# Dynamic programming



system dynamics: $\quad x_{k+1}=f_k(x_k,u_k,w_k)$

$\left(\begin{array}{lll} k\cdot\text{time index} & x_k: \text{system state} \in S_k & X=(x_1,...,x_N) \\ N\cdot\text{time horizon} & u_k: \text{control input} \in \mathcal{U}_k(x_k) & U=(u_0,...,u_{N-1}) \\ (\text{dep on } x_k, u_k) \to w_k: \text{disturbances} & W=(w_0,...,w_{N-1}) \end{array}\right)$

objective: find optimal input policies $\pi=(\mu_1,...,\mu_{N-1})$ to minimize cost $\quad$ ($g_k$: stage cost ; $g_N$: terminal cost)

$$J^*(x_0)=\min_\pi J(x_0,\pi)=\min_\pi E\{g_N(x_N)+\sum_{k=0}^{N-1}g_k(x_k,u_k=\mu_k(x_k),w_k)\} \quad \text{subj. to} \quad x_{k+1}=f_k(x_k,u_k=\mu_k(x_k),w_k)$$

[ using expected value in cost is not ideal, as variance could be high, but having variance in cost ]
[ makes dynamic programing no longer applicable! ]

## dynamic programing algorithm (DPA)

optimized tail of trajectory is also part of optimal trajectory, so: start from tail and iteratively find one input/policy to append that minimizes cost of current trajectory (=cost-to-go)

optimal cost-to-go: $J_k(x_k)$ optimal cost to go from state $x_k$ to end of trajectory

recursion: $J_k(x_k)=\min_{\mu_k} E_{(w_k|x_k,\mu_k)}\{g_k(x_k,u_k=\mu_k(x_k),w_k)+J_{k+1}(f_k(x_k,u_k=\mu_k(x_k),w_k))\}$ initial cond: $J_N(x_N)=g_N(x_N)$

( in practice: minimize over $u_k$ for all possible states $x_k$ at time k $\to$ collection of $u_k \triangleq \mu_k(x_k)$ )
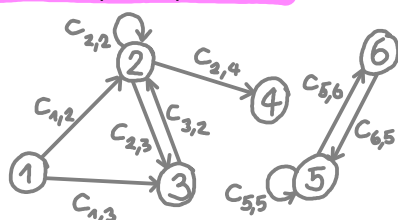
⊙ conversions to standard form:

- time lag: $x_{k+1}=f'(x_k,x_{k-1},u_k,u_{k-1},w_k) \to \tilde{x}_k=\begin{bmatrix}x_k=x_k\\y_k=x_{k-1}\\s_k=u_{k-1}\end{bmatrix}$ ; $\tilde{x}_{k+1}=\begin{bmatrix}x_{k+1}\\y_{k+1}\\s_{k+1}\end{bmatrix}=\begin{bmatrix}f_k(x_k,y_k,u_k,s_k,w_k)\\x_k\\u_k\end{bmatrix}=\tilde{f}_k(\tilde{x}_k,u_k,w_k)$

- correlated disturbance: $w_k=C_ky_k$, $y_{k+1}=A_ky_k+\xi_k \to \hat{x}_k=\begin{bmatrix}x_k\\y_k\end{bmatrix}$ ; $\hat{x}_{k+1}=\begin{bmatrix}x_{k+1}\\y_{k+1}\end{bmatrix}=\begin{bmatrix}f_k(x_k,u_k,C_k(A_ky_k+\xi_k))\\A_ky_k+\xi_k\end{bmatrix}=\tilde{f}_k(\tilde{x}_k,u_k,\xi_k)$

- forecast: at the start of each step k we get a forecast $y_k$ with which $w_k$ correlates. the a priori prob of $y_k$ is known.

