

4F3 cribs

Question 1

- (a) Let us denote x_0, x_1, \dots, x_N the sequence of nodes along the optimal path from node x_0 to node x_N . By Bellman's Principle of Optimality, the truncated sequence x_0, x_1, \dots, x_{N-1} is the optimal path from x_0 to x_{N-1} as well (truncated problem). This principle allows to take advantage of recursive algorithms in optimization.
- (b) The sequence of graphs in Figure 1 represents the main steps of the dynamic programming computation. At each step we update the cost to go, whose values are represented within each node.

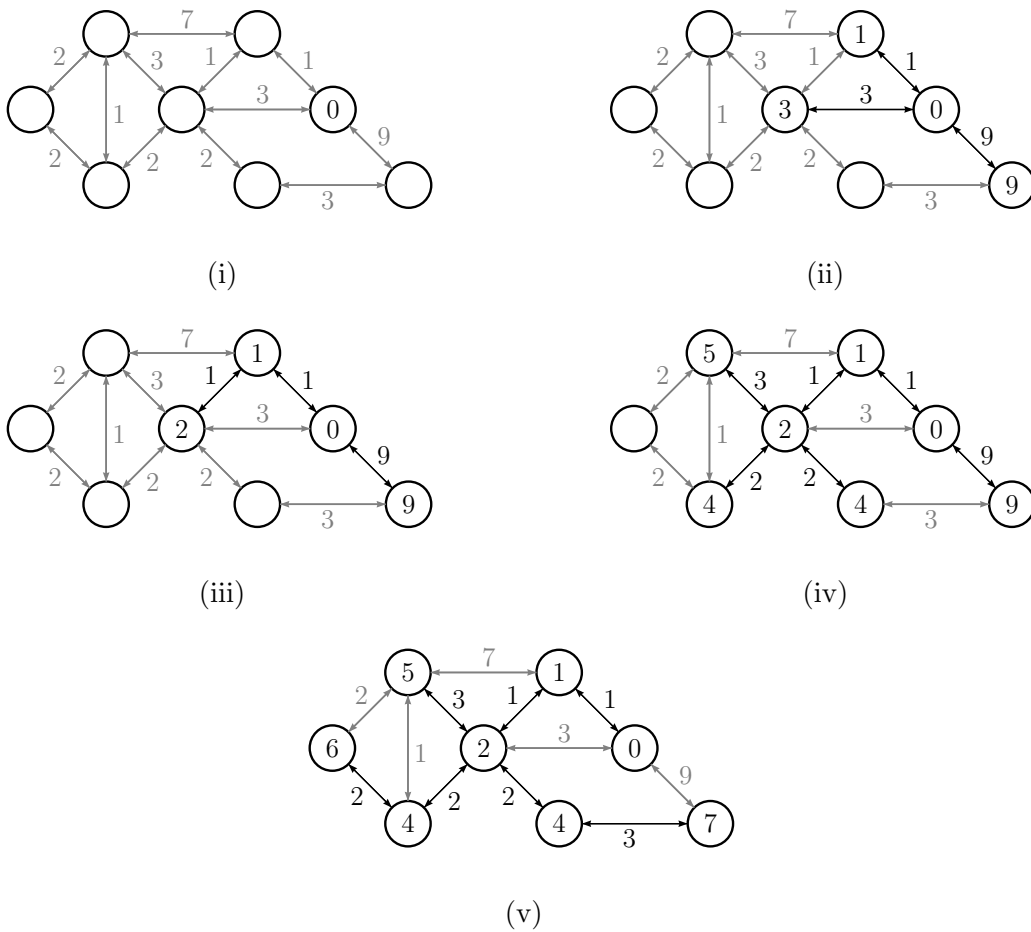


Figure 1

- (c) The minimal-time routing paths needs to be recalculated, with the delay per edge increased by 1, as paths with more hops (intermediate nodes) will be penalized more. After recalculating there are some nodes with two paths with the same minimal cost.
- (d) Running back and forth on the edge with negative weight makes the cost to go smaller and smaller. That is, with a negative edge the path length cost has no minimum.
- (e.i) Define the index set $I = \{1, 2, \dots, 8\}$ associated with each node. For a path from node 1 to node 7, the indicator variable x_{ij} is defined as $x_{ij} = 1$ if there exists an edge from i to j that belongs to the path and 0 otherwise, i.e. the indicator set identifies the path. Every path from node 1 to node 7 satisfies the following three conditions

