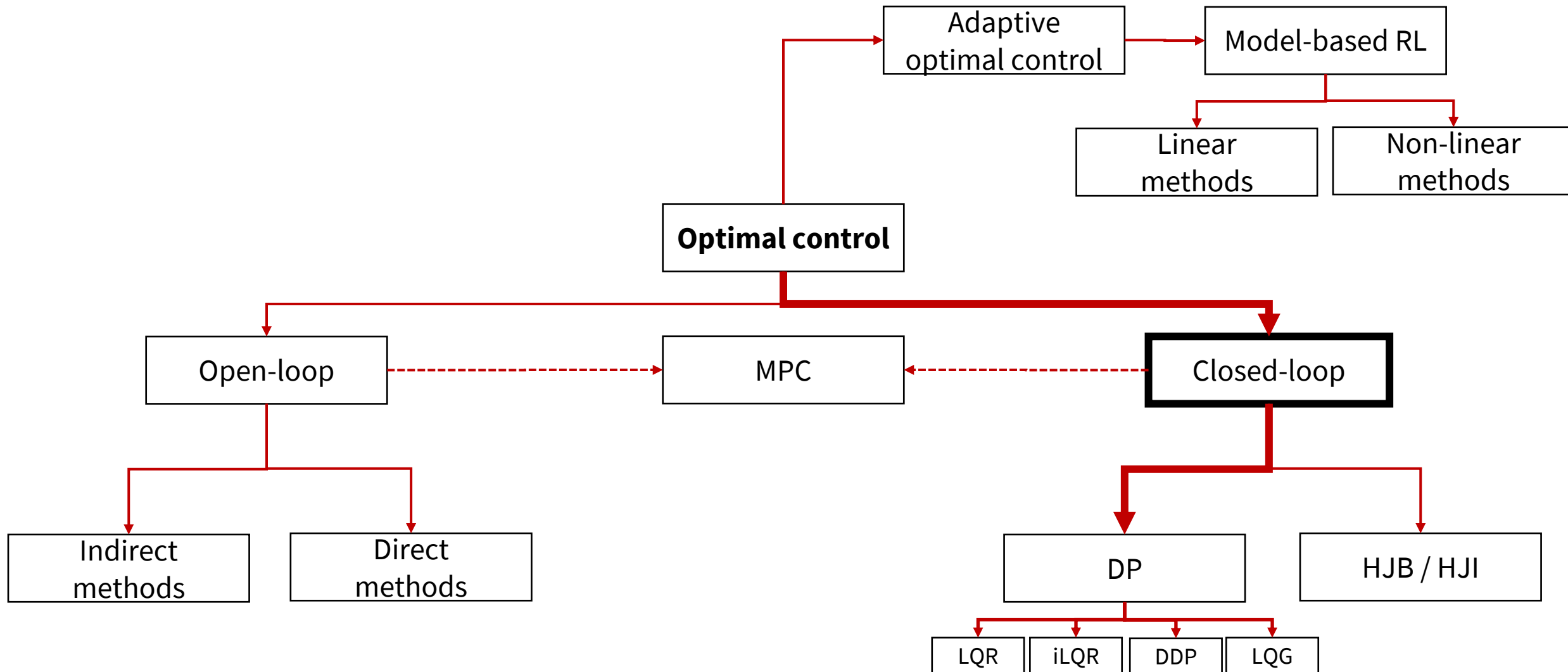


AA203

Optimal and Learning-based Control

iLQR, stochastic optimal control, LQG

Roadmap



LQR extensions

- Linear tracking problems
- LQR with cross-quadratic cost and affine dynamics
- Non-linear tracking problems
- Using LQR to solve non-linear optimal control problems
 - Iterative LQR
 - Differential dynamic programming

Linear tracking problems

- Imagine you are given a *nominal trajectory*

$$(\bar{\mathbf{x}}_0, \dots, \bar{\mathbf{x}}_N), (\bar{\mathbf{u}}_0, \dots, \bar{\mathbf{u}}_{N-1})$$