  $\quad\hookrightarrow p(w_k|y_k)=...$
  $\quad\hookrightarrow p(\xi_k)=... ; y_{k+1}=\xi_k$ $\quad \to \quad \tilde{x}_k=\begin{bmatrix}x_k\\y_k\end{bmatrix}$ ; $\tilde{w}_k=\begin{bmatrix}w_k\\\xi_k\end{bmatrix}$ ; $\tilde{x}_{k+1}=\begin{bmatrix}x_{k+1}\\y_{k+1}\end{bmatrix}=\begin{bmatrix}f_k(x_k,u_k,w_k)\\\xi_k\end{bmatrix}=\tilde{f}_k(\tilde{x}_k,u_k,\tilde{w}_k)$

  $p(\tilde{w}_k|\tilde{x}_k,u_k)=...=p(w_k|y_k)p(\xi_k)$ $\quad J_k(\tilde{x}_k)=\min_{\mu_k} E_{(w_k|y_k)}\{g_k(x_k,u_k=\mu_k(x_k),w_k)+E_{\xi_k}\{J_{k+1}(\tilde{f}_k(\tilde{x}_k,u_k=\mu_k(x_k),\tilde{w}_k))\}\}$

## infinite horizon (BE)

for a time invariant system+cost, let $N\to\infty$. Then $J_k=J_{k+1}=J_\infty$ and $\mu_k=\mu_{k+1}=\mu_\infty$ as 1 timestep is small compared to $\infty$

$\hookrightarrow$ Bellmann Equation. $J_\infty(x)=\min_{\mu_\infty} E_{(w|x,\mu)}\{g(x,u,w)+J_\infty(f(x,u=\mu_\infty(x),w))\} \to u=\mu(x)$: optimal input

## shortest path problem (SP)



$\mathcal{V}$: vertex space ①,...
$C$: edge space
$C:=\{(i,j,c_{i,j})\in\mathcal{V}\times\mathcal{V}\times\mathbb{R}\}$
$c_{i,j}$: length from vertex i to j

$Q$: path = ordered list of nodes
$Q:=(i_1,...,i_q)$ $\quad$ ($q$: path length)
$Q_{s,\tau}$: set of all paths starting at vertex $s\in\mathcal{V}$ and ending at $\tau\in\mathcal{V}$

objective: $Q^*=\underset{Q\in Q_{s,\tau}}{\text{argmin}} J_Q$ (shortest path from vertex s to $\tau$) $\leftarrow J_Q=\sum_{h=1}^q c_{i_h,i_{h+1}}$ (length of path Q)

assumption: no negative cycles! $\forall i\in\mathcal{V}$ and $\forall Q\in Q_{i,i}$ $\nexists J_Q<0$

# <u>stochastic shortest path problem</u> (SSP) *(discrete system defined by transition probabilities)*

system: $x_{k+1} = w_k$      $x_k \in S_k, u_k \in U_k(x_k)$    ( j: discr. state $x_{k+1}$ #i )  ( convert system of standard formulation: )

$P_{x_{k+1}|x_k u_k}(j|i,u) = P_{w_k|x_k u_k}(j|i,u) = P_{ij}(u,k)$    ( i: discr. state $x_k$ #j )   ( $P_{ij}(u) = \sum_{\{w_k|f(i,u,w_k)=j\}} P_{w_k|x_k,u_k}(\bar{w}_k|i,u)$ )

             ( u: discr. input $u_k$ #U )

objective: $J^*(x_0) = \min_\pi E\{g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k = \mu_k(x_k), w_k)\}$   ( $g_k$: stage cost ; $g_N$: terminal cost )

- **DPA + SSP problem:**   $J_k(i) = \min_{\mu_k} q_k(i, u=\mu_k(i)) + \sum_{j=1}^n P_{ij}(u=\mu_k(i),k) \cdot J_{k+1}(j)$    ( $q_k(i,u) = E_{(j|i,u)}\{g_k(i,u,j)\}$ )

- **BE + SSP problem:**

assumptions:   • time invariant SSP: $P_{ij}(u,k) \to P_{ij}(u)$ ; $g_k \to g$ ; $S_k \to S$ ; $U_k(x_k) \to U(x_k)$

             • cost free termination state: state with i=0 is such that $P_{00}(u)=1$ and $g(0,u,0)=0$   $\forall u$

             • $\exists$ a policy $\mu$ such that i=0 is eventually reached from all $x_0$   (propper policy)

$\hookrightarrow$ Bellmann equation: $J_\infty(i) = \min_{\mu_\infty}[q(i, u=\mu_\infty(i)) + \sum_{j=1}^n P_{ij}(u=\mu_\infty(i)) \cdot J_\infty(j)]$    $q(i,u) = E_{(j|i,u)}\{g(i,u,j)\} = \sum_{j=1}^n P_{ij}(u) g(i,u,j)$

