

Bootstrapped MPC

CS 294-112

Vladimir Feinberg, Samvit Jain, Michael Whittaker

10 October 2017

1 Introduction

We changed our research direction from the encoded states research to exploring model predictive control (MPC), namely a method we call bootstrapped MPC (BMPC). BMPC is a model-based algorithm that can be viewed as instantiation of the meta-MPC algorithm, which also explains the motivation for bootstrapped MPC (Algorithm 1). Note the (continuous) domains for actions and states are \mathcal{A}, \mathcal{S} and our restriction to deterministic dynamics and rewards for simplicity (though extensions would not require much modification).

Algorithm 1 The METAMPC algorithm is a template for MPC-based algorithms. As template parameters, it accepts a number of simulations to perform K and the simulation horizon H . The critical template parameter of interest is the time-dependent action-sampling distribution A_t . This is a stochastic policy, returning an action a for a provided state s .

```
1: procedure METAMPC(state  $s$ , dynamics  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ , reward  $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ )
2:    $\left\{ s_1^{(i)} \leftarrow s \right\}_{i=1}^K$ .
3:   for  $i \leftarrow 1, \dots, K$  do
4:     for  $t \leftarrow 1 \dots, H$  do
5:       sample  $a_t^{(i)} \sim A_t(s_t)$ 
6:        $s_{t+1}^{(i)} \leftarrow f(s_t^{(i)}, a_t^{(i)})$ 
7:        $R_i \leftarrow \sum_{t=1}^H r(s_t^{(i)}, a_t^{(i)}, s_{t+1}^{(i)})$ 
8:    $i_* = \operatorname{argmax}_i R_i$ 
9:   return  $a_1^{(i_*)}$ 
```

We believe METAMPC works when R_i is a decent estimator for the true H -step reward of the agent under the METAMPC policy, supposing the agent takes $a_1^{(i)}$ now and continues with METAMPC later. Then, choosing the first action $a_1^{(i_*)}$ that results in the best R_i would be the smartest move one could make. Regular MPC specifies $A_t(s) = \text{Uniform}(\mathcal{A})$ for a rectangle \mathcal{A} , independent of t, s . But a uniform distribution is a poor representation of future METAMPC steps: it stands to reason that a more intelligent selection of A_t might improve performance.

In particular, we consider learning a (mostly) stationary action distribution conditional on the state, $\pi_\theta(s)$, which is trained to mimic the result of the MPC algorithm. In particular, every iteration, we roll out METAMPC with $A_t = \pi_\theta^{(t)}$ (with θ held fixed during both the real and simulated rollouts). After collecting some rollouts, resulting in a large dataset of state-action pairs $(s, a, s') \in \mathcal{D}$, we train $\pi_\theta^{(t)}$ to replicate METAMPC's behavior with the expectation that the reward estimates improve (let this be BMPC). We can view this as bootstrapping performance because as π_θ gets more accurate, the BMPC policy is able to take better actions, which are then learned by π , resulting in a virtuous cycle. The dynamics are learned in a similar manner with the dataset of transitions.

2 Related Work

To some degree, this approach is similar to other continuous-action methods that seek to optimize the reward function directly, since in effect we are proposing a method for stochastic optimization of the rollout reward. One difference than, say, CMA-ES, would be that we explicitly model the reinforcement learning setting by modelling the BMPC agent as a stationary policy.

To some degree this approach is similar to Alpha-Go style sampling, where instead of replicating agent behavior the goal is to have intelligent MCTS pruning for a UCT. If UCT pruning is guided by a neural network, then one might view UCT as a tree-version of this BMPC approach (where the space for exploration, being discrete, is much more tractable to explore in the exhaustive tree manner).

3 Easy Examples

We consider learning the dynamics for the **HalfCheetah-v1** environment and using those learned dynamics for BMPC planning. We note two important difference from the regular **gym** environment: the full observation is provided so the dynamics and rewards are true functions of the previous state (the usual observation excerpts the first coordinate of the agent).

