# Linear System
(type of state transitions for continuos MDPs)

$$S_{t+1} = AS_t + Ba_t + W_t$$
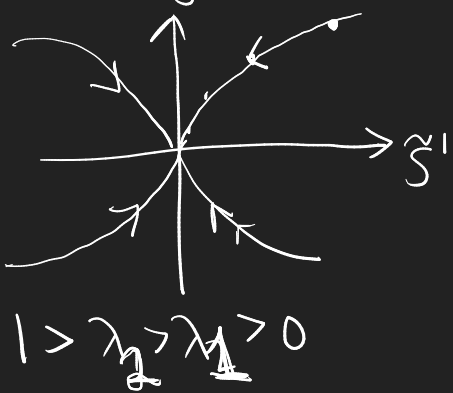$$\in \mathbb{R}^{n_s} \qquad \in \mathbb{R}^{n_a} \quad \sim \mathcal{N}(0, \sigma^2 I)$$

Last time:

$\rho(A)$ spectral radius
determining stability when $\underline{a_t = 0}$ and $\underline{W_t = 0}$

Key idea: $S_{t+1} = \underline{A} S_t$ $\qquad$ $A = \underline{VDV^{-1}}$ (if diagonalizable)

Then $\tilde{S}_t = V^{-1} S_t$ $\qquad\qquad$ diagonal
$$\begin{bmatrix} \lambda_1 & \lambda_2 \cdots \lambda_{ns} \end{bmatrix}$$

$\tilde{S}_t = \begin{bmatrix} \lambda_1^t & \\ & \ddots \\ & & \lambda_{ns}^t \end{bmatrix} \tilde{S}_0$

Alternative:
$$S_{t+1} = AS_t$$
$$S_0 = \boxed{V_i} \quad \text{(eigenvector, so } Av_i = \lambda_i v_i)$$
$$S_1 = AS_0 = Av_i = \lambda v_i$$
$$S_t = \boxed{\lambda^t v_i}$$



$1 > \lambda_2 > \lambda_1 > 0$

# 1) LQR (Linear Quadratic Regulator)

$$S_{t+1} = \underline{A} S_t + \underline{B} a_t + W_t \qquad W_t \sim \mathcal{N}(0, \sigma^2 I)$$

$$C(S, u) = S^T Q S + a^T R a$$

$Q, R$ symetric $Q^T = Q$
and positive definite (positive eigenvalues)

## OCP

$$\min_{\Pi} \; \mathbb{E}\left[ S_H^T Q S_H + \sum_{t=0}^{H-1} S_t^T Q S_t + a_t R a_t \; \middle| \; \begin{array}{l} S_{t+1} = A S_t + B a_t + w_t, \quad w_t \sim \mathcal{N}(0, \sigma^2 I) \\ a_t = \Pi_t(S_t), \quad S_0 \sim \mu_0 \end{array} \right]$$

## ex  1D robot

$$S_{t+1} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} S_t + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} a_t \qquad S_t = \begin{bmatrix} p_t \\ v_t \end{bmatrix}$$

$$C(S, a) = \gamma_p \, p_t^2 + \gamma_v \, v_t^2 + \gamma_a \, a_t$$

## Value & Q functions

"cost-to-go"

$$\rightarrow V_t^{\Pi}(S) = \mathbb{E}\left[ S_H^T Q S_H + \sum_{k=t}^{H-1} S_k^T Q S_k + a_k^T R a_k \; \middle| \; \begin{array}{l} \text{dynamics} \\ a_k = \boxed{\Pi_k(S_k)} \\ \boxed{S_t = S} \end{array} \right]$$

$$Q_t^{\Pi}(S, a) \quad " \qquad\qquad\qquad\qquad\qquad\qquad\qquad " \quad a_t = a$$

Because of terminal cost,

$$- V_H^{\Pi}(S) = S^T Q S$$

## 2) Optimal LQR Policy

Dynamic Programing for OCP:

Start: $V_H^{\Pi}(S) = C_H(S)$

for $t = H-1, H-2, \ldots, 0$:

$$Q_t^*(S, a) = C_t(S, a) + \underbrace{\mathbb{E}_{s' \sim P(S,a)}\left[ V_{t+1}^*(s') \right]}_{\substack{s' = f(s, a, w) \\ w \sim \emptyset}}$$

$$\Pi_t^*(S) = \underset{a}{\text{argmin}} \; Q_t^*(S, a)$$

$$V_t^*(S) = Q_t^*(S, \Pi_t^*(S))$$

Theorem (LQR optimal value function & policy)