---

| **value iteration**: do DPA, until cost-to-go stops changing | **policy iteration**: improve $\mu_\infty(i)$ guess until it converges |
|---|---|
| $V_l \triangleq J_{N-l}$ cost-to-go $l$ steps from end. $\to$ if $l$ is large, then $V_l \approx J_\infty$ | $J_{\mu_n}(i)$: cost to $\infty$ with policy $\mu_n(i)$ |
| initial cond:   $V_0(i) =$ some arbitrary cost | initial cond:   $\mu_0(i) =$ some arbitrary <u>propper</u> policy |
| policy update: $\mu_l(i) = \arg\min_\mu [q(i, u=\mu(i)) + \sum_{j=1}^n P_{ij}(u=\mu(i)) V_l(j)]$ | policy eval:   solve $J_{\mu_n}(i) = q(i, u=\mu_n(i)) + \sum_{j=1}^n P_{ij}(u=\mu_n(i)) J_{\mu_n}(j)$ |
| value update: $V_{l+1}(i) = q(i, u=\mu_l(i)) + \sum_{j=1}^n P_{ij}(u=\mu_l(i)) V_l(j)$ | policy update: $\mu_{n+1}(i) = \arg\min_\mu [q(i, u=\mu(i)) + \sum_{j=1}^n P_{ij}(u=\mu(i)) J_{\mu_n}(j)]$ |
| break cond:   $\|V_{l+1}(i) - V_l(i)\| < tol$ | break cond:   $\|J_{\mu_{n+1}}(i) - J_{\mu_n}(i)\| < tol$ |
| ( not guaranteed to converge in finite iterations ) ( less comp expensive, needs many iterations ) | ( guaranteed to converge in finite iterations ) ( more comp expensive, needs few iterations ) |

---

$\to$ variations: - Gauss-Seidel update: in value update use $V_{l+1}(i)$ instead of $V_l(i)$ if it was already computed

           - run multiple value updates with same policy before policy update ($\triangleq$ policy eval for inf.)

           - update only some states in policy/value update (e.g. half in one iter., half in the following)

$\to$ policy eval in matrix form: $\underline{J_{\mu_n}} = q + \underline{P} \underline{J_{\mu_n}}$   ( $\underline{J_{\mu_n}} = \begin{bmatrix} J_{\mu_n}(1) \\ \vdots \\ J_{\mu_n}(n) \end{bmatrix}$   $q = \begin{bmatrix} q(1,\mu_n(1)) \\ \vdots \\ q(n,\mu_n(n)) \end{bmatrix}$   $\underline{P} = \begin{bmatrix} P_{11}(\mu_n(1)) \cdots P_{1n}(\mu_n(1)) \\ \vdots \\ P_{n1}(\mu_n(n)) \cdots P_{nn}(\mu_n(n)) \end{bmatrix}$ )

           $\hookrightarrow \underline{J_{\mu_n}} = (\underline{\mathbb{I}} - \underline{P})^{-1} q$

$\to$ linear program equivalent to value iteration: $J^* = \max_V \sum_{i=0}^n V(i)$   subj.to: $V(i) \leq (q(i,u) + \sum_{j=0}^n P_{ij}(u) V(j))$   $\forall u, \forall i$

$\to$ discounted problem: ( solve an auxiliary prob. to get solution to discounted prob. )

discounted SSP + BE:          $\to$ auxiliary SSP + BE:

$\tilde{x}_{k+1} = \tilde{w}_k$   ( $\tilde{x}_k \in S^+ = \{1,...,n\}$ , $\tilde{u}_k \in \tilde{U}(\tilde{x}_k)$ )    $\to$ $x_k \in S = S^+ \cup \{0\}$ , $U(0) = \{stay\}$    add aux terminal state

$P_{\tilde{w}_k|\tilde{x}_k \tilde{u}_k}(j|i,u) = \tilde{P}_{ij}(u)$  $\leftarrow$ terminal state not needed     $\to$ $P_{ij}(u) = \alpha \tilde{P}_{ij}(u)$ , $P_{i0}(u) = 1-\alpha$ , $P_{0j}(u=stay) = 0$ , $P_{00}(u=stay) = 1$

