시간 스텝 t 에서 최종 시간스텝 T 까지 최소의 비용함수를 실현 하는 현재 시간스텝 t 의 최적제어는 다음 식을 만족해야 한다.

$$V(x_t) = \min_{\pi} E_{u_t \sim \pi(u_t|x_t)} [\log \pi(u_t|x_t) + c(x_t, u_t) + E_{u_{t+1} \sim p(x_{t+1}|x_t, u_t)} [V(x_{t+1})]]$$

최종 시간인 $t = T$ 에서는 소프트 상태 가치 함수는 다음과 같다.

$$\begin{aligned} V(x_T) &= E_{u_T \sim p(u_T|x_T)} \left[\log \pi(u_T|x_T) + \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T C_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T c_T \right] \\ &= E_{u_T \sim p(u_T|x_T)} \left[\log \pi(u_T|x_T) + \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T Q_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T q_T \right] \\ &= E_{u_T \sim p(u_T|x_T)} [\log \pi(u_T|x_T) + Q(x_T, u_T)] \end{aligned} \quad [8.83]$$

여기서 $Q_T = C_T$, $q_T = c_T$ 로 놓았다. 최적정책을 가우시안 확률밀도 함수로 가정했으므로 소프트 상태 가치 함수는 다음과 같이 전개 된다.

$$\begin{aligned} V(x_T) &= E_{u_T \sim \pi(u_T|x_T)} \left[-\frac{m}{2} \log(2\pi) - \frac{1}{2} \log(\det S_T) \right. \\ &\quad \left. - \frac{1}{2} (u_T - \mu_T)^T S_T^{-1} (u_T - \mu_T) \right. \\ &\quad \left. + \frac{1}{2} x_T^T Q_{xxT} x_T + x_T^T Q_{xuT} u_T + \frac{1}{2} u_T^T Q_{uuT} u_T \right. \\ &\quad \left. + x_T^T Q_{xT} + u_T^T Q_{uT} \right] \\ &= -\frac{m}{2} \log(2\pi) - \frac{1}{2} \log(\det S_T) - \frac{1}{2} \text{tr}(I) \\ &\quad + \frac{1}{2} x_T^T Q_{xxT} x_T + x_T^T Q_{xuT} \mu_T \\ &\quad + \frac{1}{2} \mu_T^T Q_{uuT} \mu_T + \frac{1}{2} \text{tr}(Q_{uuT} S_T) \\ &\quad + x_T^T Q_{xT} + \mu_T^T Q_{uT} \end{aligned} \quad [8.84]$$

여기서 m 은 u_T 의 차원이고 I 는 단위행렬이며

$$\begin{aligned} Q_T &= \begin{bmatrix} Q_{xxT} & Q_{xuT} \\ Q_{uxT} & Q_{uuT} \end{bmatrix} \\ q_T &= \begin{bmatrix} Q_{xT} \\ Q_{uT} \end{bmatrix} \end{aligned}$$

소프트 상태가치 함수를 최소로 만드는 최적 정책은 다음과 같이 가우시안 분포의 평균과 공분산에 대해서 각각 미분하여 구할 수 있다.

$$\frac{\partial V(x_T)}{\partial \mu_T} = Q_{uuT} \mu_T + Q_{uxT} x_T + Q_{uT} = 0 \quad [8.86]$$

$$\frac{\partial V(x_T)}{\partial S_T} = -\frac{1}{2} S_T^{-1} + \frac{1}{2} Q_{uuT} = 0$$

그러면 μ_T^* 와 S_T^* 는 다음과 같이 구해진다.

$$\begin{aligned}\dot{\mu}_T &= -Q_{uuT}^{-1} Q_{uxT} x_T - Q_{uuT}^{-1} Q_{uT} \\ \dot{S}_T &= Q_{uuT}^{-1}\end{aligned}\tag{8.87}$$

칼만 게인을 다음과 같이 정의를 하면

$$\begin{aligned}K_T &= -Q_{uuT}^{-1} Q_{uxT} \\ k_T &= -Q_{uuT}^{-1} Q_{uT}\end{aligned}\tag{8.88}$$

t=T에서의 최적 가우시안 LQR의 평균제어는 다음과 같이 쓸 수 있다.

