
Supplementary Material: Probabilistic Differential Dynamic Programming

Yunpeng Pan and **Evangelos A. Theodorou**
Daniel Guggenheim School of Aerospace Engineering
Institute for Robotics and Intelligent Machines
Georgia Institute of Technology
Atlanta, GA 30332

ypan37@gatech.edu, evangelos.theodorou@ae.gatech.edu

Abstract

This is the supplementary document for the paper on Probabilistic Differential Dynamic Programming (PDDP). It includes derivations for the probabilistic representation of the stochastic dynamics, the linearization of the dynamics model and the cost function formulation.

1 Problem formulation

We consider a general unknown stochastic system described by the following differential equation

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, \mathbf{u})dt + \mathbf{C}(\mathbf{x}, \mathbf{u})d\omega, \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad d\omega \sim \mathcal{N}(0, \Sigma_\omega), \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the state, $\mathbf{u} \in \mathbb{R}^m$ is the control, t is time and $\omega \in \mathbb{R}^p$ is standard Brownian motion noise. The trajectory optimization problem is defined as finding a sequence of state and controls that minimize the expected cost

$$J^\pi(\mathbf{x}(t_0)) = \mathbb{E} \left[h(\mathbf{x}(T)) + \int_{t_0}^T \mathcal{L}(\mathbf{x}(t), \pi(\mathbf{x}(t)), t) dt \right], \quad (2)$$

where $h(\mathbf{x}(T))$ is the terminal cost, $\mathcal{L}(\mathbf{x}(t), \pi(\mathbf{x}(t)), t)$ is the instantaneous cost rate, $\mathbf{u}(t) = \pi(\mathbf{x}(t))$ is the control policy. The cost $J^\pi(\mathbf{x}(t_0))$ is defined as the expectation of the total cost accumulated from t_0 to T . For the rest of our analysis, we denote $\mathbf{x}_k = \mathbf{x}(k)$ in discrete-time where $k = 0, 1, \dots, H$ is the time step, we use this subscript rule for other variables as well.

2 Probabilistic model learning

The continuous functional mapping from state-control pair $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{u}) \in \mathbb{R}^{n+m}$ to state transition $d\mathbf{x}$ can be viewed as an inference with the goal of inferring $d\mathbf{x}$ given $\tilde{\mathbf{x}}$. We view this inference as a nonlinear regression problem. In this subsection, we introduce the Gaussian processes (GP) approach to learning the dynamics model in (1). A GP is defined as a collection of random variables, any finite number subset of which have a joint Gaussian distribution. Given a sequence of state-control pair $\tilde{\mathbf{X}} = \{(\mathbf{x}_0, \mathbf{u}_0), \dots, (\mathbf{x}_H, \mathbf{u}_H)\}$, and the corresponding state transition $d\tilde{\mathbf{X}} = \{d\mathbf{x}_0, \dots, d\mathbf{x}_H\}$, a GP is completely defined by a mean function and a covariance function. The joint distribution of the observed output and the output corresponding to a given test state-control pair $\tilde{\mathbf{X}}^* = (\mathbf{x}^*, \mathbf{u}^*)$ can be written as

$$\mathbf{p} \left(\begin{array}{c} d\mathbf{X} \\ d\mathbf{x}^* \end{array} \right) \sim \mathcal{N} \left(0, \begin{bmatrix} \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n \mathbf{I} & \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}^*) \\ \mathbf{K}(\tilde{\mathbf{x}}^*, \tilde{\mathbf{X}}) & \mathbf{K}(\tilde{\mathbf{x}}^*, \tilde{\mathbf{x}}^*) \end{bmatrix} \right). \quad (3)$$

The covariance of this multivariate Gaussian distribution is defined via a kernel matrix $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$. In particular, in this paper we consider the Gaussian kernel

$$\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) = \sigma_s^2 \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{W}(\mathbf{x}_i - \mathbf{x}_j)\right) + \sigma_n^2, \quad (4)$$

with $\sigma_s, \sigma_n, \mathbf{W}$ the hyper-parameters of the GP. The kernel function can be interpreted as a similarity measure of random variables. More specifically, if the training pairs $\tilde{\mathbf{X}}_i$ and $\tilde{\mathbf{X}}_j$ are close to each other in the kernel space, their output $d\mathbf{x}_i$ and $d\mathbf{x}_j$ are highly correlated.