C₁: there is a single edge entering the final node:

$$\sum_{i \in I} x_{i7} = 1 ;$$

C₂: there is a single edge leaving the initial node:

$$\sum_{i \in I} x_{1i} = 1 ;$$

C₃: the number of edges entering the node j and the number of edges leaving the node j are the same: for all $j \in I$ such that $j \neq 0$ and $j \neq 7$,

$$\sum_{i \in I} x_{ij} - \sum_{i \in I} x_{ji} = 0 .$$

Subject to these constraints, the suggested cost associated to the path from node 1 to node 7 is given by

$$J = \sum_{i \in I, j \in I} w_{ij} x_{ij} .$$

Summarizing, the mathematical optimisation problem reads

$$\min \sum_{i \in I, j \in I} w_{ij} x_{ij} \text{ subject to } x_{ij} \in \{0, 1\} \text{ and } \mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3 .$$

- (e.ii) For $x_{i,j} \geq 0$ and continuous, this is a linear program, since cost function and constraints are all linear.

Question 2

(a) The system can be written as follows

$$\dot{x} = Ax + B_1 w_1 + B_2 u \quad y = Cx + w_2$$

where

$$A = \begin{bmatrix} -k_{12} - d & k_{21} \\ k_{12} & -k_{21} - d \end{bmatrix} = \begin{bmatrix} -11 & 20 \\ 10 & -21 \end{bmatrix}$$

and

$$B_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad B_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad C = [0 \quad 1].$$

We can now derive the transfer function as

$$T_{w_1 \rightarrow y} = C(sI - A)^{-1} B_1 = \frac{s + 11}{s^2 + 32s + 31}$$

(b) The \mathcal{H}_2 norm is defined as follows

$$\|T_{w_1 \rightarrow y}\|_2^2 = \int_{-\infty}^{\infty} \text{trace} \{ \bar{T}_{w_1 \rightarrow y}(j\omega)^T T_{w_1 \rightarrow y}(j\omega) \} d\omega$$

Its physical meaning in terms of system performance is given by the expression

$$\|y\|_{\infty} \leq \frac{1}{\sqrt{2\pi}} \|T_{w_1 \rightarrow y}\|_2 \|w_1\|_2.$$

Thus, $\|T_{w_1 \rightarrow y}\|_2$ provides a bound on the largest amplification between the energy of the disturbance w_1 and the magnitude of the output y .

The \mathcal{H}_{∞} norm satisfies

$$\|T_{w_1 \rightarrow y}\|_{\infty} = \sup_{\omega} \bar{\sigma}(T_{w_1 \rightarrow y}(j\omega)).$$

Its physical meaning is well described by the expression

$$\|T_{w_1 \rightarrow y}\|_{\infty} = \sup_{w_1} \frac{\|T_{w_1 \rightarrow y} w_1\|_2}{\|w_1\|_2} = \sup_{w_1} \frac{\|y\|_2}{\|w_1\|_2}.$$

Thus, $\|T_{w_1 \rightarrow y}\|_{\infty}$ provides a bound on the largest amplification between the energies of the disturbance w_1 and the output y .