$$\mu_T^* = K_T x_T + k_T$$

위 식을 이용해 $t = T$ 에서의 최소 소프트 상태가치 값을 구하면 다음과 같다.

$$\begin{aligned}V(x_T) &= \frac{1}{2} x_T^T Q_{xxT} x_T + x_T^T Q_{xuT} K_T x_T + x_T^T Q_{xuT} k_T \\ &\quad + \frac{1}{2} x_T^T K_T^T Q_{uuT} K_T x_T + x_T^T K_T^T Q_{uuT} k_T \\ &\quad + x_T^T Q_{xT} + x_T^T K_T^T Q_{uT} + const \\ &= \frac{1}{2} x_T^T V_T x_T + x_T^T v_T + const\end{aligned}\tag{8.90}$$

여기서

$$\begin{aligned}V_T &= Q_{xxT} + Q_{xuT} K_T + K_T^T Q_{uxT} + K_T^T Q_{uuT} K_T \\ v_T &= Q_{xT} + K_T^T Q_{uT} + Q_{xuT} k_T + K_T^T Q_{uuT} k_T\end{aligned}\tag{8.91}$$

다음으로 시간의 다음 역방향 단계인 $t = T - 1$ 에서는 소프트 상태가치 함수가 다음과 같다.

$$\begin{aligned}V(x_{T-1}) &= E_{u_{T-1} \sim \pi(x_{T-1}|x_{T-2})} \left[\log \pi(u_{T-1}|x_{T-1}) + \right. \\ &\quad \left. + \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T C_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T c_T \right. \\ &\quad \left. + \frac{1}{2} x_T^T V_T x_T + x_T^T v_T + const \right] \\ &= E_{u_{T-1} \sim \pi(x_{T-1}|x_{T-2})} \left[\log \pi(u_{T-1}|x_{T-1}) + \right. \\ &\quad \left. + \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T Q_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T q_{T-1} + const \right] \\ &= E_{u_{T-1} \sim \pi(x_{T-1}|x_{T-2})} \left[\log \pi(u_{T-1}|x_{T-1}) + Q(x_{T-1}, u_{T-1}) \right]\end{aligned}\tag{8.92}$$

여기서

$$Q_{T-1} = C_{T-1} + F_{T-1}^T V_T F_{T-1}$$

$$q_{T-1} = c_{T-1} + F_{T-1}^T V_T F_{T-1} + F_{T-1}^T v_T$$

최적정책이 가우시안 확률밀도 함수이므로 소프트 상태가치 함수는 다음과 같이 전개된다.

$$V(x_{T-1}) = E_{u_{T-1} \sim \pi(u_T|x_T)} \left[-\frac{1}{2} \log(\det S_{T-1}) - \frac{1}{2} (u_{T-1} - \mu_{T-1})^T S_{T-1}^{-1} (u_{T-1} - \mu_{T-1}) \right. \\ \left. + \frac{1}{2} x_{T-1}^T Q_{xxT-1} x_{T-1} + x_{T-1}^T Q_{xuT-1} u_{T-1} + \frac{1}{2} u_{T-1}^T Q_{uuT-1} u_{T-1} \right. \\ \left. + x_{T-1}^T Q_{xT-1} + u_{T-1}^T Q_{uT-1} + const \right] \quad [8.94]$$

$$= -\frac{1}{2} \log(\det S_{T-1}) - \frac{1}{2} tr(I) + \frac{1}{2} x_{T-1}^T Q_{xxT-1} x_{T-1} + x_{T-1}^T Q_{xuT-1} \mu_{T-1} \\ + \frac{1}{2} \mu_{T-1}^T Q_{uuT-1} \mu_{T-1} + \frac{1}{2} tr(Q_{uuT-1} S_{T-1}) \\ + x_{T-1}^T Q_{xT-1} + \mu_{T-1}^T Q_{uT-1} + const$$