- Assume nominal trajectory satisfies linear dynamics
- Linear tracking problem: find policy to minimize cost

$$\frac{1}{2}(\mathbf{x}_N - \bar{\mathbf{x}}_N)' H(\mathbf{x}_N - \bar{\mathbf{x}}_N) + \frac{1}{2} \sum_{k=0}^{N-1} [(\mathbf{x}_k - \bar{\mathbf{x}}_k) Q (\mathbf{x}_k - \bar{\mathbf{x}}_k) + (\mathbf{u}_k - \bar{\mathbf{u}}_k) R (\mathbf{u}_k - \bar{\mathbf{u}}_k)]$$

- Then define *deviation variables*

$$\delta \mathbf{x}_k := \mathbf{x}_k - \bar{\mathbf{x}}_k \text{ and } \delta \mathbf{u}_k := \mathbf{u}_k - \bar{\mathbf{u}}_k$$

and solve standard LQR with respect to deviation variables

LQR with cross-quadratic cost & affine dynamics

- Consider the LQR problem with the generalized cost

$$\frac{1}{2} \mathbf{x}_k^T Q_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T R_k \mathbf{u}_k + \mathbf{u}_k^T H_k \mathbf{x}_k + \mathbf{q}_k^T \mathbf{x}_k + \mathbf{r}_k^T \mathbf{u}_k + c_k$$

- and dynamics

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k$$

- We can derive an *affine* optimal feedback law for this system via DP recursion

LQR with cross-quadratic cost & affine dynamics

- The cost-to-go at time k takes the form

$$\frac{1}{2} \mathbf{x}_k^T P_k \mathbf{x}_k + \mathbf{p}_k^T \mathbf{x}_k + p_k$$

- Optimal control takes the form $\mathbf{u}_k^* = \mathbf{l}_k + L_k \mathbf{x}_k$ with

$$\mathbf{l}_k = -(R_k + B_k^T P_{k+1} B_k)^{-1} (\mathbf{r}_k + \mathbf{p}_{k+1}^T B_k + \mathbf{d}_k P_{k+1} B_k)$$

$$L_k = -(R_k + B_k^T P_{k+1} B_k)^{-1} (B_k^T P_{k+1} A_k + H_k)$$

- Equations for the constant/linear/quadratic cost-to-go terms are unwieldy but not hard to derive, and are given in the lecture notes

Nonlinear tracking problems

- Imagine you are given a *feasible nominal trajectory*

$$(\bar{\mathbf{x}}_0, \dots, \bar{\mathbf{x}}_N), (\bar{\mathbf{u}}_0, \dots, \bar{\mathbf{u}}_{N-1})$$

- The tracking cost is still quadratic, but the dynamics are now nonlinear

$$\bar{\mathbf{x}}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k)$$

- To apply LQR, we can linearize around the nominal trajectory

$$\mathbf{x}_{k+1} \approx f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) + \underbrace{\frac{\partial f}{\partial \mathbf{x}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)}_A \underbrace{(\mathbf{x}_k - \bar{\mathbf{x}}_k)}_{\delta \bar{\mathbf{x}}_k} + \underbrace{\frac{\partial f}{\partial \mathbf{u}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)}_B \underbrace{(\mathbf{u}_k - \bar{\mathbf{u}}_k)}_{\delta \bar{\mathbf{u}}_k}$$

- And apply LQR to the deviation variables (with dynamics $\delta \bar{\mathbf{x}}_{k+1} = A\delta \bar{\mathbf{x}}_k + B\delta \bar{\mathbf{u}}_k$)

Non-linear optimal control problem

- Consider now non-linear optimal control problem

$$\min_{\mathbf{u}} \sum_{k=0}^{N-1} c(\mathbf{x}_k, \mathbf{u}_k)$$

subject to $\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k)$

- Can we apply LQR-techniques to approximately solve it?

Iterative LQR

- Imagine you are given a *feasible nominal trajectory*

$$(\bar{\mathbf{x}}_0, \dots, \bar{\mathbf{x}}_N), (\bar{\mathbf{u}}_0, \dots, \bar{\mathbf{u}}_{N-1})$$

- Linearize the dynamics around feasible trajectory

$$\mathbf{x}_{k+1} \approx f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) + \frac{\partial f}{\partial \mathbf{x}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)(\mathbf{x}_k - \bar{\mathbf{x}}_k) + \frac{\partial f}{\partial \mathbf{u}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)(\mathbf{u}_k - \bar{\mathbf{u}}_k)$$