Unless otherwise specified, all parameters are their defaults: 3 different seeds for each setup, 60 epochs of dynamics training on the full dataset per on-policy iteration, a dynamics batch size of 512, maximum episode length of 1000, 1000 simulated paths in the MPC controller (each of which has a simulated horizon of H 15), and 10 paths sampled by the initial random agent and then by the MPC controller that are aggregated every policy iteration. The learning rate was the default 10^{-3} . Each such policy iteration contains a round of training the learner policy π_θ in BMPC. All policies π_θ are neural networks with $\dim \mathcal{A}$ outputs and $\dim \mathcal{S}$ inputs, all depth 5, width 32, trained with 100 epochs after each policy iteration with 512-size minibatches and a 10^{-3} learning rate.

We consider a couple variants of BMPC, which are all defined by the conditional action distribution A_t . While our conditional distributions are mostly stationary, we consider exploring at the first step $t = 1$. Equation 1 describes the δ -BMPC, with a point mass describing the deterministic action in the METAMPC simulation taken by A_t when $t > 1$. δ -BMPC is trained on minimizing the loss ℓ , which is the MSE from the expert.

$$A_t(s) = \begin{cases} \text{Uniform}(\mathcal{A}) & t = 1 \\ \delta_{\pi_\theta(s)} & \text{otherwise} \end{cases} \quad \ell(\theta) = \sum_{(s,a) \in \mathcal{D}} \|\pi_\theta(s) - a\|_2^2 \quad (1)$$

Gaussian BMPC relaxes the determinism of δ -BMPC, fitting a conditional diagonal Gaussian to the METAMPC actions, as described in Equation 2.

$$A_t(s) = \begin{cases} \text{Uniform}(\mathcal{A}) & t = 1 \\ N(\pi_\theta(s), \text{diag } \sigma^2) & \text{otherwise} \end{cases} \quad \ell(\theta, \sigma) = - \sum_{(s,a) \in \mathcal{D}} \log p_{\theta, \sigma}(a|s) \quad (2)$$

No-explore BMPC reduces the exploration by keeping the policy purely conditional diagonal Gaussian and stationary, with no initial exploration, as described in Equation 3.

$$A_t(s) = N(\pi_\theta(s), \text{diag } \sigma^2) \quad \ell(\theta) = \sum_{(s,a) \in \mathcal{D}} \|\pi_\theta(s) - a\|_2^2 \quad (3)$$

We evaluate these examples with an easy, supervised reward for **HalfCheetah-v1** which was provided in HW4 (Figure 1). Note that this is different than the classical **gym** reward.

In Figure 1, we see a couple interesting things going on already. Amazingly, the improvement above happens even though the learners themselves are not near the MPC performance at any iteration (Figure 2).

The poor learner performance does not kill the original hypothesis: perhaps even weak learners can improve reward estimates R_{i_*} . But this does suggest a new hypothesis for the improvement: perhaps the

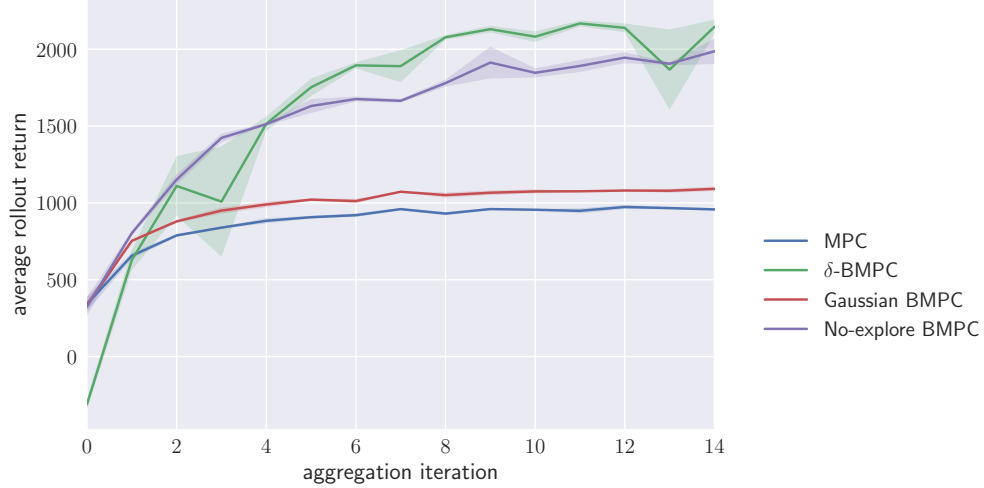


Figure 1: Average return over policy aggregation iterations for the learning-based sampling MPC agent, compared to the standard uniform sampling MPC agent.