The posterior distribution, which is also a Gaussian distribution, can be obtained by constraining the joint distribution to contain the output $d\mathbf{x}^*$ that are consistent with the observations. Assuming independent outputs (no correlation between each output dimension) and given test input $\tilde{\mathbf{x}}_k = [\mathbf{x}_k, \mathbf{u}_k]$ at time k . The one-step prediction of dynamics based on GP can be evaluated as

$$p(d\mathbf{x}_k | \tilde{\mathbf{x}}_k) \sim \mathcal{N}(d\boldsymbol{\mu}_k, d\boldsymbol{\Sigma}_k), \quad (5)$$

where the mean and variance are given by

$$\begin{aligned} d\boldsymbol{\mu}_k &= \mathbb{E}[d\mathbf{x}_k] = \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}})(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n \mathbf{I})^{-1} d\mathbf{X}, \\ d\boldsymbol{\Sigma}_k &= \text{Var}[d\mathbf{x}_k] = \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_k) - \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}})(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n \mathbf{I})^{-1} \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \end{aligned} \quad (6)$$

where $d\boldsymbol{\mu}_k$ and $d\boldsymbol{\Sigma}_k$ are predictive mean and variance of the state transition, respectively. Therefore, the state distribution at $k + 1$ would be:

$$p(\mathbf{x}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (7)$$

where the state mean and variance are

$$\boldsymbol{\mu}_{k+1} = \mathbf{x}_k + d\boldsymbol{\mu}_k, \quad \boldsymbol{\Sigma}_{k+1} = d\boldsymbol{\Sigma}_k. \quad (8)$$

When propagating the GP-based dynamics over a trajectory of time horizon H , the input state \mathbf{x}_k becomes uncertain with Gaussian distribution, where $k = 1, \dots, H$ (the initial state \mathbf{x}_0 is deterministic). Thus the distribution over state transition can be computed as:

$$p(d\mathbf{x}_k) = \int \int p(\mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) | \mathbf{x}_k, \mathbf{u}_k) p(\mathbf{x}_k, \mathbf{u}_k) d\mathbf{x}_k d\mathbf{u}_k. \quad (9)$$

Generally, the above distribution cannot be computed analytically because the nonlinear mapping of input Gaussian distributions lead to non-Gaussian predictive distributions. However, the predictive distribution can be approximated by a Gaussian. Thus the state distribution at $t + 1$ is also a Gaussian

$$\boldsymbol{\mu}_{k+1} = \boldsymbol{\mu}_k + d\boldsymbol{\mu}_k, \quad \boldsymbol{\Sigma}_{k+1} = \boldsymbol{\Sigma}_k + d\boldsymbol{\Sigma}_k + \text{Cov}[\mathbf{x}_k, d\mathbf{x}_k] + \text{Cov}[d\mathbf{x}_k, \mathbf{x}_k]. \quad (10)$$

In order to obtain the distribution over state $\mathcal{N}(\boldsymbol{\mu}_{k+1}, \boldsymbol{\Sigma}_{k+1})$. Firstly, we compute the joint distribution over state-control pair $p(\tilde{\mathbf{x}}_k) = p(\mathbf{x}_k, \mathbf{u}_k)$ as follow

$$p\left(\begin{array}{c} \mathbf{x}_k \\ \mathbf{u}_k \end{array}\right) \sim \mathcal{N}\left(\begin{array}{c} \boldsymbol{\mu}_k \\ \mathbb{E}[\mathbf{u}_k] \end{array}, \left[\begin{array}{cc} \boldsymbol{\Sigma}_k & \text{Cov}[\mathbf{x}_k, \mathbf{u}_k] \\ \text{Cov}[\mathbf{u}_k, \mathbf{x}_k] & \text{Cov}[\mathbf{u}_k] \end{array} \right]\right) \quad (11)$$

where $\mathbb{E}[\mathbf{u}_k]$ and $\text{Cov}[\mathbf{u}_k]$ are mean and covariance of the distribution over control policy $p(\mathbf{u}_k)$. To simplify notation, we denote the mean and covariance of the above distribution as $p(\tilde{\mathbf{x}}_k) \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k)$. Since the control policy is a linear function of the Gaussian belief augmented state in this paper, the control is actually deterministic.

