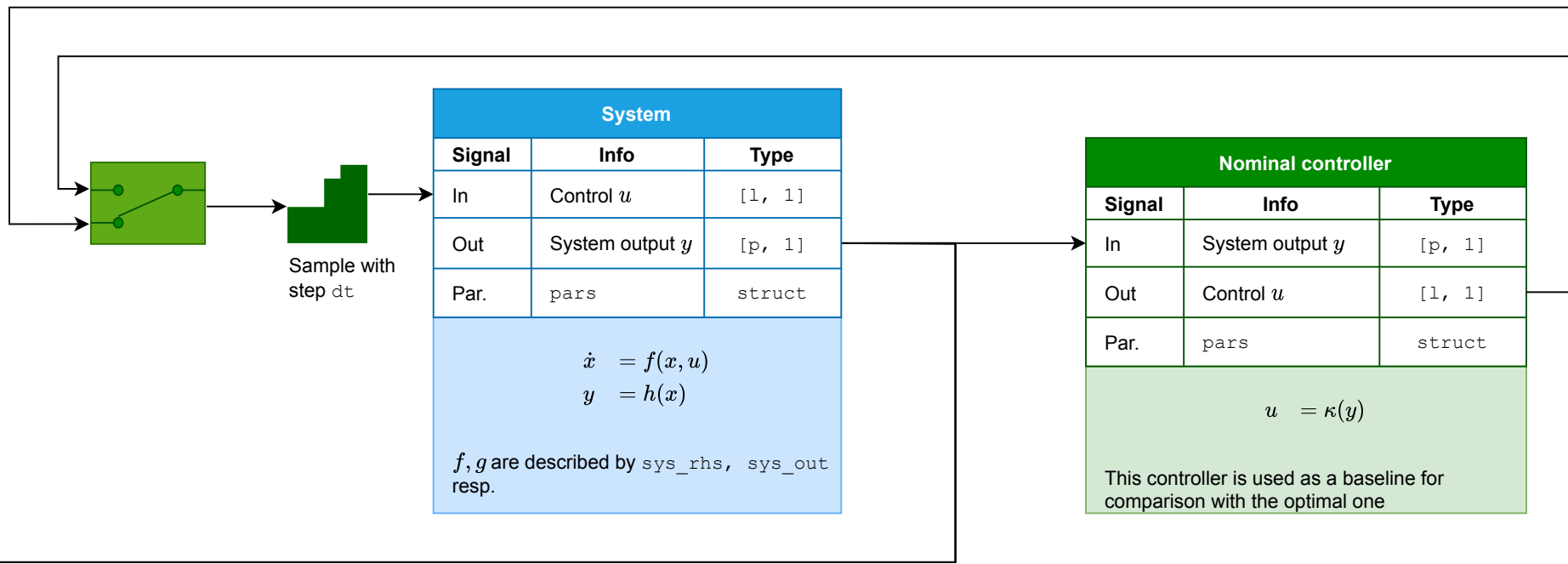


| Initialization <code>init.m</code> |  |
|------------------------------------|--|
| <code>pars</code>                  | Parameter structure for simulation environment |
| <code>Tsim</code>                  | Total simulation time for one run              |
| <code>Nruns</code>                 | Number of runs                                 |

| Key fields in <code>pars</code> |                                 |
|---------------------------------|---------------------------------|
| <code>estBufferSize</code>      | Buffer size for model estimator |
| <code>optCtrlMode</code>        | Optimal controller mode         |
| <code>ctrlStackSize</code>      | $N$                             |
| <code>criticStackSize</code>    | $M$                             |
| <code>rcostS, rcostR</code>     | $S, R$                          |

| <code>optCtrlMode</code>                                  |  |
|---|--|
| 1, 2 - model-predictive control (MPC)                     | $J(y_1, \{u\}_1^N) = \sum_{k=1}^N r(y_k, u_k)$                         |
| 3, 4 - RL/ADP via stacked Q-learning                      | $J(y_1, \{u\}_1^N) = \sum_{k=1}^N \hat{Q}(y_k, u_k)$                   |
| 5, 6 - RL/ADP via $N$ roll-outs of $r$                    | $J(y_1, \{u\}_1^N) = \sum_{k=1}^{N-1} r(y_k, u_k) + \hat{Q}(y_N, u_N)$ |
| Modes 1, 3, 5 use true model $(f, h)$ for prediction      |  |
| Modes 2, 4, 6 use an a state-space model estimated online |  |

| Legend               |   |
|----------------------|---|
| $[\bullet, \bullet]$ | a matrix  |
| <code>struct</code>  | a structure                                     |
| $r$                  | running cost                                    |
| $N$                  | control horizon                                 |
| $M$                  | critic stack size                               |
| $J$                  | controller cost function                        |
| $J_c$                | critic cost function                            |
| $Q, \hat{Q}$         | Q-function and its approximate                  |
| $e$                  | temporal difference                             |
| $W, W^-$             | current and previous critic weights             |
| $R, S$               | matrices in $r(y, u) = y^\top S y + u^\top R u$ |
| $\gamma$             | discounting factor                              |
| $L$                  | number of critic weights                        |
| $\Delta t$           | controller sampling time                        |



