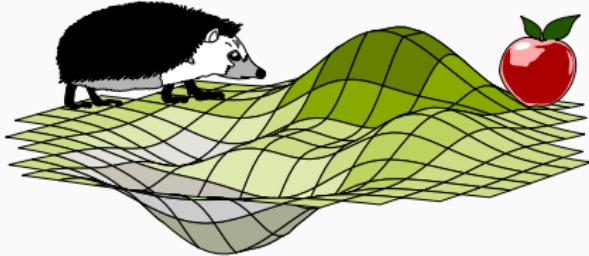


REINFORCEMENT LEARNING

PREDICTIVE ALGORITHMS

Pavel Osinenko





The picture of reinforcement learning is incomplete without **model - predictive control**, which is in most cases a synonym of online, finite-horizon optimal control. The related problem also, like the ∞ -horizon one, admits an HJB (although time-dependent). However, unlike the ∞ -horizon optimal control problem, the finite-horizon one is computationally tractable.

Model - predictive control (MPC) is listed within some popular RL libraries (such as SLM Lab) and has a number of advantages



A fairly general MPC setup reads:

$$\min_{\{u_{i|k}\}_i^N} \mathbb{E} \left[\sum_{i=1}^N p(\hat{X}_{i|k}, u_{i|k}) \mid X_{1|k} = x_k \right]$$

$$\text{s.t. } dX_t = f(X_t, u^\delta)dt + g(X_t, u^\delta)dB_t,$$

$$u^\delta(t) \equiv u_{i|k}, t \in [(k+i-1)\delta, (k+i)\delta),$$

$$\hat{X}_{i+1|k} = Q^\delta(\hat{X}_{i|k}, u_{i|k}),$$

$$\mathbb{P}[X_{s(k+i-1)} \in \mathcal{X}] \geq 1 - \varepsilon_x,$$

$$u_{i|k} \in \mathcal{U}$$



(cont.)

Description:

- sequence notation: $\{\mathbf{o}_{ilk}\}_i^N = \{\mathbf{o}_k, \dots, \mathbf{o}_{k+N-1}\}$
 $\mathbf{o}_{ilk} = \mathbf{o}_{k+i-1}$
- N : horizon (length)
- u^δ : sampled control that is piece-wise constant on δ -intervals of time
- \hat{Q}^δ : a digital model of the environment that produces estimated states
 \hat{x}_{ilk} at time nodes $(k+i-1)\delta$



(cont.)

If the environment description, i.e., the functions f, \mathcal{G} , is known, then \underline{Q}^δ may, e.g., be the Euler-Maruyama numerical solver:

$$\underline{Q}^\delta(\hat{X}_{i|k}, u_{i|k}) = \hat{X}_{i|k} + f(\hat{X}_{i|k}, u_{i|k})\delta + \mathcal{G}(\hat{X}_{i|k}, u_{i|k})\Delta B_{i|k},$$

where $\{\Delta B_{i|k}\}_i$ are the Wiener increments, i.e., i.i.d. normal zero-mean random variables with variance δ



(cont.)

Otherwise, an environment model may be learned from experience during exploration

- \mathcal{X}, \mathcal{U} are the sets of allowed states and, respectively, actions. These could, e.g., be interpreted as **safe** state and action sets. The expression $u_{ik} \in \mathcal{U}, \forall i \in [N]$ is also called **input constraint** (here, $[N]$ means $\{1, \dots, N\}$). The probabilistic type of a **state constraint** $P[X_{s(k+i-1)} \in \mathcal{X}] \geq 1 - \epsilon_x$ is also called **chance** state constraint with a user-defined tolerance ϵ_x .



MPC has a number of advantages compared to RL, especially in its online neural-network-based variants. Namely, it has an apparatus to guarantee constraint satisfaction at all times (this is related to the so-called recursive feasibility in the MPC theory). That is a strong feature that makes MPC an industrial standard. Furthermore, MPC has a strong theory of environment stabilization, whereas the latter is a major pain for online RL.