Given the input joint distribution $p(\tilde{\mathbf{x}}_k)$, we will compute the predictive distribution of state transition $p(d\mathbf{x}_k)$. The predictive mean can be computed using the law of iterated expectations (Fubini's

theorem)

$$\begin{aligned}
d\boldsymbol{\mu}_k &= \int \int \int \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) p(\mathbf{x}_k, \mathbf{u}_k) d\mathbf{f} d\mathbf{x}_k d\mathbf{u}_k \\
&= \int \int \mathbf{f}(\tilde{\mathbf{x}}_k) p(\tilde{\mathbf{x}}_k) d\mathbf{f} d\tilde{\mathbf{x}}_k dt \\
&= \mathbb{E}_{\mathbf{f}, \tilde{\mathbf{x}}_k} [\mathbf{f} | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k] \\
&= \mathbb{E}_{\tilde{\mathbf{x}}_k} \left[\mathbb{E}_{\mathbf{f}} [\mathbf{f}(\tilde{\mathbf{x}}_k) | \tilde{\mathbf{x}}_k] | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k \right] dt \\
&= \int \left(\mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) (\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} d\mathbf{X} \right) \mathcal{N}(\tilde{\mathbf{x}}_k | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) d\tilde{\mathbf{x}}_k \\
&= \underbrace{\left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right)^{-1} d\mathbf{X}}_{\boldsymbol{\Psi}} \underbrace{\int \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \mathcal{N}(\tilde{\mathbf{x}}_k | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) d\tilde{\mathbf{x}}_k}_{\mathbf{q}_k} \\
&= \boldsymbol{\Psi}^T \mathbf{q}_k
\end{aligned} \tag{12}$$

where $\boldsymbol{\Psi} \in \mathbb{R}^{N \times n}$ and $\mathbf{q}_k = [q_{k1}, \dots, q_{kn}]^T \in \mathbb{R}^N$ with each element

$$\begin{aligned}
q_{ki} &= \int \mathbf{K}(\tilde{\mathbf{X}}_i, \tilde{\mathbf{x}}_k) \mathcal{N}(\tilde{\mathbf{x}}_k | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) d\tilde{\mathbf{x}}_k \\
&= \alpha^2 |\tilde{\boldsymbol{\Sigma}}_k + \mathbf{W}|^{\frac{1}{2}} \exp \left(-\frac{1}{2} (\tilde{\mathbf{X}}_i - \tilde{\boldsymbol{\mu}}_k)^T (\tilde{\boldsymbol{\Sigma}}_k + \Lambda)^{-1} (\tilde{\mathbf{X}}_i - \tilde{\boldsymbol{\mu}}_k) \right).
\end{aligned} \tag{13}$$

Next, we compute the predictive covariance matrix

$$\text{Cov}(d\mathbf{x}_k | \tilde{\mathbf{x}}_k) = \begin{bmatrix} \text{Var}(d\mathbf{x}_{k1}) & \dots & \text{Cov}(d\mathbf{x}_{kn}, d\mathbf{x}_{k1}) \\ \vdots & \ddots & \vdots \\ \text{Cov}(d\mathbf{x}_{k1}, d\mathbf{x}_{kn}) & \dots & \text{Var}(d\mathbf{x}_{kn}) \end{bmatrix} \tag{14}$$

where the variance terms can be obtained as

$$\begin{aligned}
\text{Var}(d\mathbf{x}_k) &= \mathbb{E}_{\tilde{\mathbf{x}}_k} [\text{Var}(\mathbf{f}(\tilde{\mathbf{x}}_k) | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k)] + \text{Var}(\mathbb{E}_{\mathbf{f}} [\mathbf{f}(\tilde{\mathbf{x}}_k) | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k]) \\
&= \mathbb{E}_{\tilde{\mathbf{x}}_k} [\text{Var}(d\mathbf{x}_k)] + \left(\mathbb{E}_{\tilde{\mathbf{x}}_k} [(d\mathbf{x}_k)^2] - \mathbb{E}_{\tilde{\mathbf{x}}_k} [d\mathbf{x}_k]^2 \right) \\
&= \int \left(\mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_k) - \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) (\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \right) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k \\
&\quad + \int \left(\mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) (\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} d\mathbf{X} \right)^2 p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k \\
&\quad - \left(\left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right)^{-1} d\mathbf{X} \right)^T \int \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) \mathcal{N}(\tilde{\mathbf{x}}_k | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) d\tilde{\mathbf{x}}_k \Big)^2.
\end{aligned} \tag{15}$$

The last term in the above equation can be represented by $\boldsymbol{\Psi}$ and \mathbf{q} defined earlier, then the equation becomes

$$\begin{aligned}
\text{Var}(d\mathbf{x}_k) &= \int \left(\mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_k) - \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) (\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \right) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k \\
&\quad + \int \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) \boldsymbol{\Psi} \boldsymbol{\Psi}^T \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k - (\boldsymbol{\Psi}^T \mathbf{q}_k)^2.
\end{aligned} \tag{16}$$