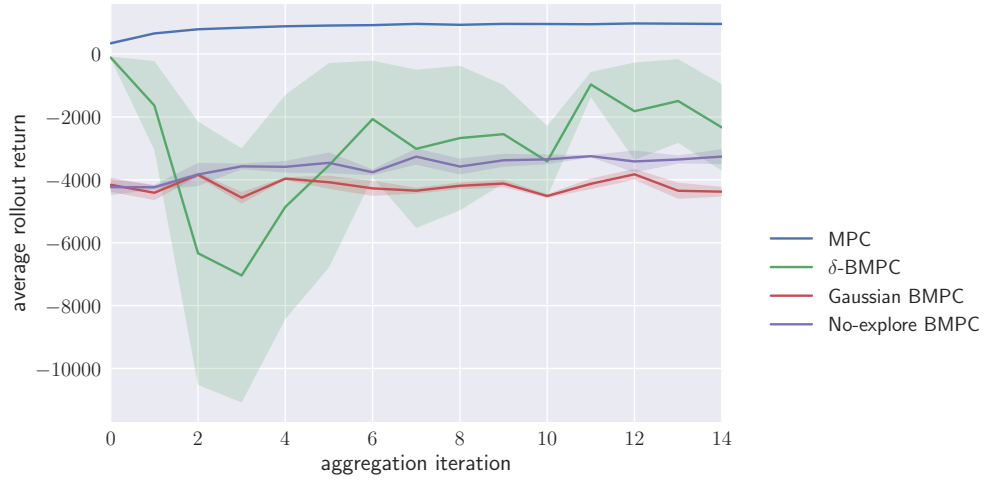


Figure 2: Average return over policy aggregation iterations for an agent sampling from A_2 for each model, along with the original MPC as a baseline.

learners prevent catastrophic changes in the policy. Maybe simply reducing the variance of estimates R_{i_*} yields the improvement. We note that, at least for this simple case, learners don't improve performance by making the dynamics more learnable: for both the original MPC and BMPC, the learned dynamics are about equally performant (Figure 3).

4 Hard Cost

If we force ourselves to only rely on the usual `HalfCheetah-v1` reward, which only rewards forward movement and discourages large-magnitude actions, we are faced with a more difficult setting.

In this case, it is unclear if bootstrapping still helps (Figure 4).



Figure 3: We evaluate dynamics in the setting of Figure 1. This evaluation is measures how predictive the dynamics were of the actual transitions in validation rollouts (it is not the training MSE). We omit the first couple of iterations because their scales are vastly different than the rest.

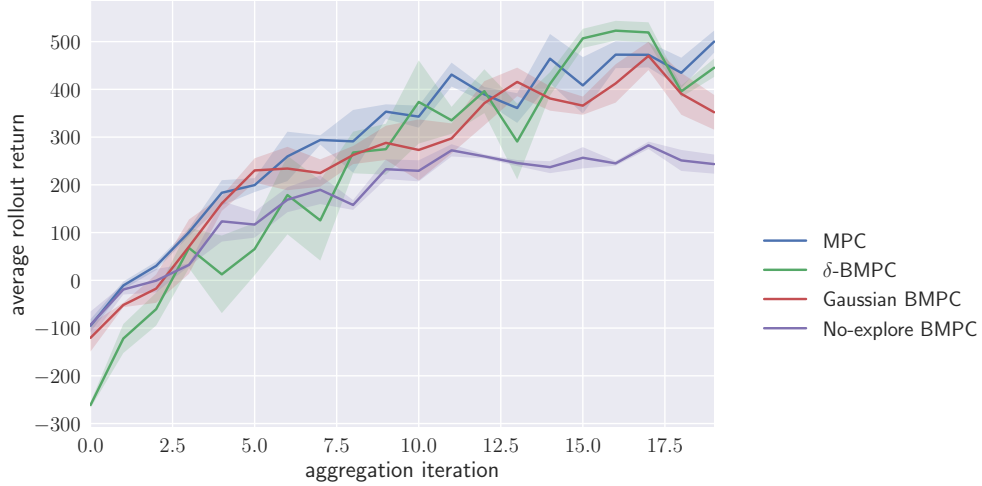


Figure 4: Average return over policy aggregation iterations for the learning-based sampling MPC agent, compared to the standard uniform sampling MPC agent, this time with a hard cost function.