In the following, we will study how such a stabilization works on the example of an MDP with a specific structure



Consider the following particular form of an MDP :

$$X_{k+1} = f(X_k, U_k, W_k),$$

$f(0, 0, 0) = 0$, $W_k \in \mathcal{W} \subseteq \mathbb{R}^d$ are i.i.d., 0-mean
with a PDF P_W

State prediction is done via a direct use of the ground truth model (although this is not the only option) :

$$X_{i+1|k} = f(X_{i|k}, U_{i|k}, W_{i|k})$$

Reference:

Lorenzen, M., Müller, M. A., & Allgöwer, F. (2019). Stochastic model predictive control without terminal constraints. International Journal of Robust and Nonlinear Control, 29(15), 4987-5001.

Predictive algorithms



Denote: $P_\infty := \bigotimes_{k=0}^{\infty} P_k$,

$$P_k(\bullet) := \mathbb{P}[\bullet | X_k = x_k],$$

$$\mathbb{E}_k[\bullet] := \mathbb{E}[\bullet | X_k = x_k].$$

$P_\infty(\bullet) = 1$ will also be abbreviated „ P_∞ a.s.”

We will assume input and state constraints in the format described above.

Also, the stage cost ρ is assumed continuous, positive-definite and radially unbounded.

Let's denote the objective as:

$$J_N(x_0 | \{u_{ik}\}_{i=1}^N) := \mathbb{E}_k \left[\sum_{i=1}^N \rho(X_{ik}, u_{ik}) \right]$$



The corresponding value function reads:

$$V_N(x_k) := \min_{\{u_{ik}\}_i^N} J_N(x_k | \{u_{ik}\}_i^N)$$

s.t. system dynamics,
state prediction,
input & state constraints

Let's denote the corresponding minimizer as $\{u_{ik}^*\}_i^N$. Accordingly, the state sequence will be denoted $\{X_{ik}^*\}_i^N$.

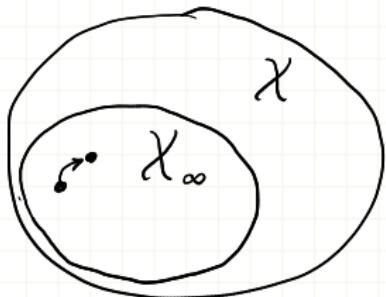


The MPC policy is $\kappa_N(x_k) := u_{1:k}^* = u_k^*$.

For the forthcoming analysis to work, we need the following **viability** assumption:

there is a set $X_\infty \subseteq X$ s.t. $\forall x_k \in X_\infty$,
 $\exists N$ the value function is attained (i.e., a solution $\kappa_N(x_k)$ to the MPC problem exists), and
 $\forall w \in W \quad f(x_k, \kappa_N(x_k), w) \in X_\infty$.

Such a property that the future state lie in the same set is called **positive invariance**. Under control and noise, it's **robust controlled invariance**





The idea of viability is that inside X_∞ , there is always a way (via MPC) to stay inside, no matter what value the noise assumes.

The next assumption describes the properties of the stage cost. First, some notation:

consider, for any N ,

$$\{v_{i|0}^*\}_{i=1}^N := \arg \min_{\{u_{i|0}\}_{i=1}^N} J(o | \{u_{i|0}\}_{i=1}^N),$$

the minimizer at zero initial state. Let the corresponding state sequence be $\{Z_{i|0}^*\}_{i=1}^N$.

Denote

$$\bar{p}_i^N := \mathbb{E}_o \left[\rho(Z_{i|0}^*, v_{i|0}^*) \right],$$

the i th expected stage cost

Predictive algorithms



(cont.)

With this at hand, we assume:

$$1. \bar{p}_1^N = 0$$

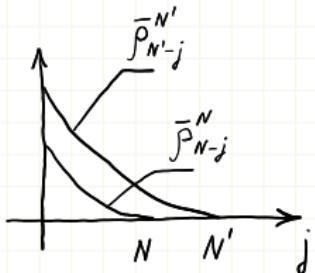
$$2. \forall N' \geq N \quad \forall j \quad \bar{p}_{N-j}^N \leq \bar{p}_{N'-j}^{N'}$$

$$3. \forall x_k \quad \mathbb{E}_k [\rho(x_{ik}^*, u_{ik}^*)] \geq \bar{p}_i^N$$

4: If p.-d., continuous, radially unbounded function $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ and $\alpha > 0$ s.t.