Re-arrange the above expressions by pulling the terms that are independent of $\tilde{\mathbf{x}}_k$ out of the integrals:

$$\begin{aligned}
\text{Var}(d\mathbf{x}_k) &= \sigma_s^2 - \text{tr} \left(\left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right)^{-1} \int \left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) \right) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k \right) \\
&\quad + \boldsymbol{\Psi}^T \left(\underbrace{\int \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k}_{\Phi_k} \right) \boldsymbol{\Psi} - (\boldsymbol{\Psi}^T \mathbf{q}_k)^2 \\
&= \sigma_s^2 - \text{tr} \left(\left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right)^{-1} \Phi_k \right) + \boldsymbol{\Psi}^T \Phi_k \boldsymbol{\Psi} - (\boldsymbol{\Psi}^T \mathbf{q}_k)^2,
\end{aligned} \tag{17}$$

where the integral terms Φ_k can be evaluated as

$$\begin{aligned}\Phi_{ij} &= \int \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{x}}_k) \mathbf{K}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{X}}) p(\tilde{\mathbf{x}}_k) d\tilde{\mathbf{x}}_k \\ &= \frac{\mathbf{K}(\tilde{\mathbf{X}}_i, \tilde{\boldsymbol{\mu}}_k) \mathbf{K}(\tilde{\mathbf{X}}_j, \tilde{\boldsymbol{\mu}}_k)}{|2\tilde{\boldsymbol{\Sigma}}_k(\mathbf{W}_i^{-1} + \mathbf{W}_j^{-1}) + \mathbf{I}|^{\frac{1}{2}}} \exp\left(\left(\frac{1}{2}\left(\frac{\mathbf{W}_j}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_i + \frac{\mathbf{W}_i}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_j\right) - \tilde{\boldsymbol{\mu}}_k\right)^{\text{T}}\right. \\ &\quad \left. \left(\tilde{\boldsymbol{\Sigma}} + \frac{1}{2}\mathbf{W}\right)^{-1} \tilde{\boldsymbol{\Sigma}}_k \mathbf{W}^{-1} \left(\frac{1}{2}\left(\frac{\mathbf{W}_j}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_i + \frac{\mathbf{W}_i}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_j\right) - \tilde{\boldsymbol{\mu}}_k\right)\right),\end{aligned}\quad (18)$$

where $\mathbf{W}_i, \mathbf{W}_j$ are the kernel parameters corresponding to output dimension i and j , respectively. The cross covariance terms can be obtained by

$$\text{Cov}(d\mathbf{x}_{ki}, d\mathbf{x}_{kj}) = \mathbb{E}_{\tilde{\mathbf{x}}_k} [d\mathbf{x}_{ki} d\mathbf{x}_{kj}] - \mathbb{E}_{\tilde{\mathbf{x}}_k} [d\mathbf{x}_{ki}] \mathbb{E}_{\tilde{\mathbf{x}}_k} [d\mathbf{x}_{kj}] \quad (19)$$

Similarly, it can be found that the first term is

$$\mathbb{E}_{\tilde{\mathbf{x}}_k} [d\mathbf{x}_{ki} d\mathbf{x}_{kj}] = \Psi_i^{\text{T}} \Phi_k \Psi_j, \quad (20)$$

where

$$\Psi_i = \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} d\mathbf{X}_i, \quad \Psi_j = \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I})^{-1} d\mathbf{X}_j. \quad (21)$$

Therefore

$$\text{Cov}[d\mathbf{x}_{ki}, d\mathbf{x}_{kj}] = \Psi_i^{\text{T}} \Phi_k \Psi_j - (\Psi_i^{\text{T}} \mathbf{q}_k)^{\text{T}} (\Psi_j^{\text{T}} \mathbf{q}_k). \quad (22)$$

The input-output cross-covariances can be obtained by

$$\text{Cov}[\mathbf{x}_k, d\mathbf{x}_k] = \mathbb{E}[\mathbf{x}_k d\mathbf{x}_k] - \mathbb{E}[\mathbf{x}_k] \mathbb{E}[d\mathbf{x}_k] = \mathbb{E}[\mathbf{x}_k \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)] - \boldsymbol{\mu}_k d\boldsymbol{\mu}_k. \quad (23)$$

The kernel or hyper-parameter $\Theta = (\sigma_n, \sigma_s, \mathbf{W})$ can be learned by maximizing the log-likelihood of the training outputs given the inputs:

$$\Theta^* = \underset{\Theta}{\text{argmax}} \left\{ \log \left(p(d\mathbf{X} | \tilde{\mathbf{X}}, \Theta) \right) \right\}. \quad (24)$$

where

$$\begin{aligned}\log \left(p(d\mathbf{X} | \tilde{\mathbf{X}}, \Theta) \right) &= -\frac{1}{2} d\mathbf{X}^{\text{T}} \left(\mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right)^{-1} d\mathbf{X} \\ &\quad - \frac{1}{2} \log \left| \mathbf{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I} \right| - \frac{H}{2} \log 2\pi.\end{aligned}\quad (25)$$

The optimization problem can be solved using numerical methods such as conjugate gradient.

