

# Meta Reinforcement Learning with Latent Variable Gaussian Processes: Supplement

## 1 Model-Based Meta Reinforcement Learning

We consider stochastic systems of the form

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{c}_t) + \boldsymbol{\epsilon} \quad (1)$$

with state variables  $\mathbf{x} \in \mathbb{R}^D$ , control signals  $\mathbf{c} \in \mathbb{R}^K$  and i.i.d. system noise  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{E})$ , where  $\mathbf{E} = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$ .

We model the distribution over systems using a latent embedding  $\mathbf{h}$  and model the dynamics using a global function conditioned on the latent embedding. For a sample  $f_p$  from the distribution we modify the dynamics from (1) to

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p) + \boldsymbol{\epsilon}, \quad (2)$$

such that the successor state depends on the latent system specification  $\mathbf{h}_p$ . This means, we explicitly model the global properties through a shared function  $f$  and the task-specific variation using a distribution over the latent variables  $p(\mathbf{h}_p)$ . Framing the meta learning problem as a hierarchical Bayesian model means that meta-training becomes inference in a meta-learning model. The meta RL procedures for training and testing are detailed in algorithms 1 and 2, respectively.

### 1.1 Meta-Learning Gaussian Process

Our meta-learning model is a GP prior on the unknown transition function in (2) with a concatenated state  $\tilde{\mathbf{x}}_t = (\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p) \in \mathbb{R}^{D+K+Q}$  as the input to the model. We define  $\mathbf{y}_t = \mathbf{x}_{t+1} - \mathbf{x}_t$  as the targets of the GP and take the mean function to be  $m(\tilde{\mathbf{x}}_t) = \mathbf{0}$ , which encodes that a priori the state does not change [2]. Each dimension of the targets  $\mathbf{y}$  is modeled by an independent GP. We use a Gaussian likelihood

$$p(\mathbf{y}_t | \tilde{\mathbf{x}}_t, \mathbf{f}(\cdot), \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}_t | \mathbf{f}(\tilde{\mathbf{x}}_t), \mathbf{E}), \quad (3)$$

where  $\boldsymbol{\theta} = \{\mathbf{E}, \mathbf{L}, \sigma_f^2, Q\}$  are the hyperparameters and  $\mathbf{f}(\cdot) = (f^1(\cdot), \dots, f^D(\cdot))$  denotes a multi-dimensional function. We place a standard-normal prior  $\mathbf{h}_p \sim$

```

Initialize dataset  $D$  and model  $M$ 
▷ Initial random rollouts

forall training tasks do
  | execute random policy
  | add observations to  $D$ 
end
▷ Meta training

while training tasks not solved do
  | —update—: train  $M$  and infer  $\mathbf{h}$  given  $D$ 
  | forall unsolved training tasks do
  |   | for each step in horizon do
  |     | | —plan—: get control sequence using (13)
  |     | | —execute—: execute first control in sequence
  |     | end
  |     | add observations to  $D$ 
  |     | check if task solved
  |   | end
  | end
end

```

**Algorithm 1:** Model-based Meta RL with MPC (Train)

$\mathcal{N}(\mathbf{0}, \mathbf{I})$  on the latent variables  $\mathbf{h}_p$ . The full specification of the model is

$$\begin{aligned}
& p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{C}) \\
&= \prod_{p=1}^P p(\mathbf{h}_p) \prod_{t=1}^{T_p} p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p, \mathbf{f}(\cdot)) p(\mathbf{f}(\cdot))
\end{aligned} \tag{4}$$

where we denote a collection of vectors in bold upper-case and we have dropped dependence on the hyperparameters for notation purposes. The corresponding graphical model is given in Fig. 1. The figure shows the dependence of individual system observations on the global GPs  $\mathbf{f}(\cdot)$  modeling each dimension of the outputs, the system-specific latent embeddings  $\mathbf{h}_p$  and the observed states and controls.

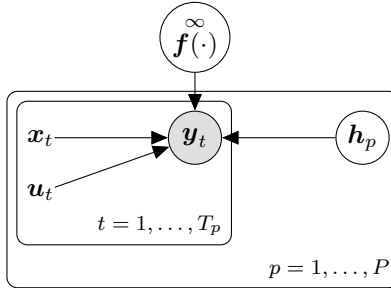


Figure 1: Graphical model for our ML-GP model.

