

Research Idea

Xu Yunpeng

November 2, 2024

Uniformly Quantized State Variables as Training Samples Let $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ be a discrete-time filtered probability space with a finite discrete index set $N = \{0, 1, \dots, T\}$.

Consider a vector of n controlled Markovian processes $X_{i,t}$ with Gaussian noise $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \Sigma)$ live in \mathbb{R}^n , where the control variable is denoted by u live in the set \mathbb{U} , and the control at a time t is $u_t \in \mathbb{U}_t$, $\mathbb{U}_t \subseteq \mathbb{U}$:

$$\mathbf{Y}_t = (X_{1,t}, \dots, X_{n,t})^\top.$$

$$\mathbf{Y}_{t+1} = g(\mathbf{Y}_t, u_t, \epsilon_{t+1})$$

Here $g : \mathbb{R}^n \times \mathbb{U} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ establishes the dynamics of \mathbf{Y}_t .

Denote the realisation of \mathbf{Y}_t as \mathbf{Y} , and consider functions that do the following mappings, reward functions $f_1 : \mathbb{R}^n \times \mathbb{U} \rightarrow \mathbb{R}$, $f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$, value function $J : N \times \mathbb{R}^n \times \mathbb{U} \rightarrow \mathbb{R}$, and optimal value function $V : N \times \mathbb{R}^n \rightarrow \mathbb{R}$,

$$J_t(\mathbf{Y}, u) = \mathbb{E} \left(\sum_{s=t}^{T-1} f_1(\mathbf{Y}_s, u) + f_2(\mathbf{Y}_T) \middle| \mathcal{F}_t \right)$$

$$V_t(\mathbf{Y}) = \sup_{\forall u \in \mathbb{U}} \{ J_t(\mathbf{Y}, u) \} \tag{1}$$

$$= \sup_{\forall u \in \mathbb{U}_t} \left\{ f_1(\mathbf{Y}, u) + \mathbb{E} \left(V_{t+1} \left(g(\mathbf{Y}, u, \epsilon_{t+1}) \right) \middle| \mathcal{F}_t \right) \right\}, t \leq T-1 \tag{2}$$

$$V_T(\mathbf{Y}) = f_2(\mathbf{Y}) \tag{3}$$

where the closed-form expressions of f_1, f_2, g are known. The aim is to find the surrogate of V_t using the above dynamic programming set up, so that given a sample path, one can use the surrogate to find the optimal control according to Equation 1. The closed-form expression of V_{T-1} can be very challenging to find, but we can input a training sample $(\mathbf{y}_1, \dots, \mathbf{y}_p)$ generated by random uniform with suitable ranges, optimize $f_1(\mathbf{y}_i, u) + \mathbb{E}(V_T(g(\mathbf{y}_i, u, \epsilon_T)))$ using existing algorithms, denote the optimized value of it as \hat{v}_i , and gather p number of pairs $(\mathbf{y}_i, \hat{v}_i), i = 1, \dots, p$. We can then do a feedforward neural network training to find an approximated mapping (a surrogate) of V_{T-1} , denoted as \hat{V}_{T-1} , and we can repeat such procedure to find $V_t, t \leq T-2$ by replacing V_{t+1} in Equation 1 by \hat{V}_{t+1} .

The problem is whether if we can find the optimal training sample such that \hat{V} is consistent for V . Standing at time $T - 1$, if one can find a m level optimal uniform quantization of \mathbf{Y}_{T-1}

$$\tilde{\mathbf{y}}_{T-1}(m) = \begin{pmatrix} \tilde{\mathbf{y}}_{1,T-1} \\ \tilde{\mathbf{y}}_{2,T-1} \\ \vdots \\ \tilde{\mathbf{y}}_{m,T-1} \end{pmatrix} \text{ with } \begin{pmatrix} \frac{1}{m} \\ \frac{1}{m} \\ \vdots \\ \frac{1}{m} \end{pmatrix}, \quad (4)$$

such that $\tilde{\mathbf{y}}_{T-1}(m) \xrightarrow{d} \mathbf{Y}_{T-1}$ as $m \rightarrow \infty$, then we shall treat the $\tilde{\mathbf{y}}_{T-1}(m)$ as training samples, and at the same time a discrete uniform random vector.

Assuming there exist an optimal feedforwork neural network parameter (node number, number of layers, weights and biases) such that the MSE can equal 0 at least in 2 decimal places, standing at $T - 1$, we can use a similar method to optimize $f_1(\tilde{\mathbf{y}}_{i,T-1}, u) + \mathbb{E}(V_T(g(\tilde{\mathbf{y}}_{i,T-1}, u, \epsilon_T)))$ using existing algorithms, denote the optimized value of it as \hat{v}_i , and gather m number of pairs $(\tilde{\mathbf{y}}_{i,T-1}, \hat{v}_i), i = 1, \dots, m$. We can then do a feedforward neural network training to find an approximated mapping (a surrogate) of V_{T-1} , denoted as \hat{V}_{T-1} . By the continuous mapping theorem, we have

$$\lim_{m \rightarrow \infty} \mathbb{E} \left(\left(\hat{V}_{T-1}(\tilde{\mathbf{y}}_{T-1}(m)) - V_{T-1}(\tilde{\mathbf{y}}_{T-1}(m)) \right)^2 \right) \xrightarrow{d} \mathbb{E} \left((\hat{V}_{T-1}(\mathbf{Y}_{T-1}) - V_{T-1}(\mathbf{Y}_{T-1}))^2 \right)$$

by passing the limit inside the expectatoin. Thus, based on the below

$$\mathbb{E} \left((\hat{V}_{T-1}(\tilde{\mathbf{y}}_{T-1}(m)) - V_{T-1}(\tilde{\mathbf{y}}_{T-1}(m)))^2 \right) = \sum_{j=1}^m \frac{1}{m} \cdot (\hat{V}_{T-1}(\tilde{\mathbf{y}}_{j,T-1}) - v_j)^2 = 0$$

we can conclude $\text{Var}(\hat{V}_{T-1}(\mathbf{Y}_{T-1}) - V_{T-1}(\mathbf{Y}_{T-1}))$ and $\mathbb{E}(\hat{V}_{T-1}(\mathbf{Y}_{T-1}) - V_{T-1}(\mathbf{Y}_{T-1}))$ are “asymptotically” 0, therefore we can say \hat{V}_{T-1} is an “asymptotically” consistent esitmate of V_{T-1} . The procedure can then be repeated to find $V_t, t \leq T - 2$ by replacing V_{t+1} in Equation 1 by \hat{V}_{t+1} .