| Optimal controller <code>dataDrivOptCtrl</code> |                   |          |
|---|-------------------|----------|
| Signal  | Info              | Type     |
| In  | System output $y$ | $[p, 1]$ |
| Out   | Control $u$       | $[1, 1]$ |
| Par.  | <code>pars</code> | struct   |

| Model estimator <code>mySSest.m</code>   |   |   |
|--|---|---|
| Signal   | Info  | Type  |
| In   | Buffers of previous controls and measured outputs $u_s, y_s$  | $[1, \text{estBufferSize}]$<br>$[p, \text{estBufferSize}]$  |
| In   | Sampling time $\Delta t$  | scalar $> 0$  |
| In   | <code>modelOrder</code>   | natural   |
| In   | Initial state-space model ( $A_{\text{init}}, B_{\text{init}}, C_{\text{init}}, D_{\text{init}}$ )                  | $[\text{modelOrder}, \text{modelOrder}]$<br>$[\text{modelOrder}, 1]$<br>$[p, \text{modelOrder}, 1]$<br>$[p, 1]$                             |
| Out  | Par. of estimated state-space model $(A, B, C, D)$ of desired order <code>modelOrder</code> and initial state $x_0$ | $[\text{modelOrder}, \text{modelOrder}]$<br>$[\text{modelOrder}, 1]$<br>$[p, \text{modelOrder}, 1]$<br>$[p, 1]$<br>$[\text{modelOrder}, 1]$ |
| Fits a state-space model   |   |   |
| $\hat{x}^+ = A\hat{x} + Bu$ $y^+ = C\hat{x} + Du,$   |   |   |
| from data $u_s, y_s$ that are organized in buffers of previous controls and measured outputs |   |   |

| Critic <code>critic.m</code>  |  |                      |
|---|--|----------------------|
| Signal  | Info   | Type                 |
| In  | Buffers of previous controls and measured outputs $u, y$ | $[1, M]$<br>$[p, M]$ |
| In  | Running cost par. $S, R$                                 | $[p, p]$<br>$[1, 1]$ |
| In  | Discounting factor $\gamma$                              | scalar in $(0, 1]$   |
| Out   | Critic weights and cost $W, J_c$                         | $[L, 1]$<br>scalar   |
| Critic solves the problem   |  |                      |
| $\min_W J_c = \frac{1}{2} \sum_{k=1}^M e^2$   |  |                      |
| where $e = W^\top \varphi(y^-, u^-) - \gamma W^- \varphi(y, u) - r(y^-, u^-)$ and $(y, u), (y^-, u^-)$ are the current and, resp., previous output and control taken from the buffers $u, y$ . The Q-function approximate is described by |  |                      |
| $\hat{Q}(y, u) = W^\top \varphi(y, u)$  |  |                      |

| Actor <code>optCtrl.m</code>   |   |   |
|--|---|---|
| Signal   | Info  | Type  |
| In   | System output $y$   | $[p, 1]$  |
| In   | Running cost par. $S, R$  | $[p, p]$<br>$[1, 1]$  |
| In   | Critic weights $W$  | $[1, 1]$  |
| In   | Sampling time $\Delta t$  | scalar $> 0$  |
| In   | Par. of estimated state-space model $(A, B, C, D)$ of desired order <code>modelOrder</code> and initial state $x_0$ | $[\text{modelOrder}, \text{modelOrder}]$<br>$[\text{modelOrder}, 1]$<br>$[p, \text{modelOrder}, 1]$<br>$[p, 1]$ |
| Out  | Control $u$ and actor cost function $J$   | $[1, 1]$<br>scalar  |
| Actor solves the problem   |   |   |
| $\min_U J(y, U)$ $\text{s.t. } y^+ = \mathcal{M}(\bullet, u)$  |   |   |
| w.r.t. to a <b>sequence</b> of controls $U = \{u\}_1^N$ over a horizon of $N$ (see definition of $J$ in <code>optCtrlMode</code> ) starting with the current output $y$ . Here, $\mathcal{M}$ is a model that predicts each next output from the current one along with the control. According to the mode <code>optCtrlMode</code> , $\mathcal{M}$ can be either (a discretized) version of $(f, h)$ , i. e., |   |   |
| $x^+ = x + \delta f(x, u)$ $y^+ = h(x),$   |   |   |
| where $\delta$ is $\Delta t$ ; or an estimated state-space model   |   |   |
| $\hat{x}^+ = A\hat{x} + Bu$ $y^+ = C\hat{x} + Du,$   |   |   |