$\tilde{g}(\tilde{x}_k, \tilde{u}_k, \tilde{w}_k)$                        $\to$ $g(x_k, u_k, w_k) = \alpha^{-1} \tilde{g}(x_k, u_k, w_k)$ , $g(x_k, u_k, 0) = 0$ , $g(0, stay, 0) = 0$

$\tilde{q}(\tilde{x}_k, \tilde{u}_k) = E\{\tilde{g}\} = \sum_{j=1}^n \tilde{P}_{ij}(u) \tilde{g}(i,u,j)$      $\to$ $q(x_k, u_k) = E\{g\} = \sum_{j=0}^n P_{ij}(u) g(i,u,j) = ... = \tilde{q}(x_k, u_k)$ , $q(0, u_k) = 0$

$\tilde{J}^*(\tilde{x}_0) = \min_\pi E\{\sum_{k=0}^{N-1} \alpha^k \tilde{g}(\tilde{x}_k, \tilde{u}_k = \tilde{\mu}_k(\tilde{x}_k), \tilde{w}_k)\}$   $\to$ $J^*(x_0) = \min_\pi E\{\sum_{k=0}^{N-1} g(x_k, u_k = \mu_k(x_k), w_k)\}$

$\tilde{J}_\infty(i) = \min_{\tilde{\mu}_\infty}[\tilde{q}(i, u=\tilde{\mu}_\infty(i)) + \alpha \sum_{j=1}^n \tilde{P}_{ij}(u=\tilde{\mu}_\infty(i)) \cdot \tilde{J}_\infty(j)]$   $\to$ $J_\infty(i) = \min_{\mu_\infty}[q(i, u=\mu_\infty(i)) + \sum_{j=0}^n P_{ij}(u=\mu_\infty(i)) \cdot J_\infty(j)]$

                                        $\hookrightarrow \underline{\mu_\infty(i) = \tilde{\mu}_\infty(i)}$   $\forall i \neq 0$
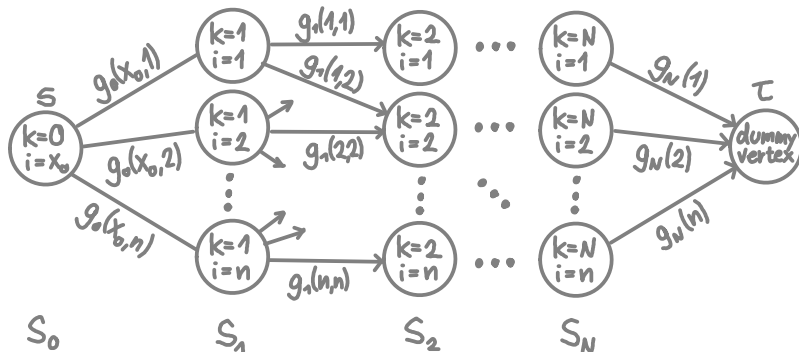
# deterministic finite state problem (DFS)

system: $x_{k+1} = u_k$    $x_k \in S_k, u_k \in U_k(x_k)$

objective: $J^*(x_0) = \min_\pi g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k = \mu_k(x_k))$

$\begin{pmatrix} j: discr.\ state\ x_{k+1}\ \#i \\ i: discr.\ state\ x_k\ \ \ \#j \\ U: discr.\ input\ u_k\ \#U \end{pmatrix}$

- ## DFS to equivalent SP:



$S_0$      $S_1$      $S_2$      $S_N$

- ## SP to equivalent DFS:

set $c_{i,i} = 0\ \forall i \in V$, search for paths of len. $N = |V|$:

$S_0 = \{s\}$ ; $S_k = V \setminus \{\tau\}$ for $k = 1, .., N-1$ ; $S_N = \{\tau\}$

$U_k = V \setminus \{\tau\}$ for $k = 0, .., N-2$ ; $U_{N-1} = \{\tau\}$

$g_N(\tau) = 0$ ; $g(x_k, u_k) = \begin{cases} 0 & \text{if } x_k = u_k \\ \infty & \text{if } \nexists c_{x_k, u_k} \\ c_{x_k, u_k} & \text{otherwise} \end{cases}$