The dynamics and learning returns for the runs in Figure 4 paint a similar picture as for easy cost. In case hard cost required additional simulation, we changed H, K from 15, 1000 to 50, 1000. Performance was generally better but more chaotic, and still with no clear advantage to BMPC (Figure 5).

Nonetheless, the trend in Figure 4 deserves elaboration—perhaps at later iterations BMPC overtakes MPC.

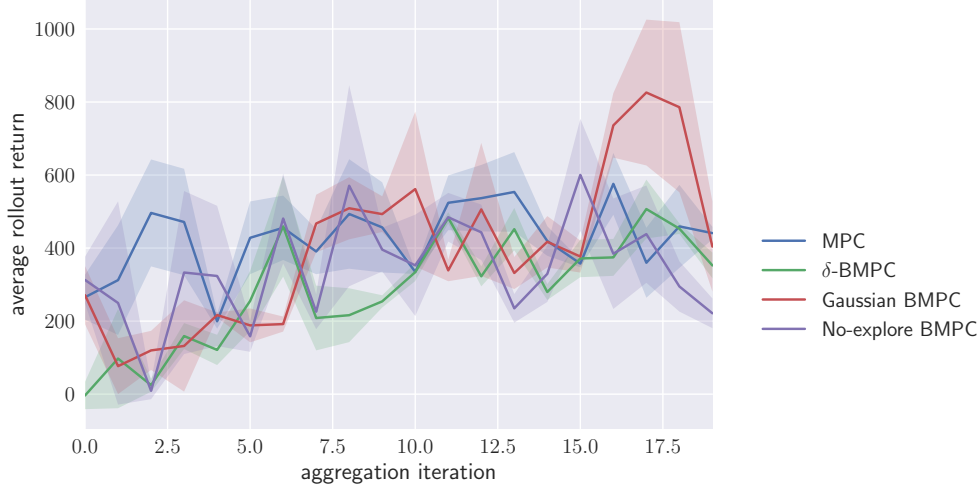


Figure 5: The same setting as Figure 5 but with longer simulated horizons.

5 Planned Work

The proposition being investigated (more accurate or lower-variance estimates of reward with R_{i_*}) can be directly sampled and tested to be the case! After simulating several rollouts for any METAMPC algorithm, we can sample initial state s from the empirical distribution of visited states. For each such state, we *re-simulate* H real horizon steps with independent runs (the samples A_t will differ, so the results will not be the same as the first run, even for a deterministic environment). These give us several observations of the true H -step reward \tilde{R}_H . For each such observation we compute the observed bias $\tilde{R}_H - R_{i_*}$ and squared error $(\tilde{R}_H - R_{i_*})^2$. Over several samples we can estimate variance $\text{var } R_{i_*}$ as well. Then across the samples s we can estimate average values across the entire distribution for these reward statistics.

This sampling of s is not iid but the resulting estimator is still unbiased, so we will have to use the bootstrap (haha) to calculate confidence intervals.

Measuring the explicit action distribution shift between iterations boils down to measuring distance between conditional action distributions resulting from the METAMPC algorithm, which is only a black-box generative distribution. We have a plan for this but excerpt it to keep the report shorter. Graders can contact the authors for more details.

In addition, we hope to expand to the harder **Ant** problem to prevent overfitting to **HalfCheetah**.

Finally, the chaotic performance but improving learner in the second hard setting suggest that we should explore various hyperparameter and template settings; but we should save this for the end.

6 Further Exploration

We also explored changing the initial exploration policy A_1 from $\text{Uniform}(\mathcal{A})$ to $N(\pi_\theta, \lambda I)$. This resulted in some improvements, but we leave these out of the report for simplicity (graders may contact the authors directly for these results). Further, to improve how quickly π_θ learns expert actions, we consider DAgger iteration with π_θ as the learner and METAMPC as the expert, but managing the learned dynamics for both the learner and expert will require additional work.