3 Local dynamics models

In DDP related algorithms, a local model along a nominal trajectory $(\tilde{\mathbf{x}}_k, \tilde{\mathbf{u}}_k)$, where $k = 0, \dots, H$, is created based on: i) a first or second-order local approximation of the dynamics model; ii) a second-order local approximation of the value function. In our proposed PDDP framework, we will create a local model along a trajectory of state distribution-control pair $(p(\tilde{\mathbf{x}}_k), \tilde{\mathbf{u}}_k)$. We introduce the Gaussian augmented state vector $\mathbf{z}_k^x = [\boldsymbol{\mu}_k \text{vec}(\boldsymbol{\Sigma}_k)]^{\text{T}} \in \mathbb{R}^{n+n \times n}$ where $\text{vec}(\boldsymbol{\Sigma}_k)$ is the vectorization of $\boldsymbol{\Sigma}_k$. First, we create a local linear model of the dynamics. Based on eq.(10), the dynamics model with the augmented state can be written as

$$\mathbf{z}_{k+1}^x = \mathcal{F}(\mathbf{z}_k^x, \mathbf{u}_k). \quad (26)$$

Define the control and state variations $\delta \mathbf{z}_k^x = \mathbf{z}_k^x - \bar{\mathbf{z}}_k^x$ and $\delta \mathbf{u}_k = \mathbf{u}_k - \bar{\mathbf{u}}_k$. In this work we consider the first order expansion of the dynamics. More precisely we have

$$\delta \mathbf{z}_{k+1}^x = \mathcal{F}_k^x \delta \mathbf{z}_k^x + \mathcal{F}_k^u \delta \mathbf{u}_k, \quad (27)$$

where the Jacobians \mathcal{F}_k^x and \mathcal{F}_k^u are specified as

$$\begin{aligned}\mathcal{F}_k^x &= \nabla_{\mathbf{x}_k} \mathcal{F} = \begin{bmatrix} \frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\mu}_k} & \frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\Sigma}_k} \\ \frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \boldsymbol{\mu}_k} & \frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \boldsymbol{\Sigma}_k} \end{bmatrix} \in \mathbb{R}^{(n+n^2) \times (n+n^2)}, \\ \mathcal{F}_k^u &= \nabla_{\mathbf{u}_k} \mathcal{F} = \begin{bmatrix} \frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \mathbf{u}_k} \\ \frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \mathbf{u}_k} \end{bmatrix} \in \mathbb{R}^{(n+n^2) \times m}.\end{aligned}\quad (28)$$

where the partial derivatives can be evaluated as

$$\begin{aligned}\frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\mu}_k} &= \mathbf{I} + \frac{\partial \Psi^T \mathbf{q}_k}{\partial \boldsymbol{\mu}_k} \\ &= \mathbf{I} + \frac{\partial \Psi^T \mathbf{q}_k}{\partial \tilde{\boldsymbol{\mu}}_k} \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k}\end{aligned}\quad (29)$$

For each output dimension, the partial derivative $\frac{\partial \Psi_i^T \mathbf{q}_{ki}}{\partial \boldsymbol{\mu}_k} \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k}$ can be obtain as

$$\begin{aligned}\frac{\partial \Psi_i^T \mathbf{q}_{ki}}{\partial \tilde{\boldsymbol{\mu}}_k} \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^N \Psi_{ij} \frac{\partial \mathbf{q}_{kij}}{\partial \tilde{\boldsymbol{\mu}}_k} \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k} \\ &= \left(\sum_{j=1}^N \Psi_{ij} \mathbf{q}_{kij} (\tilde{\mathbf{x}}_k - \tilde{\boldsymbol{\mu}}_k)^T (\tilde{\boldsymbol{\Sigma}}_k + \mathbf{W}_j)^{-1} \right)^T \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k}.\end{aligned}\quad (30)$$

where $\frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k}$ can be easily obtained. Similarly, the partial derivatives of predictive mean with respect to state covariance for each output dimension can be found as