$\boxed{\text{solve with DPA} \to \text{skip degenerate moves}}$

- ## DPA + DFS problem:
$J_k(x_k) = \min_{\mu_k} g_k(x_k, u_k = \mu_k(x_k)) + J_{k+1}(x_{k+1} = \mu_k(x_k))$   initial cond: $J_N(x_N) = g_N(x_N)$

- ## forward DPA for DFS derived from SP (works for all DFS)

optimal path from $s$ to $\tau$ is also optimal path from $\tau$ to $s$ if all edges are flipped → formulate aux. SP:

$\tilde{c}_{j,i} = c_{i,j}\ \ \forall (i,j,c_{i,j}) \in C \longrightarrow \tilde{J}_N(s) = 0 \to \tilde{J}_{N-1}(j) = \tilde{c}_{j,s} \to ... \to \tilde{J}_k(j) = \min_i \tilde{c}_{j,i} + \tilde{J}_{k+1}(i) \to ... \to \tilde{J}_0(\tau) = \min_i \tilde{c}_{\tau,i} + \tilde{J}_1(i)$

$J_l^F = \tilde{J}_{N-l}$ ; $l = N-k$    $\tilde{J}_N(s) = 0 \to \tilde{J}_{N-1}(j) = c_{s,j} \to ... \to \tilde{J}_k(j) = \min_i c_{i,j} + \tilde{J}_{k+1}(i) \to ... \to \tilde{J}_0(\tau) = \min_i c_{i,\tau} + \tilde{J}_1(i)$

$\boxed{J_l^F(j): \text{optimal cost-to-arrive}}$    $J_0^F(s) = 0 \to J_1^F(j) = c_{s,j} \to ... \to J_l^F(j) = \min_i c_{i,j} + J_{l-1}^F(i) \to ... \to J_N^F(\tau) = \min_i c_{i,\tau} + J_{N-1}^F(i)$

# hidden Markov model + Viterbi algorithm

hidden Markov model: ————

↳ given measurements $Z = (z_1, ..., z_N)$, find most likely state trajectory $X = (x_0, ..., x_N)$:

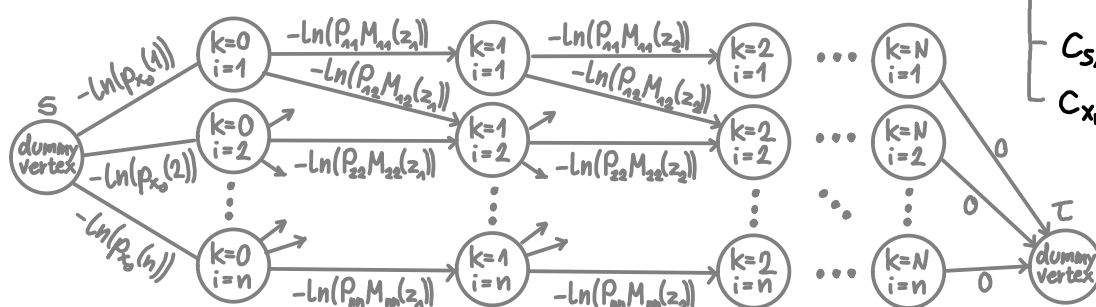$\begin{cases} x_{k+1} = w_k & x_k \in S \text{ discrete states} \\ P_{ij} = p_{w|x}(j|i) & \text{prob. of moving from } i \text{ to } j \\ M_{ij}(z) = p_{z|x,w}(z|i,j) & \text{prob. of measuring } z \text{ when moving from } i \text{ to } j \end{cases}$

$\hat{X} = \operatorname*{argmax}_X p(X|Z) = \operatorname*{argmax}_X p(X,Z) = \operatorname*{argmax}_X p(x_0) \cdot \prod_{k=1}^N P_{x_{k-1}, x_k} M_{x_{k-1}, x_k}(z_k) = \operatorname*{argmin}_X c_{s,x_0} + \sum_{k=1}^N c_{x_{k-1}, x_k}$

$c_{s,x_0} = -\ln(p(x_0))$

$c_{x_{k-1}, x_k} = -\ln(P_{x_{k-1}, x_k} M_{x_{k-1}, x_k}(z_k))$

# shortest path algorithms

- dynamic programming algorithm (DPA): convert SP to DFS and solve with DPA  (see DFS chapter for details)
  ↳ inefficient because DPA finds shortest path from _any_ vertex to $\tau$, not just from s.