$$\sqrt{N}(x) \leq \alpha \varphi(x) + \sum_{i=1}^N \bar{p}_i^N,$$

$$\rho(x, u) \geq \varphi(x), \quad \forall x \in \mathbb{R}^n, u \in \mathcal{U}$$





We will utilize the following variant of the DPP:

let $V_0(x) = 0$, then, for $l \in \mathbb{N}$, $x \in \mathcal{X}_\infty$, it holds that

$$V_l(x) = \min_u \left\{ g(x, u) + \mathbb{E}_{P_W} [V_{l-1}(X_+^u)] \right\}$$

s.t. $u \in \mathcal{U}$

$$\mathbb{P}[X_+^u \in \mathcal{X}] \geq 1 - \varepsilon$$

This can be shown by recursion. Recall we are assuming attainability of the value functions.

So, for brevity, we will omit the constraints.

First, for $l=1$, we have the relation trivially

Predictive algorithms



(cont.)

Take $l=2$. We have:

$$\bar{V}_1(x) = \min_u p(x, u) ,$$

$$\bar{V}_2(x) = \min_{u_0, u_1} \left\{ p(x, u_0) + \mathbb{E}_{P_w} \left[p(X_+^{u_0}, u_1) \mid X=x \right] \right\} .$$

By optimality of \bar{V}_1 ,

$$\min_{u_0, u_1} \left\{ p(x, u_0) + \mathbb{E}_{P_w} \left[p(X_+^{u_0}, u_1) \mid X=x \right] \right\} \geq$$

$$\min_{u_0} \left\{ p(x, u_0) + \mathbb{E}_{P_w} \left[\bar{V}_1(X_+^{u_0}) \mid X=x \right] \right\}$$

Predictive algorithms



(cont.)

On the other hand, by optimality of u_0^*, u_1^* , no matter what tail action we take — be it the minimizer corresponding to V_1 — it must hold that

$$\min_{u_0, u_1} \left\{ p(x, u_0) + \mathbb{E}_{P_w} \left[p(X_+^{u_0}, u_1) \mid X=x \right] \right\} \leq \min_{u_0} \left\{ p(x, u_0) + \mathbb{E}_{P_w} \left[V_1(X_+^{u_0}) \mid X=x \right] \right\},$$

whence the required identity

Predictive algorithms



(cont.)

Suppose the identity holds for $l-1$, so

$$V_{l-1}(x) = \min_u \left\{ p(x, u) + \mathbb{E}_{P_w} \left[V_{l-2}(X^u_+) \right] \right\}.$$

Unwrap the next one:

$$\begin{aligned} V_l(x) &= \min_{\{u_i\}_{i=0}^{l-1}} \left\{ \mathbb{E} \left[\sum_{i=0}^{l-1} p(X_{+i}, u_i) \mid X = x \right] \right\} \\ &= \min_{\{u_i\}_{i=0}^{l-1}} \left\{ p(x, u_0) + \mathbb{E} \left[\sum_{i=1}^{l-1} p(X_{+i}, u_i) \mid X = x \right] \right\}, \end{aligned}$$

and apply the same optimality argument observing that

$$V_{l-1}(X^u_+) = \min_{\{u_i\}_{i=1}^{l-1}} \mathbb{E} \left[\sum_{i=1}^{l-1} p(X_{+i}, u_i) \mid X = x \right]$$

Predictive algorithms



(cont.)

As a corollary, we may state the DPP with an arbitrary head:

$$\bar{V}_N(x_k) = \min_{\{u_{i|k}\}_{i=1}^{N-l}} \left\{ \mathbb{E} \left[\sum_{i=1}^{N-l} p(x_{i|k}, u_{i|k}) + \bar{V}_l(x_{N-l+1|k}) \right] \right\},$$

for any $l \in \mathbb{N}$, $l < N$, $x_k \in \mathcal{X}_\infty$.