$$\begin{aligned}\frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\Sigma}_k} &= \frac{\partial \Psi_i^T \mathbf{q}_{ki}}{\partial \tilde{\boldsymbol{\Sigma}}_k} \frac{\partial \tilde{\boldsymbol{\Sigma}}_k}{\partial \boldsymbol{\Sigma}_k} = \sum_{j=1}^N \Psi_{ij} \frac{\partial \mathbf{q}_{kij}}{\partial \tilde{\boldsymbol{\Sigma}}_k} \frac{\partial \tilde{\boldsymbol{\Sigma}}_k}{\partial \boldsymbol{\Sigma}_k} \\ &= \sum_{j=1}^N \Psi_{ij} \mathbf{q}_{kij} \left(-\frac{1}{2} \left((\mathbf{W}_i^{-1} \tilde{\boldsymbol{\Sigma}}_k + \mathbf{I})^{-1} \mathbf{W}_i^{-1} \right)^T - \frac{1}{2} (\tilde{\mathbf{x}}_{ki} - \tilde{\boldsymbol{\mu}}_k)^T \right. \\ &\quad \left. \frac{\partial (\mathbf{W}_j + \tilde{\boldsymbol{\Sigma}}_k)^{-1}}{\partial \tilde{\boldsymbol{\Sigma}}_k} (\tilde{\mathbf{x}}_{ki} - \tilde{\boldsymbol{\mu}}_k) \right) \frac{\partial \tilde{\boldsymbol{\Sigma}}_k}{\partial \boldsymbol{\Sigma}_k}.\end{aligned}\quad (31)$$

where $\frac{\partial \tilde{\boldsymbol{\Sigma}}_k}{\partial \boldsymbol{\Sigma}_k}$ can be easily obtained. The partial derivatives of covariance with respect to input mean for each output dimension can be evaluated as

$$\frac{\partial \boldsymbol{\Sigma}_{(k+1)ij}}{\partial \boldsymbol{\mu}_k} = \left(\frac{\partial d\boldsymbol{\Sigma}_{kij}}{\partial \tilde{\boldsymbol{\mu}}_k} + \frac{\partial \text{Cov}[\mathbf{x}_{ki}, \mathbf{d}\mathbf{x}_{kj}]}{\partial \tilde{\boldsymbol{\mu}}_k} + \frac{\partial \text{Cov}[\mathbf{d}\mathbf{x}_{ki}, \mathbf{x}_{kj}]}{\partial \tilde{\boldsymbol{\mu}}_k} \right) \frac{\partial \tilde{\boldsymbol{\mu}}_k}{\partial \boldsymbol{\mu}_k} \quad (32)$$

where

$$\begin{aligned}\frac{\partial d\boldsymbol{\Sigma}_{kij}}{\partial \tilde{\boldsymbol{\mu}}_k} &= \Psi_i^T \left(\frac{\partial \Phi_k}{\partial \tilde{\boldsymbol{\mu}}_k} - \frac{\partial \mathbf{q}_i}{\tilde{\boldsymbol{\mu}}_k} \mathbf{q}_j^T - \mathbf{q}_i \frac{\partial \mathbf{q}_j}{\tilde{\boldsymbol{\mu}}_k} \right) \Psi_j + \left(-(\mathbf{K} + \sigma_n \mathbf{I})^{-1} \frac{\partial \Phi_k}{\partial \tilde{\boldsymbol{\mu}}_k} \right) \\ \frac{\partial \Phi_{kij}}{\partial \tilde{\boldsymbol{\mu}}_k} &= \Phi_{kij} \left(\frac{\mathbf{W}_j}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_i + \frac{\mathbf{W}_i}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_j - \tilde{\boldsymbol{\mu}}_k \right)^T \left(\frac{1}{\mathbf{W}_i^{-1} + \mathbf{W}_j^{-1} + \tilde{\boldsymbol{\Sigma}}_k} \right)^{-1} \\ \frac{\partial \text{Cov}[\mathbf{d}\mathbf{x}_k, \mathbf{x}_k]}{\partial \tilde{\boldsymbol{\mu}}_k} &= \frac{\tilde{\boldsymbol{\Sigma}}_k}{\tilde{\boldsymbol{\Sigma}}_k + \mathbf{W}} \sum_{i=0}^{n+m} \Psi \left((\tilde{\mathbf{x}}_{ki} - \tilde{\boldsymbol{\mu}}_k) \frac{\partial \mathbf{q}_{ki}}{\partial \tilde{\boldsymbol{\mu}}_k} + \mathbf{q}_{ki} \mathbf{I} \right).\end{aligned}\quad (33)$$