- **label correcting algorithm (LCA):**   (same as forward DPA with DFS, but ignore branches with higher cost-to-arrive than current guess of cost-to-arrive at $\tau$.)

$d_s = 0$ ; $d_j = \infty$ $\forall j \in V \backslash \{s\}$ ; OPEN $= \{s\}$

while OPEN $\neq \{\}$
| remove a node i from OPEN ⟵
| for all children j of i  $(c_{i,j} \neq \infty)$
| | if $(d_i + c_{i,j}) < d_j$ and $(d_i + c_{i,j}) < d_\tau$
| | | set $d_j = d_i + c_{i,j}$
| | | if $j \neq \tau$, then add j to OPEN

multiple ways to pick node i from OPEN:
- depth-first search: last in, first out
- breath-first search: first in, first out
- best-first search: pick i with smallest $d_i$  (Dijkstras Algorithm)

result: $L_{Q^*} = d_\tau$   $Q^* =$ parent of $\tau$ that last was in OPEN
→ recursion until s

- **A*- algorithm:** if some positive lower bound $h_j$ is known for the optimal path-lenght from j to $\tau$, then replace $d_i + c_{i,j} < d_\tau$ with $d_i + c_{i,j} + h_j < d_\tau$ in LCA.  (e.g. $h_j =$ "air-line distance")

# deterministic continuous time optimal control

system: $\dot{x}(t) = f(x(t), u(t))$   $0 \leq t \leq T$     cost: $h(x(T)) + \int_0^T g(x(\tau), u(\tau)) d\tau$

optimal cost to go: $J^*(t, x(t)) = \min_\mu h(x(T)) + \int_t^T g(x(\tau), u = \mu(\tau, x))$  subj. to: $\dot{x}(t) = f(x(t), u(t))$     $u \in \mathcal{U}, x \in \mathcal{X}$
  ↳ optimal cost: $J^*(0, x_0)$  ; optimal policy: $u = \mu^*(t, x)$

- **Hamilton-Jacobi-Bellman Equation (HJB):**     (can be derived directly from DPA with step size $T_s \to 0$)

$\min_\mu g(x(t), u = \mu(t, x)) + \partial J^* / \partial t |_{x(t)} + \partial J^* / \partial x |_{x(t)} \cdot f(x(t), u = \mu(t, x(t))) = 0$  $\forall t \in [0, T], x$   and   $J^*(T, x(T)) = h(x(T))$ $\forall x$
  ↳ hard to solve. easier:  guess $\mu(t, x)$ → calculate $J_\mu(t, x)$ → if HJB eq. is fulfilled: $\mu^* = \mu$ $J^* = J_\mu$ optimal!

- **Pontryagin's minimum Principle:** (easier to solve than HJB, but only necessary cond. on $\mu$, not sufficient)
  (Hamiltonian function: $H(x, u, p) = g(x, u) + p^T f(x, u)$ with p: some auxiliary variable)

$u(t) = \underset{u}{\text{argmin}} H(x(t), u, p(t))$       ODE 1: $\dot{p}(t) = -\partial H / \partial x^T |_{x(t), u(t), p(t)}$       BC 1: $p(T) = \partial h / \partial x^T |_{x(T)}$

$H(x(t), u(t), p(t)) = $ const. $\forall t \in [0, T]$       ODE 2: $\dot{x}(t) = f(x(t), u(t))$       BC 2: $x(0) = x_0$

usual approach: solve ODE1, ODE2, argmin assuming x, u, p are known → reformulate such that each eq. only depends on one of x, u, p → apply BC1, BC2 to get rid of integration constants

$\Bigl[$ extensions:  – fixed terminal state: replace BC1 with $x(T) = x_T$     ($x_T$: terminal state)
  – free initial state: replace BC2 with $p(0) = -\partial L / \partial x^T |_{x(0)}$   ($L(x)$: initial cost)
  – free terminal time: solve for T with cond $H(x(t), u(t), p(t)) = 0$
  – time varying f, g: H becomes func. of t. <u>drop</u> cond. $H(x(t), u(t), p(t), t) = $ const.
  – singular problem: if $\underset{u}{\text{argmin}} H$ is undefined, assume $\dot{p} = 0$ over that interval $\Bigr]$