We will use this version in what follows



Now, we can state the **main stability result**:

under the stated assumptions, let $x_0 \in \mathcal{X}_\infty$. Then:

- $\forall k$:
1. the constraints are satisfied
 2. $\exists C_1 \geq 0$ s.t.

$$E_k [\sqrt[N]{(X_{k+1}^*)}] - \sqrt[N]{(x_k)} \leq -\rho(x_k, u_k^*) + \frac{\alpha^2}{N-1} \nu(x) + C_1.$$

Furthermore, $\exists N \in \mathbb{N}$ s.t. $\forall x_k \in \mathcal{X}_\infty$

$$E_k [\sqrt[N]{(X_{k+1}^*)}] \leq \lambda \sqrt[N]{(x_k)} + C_2$$

with $\lambda \in (0, 1)$, $C_2 \geq 0$ and

$$\sqrt[N]{(x)} := \sqrt[N]{(x)} - \sum_{i=1}^N \bar{\rho}_i^N$$



To prove this claim, let the current state to have assumed a value x_k and consider the expected difference in value:

$$\mathbb{E}_k[V_N(X_{k+1}^*)] - V_N(x_k).$$

Using the non-trivial head version of the DPP, we may write (for $j < N$):

$$\mathbb{E}_k[V_N(X_{k+1}^*)] = \mathbb{E}_k\left[\mathbb{E}_{k+1}\left[\sum_{i=1}^{N-j-1} p(X_{i|k+1}^*, u_{i|k+1}^*) + V_j(X_{N-j|k+1}^*)\right]\right]$$

Plugging in the optimal action sequence $\{u_{i|k+1}^*\}_i^N$ because it's the minimizer corresponding to $V_N(X_{k+1}^*)$

Predictive algorithms



(cont.)

Next,

$$\begin{aligned}\mathbb{E}_k[V_N(X_{k+1}^*)] &= \mathbb{E}_k\left[\mathbb{E}_{k+1}\left[\sum_{i=1}^{N-j-1} \rho(X_{i|k+1}^*, u_{i|k+1}^*) + V_j(X_{N-j|k+1}^*)\right]\right] \\ &= \mathbb{E}_k\left[\sum_{i=1}^{N-j-1} \rho(X_{i|k+1}^*, u_{i|k+1}^*) + V_j(X_{N-j|k+1}^*)\right]\end{aligned}$$

since the condition of the inner expectation is already determined by the condition of the outer one

Predictive algorithms



(cont.)

Now, we consider, at step $k+1$, a feasible candidate for the head $\{U_{i+1|k}^*\}_{i=1}^{N-j-1}$.

Why is it feasible?

Well, $\{U_{i|k}^*\}_{i=1}^N$ was feasible at step k and we just make a cyclic shift at step $k+1$ remembering that we're inside X_∞ which „neutralizes“ the noise.

With this feasible candidate, we may write:

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} p(X_{i|k+1}^*, U_{i|k+1}^*) + V_j(X_{N-j|k+1}^*) \right] \leq$$

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} p(X_{i+1|k}^*, U_{i+1|k}^*) + V_j(X_{N-j+1|k}^*) \right]$$

(be aware of the difference in the state sequences!)

Predictive algorithms



(cont.)

Now, since $\sum_{i=N-j+1}^N \bar{p}_i^N - \sum_{i=1}^j \bar{p}_i^i$ is non-negative*, we may

add it on the right-hand side of the last inequality:

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} g(X_{i|k+1}^*, u_{i|k+1}^*) + V_j^*(X_{N-j|k+1}^*) \right] \leq$$

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} g(X_{i+1|k}^*, u_{i+1|k}^*) + V_j^*(X_{N-j+1|k}^*) \right] +$$

$$\sum_{i=N-j+1}^N \bar{p}_i^N - \sum_{i=1}^j \bar{p}_i^i$$

* By item 2. of the standing assumption, each $\bar{p}_{N-i}^N \geq \bar{p}_{j-i}^i$

Predictive algorithms



(cont.)