Given dataset  $D$  and model  $M$  from training

▷ Single shot performance

```

forall test tasks do
  for each step in horizon do
    —plan—: get control sequence using (13)
    —execute—: execute first control in sequence
    —inference—: infer the value of  $\mathbf{h}_*$  given observations so far
  end
  add observations to  $D$ 
  check if task solved
end

```

▷ Meta test

```

while test tasks not solved do
  —update—: train  $M$  and infer  $\mathbf{h}$  given  $D$ 
  forall unsolved test tasks do
    for each step in horizon do
      —plan—: get control sequence using (13)
      —execute—: execute first control in sequence
    end
    add observations to  $D$ 
    check if task solved
  end
end

```

**Algorithm 2:** Model-based Meta RL with MPC (Test)

## 1.2 Inference

To learn the dynamics model we seek to optimize the hyperparameters  $\theta$  w.r.t. the log-marginal likelihood, which involves marginalization of the latent variables in (4). For predictions of the evolution of a system we also need to infer the posterior GP and the posterior distribution of the latent variables  $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_P)$ . We approach this problem with approximate variational inference. We posit a variational distribution that assumes independence between the latent functions of the GP and the latent task variables

$$Q(\mathbf{f}(\cdot), \mathbf{H}) = q(\mathbf{f}(\cdot))q(\mathbf{H}) \quad (5)$$

and minimize the Kullback-Leibler divergence between the approximate and true posterior distributions. Equivalently we can maximize the evidence lower bound

$$\mathcal{L} = \mathbb{E}_{Q(\mathbf{f}(\cdot), \mathbf{H})} \left[ \log \frac{p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{C})}{Q(\mathbf{f}(\cdot), \mathbf{H})} \right], \quad (6)$$

which lower-bounds the log-marginal likelihood [5]. We parameterize our variational distribution such that we can compute the lower bound in (6). We then jointly optimize  $\mathcal{L}$  with respect to the model hyperparameters and the variational parameters.

**Sparse Gaussian Processes** It is important to account for the fact that training a GP on a joint data set of  $P$  different systems quickly becomes infeasible due to the  $\mathcal{O}(T^3)$  computational complexity for training and  $\mathcal{O}(T^2)$  for predictions where  $T$  is the total number of observations. To address this we turn to the variational sparse GP approximation [9] and approximate the posterior GP with a variational distribution  $q(\mathbf{f}(\cdot))$  that depends on a small set of  $M \ll T$  *inducing points*. We introduce a set of  $M$  inducing inputs  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_M) \in \mathbb{R}^{M \times (D+K+Q)}$ , which live in the same space as  $\tilde{\mathbf{x}}$ , with corresponding GP function values  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_M) \in \mathbb{R}^{M \times D}$ . We follow [4] and specify the variational approximation as a combination of the conditional GP prior and a variational distribution over the inducing function values, independent across output dimensions

$$q(f^d(\cdot)) = \int p(f^d(\cdot) | \mathbf{u}^d) q(\mathbf{u}^d) d\mathbf{u}^d. \quad (7)$$

where  $q(\mathbf{u}^d) = \mathcal{N}(\mathbf{u}^d | \mathbf{m}^d, \mathbf{S}^d)$  is a full rank Gaussian. The integral in (7) can be computed in closed form since both terms are Gaussian, resulting in a GP with mean and covariance functions given by

$$m_q(\cdot) = \mathbf{k}_Z^T(\cdot) \mathbf{K}_{ZZ}^{-1} \mathbf{m}^d \quad (8)$$

$$k_q(\cdot, \cdot) = k(\cdot, \cdot) - \mathbf{k}_Z^T(\cdot) \mathbf{K}_{ZZ}^{-1} (\mathbf{K}_{ZZ} - \mathbf{S}^d) \mathbf{K}_{ZZ}^{-1} \mathbf{k}_Z(\cdot) \quad (9)$$

where  $[\mathbf{k}_Z(\cdot)]_i = k(\cdot, \mathbf{z}_i)$  and  $[\mathbf{K}_{ZZ}]_{ij} = k(\mathbf{z}_i, \mathbf{z}_j)$ . Here, the variational approach has two main benefits: a) it reduces the complexity of training to  $\mathcal{O}(TM^2)$  and predictions to  $\mathcal{O}(TM)$ , b) it enables mini-batch training for further improvement in computational efficiency.