(c) We have that $\frac{1}{\sqrt{2\pi}} \|T_{w_1 \rightarrow y}\|_2 = \sqrt{B_1^T L B_1}$ where $L = L^T$ solves

$$A^T L + L A + C^T C = 0 \quad (L \text{ is the observability Gramian}).$$

From the latter, taking $L = \begin{bmatrix} L_1 & L_2 \\ L_2 & L_3 \end{bmatrix}$ we have

$$\begin{bmatrix} -11 & 10 \\ 20 & -21 \end{bmatrix} \begin{bmatrix} L_1 & L_2 \\ L_2 & L_3 \end{bmatrix} + \begin{bmatrix} L_1 & L_2 \\ L_2 & L_3 \end{bmatrix} \begin{bmatrix} -11 & 20 \\ 10 & -21 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

which gives

$$L = \begin{bmatrix} 0.0504 & 0.0554 \\ 0.0554 & 0.0766 \end{bmatrix}.$$

Thus,

$$\sqrt{B_1^T L B_1} = \sqrt{0.0766} \simeq 0.2768$$

(d.i) Let $P(s)$ be a generalized plant with performance input/output pair w and z and control input/output pair u and y . The \mathcal{H}_2 optimal control problem corresponds to finding the stabilizing controller $u = K(s)y$ that minimizes $\|T_{w \rightarrow z}(s)\|_2$ in closed loop. Using linear fractional transformations, $T_{w \rightarrow z}(s) = \mathcal{F}_l(P(s), K(s))$ therefore the \mathcal{H}_2 optimal control problem corresponds to

$$\min_{K(s) \text{ stabilising}} \|\mathcal{F}_l(P(s), K(s))\|_2$$

(d.ii) Given the output performance z , the generalized plant is given by

$$\begin{bmatrix} \dot{x} \\ z \\ y \end{bmatrix} = \left[\begin{array}{c|cc} A & [B_1 \ 0] & B_2 \\ \hline [C_1] & 0 & [0] \\ 0 & [0 \ I] & [I] \\ C_2 & & 0 \end{array} \right] \begin{bmatrix} x \\ w \\ u \end{bmatrix}$$

where $C_2 = C_1 = [0 \ 1]$. All other matrices have been already defined above.

(d.ii) From the datasheet, we need to compute

$$A_k = A - B_2 F - H C_2 \quad B_k = -H = Y C_2^T \quad C_k = F = B_2^T X$$

Thus,

$$A_k = \begin{bmatrix} -11.0491 & 19.8387 \\ 10.0000 & -21.0747 \end{bmatrix} \quad B_k = \begin{bmatrix} -0.1072 \\ -0.0747 \end{bmatrix} \quad C_k = [0.0491 \ 0.0542]$$

Question 3

(a) $\min_{\theta} \frac{1}{2} \theta^T Q \theta + c^T \theta$ subject to $A \theta \leq b$

(b) (i) Let $P \geq 0$ be steady state solution to the Ricatti equation in datasheet, ie

$$P = Q + A^T P A - A^T P B (R + B^T P B)^{-1} B^T P A$$

and put $K = -(R + B^T P B)^{-1} B^T P A$, which is the optimal solution, a state feedback gain, for the unconstrained infinite horizon case.

Now find M, b s.t. the set $\{x : Mx \leq b\}$ is *constraint admissible* and *invariant* with respect to the control $u = Kx$ (i.e. $Mx \leq b \implies |Kx| \leq U$ and $Mx \leq b \implies M(Ax + BKx) \leq b$).

Now solve the MPC problem

$$\min_{u_0, u_1} \sum_{k=0}^1 (x_k^T Q x_k + u_k^T R u_k) + x_2^T P x_2$$

subject to $x_1 = Ax_0 + Bu_0$, $x_2 = Ax_1 + Bu_1 = A^2x_0 + ABu_0 + Bu_1$

i.e.