The partial derivatives of covariance with respect to input covariance for each output dimension can be evaluated as

$$\frac{\partial \boldsymbol{\Sigma}_{(k+1)ij}}{\partial \boldsymbol{\Sigma}_k} = \mathbf{I} + \left(\frac{\partial d\boldsymbol{\Sigma}_{kij}}{\partial \tilde{\boldsymbol{\Sigma}}_k} + \frac{\partial \text{Cov}[\mathbf{x}_{ki}, \mathbf{d}\mathbf{x}_{kj}]}{\partial \tilde{\boldsymbol{\Sigma}}_k} + \frac{\partial \text{Cov}[\mathbf{d}\mathbf{x}_{ki}, \mathbf{x}_{kj}]}{\partial \tilde{\boldsymbol{\Sigma}}_k} \right) \frac{\partial \tilde{\boldsymbol{\Sigma}}_k}{\partial \boldsymbol{\Sigma}_k}. \quad (34)$$

where

$$\begin{aligned}
\frac{\partial d\boldsymbol{\Sigma}_{kij}}{\partial \tilde{\boldsymbol{\Sigma}}} &= \Psi_i^T \left(\frac{\partial \Phi_k}{\partial \boldsymbol{\Sigma}_k} - \frac{\partial \mathbf{q}_i}{\boldsymbol{\Sigma}_k} \mathbf{q}_j^T - \mathbf{q}_i \frac{\partial \mathbf{q}_j}{\boldsymbol{\Sigma}_k} \right) \Psi_j + \left(-(\mathbf{K}^+ \sigma_n \mathbf{I})^{-1} \frac{\partial \Phi_k}{\partial \boldsymbol{\Sigma}_k} \right) \\
\frac{\partial \Phi_{kij}}{\partial \tilde{\boldsymbol{\Sigma}_k}} &= -\frac{1}{2} \Phi_{kij} \left[\left(\left(\frac{\mathbf{W}_i + \mathbf{W}_j}{\mathbf{W}_i \mathbf{W}_j} \boldsymbol{\Sigma}_k + \mathbf{I} \right)^{-1} \left(\frac{\mathbf{W}_i + \mathbf{W}_j}{\mathbf{W}_i \mathbf{W}_j} \right) \right)^T \right. \\
&\quad \left. - \left(\frac{\mathbf{W}_j}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_i + \frac{\mathbf{W}_i}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_j - \tilde{\boldsymbol{\mu}}_k \right)^T \left(\frac{1}{\mathbf{W}_i^{-1} + \mathbf{W}_j^{-1}} + \boldsymbol{\Sigma}_k \right)^{-1} \right. \\
&\quad \left. \left(\frac{\mathbf{W}_j}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_i + \frac{\mathbf{W}_i}{\mathbf{W}_i + \mathbf{W}_j} \tilde{\mathbf{x}}_j - \tilde{\boldsymbol{\mu}}_k \right) \right) \\
\frac{\partial \text{Cov}[\mathbf{d}\mathbf{x}_k, \mathbf{x}_k]}{\partial \tilde{\boldsymbol{\Sigma}_k}} &= \left(\frac{1}{\tilde{\boldsymbol{\Sigma}_k} + \mathbf{W}} + \frac{\partial((\tilde{\boldsymbol{\Sigma}} + \mathbf{W})^{-1})}{\partial \tilde{\boldsymbol{\Sigma}_k}} \right) \sum_{i=1}^{m+n} \Psi_{ki} \mathbf{q}_{ki} (\tilde{\mathbf{x}} - t_i - \tilde{\boldsymbol{\mu}}_{ki}) + \\
&\quad \tilde{\boldsymbol{\Sigma}}_k \left(\frac{1}{\tilde{\boldsymbol{\Sigma}_k} + \mathbf{W}} \right) \sum_{i=1}^{n+m} \Psi_{ki} (\tilde{\mathbf{x}}_{ki} - \tilde{\boldsymbol{\mu}}_{ki}) \frac{\partial \mathbf{q}_{ki}}{\partial \tilde{\boldsymbol{\Sigma}_k}}. \tag{35}
\end{aligned}$$

We have found the expression of $\frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\mu}_k}$, $\frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \boldsymbol{\Sigma}_k}$, $\frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \boldsymbol{\mu}_k}$, $\frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \boldsymbol{\Sigma}_k}$ analytically. The partial derivatives with respect to control $\frac{\partial \boldsymbol{\mu}_{k+1}}{\partial \mathbf{u}_k}$, $\frac{\partial \boldsymbol{\Sigma}_{k+1}}{\partial \mathbf{u}_k}$ can be found similarly.