**Latent Variables** For the latent variables  $\mathbf{H}$  we assume a Gaussian variational posterior

$$q(\mathbf{H}) = \prod_{p=1}^P \mathcal{N}(\mathbf{h}_p | \mathbf{n}_p, \mathbf{T}_p) \quad (10)$$

where  $\mathbf{T}_p$  is in general a full rank covariance matrix. We use a diagonal covariance in practice for more efficient computation of the ELBO (6).

**Evidence Lower Bound (ELBO)** Plugging in the definition of the model in eq. (3) and the specified variational distributions in (7), (10) into the definition of the ELBO in (6) we find

$$\begin{aligned} \mathcal{L} &= \mathbb{E}_{Q(\mathbf{f}(\cdot), \mathbf{H})} \left[ \log \frac{p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{C})}{Q(\mathbf{f}(\cdot), \mathbf{H})} \right] \\ &= \mathbb{E}_{q(\mathbf{f}(\cdot))q(\mathbf{H})} \left[ \log \frac{\prod_{p=1}^P p(\mathbf{h}_p) \prod_{t=1}^{T_p} p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p, \mathbf{f}(\cdot)) p(\mathbf{f}(\cdot))}{q(\mathbf{f}(\cdot))q(\mathbf{H})} \right] \\ &= \sum_{p=1}^P \sum_{t=1}^{T_p} \mathbb{E}_{q(\mathbf{f}_t | \mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p)q(\mathbf{h}_p)} [\log p(\mathbf{y}_t | \mathbf{f}_t)] \\ &\quad - \text{KL}[q(\mathbf{H}) || p(\mathbf{H})] - \text{KL}[q(\mathbf{f}(\cdot)) || p(\mathbf{f}(\cdot))] \end{aligned}$$

where we emphasize that  $q(\mathbf{f}_t|\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p) = q(\mathbf{f}(\tilde{\mathbf{x}}_t)|\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p)$  is the marginal of the GP evaluated at the inputs  $\tilde{\mathbf{x}}_t$ . The KL term for the latent variables  $\mathbf{H}$  is analytically tractable since both terms are Gaussian. The KL term between the GPs has been shown to simplify to  $\text{KL}[q(\mathbf{U})||p(\mathbf{U})]$  [7], which again is analytically tractable since both terms are Gaussian. Thus the ELBO can be written

$$\begin{aligned} \mathcal{L} = & \sum_{p=1}^P \sum_{t=1}^T \mathbb{E}_{q(\mathbf{f}_t|\mathbf{x}_t, \mathbf{c}_t)} [\log p(\mathbf{y}_t|\mathbf{f}_t)] \\ & - \text{KL}[q(\mathbf{H})||p(\mathbf{H})] - \text{KL}[q(\mathbf{U})||p(\mathbf{U})] \end{aligned} \quad (11)$$

The expected log-likelihood term in (11) needs further consideration: we would like to integrate out the latent variable  $\mathbf{h}_p$  to obtain

$$q(\mathbf{f}_t|\mathbf{x}_t, \mathbf{c}_t) = \int q(\mathbf{f}_t|\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p)q(\mathbf{h}_p)d\mathbf{h}_p \quad (12)$$

The integral in (12) is intractable due to the non-linear dependence on  $\mathbf{h}_p$  in (8) and (9). Given our choice of kernel (RBF) and Gaussian variational distribution  $q(\mathbf{h}_p)$  the first and second moments can be computed in closed form. We could use these terms to compute the log-likelihood term in closed form since the likelihood is Gaussian but in practice this can be prohibitively expensive since it requires the evaluation of a  $TM^2D$  tensor. Instead we avoid computing the moments by approximately integrating out the latent variable using Monte Carlo sampling.

$$\begin{aligned} \text{Hyperparameters: } & \boldsymbol{\theta} = \{\mathbf{E}, \mathbf{L}, \sigma_f^2, Q\} \\ \text{Variational parameters: } & \boldsymbol{\phi} = \{\mathbf{Z}, M \{\mathbf{m}^d, \mathbf{S}^d\}_{d=1}^D, \{\mathbf{n}_p, \mathbf{T}_p\}_{p=1}^P\} \end{aligned}$$

**Training** For the update steps in algorithms 1 and 2 we jointly optimize the GP hyperparameters and the variational parameters w.r.t. the ELBO. For the inference step in algorithm 2, we optimize only the variational parameters for the latent variables  $\mathbf{h}$ , i.e.  $\boldsymbol{\phi}_h = \{\mathbf{n}_p, \mathbf{T}_p\}_{p=1}^P$ .