Plugging this back into the expected difference in value and unwrapping the pivot yields:

$$\mathbb{E}_k \left[V_N(X_{k+1}^*) \right] - V_N(x_k) \leq \mathbb{E}_k \left[\sum_{i=1}^{N-j-1} p(X_{i+1|k}^*, u_{i+1|k}^*) + V_j(X_{N-j+1|k}^*) \right] + \\ \sum_{i=N-j+1}^N \bar{p}_i^N - \sum_{i=1}^j \bar{p}_i^i - \mathbb{E}_k \left[\sum_{i=1}^N p(X_{i|k}^*, u_{i|k}^*) \right]$$

By the standing assumption (item 3.),

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} p(X_{i+1|k}^*, u_{i+1|k}^*) \right] + \sum_{i=N-j+1}^N \bar{p}_i^N \leq \mathbb{E}_k \left[\sum_{i=1}^N p(X_{i+1|k}^*, u_{i+1|k}^*) \right]$$

Predictive algorithms

(cont.)

Therefore,

$$\mathbb{E}_k \left[\sum_{i=1}^{N-j-1} \rho(X_{i+1|k}^*, u_{i+1|k}^*) \right] + \sum_{i=N-j+1}^N \bar{\rho}_i^N - \mathbb{E}_k \left[\sum_{i=1}^N \rho(X_{i|k}^*, u_{i|k}^*) \right] \leq -\rho(x_k, u_k^*),$$

whence

$$\mathbb{E}_k \left[V_N(X_{k+1}^*) \right] - V_N(x_k) \leq -\rho(x_k, u_k^*) - \sum_{i=1}^j \bar{\rho}_i^i + \mathbb{E}_k \left[V_j(X_{N-j+1|k}^*) \right]$$

Predictive algorithms



(cont.)

Again, by the standing assumption (item 4), we may write :

$$\begin{aligned} \mathbb{E}_k[V_N(X_{k+1}^*)] - V_N(x_k) &\leq -\rho(x_k, u_k^*) - \sum_{i=1}^j \bar{\rho}_i^i + \\ &\quad \mathbb{E}_k[V_j^*(X_{N-j+1|k}^*)] \\ &\leq -\rho(x_k, u_k^*) + \alpha \mathbb{E}_k[V^*(X_{N-j+1|k}^*)]. \end{aligned}$$

Predictive algorithms



(cont.)

Now, by the same item 4. of the assumption, we have

$$\mathbb{E}_k \left[\sum_{i=2}^N p(X_{i|k}^*, u_{i|k}^*) \right] * \leq V_N(x_k) \leq \alpha V(x_k) + \sum_{i=1}^N \bar{p}_i^N.$$

From this, we may deduce, using $\forall x, u \in \mathcal{U} \quad p(x, u) \geq V(x)$, that there is an $l \in [2:N]$ s.t.

$$\mathbb{E}_k [V(X_{l|k}^*)] \leq \mathbb{E}_k [p(X_{l|k}^*, u_{l|k}^*)] \leq \frac{\alpha}{N-1} V(x_k) + \frac{1}{N-1} \sum_{i=1}^N \bar{p}_i^N$$

(there is a term that is not greater than the average among the $N-1$ ones)

* This is just dropping the first term

Predictive algorithms



(cont.)

Going back to the expected difference in value :

$$\mathbb{E}_k[V_N(X_{k+1}^*)] - V_N(x_k) \leq -\rho(x_k, u_k^*) + \alpha \mathbb{E}_k[V(X_{N-j+2|k})],$$

and using the last derived inequality (from the previous slide), we may write :

$$\mathbb{E}_k[V_N(X_{k+1}^*)] - V_N(x_k) \leq -\rho(x_k, u_k^*) + \frac{\alpha^2}{N-1} V(x_k) + \frac{\alpha}{N-1} \sum_{i=1}^N \bar{\rho}_i^N$$

(assuming we took $j = N-l+1$ from the start).

Setting $c_1 := \frac{\alpha}{N-1} \sum_{i=1}^N \bar{\rho}_i^N$ validates the first part of the original claim

Predictive algorithms



(cont.)

As for the second, observe that if $N \geq \alpha^2 + 1$,

$$\frac{\alpha^2}{N-1} \nu(x_k) \leq \nu(x_k)$$

and so

$$\frac{\alpha^2}{N-1} \nu(x_k) \leq p(x_k, u), \forall u \in U,$$

whence

$$-p(x_k, u_k^*) + \frac{\alpha^2}{N-1} \nu(x_k) + c_1 \leq \left(\frac{\alpha^2}{N-1} - 1 \right) \nu(x_k) + c_1$$

Predictive algorithms



(cont.)