4 Cost function

In classic DDP/LQG and most optimal control problems, the following quadratic cost function was used:

$$\mathcal{L}(\mathbf{x}_k, \mathbf{u}_k) = (\mathbf{x}_k - \mathbf{x}_k^{goal})^T \mathbf{Q} (\mathbf{x}_k - \mathbf{x}_k^{goal}) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k, \tag{36}$$

where \mathbf{x}_k^{goal} is the target state. In probabilistic DDP, given the distribution $p(\mathbf{x}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. Let $\sigma_{kij} = [\boldsymbol{\Sigma}_k]_{ij}$ and $q_{ij} = [\mathbf{Q}]_{ij}$. The expectation of original quadratic cost function can be obtained as:

$$\begin{aligned}
\mathbb{E}[\mathcal{L}(\mathbf{x}_k, \mathbf{u}_k)] &= \mathbb{E}[(\mathbf{x}_k - \mathbf{x}_k^{goal})^T \mathbf{Q} (\mathbf{x}_k - \mathbf{x}_k^{goal}) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k] \\
&= \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n q_{ij} (x_{ki} - x_i^{goal})(x_{kj} - x_j^{goal}) \right] + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \sum_{i=1}^n \sum_{j=1}^n q_{ij} \mathbb{E} \left[(x_{ki} - x_i^{goal})(x_{kj} - x_j^{goal}) \right] + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \sum_{i=1}^n \sum_{j=1}^n q_{ij} \left(\text{Cov}((x_{ki} - x_i^{goal}), (x_{kj} - x_j^{goal})) + \right. \\
&\quad \left. \mathbb{E}[x_{ki} - x_i^{goal}] \mathbb{E}[x_{kj} - x_j^{goal}] \right) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \sum_{i=1}^n \sum_{j=1}^n q_{ij} (\sigma_{kij} + (\mu_{ki} - x_i^{goal})(\mu_{kj} - x_j^{goal})) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \sum_{i=1}^n \sum_{j=1}^n q_{ij} \sigma_{tji} + \sum_{i=1}^n \sum_{j=1}^n q_{ij} (\mu_{ki} - x_i^{goal})(\mu_{kj} - x_j^{goal}) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \sum_{i=1}^n [\mathbf{Q} \boldsymbol{\Sigma}_k]_{ii} + (\boldsymbol{\mu}_k - \mathbf{x}_k^{goal})^T \mathbf{Q} (\boldsymbol{\mu}_k - \mathbf{x}_k^{goal}) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k \\
&= \text{tr}(\mathbf{Q} \boldsymbol{\Sigma}_k) + (\boldsymbol{\mu}_k - \mathbf{x}_k^{goal})^T \mathbf{Q} (\boldsymbol{\mu}_k - \mathbf{x}_k^{goal}) + \mathbf{u}_k^T \mathbf{R} \mathbf{u}_k
\end{aligned}$$

Therefore, in this paper we use the cost function with the augmented state

$$\mathcal{L}(\mathbf{z}_k^x, \mathbf{u}_k) = \text{tr}(\mathbf{Q}\boldsymbol{\Sigma}_k) + (\boldsymbol{\mu}_k - \mathbf{x}_k^{goal})^\top \mathbf{Q}(\boldsymbol{\mu}_k - \mathbf{x}_k^{goal}) + (\mathbf{u}_k)^\top \mathbf{R}\mathbf{u}_k. \quad (37)$$

The partial derivatives of the above cost function with respect to $(\mathbf{z}_k^x, \mathbf{u}_k)$ can be easily obtained by

$$\begin{aligned} \frac{\partial}{\partial \mathbf{z}_k^x} \mathcal{L}(\mathbf{z}_k^x, \mathbf{u}_k) &= \left[\frac{\partial}{\partial \boldsymbol{\mu}_k} \mathcal{L}(\mathbf{z}_k^x, \mathbf{u}_k) \quad \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \mathcal{L}(\mathbf{z}_k^x, \mathbf{u}_k) \right]^\top, \\ &= \left[2(\boldsymbol{\mu}_k - \mathbf{x}_k^{goal})^\top \mathbf{Q} \quad \mathbf{Q} \right]^\top, \\ \frac{\partial}{\partial \mathbf{u}_k} \mathcal{L}(\mathbf{z}_k^x, \mathbf{u}_k) &= 2(\mathbf{u}_k)^\top \mathbf{R}. \end{aligned}$$

The cost function scales linearly with the state covariance, therefore the exploration strategy of PDDP is balanced between the distance from the target and the variance of the state and avoids high risk explorations.