$$\begin{aligned} & \min_{u_0, u_1} x_0^T Q x_0 + (Ax_0 + Bu_0)^T Q (Ax_0 + Bu_0) + \\ & (A^2x_0 + ABu_0 + Bu_1)^T P (A^2x_0 + ABu_0 + Bu_1) + u_0^T R u_0 + u_1^T R u_1 \\ = & \text{const} + \min_{u_0, u_1} \begin{bmatrix} u_0 \\ u_1 \end{bmatrix}^T \begin{bmatrix} B^T Q B + B^T A^T P A B + R & B^T A^T P B \\ B^T P A B & B^T P B + R \end{bmatrix} \begin{bmatrix} u_0 \\ u_1 \end{bmatrix} \\ & + 2 \begin{bmatrix} x_0^T A^T Q B + x_0^T A^2 P A B & x_0^T A^2 P B \end{bmatrix} \begin{bmatrix} u_0 \\ u_1 \end{bmatrix} \end{aligned}$$

s.t

$$\begin{bmatrix} I \\ -I \\ MAB & B \end{bmatrix} \begin{bmatrix} u_0 \\ u_1 \end{bmatrix} \leq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ b - A^2x_0 \end{bmatrix}$$

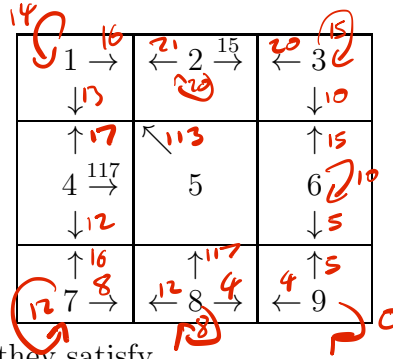
(ii) The control above is suboptimal because of the added constraint that $Mx \leq b$ at step 2. Increasing N pushes this constraint further into the future, and the resulting controller will be closer to optimal, and eventually optimal.

(iii) For large x_0 it may be impossible to get inside the control invariant set in 2 time steps.

Question 4

(a) $Q(s, a)$ is the cost of taking action a from state s and taking optimal actions thereafter. It enables you to find the optimal action, by enumeration, without knowing the model or reward.

(b) (i)



These are optimal, since they satisfy

$$Q(s, a) = c(s) + \min_{a'} Q(s', a')$$

in each case (where s' is that state reached by taking action a at state s .)

- (ii) The optimal cost is 13, associated with the path $1 \rightarrow 4 \rightarrow 7 \rightarrow 8 \rightarrow 9$ (and remain at 9 thereafter).
- (iii) After each state action pair has been visited a sufficiently large number of times Q will have converged to the optimal value. There is then a 5% chance of taking the suboptimal action $4 \rightarrow 5$, and the same for $8 \rightarrow 5$. There is thus a 10% chance of an extra cost of at least 100. Hence the cost $> 13 + 0.1 * 100 = 23$.
- (iv) For SARSA the iteration is instead

$$Q(s, a) \leftarrow c(s) + Q(s', a')$$

where a' is selected according to the current policy. If this policy is ϵ -greedy, then actions to states 4 and 8 will be marked with a higher cost and the path found will be $1 \rightarrow 2 \rightarrow 3 \rightarrow 6 \rightarrow 9$ with an optimal cost of slightly greater than 16.

Q1. Dynamic programming.

A quite popular question that was attempted by most candidates. Parts (a)–(d) associated with finding the minimal-time routing path was generally well answered. Part (e) was more challenging and many candidates had difficulties formulating the constraints associated with the optimization problem.

Q2. H2 optimal control.

The question was attempted by about two third of the candidates. Part (a) associated with providing a state space representation and finding a transfer function was generally well answered. Many candidates had difficulties in part (c) associated with the computation of the H2 norm. There were also good attempts for part (d) where candidates were asked to formulate and solve the H2 optimal control problem.

Q3. Model Predictive Control.

A popular question on using MPC to solve the constrained LQR problem. Attempted by most candidates and with many excellent answers. Most marks were lost by a lack of precision in answering a).

Q4 Reinforcement learning

A less popular question, attempted by about half the candidates, comparing Q-learning and SARSA on a cliff edge walking problem. The basics were well done, but ost marks were lost on b) parts iii) and iv), with many candidates not grasping the effect of using epsilon-greedy action selection.