## 2 Expected Long Term Cost

**RL with MPC** Our objective is to find a sequence of optimal controls  $\mathbf{c}_0^*, \dots, \mathbf{c}_{H-1}^*$  that minimizes the expected finite-horizon cost

$$J = \mathbb{E} \left[ \sum_{t=1}^H \ell(\mathbf{x}_t) \right], \quad (13)$$

where  $\mathbf{x}_t$  is the state of the system at time  $t$  and  $\ell$  is a known immediate/instantaneous cost function that encodes the task objective. We consider an episodic setting. Initial states  $\mathbf{x}_0$  are sampled from  $p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ .

To find the optimal open-loop sequence  $\mathbf{c}_0^*, \dots, \mathbf{c}_{H-1}^*$ , we compute the expected long-term cost  $J$  in (13) using Gaussian approximations  $p(\mathbf{x}_1), \dots, p(\mathbf{x}_H)$  for a given control sequence  $\mathbf{c}_0, \dots, \mathbf{c}_{H-1}$ . Then, we find an open-loop control sequence that minimizes the expected long-term cost and apply the first control signal  $\mathbf{c}_0^*$  to the system, which transitions into the next state. Next we re-plan, i.e., we determine the next open-loop control sequence  $\mathbf{c}_0^*, \dots, \mathbf{c}_{H-1}^*$  from the new state. This iterative MPC approach turns an open-loop controller into a closed-loop controller. Combining MPC with learned GP models for the underlying dynamics increases the robustness to model errors and has shown improved data efficiency in RL [6].

**Computing the Expected Long-Term Cost** To compute the expected long-term cost in (13), we sum up the expected immediate costs  $\mathbb{E}[\ell(\mathbf{x}_t)] = \int \ell(\mathbf{x}_t)p(\mathbf{x}_t)d\mathbf{x}_t$  for  $t = 0, \dots, H - 1$ . We assume Gaussian approximations  $p(\mathbf{x}_t)$  of the state distribution and choose  $\ell$  so that this expectation can be computed analytically. To obtain the marginals  $p(\mathbf{x}_t)$  we iteratively predict the state evolution using the GP dynamics model. Section 1.2 details the long-term predictions the Gaussian approximations.

**Long-term Predictions** Key to computing the expected long-term cost in (13) is to compute long-term predictions  $p(\mathbf{x}_1), \dots, p(\mathbf{x}_{H-1})$  of the state evolution. To accomplish this we iteratively predict

$$p(\mathbf{x}_{t+1}|\mathbf{c}_t) = \iiint p(\mathbf{x}_{t+1}|\mathbf{f}_t, \mathbf{x}_t)q(\mathbf{f}_t|\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p) \times p(\mathbf{x}_t)q(\mathbf{h}_p)d\mathbf{f}_td\mathbf{x}_td\mathbf{h}_p, \quad (14)$$

which we approximate as a Gaussian using GP moment matching [8, 3].<sup>1</sup> We first define the distribution

$$q(\tilde{\mathbf{x}}_t|\mathbf{c}_t) = \mathcal{N}(\tilde{\mathbf{x}}_t|\boldsymbol{\mu}_{\tilde{\mathbf{x}}_t}, \boldsymbol{\Sigma}_{\tilde{\mathbf{x}}_t}) \quad (15)$$

over the concatenated state  $\tilde{\mathbf{x}}_t = (\mathbf{x}_t, \mathbf{c}_t, \mathbf{h}_p)$ , where the Gaussianity of  $q(\tilde{\mathbf{x}}_t|\mathbf{c}_t)$  follows from the fact that  $p(\mathbf{x}_t)$  and  $q(\mathbf{h}_p)$  are both multivariate Gaussians and independent of each other [1]. The mean and covariance of  $q(\tilde{\mathbf{x}}_t|\mathbf{c}_t)$  are

$$\boldsymbol{\mu}_{\tilde{\mathbf{x}}_t} = [\boldsymbol{\mu}_{\mathbf{x}_t}, \mathbf{c}_t, \mathbf{n}_p] \quad (16)$$

$$\boldsymbol{\Sigma}_{\tilde{\mathbf{x}}_t} = \text{blockdiagonal}([\boldsymbol{\Sigma}_{\mathbf{x}_t}, \mathbf{0}, \mathbf{T}_p]). \quad (17)$$