Recall

$$\sqrt{N}^o(x_k) = \sqrt{N}(x_k) - \sum_{i=1}^N \bar{\rho}_i^N \quad \text{and} \quad c_1 = \frac{\alpha}{N-1} \sum_{i=1}^N \bar{\rho}_i^N$$

Also, $\sqrt{N}(x_k) \leq \alpha \mathcal{V}(x_k) + \frac{N-1}{\alpha} c_1$,

whence $-\mathcal{V}(x_k) \leq -\frac{1}{\alpha} \sqrt{N}(x_k) + \frac{N-1}{\alpha^2} c_1 = -\frac{1}{\alpha} \sqrt{N}^o(x_k)$

and so

$$\left(\frac{\alpha^2}{N-1} - 1 \right) \mathcal{V}(x_k) + c_1 \leq \left(\frac{\alpha}{N-1} - \frac{1}{\alpha} \right) \sqrt{N}^o(x_k) + c_1 .$$

Predictive algorithms



(cont.)

$$\text{Now, since } E_k[V_N^*(X_{k+1}^*)] - V_N^*(x_k) = E_k[V_N^o(X_{k+1}^*)] - V_N^o(x_k),$$

we have:

$$E_k[V_N^o(X_{k+1}^*)] - V_N^o(x_k) \leq \left(\frac{\alpha}{N-1} - \frac{1}{\alpha}\right)V_N^o(x_k) + C_1,$$

whence

$$E_k[V_N^o(X_{k+1}^*)] \leq \left(\frac{\alpha}{N-1} - \frac{1}{\alpha} + 1\right)V_N^o(x_k) + C_1$$

which validates the second part of the claim
with

$$\lambda = \frac{\alpha}{N-1} - \frac{1}{\alpha} + 1 \quad \text{and} \quad C_2 = C_L = \frac{\alpha}{N-1} \sum_{i=1}^N \bar{\rho}_i^N$$

Predictive algorithms



(cont.)

If $N > \alpha^2 + 1$, then

$$\lambda = \frac{\alpha}{N-1} - \frac{1}{\alpha} + 1 = \frac{\alpha^2 - N + 1 + (N-1)\alpha}{(N-1)\alpha} = \frac{\alpha^2 + (\alpha-1)(N-1)}{(N-1)\alpha} < 1$$

which shows that $E_k[V_N^*(X_{k+1}^*)]$ stays bounded for any k (in fact, it'd decay asymptotically to zero if c_1 were zero — so is the case if $f(0, 0, w) = 0$ for any $w \in W$).

Predictive algorithms



To actually claim stochastic stability (in mean), we need to relate

$$\mathbb{E}_k[V_N^o(X_{k+1}^*)]$$

to the mean state.

Assume, for simplicity, that indeed $f(o, o, w) = o$ for any $w \in W$. Then, zero control at zero initial state yields zero-mean cost, i.e., $\bar{\rho}_i^N = 0 \forall i$, whence $V_N^o = V_N$.

We have $\mathbb{E}_k[V_N(X_{k+1}^*)] \leq \lambda^k V_N(x_o) \leq \lambda^k \alpha \nu(x_o)$

Predictive algorithms

REINFORCEMENT LEARNING



(cont.)

So,

$$E_k[V_N(X_{k+1}^*)] \leq \lambda^k \alpha V(x_0).$$

Since $\forall x, u \in \mathcal{U} \quad V(x) \leq p(x, u)$, deducting all but the first stage cost from the value function on the left yields :

$$E_k[V(X_{k+1}^*)] \leq \lambda^k \alpha V(x_0).$$

So, if, say $V(x) = b \|x\|^2$ for some $b > 0$, we may claim mean-square stability of the (norm) state and, moreover, its asymptotic convergence to zero



Discussion

- We considered environment stabilization for an MDP, but the results could be extended to an SDE in the sample-and-hold setup (as described at the beginning). To this end, one may employ moment-based accuracy bounds of the Euler-Maruyama or some other numerical integration method
- In the studied result, stabilization was achieved for all sufficiently long horizons. MPC theory actually offers further machinery for stabilization, such as **terminal constraints** and **terminal costs**