여기서

$$Q_{T-1} = \begin{bmatrix} Q_{xxT-1} & Q_{xuT-1} \\ Q_{uxT-1} & Q_{uuT-1} \end{bmatrix}$$

$$q_{T-1} = \begin{bmatrix} Q_{xT-1} \\ Q_{uT-1} \end{bmatrix}$$

평균과 공분산에 대해서 각각 미분하여 계산하면

$$\frac{\partial V(x_{T-1})}{\partial \mu_{T-1}} = Q_{uuT-1} \mu_{T-1} + Q_{uxT-1} x_{T-1} + Q_{uT-1} = 0 \quad [8.96]$$

$$\frac{\partial V(x_{T-1})}{\partial S_{T-1}} = -\frac{1}{2} S_{T-1}^{-1} + \frac{1}{2} Q_{uuT-1} = 0$$

그러면 μ_{T-1}^* 와 S_{T-1}^* 는 다음과 같이 구해진다.

$$\dot{\mu}_{T-1} = -Q_{uuT-1}^{-1} Q_{uxT-1} x_{T-1} - Q_{uuT-1}^{-1} Q_{uT-1}$$

$$\dot{S}_{T-1} = Q_{uuT-1}^{-1} \quad [8.97]$$

칼만 게인을 다음과 같이 정의를 하면

$$K_{T-1} = -Q_{uuT-1}^{-1} Q_{uxT-1}$$

$$k_{T-1} = -Q_{uuT-1}^{-1} Q_{uT-1} \quad [8.98]$$

t=T-1에서의 최적 가우시안 LQR의 평균제어는 다음과 같이 쓸 수 있다.

$$\mu_{T-1}^* = K_{T-1} x_{T-1} + k_{T-1}$$

위의 식을 이용해 $t = T - 1$ 에서의 최소 소프트 상태가치 값을 구하면 다음과 같다.

$$\begin{aligned}
V'(x_{T-1}) &= \frac{1}{2} x_{T-1}^T Q_{xxT-1} x_{T-1} + x_{T-1}^T Q_{xuT-1} K_{T-1} x_{T-1} + x_{T-1}^T Q_{uxT-1} k_{T-1} \\
&\quad + \frac{1}{2} x_{T-1}^T K_{T-1}^T Q_{uuT-1} K_{T-1} x_{T-1} + x_{T-1}^T K_{T-1}^T Q_{ouT-1} k_{T-1} \\
&\quad + x_{T-1}^T Q_{sxT-1} + x_{T-1}^T K_{T-1}^T Q_{osT-1} + const \\
&= \frac{1}{2} x_{T-1}^T V_{T-1} x_{T-1} + x_{T-1}^T v_{T-1} + const
\end{aligned} \tag{8.100}$$

여기서

$$\begin{aligned}
V_{T-1} &= Q_{xxT-1} + Q_{xuT-1} K_{T-1} + K_{T-1}^T Q_{uxT-1} + K_{T-1}^T Q_{ouT-1} K_{T-1} \\
v_{T-1} &= Q_{sxT-1} + K_{T-1}^T Q_{ouT-1} + Q_{xuT-1} k_{T-1} + K_{T-1}^T Q_{ouT-1} k_{T-1}
\end{aligned} \tag{8.101}$$

위의 전개식을 확정적 LQR과 비교를 해보면

- 가우시안 정책이 산출하는 평균은 확정적 LQR의 제어값과 동일하고
- 공분산은 $S_t = Q^{-1} u u^T$ 임을 알수 있다.

$$\pi(u_t | x_t) = N(K_t x_t + k_t, Q^{-1} u u^T)$$

8.2.4 반복적 LQR (iLQR)

- 비선형 시스템에 LQR를 적용한 것
- LQR은 선형 시스템과 2차 함수로 된 비용함수에만 적용할 수 있다.
- iLQR은 현재의 궤적을 기준으로 비선형 시스템을 1차 시스템으로 근사하고 비용함수를 2차 함수로 근사한 후에 LQR을 적용하는 방법
- 궤적이 수렴할 때 까지 반복한다.