- And Taylor expand cost function around feasible trajectory

$$c(\delta \mathbf{x}_k, \delta \mathbf{u}_k) = c_k + \underbrace{c_{\mathbf{x},k}^T}_{\mathbf{q}_k} \delta \mathbf{x}_k + \underbrace{c_{\mathbf{u},k}^T}_{\mathbf{r}_k} \delta \mathbf{u}_k + \frac{1}{2} \delta \mathbf{u}_k^T \underbrace{c_{\mathbf{u}\mathbf{u},k}^T}_{R_k} \delta \mathbf{u}_k + \frac{1}{2} \delta \mathbf{x}_k^T \underbrace{c_{\mathbf{x}\mathbf{x},k}^T}_{Q_k} \delta \mathbf{x}_k + \delta \mathbf{u}_k^T \underbrace{c_{\mathbf{u}\mathbf{x},k}^T}_{H_k} \delta \mathbf{x}_k$$

Iterative LQR

- By optimizing over deviation variables (using results for LQR with cross-quadratic cost & affine dynamics), we obtain new solution:

$$\{\bar{\mathbf{x}}_k + \delta \mathbf{x}_k^*\} \text{ and } \{\bar{\mathbf{u}}_k + \delta \mathbf{u}_k^*\}$$

- We can then re-linearize and Taylor expand around this new trajectory, and iterate!

Iterative LQR

- Backward pass ($k = N$ to 0):
 - Compute locally linear dynamics, locally quadratic cost around nominal trajectory
 - Solve local approximation of DP recursion to compute control law
 - Compute cost-to-go
- Forward pass ($k = 0$ to N):
 - Use control law to update nominal trajectory
- Iterate until convergence

Differential Dynamic Programming (DDP)

- iLQR first approximates dynamics and cost, then performs exact DP recursion
- DDP instead approximates DP recursion directly
 - Define change in cost-to-go J_k under perturbation $(\delta \mathbf{x}_k, \delta \mathbf{u}_k)$ as

$$Q(\delta \mathbf{x}_k, \delta \mathbf{u}_k) := c(\bar{\mathbf{x}}_k + \delta \mathbf{x}_k, \bar{\mathbf{u}}_k + \delta \mathbf{u}_k) + J_{k+1}(f(\bar{\mathbf{x}}_k + \delta \mathbf{x}_k, \bar{\mathbf{u}}_k + \delta \mathbf{u}_k))$$

- Then, second order expansion

$$Q(\delta \mathbf{x}_k, \delta \mathbf{u}_k) \approx \frac{1}{2} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix}^T \begin{bmatrix} Q_k & Q_{\mathbf{x},k}^T & Q_{\mathbf{u},k}^T \\ Q_{\mathbf{x},k} & Q_{\mathbf{x}\mathbf{x},k} & Q_{\mathbf{u}\mathbf{x},k}^T \\ Q_{\mathbf{u},k} & Q_{\mathbf{u}\mathbf{x},k} & Q_{\mathbf{u}\mathbf{u},k} \end{bmatrix} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix}$$

Differential Dynamic Programming (DDP)

- The optimal control perturbation is

$$\delta \mathbf{u}_k^* = \operatorname{argmin}_{\delta \mathbf{u}} Q(\delta \mathbf{x}_k, \delta \mathbf{u})$$

- Leveraging the approximation, one can re-use LQR results and find that the optimal deviation is

$$\delta \mathbf{u}_k^* = \mathbf{l}_k + \mathbf{L}_k \delta \mathbf{x}_k$$

- Algorithm proceeds via same forward/backward passes as iLQR

Stochastic optimal control problem (discrete time)

- **System:** $\mathbf{x}_{k+1} = f_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), k = 0, \dots, N$
- **Control constraints:** $\mathbf{u}_k \in U(\mathbf{x}_k)$
- **Probability distribution:** $P_k(\cdot | \mathbf{x}_k, \mathbf{u}_k)$ of \mathbf{w}_k
- **Policies:** $\pi = \{\pi_0, \dots, \pi_{N-1}\}$, where $\mathbf{u}_k = \pi_k(\mathbf{x}_k)$
- **Expected Cost:**

$$J_\pi(\mathbf{x}_0) = E_{\mathbf{w}_k, k=0, \dots, N-1} \left[g_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} g_k(\mathbf{x}_k, \pi_k(\mathbf{x}_k), \mathbf{w}_k) \right]$$