Given $(A, B, Q, R, \sigma^2)$:

$$V_t^*(s) = s^\top P_t s + p_t$$

$$\pi_t^*(s) = -K_t^* s \qquad K_t^* \in \mathbb{R}^{n_s \times n_a}$$

where $P_t, K_t, p_t$ depend on $(A, B, Q, R, \sigma^2)$

Proof: by induction

claim 1: (Base case) $V_H^*(s) = s^\top P_H s + p_H$ is quadratic.

claim 2: (induction) if $V_{t+1}^*(s) = s^\top P_{t+1} s + p_{t+1}$

Then
1) $Q_t^*(s,a)$ is quadratic in $s, a$
2) $\pi_t^*(a) = \underset{a}{\text{argmin}} \; Q_t^*(s,a)$ is linear in $s$

Thus $V_t^*(s) = s^\top P_t s + p_t$ is quadratic.

$$P_H = Q \quad \text{and} \quad p_H = 0$$

$$Q_t^*(s,a) = \boxed{c(s,a)} + \underset{s'}{\mathbb{E}}\left[V_{t+1}^*(s')\right]$$

$$V_{t+1}^*(s) = s^\top P_{t+1} s + p_{t+1}$$

$$\underset{w_t \sim N(0,\sigma^2 I)}{\mathbb{E}}\left[V_{t+1}^*(As_t + Ba_t + w_t)\right]$$

$$= (As)^\top P_{t+1}(As) + (As)^\top P_{t+1} Ba + 0$$

$$(Ba)^\top P_{t+1} As + (Ba)^\top P_{t+1} Ba + 0$$

$$0 \qquad + \qquad 0 \qquad + \underset{w}{\mathbb{E}}\left[w^\top P_{t+1} w\right] + p_{t+1}$$

$$\Rightarrow \underset{w}{\mathbb{E}}\left[w^\top P w\right] = \sigma^2 \text{Tr}(P)$$

$$\mathbb{E}(\text{Tr}(w^\top P w)) = \text{Tr}(P \mathbb{E}[w w^\top]) = \sigma^2 I$$

$$\mathbb{E}[MW] = M\mathbb{E}[W]$$
$$w \sim N(0, \sigma^2 I) \qquad = M 0 = 0$$
$$\mathbb{E}[w^\top w] \neq \mathbb{E} w^\top \mathbb{E} w$$

$$Q_t^*(s,a) = s^T \underbrace{(Q + A^T P_{t+1} A)}_{M_1} s + a^T \underbrace{(R + B^T P_{t+1} B)}_{M_2} a$$
$$+ 2 \underbrace{s^T A^T P_{t+1} B a}_{M_3} + \underbrace{\sigma^2 tr(P_{t+1}) + p_{t+1}}_{C}$$

$$\pi_t^*(s) = \underset{a}{argmin}\ Q_t^*(s,a) \qquad M_1, M_2\ symetric \qquad \nabla_a(\omega^T a) = \omega$$

$$\rightarrow Q(s,a) = \underbrace{s^T M_1 s} + \boxed{a^T M_2\, a} + 2 \underbrace{s^T M_3 a} + C$$

minimun occurs $\nabla_a\ Q(s,a)$

$$\nabla_a Q(s,a) = 0 + 2 M_2 a + 2 M_3^T s + \bar{0} = 0$$

$$\pi_t^*(s) = - \underbrace{(R + B^T P_{t+1} B)^{-1} B^T P_{t+1} A}\, s \qquad a^* = -M_2^{-1} M_3^T s$$

$$(x^T M x) = x^T \left(\frac{M}{2} + \frac{M^T}{2}\right) x \qquad \underline{K_t^*} \qquad\qquad (AB)^T = B^T A^T$$
$$= (x^T M x)^T = x M^T x$$
$$V_t^*(s) = Q_t^* = (s, \pi^*(s)) \qquad\qquad Q(s, a^*) = s^T(M_1 - M_3 M_2^{-1} M_3^T) s + C$$