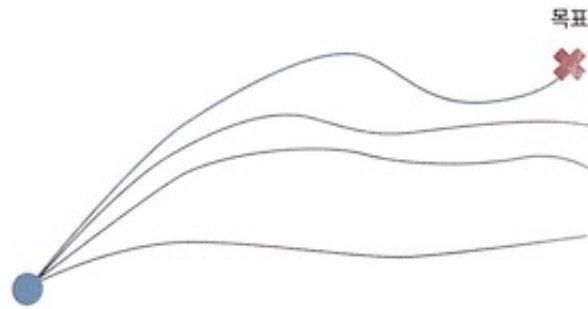


그림 8.3 반복적 LQR의 아이디어

다음과 같은 비선형 이산시간 시스템을 생각해 보자

$$x_{t+1} = f(x_t, u_t)$$

다음과 같은 목적 함수를 가정한다.

$$J_0 = \sum_{t=0}^T c(x_t, u_t)$$

테일러 시리즈를 이용해 명목(norminal) 궤적 (\hat{x}_t, \hat{u}_t) 를 기준으로 비선형 시스템을 선형화해 보자

$$\begin{aligned} x_{t+1} &= f(x_t, u_t) \\ &\approx f(\hat{x}_t, \hat{u}_t) + f_{x_t}(x_t - \hat{x}_t) + f_{u_t}(u_t - \hat{u}_t) \\ &= f_{x_t} x_t + f_{u_t} u_t + f(\hat{x}_t, \hat{u}_t) - f_{x_t} \hat{x}_t - f_{u_t} \hat{u}_t \\ &= f_{x_t} x_t + f_{u_t} u_t + f_{c_t} \end{aligned}$$

여기서 $f_{x_t} = \nabla_{x_t} f(\hat{x}_t, \hat{u}_t)$, $f_{u_t} = \nabla_{u_t} f(\hat{x}_t, \hat{u}_t)$, $f_{c_t} = f(\hat{x}_t, \hat{u}_t) - f_{x_t} \hat{x}_t - f_{u_t} \hat{u}_t$ 이다.

용어 설명 비선형 시스템을 테일러 시리즈를 이용해 선형화할 때 기준이 되는 궤적을 명목궤적이라고 한다. 실제 궤적은 명목궤적과 큰 차이가 나지 않는 주변에 있다고 가정해 실제 궤적 x_t 를 명목궤적 \bar{x}_t 와 섭동(perturbation) Δx_t 의 합으로 표현한다. 즉, 비선형 동적 시스템을 $y_t = h(x_t)$ 라고 할 때 명목궤적 \bar{x}_t 를 기준으로 전개하면 다음과 같다.

$$\begin{aligned} y_t &= h(x_t) = h(\bar{x}_t + \Delta x_t) \\ &= h(\bar{x}_t) + \left. \frac{dh}{dx} \right|_{x_t=\bar{x}_t} \Delta x_t + H.O.T. \\ &\approx h(\bar{x}_t) + \left. \frac{dh}{dx} \right|_{x_t=\bar{x}_t} \Delta x_t \end{aligned}$$

여기서 $H.O.T.$ 는 Δx_t 의 고차항을 나타내며, Δx_t 가 작다고 가정하고 무시한다. 또한 $\left. \frac{dh}{dx} \right|_{x_t=\bar{x}_t} = \nabla_{x_t} h(\bar{x}_t)$ 를 자코비안(Jacobian) 행렬이라고 한다.

비용함수도 2차함수로 근사해보자.

$$\begin{aligned} J_0 &\approx \sum_{t=0}^T \left(c_t(\hat{x}_t, \hat{u}_t) + b_t^T \begin{bmatrix} x_t - \hat{x}_t \\ u_t - \hat{u}_t \end{bmatrix} + \frac{1}{2} \begin{bmatrix} x_t - \hat{x}_t \\ u_t - \hat{u}_t \end{bmatrix}^T C_t \begin{bmatrix} x_t - \hat{x}_t \\ u_t - \hat{u}_t \end{bmatrix} \right) \\ &= \sum_{t=0}^T \left(\begin{bmatrix} x_t \\ u_t \end{bmatrix}^T b_t + \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T C_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} - \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix} \right) \\ &\quad + \sum_{t=0}^T \left(c_t(\hat{x}_t, \hat{u}_t) + \frac{1}{2} \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix} - \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T b_t \right) \\ &= \sum_{t=0}^T \left(\frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T C_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T c_t + d_t \right) \end{aligned}$$