- **Stochastic optimal control problem**

$$J^*(\mathbf{x}_0) = \min_{\pi} J_\pi(\mathbf{x}_0)$$

Key points

- Discrete-time model
- Markovian model
- Objective: find optimal **closed-loop**
- Additive cost (central assumption)
- Risk-neutral formulation

Other communities use different notation: Powell, W. B. *AI, OR and control theory: A Rosetta Stone for stochastic optimization*. Princeton University, 2012.

http://castlelab.princeton.edu/Papers/AIOR_July2012.pdf

Principle of optimality

- Let $\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{N-1}^*\}$ be an optimal policy
- Consider **tail subproblem**

$$E \left[g_N(\mathbf{x}_N) + \sum_{k=i}^{N-1} g_k(\mathbf{x}_k, \pi_k(\mathbf{x}_k), \mathbf{w}_k) \right]$$

and the **tail policy** $\{\pi_i^*, \dots, \pi_{N-1}^*\}$

Principle of optimality: The tail policy is optimal for the tail subproblem

The DP algorithm (stochastic case)

Intuition

- DP first solves ALL tail subproblems at the final stage
- At generic step, it solves ALL tail subproblems of a given time length, using solution of tail subproblems of shorter length

The DP algorithm (stochastic case)

The DP algorithm

- Start with

$$J_N(\mathbf{x}_N) = g_N(\mathbf{x}_N)$$

and go backwards using

$$J_k(\mathbf{x}_k) = \min_{\mathbf{u}_k \in U(\mathbf{x}_k)} E_{\mathbf{w}_k} [g_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) + J_{k+1}(f(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k))]$$

for $k = 0, 1, \dots, N - 1$

- Then $J^*(\mathbf{x}_0) = J_0(\mathbf{x}_0)$ and optimal policy is constructed by setting $\pi_k^*(\mathbf{x}_k) = \mathbf{u}_k^*$

Example: Inventory Control Problem (1/3)

- Stock available $x_k \in \mathbb{N}$, inventory $u_k \in \mathbb{N}$, and demand $w_k \in \mathbb{N}$
- Dynamics: $x_{k+1} = \max(0, x_k + u_k - w_k)$
- Constraints: $x_k + u_k \leq 2$
- Probabilistic structure: $p(w_k = 0) = 0.1$, $p(w_k = 1) = 0.7$, and $p(w_k = 2) = 0.2$
- Cost

$$E \left[\underbrace{0}_{g_3(x_3)} + \sum_{k=0}^2 \underbrace{(u_k + (x_k + u_k - w_k)^2)}_{g_k(x_k, u_k, w_k)} \right]$$

Example: Inventory Control Problem (2/3)

- Algorithm takes form

$$J_k(x_k) = \min_{0 \leq u_k \leq 2-x_k} E_{w_k} [u_k + (x_k + u_k - w_k)^2 + J_{k+1}(\max(0, x_k + u_k - w_k))]$$

for $k = 0, 1, 2$

- For example

$$J_2(0) = \min_{u_2=0,1,2} E_{w_2} [u_2 + (u_2 - w_2)^2] = \min_{u_2=0,1,2} E_{w_2} [u_2 + 0.1(u_2)^2 + 0.7(u_2 - 1)^2 + 0.2(u_2 - 2)^2]$$

which yields $J_2(0) = 1.3$, and $\pi_2^*(0) = 1$

Example: Inventory Control Problem (3/3)

Final solution:

- $J_0(0) = 3.7$,
- $J_0(1) = 2.7$, and
- $J_0(2) = 2.818$

Problems with imperfect state information

- Now the controller, instead of having perfect knowledge of the state, has access to observations \mathbf{z}_k of the form

$$\mathbf{z}_0 = h_0(\mathbf{x}_0, \mathbf{v}_0), \quad \mathbf{z}_k = h_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{v}_k), \quad k = 1, 2, \dots, N - 1$$