Next we predict the change in state given the control  $\mathbf{c}_t$

$$p(\mathbf{y}_t|\mathbf{c}_t) = \iint p(\mathbf{y}_t|\mathbf{f}_t)q(\mathbf{f}_t|\tilde{\mathbf{x}}_t)q(\tilde{\mathbf{x}}_t|\mathbf{c}_t)d\mathbf{f}_td\tilde{\mathbf{x}}_t \quad (18)$$

---

<sup>1</sup>Note that we integrate out the distribution of the state, the GP dynamics model and the latent variable encoding the task. The control has no distribution as it is treated as a free parameter within an MPC context.

which is intractable due to the integral over the state  $q(\tilde{\mathbf{x}}_t|\mathbf{c}_t)$  through the GP. However, given that the state is a Gaussian and the kernel is the RBF kernel, we can analytically compute the first and second moments of the distribution  $p(\mathbf{y}_t|\mathbf{c}_t)$  [8, 3, 2]. We then make a Gaussian approximation using those moments

$$q(\mathbf{y}_t|\mathbf{c}_t) = \mathcal{N}(\mathbf{y}_t|\boldsymbol{\mu}_{\mathbf{y}_t}, \boldsymbol{\Sigma}_{\mathbf{y}_t}) \quad (19)$$

We are really interested in

$$p(\mathbf{x}_{t+1}|\mathbf{c}_t) = \iint p(\mathbf{x}_{t+1}|\mathbf{y}_t, \mathbf{x}_t)q(\mathbf{y}_t|\mathbf{c}_t)p(\mathbf{x}_t|\mathbf{c}_t)d\mathbf{y}_td\mathbf{x}_t \quad (20)$$

which gives

$$\mathbb{E}[\mathbf{x}_{t+1}] = \mathbb{E}[\mathbf{y}_t] + \boldsymbol{\mu}_{\mathbf{x}_t} = \boldsymbol{\mu}_{\mathbf{y}_t} + \boldsymbol{\mu}_{\mathbf{x}_t} \quad (21)$$

$$\text{var}[\mathbf{x}_{t+1}] = \text{var}[\mathbf{x}_t] + \text{var}[\mathbf{y}_t] + \text{cov}[\mathbf{x}_t, \mathbf{y}_t] + \text{cov}[\mathbf{y}_t, \mathbf{x}_t] \quad (22)$$

where the computation of the covariance terms in (22) is detailed in [2]. By iterating this procedure for  $t = 1, \dots, H - 1$ , we obtain the desired sequence of approximate Gaussian state distributions  $p(\mathbf{x}_1), \dots, p(\mathbf{x}_{H-1})$ .

## References

- [1] J. Blitzstein and J. Hwang. *Introduction to Probability*. CRC Press, 2015.
- [2] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian processes for data-efficient learning in robotics and control. *PAMI*, 37(2):408–423, 2015.
- [3] M. P. Deisenroth and C. E. Rasmussen. PILCO: A model-based and data-efficient approach to policy search. In *ICML*, 2011.
- [4] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. In *UAI*, 2013.
- [5] M. Hoffman, D. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *JMLR*, 2013.
- [6] S. Kamthe and M. P. Deisenroth. Data-efficient reinforcement learning with probabilistic model predictive control. In *AISTATS*, 2018.
- [7] A. Matthews, J. Hensman, R. Turner, and G. Zoubin. On sparse variational methods and the Kullback-Leibler divergence between stochastic processes. In *AISTATS*, 2016.
- [8] J. Quiñonero-Candela, A. Girard, J. Larsen, and C. E. Rasmussen. Propagation of uncertainty in Bayesian kernel models—application to multiple-step ahead forecasting. In *ICASSP*, 2003.
- [9] M. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *AISTATS*, 2009.