여기서

$$\begin{aligned}
C_t &= \nabla_{\hat{x}_t, \hat{u}_t}^2 c_t(\hat{x}_t, \hat{u}_t), \quad b_t = \nabla_{\hat{x}_t, \hat{u}_t} c_t(\hat{x}_t, \hat{u}_t) \\
c_t &= b_t - C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}, \\
d_t &= c_t(\hat{x}_t, \hat{u}_t) + \frac{1}{2} \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix} - \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T b_t \\
&= c_t(\hat{x}_t, \hat{u}_t) - \frac{1}{2} \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix} - \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}^T c_t
\end{aligned}$$

d_t 는 상수 항으로서 2차 함수로 근사된 목적함수의 최적화에 영향을 미치지 않는다.

이제 선형화된 시스템과 2차 함수로 근사된 목적함수에 LQR 알고리즘을 적용할 수 있다.

- LQR의 역방향 패스를 적용하면 칼만 게인 시퀀스 (K_0, K_1, \dots, K_T) 와 (k_0, k_1, \dots, k_T) 를 계산할 수 있다.
- 그런 다음, 순방향 패스를 계산하면 새로운 최적 제어 (u_0, u_1, \dots, u_t) 와 같은 상태 시퀀스 (x_0, x_1, \dots, x_T) 를 얻을 수 있다.
- 그리고 명목 제어 시퀀스와 상태 시퀀스를 다음과 같이 업데이트 한 후, 위 과정을 수렴할 때 까지 반복하면 된다.

$$\begin{aligned}
(\hat{u}_0, \hat{u}_1, \dots, \hat{u}_T) &\leftarrow (u_0, u_1, \dots, u_T) \\
(\hat{x}_0, \hat{x}_1, \dots, \hat{x}_T) &\leftarrow (x_0, x_1, \dots, x_T)
\end{aligned}$$

1. 초기 상태 x_0 에 대해서 명목 제어 $(\hat{u}_0, \hat{u}_1, \dots, \hat{u}_T)$ 와 명목 상태 $(\hat{x}_0, \hat{x}_1, \dots, \hat{x}_T)$ 초기화

2. Repeat {

[1] $\hat{J}_0 = \sum_{t=0}^T c(\hat{x}_t, \hat{u}_t)$ 계산

[2] $f_{xt} = \nabla_x f(\hat{x}_t, \hat{u}_t)$, $f_{ut} = \nabla_u f(\hat{x}_t, \hat{u}_t)$, $f_{ct} = f(\hat{x}_t, \hat{u}_t) - f_{xt}\hat{x}_t - f_{ut}\hat{u}_t$ 계산

[3] $C_t = \nabla_{x_t, u_t}^2 c_t(\hat{x}_t, \hat{u}_t)$, $b_t = \nabla_{x_t, u_t} c_t(\hat{x}_t, \hat{u}_t)$ 계산

[4] $c_t = b_t - C_t \begin{bmatrix} \hat{x}_t \\ \hat{u}_t \end{bmatrix}$ 계산

[5] LQR 역방향 패스

[6] LQR 순방향 패스

for $t = 0 : T$ {

[1] $u_t = K_t x_t + k_t$ 계산

[2] $x_{t+1} = f(x_t, u_t)$ 계산

} end

[7] 명목 궤적과 제어 업데이트

$$\begin{aligned} (\hat{u}_0, \hat{u}_1, \dots, \hat{u}_T) &\leftarrow (u_0, u_1, \dots, u_T) \\ (\hat{x}_0, \hat{x}_1, \dots, \hat{x}_T) &\leftarrow (x_0, x_1, \dots, x_T) \end{aligned}$$

[8] 목적함수를 계산한다.