- The random observation disturbance is characterized by a given probability distribution

$$P_{\mathbf{v}_k}(\cdot | \mathbf{x}_k, \dots, \mathbf{x}_0, \mathbf{u}_{k-1}, \dots, \mathbf{u}_0, \mathbf{w}_{k-1}, \dots, \mathbf{w}_0, \mathbf{v}_{k-1}, \dots, \mathbf{v}_0)$$

- The initial state \mathbf{x}_0 is also random and characterized by given $P_{\mathbf{x}_0}$

Control policies

- Define the *information vector* as

$$\mathbf{I}_k = (\mathbf{z}_0, \dots, \mathbf{z}_k, \mathbf{u}_0, \dots, \mathbf{u}_{k-1}), \quad \mathbf{I}_0 = \mathbf{z}_0$$

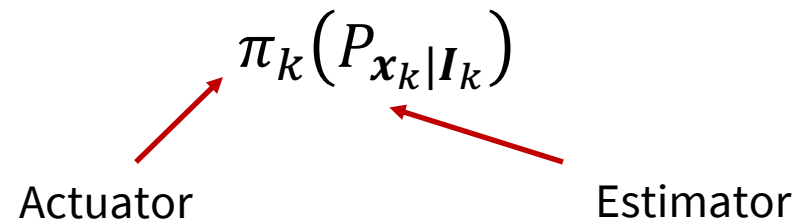
- Focus is now on *admissible* policies $\pi_k(\mathbf{I}_k) \in U_k$

- We want then to find an admissible policy that minimizes

$$J_\pi = E_{\mathbf{x}_0, \mathbf{w}_k, \mathbf{v}_k} \left[g_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} g_k(\mathbf{x}_k, \pi_k(\mathbf{I}_k), \mathbf{w}_k) \right]$$

Solution strategies

1. Reformulation as a perfect state information problem (main idea: make the information vector the state of the system)
 - Main drawback: state has *expanding* dimension!
2. Reason in terms of sufficient statistics, i.e., quantities that ideally are smaller than \mathbf{I}_k and yet summarize all its essential content
 - Main example: conditional probability distribution $P_{\mathbf{x}_k|\mathbf{I}_k}$
 - Condition probability distribution leads to a decomposition of the optimal controller in two parts:



LQG

Discrete LQG: find admissible control policy that minimizes

$$E \left[\mathbf{x}'_N Q \mathbf{x}_N + \sum_{k=0}^{N-1} (\mathbf{x}'_k Q_k \mathbf{x}_k + \mathbf{u}'_k R_k \mathbf{u}_k) \right]$$

subject to

- the dynamics $\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{w}_k$
- the measurement equation $\mathbf{z}_k = C_k \mathbf{x}_k + \mathbf{v}_k$

and with $\mathbf{x}_0, \{\mathbf{w}_k\}, \{\mathbf{v}_k\}$, independent and Gaussian vectors (and in addition $\{\mathbf{w}_k\}, \{\mathbf{v}_k\}$ zero mean)

LQG – solution

Let

- $M_k := E[\mathbf{w}_k \mathbf{w}_k']$
- $N_k := E[\mathbf{v}_k \mathbf{v}_k']$
- $S := E[(\mathbf{x}_0 - E[\mathbf{x}_0])(\mathbf{x}_0 - E[\mathbf{x}_0])']$

LQG – solution

The optimal controller is $\mathbf{u}_k = F_k \hat{\mathbf{x}}_k$, where

- F_k is the LQR gain
- $\hat{\mathbf{x}}_{k+1} = A_k \hat{\mathbf{x}}_k + B_k \mathbf{u}_k + \Sigma_{k+1|k+1} C'_{k+1} N_{k+1}^{-1} (\mathbf{z}_{k+1} - C_{k+1} (A_k \hat{\mathbf{x}}_k + B_k \mathbf{u}_k))$
- $\hat{\mathbf{x}}_0 = E[\mathbf{x}_0] + \Sigma_{0|0} C'_0 N_0^{-1} (\mathbf{z}_0 - C_0 E[\mathbf{x}_0])$
- and matrices $\Sigma_{k|k}$ are *precomputable* (given in the lecture notes)
- Key property: the estimation portion of the optimal controller is an optimal solution of the problem of estimating the state x_k assuming no control takes place, while the actuator portion is an optimal solution of the control problem assuming perfect state information
→ **separation principle**

Next time

- HJB and continuous-time LQR