$$J_0 = \sum_{t=0}^T c_t(x_t, u_t)$$

} 수렴할 때 $(|\hat{J}_0 - J_0| \leq \epsilon)$ 까지 반복

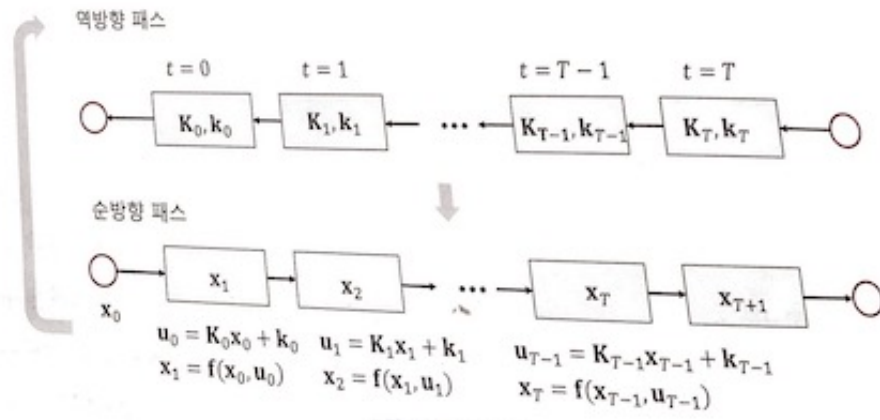


그림 8.4 반복적 LQR

8.3 모델 학습 방법

모델 강화 학습에서는 에이전트 환경과 상호 작용 하면서 얻은 샘플로 시스템의 동역학 모델을 지도학습 한다.

가장 간단한 알고리즘

1. 랜덤 정책이나 기본 정책을 실행해 상태천이 데이터셋 $D = \{(X_t, u_t, x_{t+1})_j\}$ 를 수집한다.
2. $x_{t+1} \approx f(x_t, u_t)$ 가 되도록 함수 $f(x_t, u_t)$ 를 학습한다.
3. 동역학 모델을 이용해 정책을 계산한다.

- 시스템의 동역학 구조는 알지만 일부 파라미터 값이 불확실한 경우에 유용
- 데이터를 수집하는 데 사용한 정책과 추정된 모델을 이용해 계산하는 정책이 다르다.

이 책에서 사용할 모델 학습 방법

1. 랜덤 정책이나 기본 정책을 실행해 상태천이 데이터셋 $D = \{(X_t, u_t, x_{t+1})_j\}$ 를 수집한다.
2. $x_{t+1} \approx f(x_t, u_t)$ 가 되도록 함수 $f(x_t, u_t)$ 를 학습한다.
3. 동역학 모델을 이용해 정책을 계산한다.
4. 계산된 정책을 실행해 새로운 궤적을 발생시키고 상태 천이 데이터셋 D에 추가한다.
5. 2번으로 돌아가 절차를 반복한다.

시스템의 동역학 모델 $x_{t+1} \approx f(x_t, u_t)$ 를 사용할 수 있지만, 가우시안 프로세스, 가우시안 혼합모델 (GMM), 신경망과 같은 일반적인 모델로도 표현할 수 있다.

글로벌 모델

- 모든 상태 공간에서 작동하는 단일 모델
- 상태 공간 전체에서 연속성을 갖고 있다는 장점이 있다.
- **시스템의 운동이 매우 복잡하다면** 전체 상태 공간에서 매우 복잡한 운동 모델을 고려해야 하고, 이 모델을 학습하기 위해서는 **많은 데이터**가 필요할 것이다.
- 가우시안 프로세스나 신경망 등으로 모델링

로컬 모델

- 상태 공간 일부에서만 작동하는 모델
- 단순하기 때문에 적은 수의 데이터셋 만으로도 추정하기 쉽다는 장점이 있다.
- 로컬 모델 기반으로 계산한 정책이 업데이트가 되면 해당 모델이 부정확해진다는 단점이 있다. → 명목 궤적 근처에서만 국지적으로 유효하기 때문

- 이 단점을 극복한다면 로컬 모델은 모델 기반 강화학습 기법을 실제 문제에 적용하는데 있어서 매우 유용한 수단이 될 수 있다.