

Differentiable Optimization

Chen Wang, Krishna Murthy Jatavallabhula, and Mustafa Mukadam

4.1 Introduction

As presented in Chapter 1, the design of a contemporary SLAM system generally adheres to a front-end and back-end architecture. In this structure, the front-end is typically responsible for pre-processing sensor data and generating an initial estimate of the robot’s trajectory and the map of the environment, while the back-end refines these initial estimates to improve overall accuracy. Recent advances in machine learning have provided new approaches, based on deep neural networks, that have the potential to enhance some of the functionalities in the SLAM front-end. For instance, deep learning-based methods can exhibit impressive performance in feature detection and matching [976, 275, 1262] and front-end motion estimation [1165, 1086]. These methods train a neural network from a large dataset of examples, and then make estimations without being explicitly programmed to perform the task. Meanwhile, geometry-based techniques persist as an essential element for the SLAM back-end, primarily due to their generality and effectiveness in producing a globally consistent estimate by solving an optimization problem.

While in principle one could just “plug” a learning-based SLAM front-end in the SLAM architecture and feed the corresponding outputs to the back-end, the use of learning-based techniques opens the door for a less unidirectional information exchange. In particular, the back-end can now provide feedback to the front-end, enabling it to learn directly from the back-end estimates in a way that the two modules can more harmoniously cooperate to reduce the estimation errors. Reconciling geometric approaches with deep learning to leverage their complementary strengths is a common thread in a large body of recent work in SLAM. In particular, an emerging trend is to differentiate through geometry-based optimization problems arising in the SLAM back-end. Intuitively, differentiating through an optimization problem allows understanding how the optimal solution of that problem (*e.g.*, our SLAM estimate) depends on the parameters of that problem — in our case, the measurements produced by a learning-based front-end; this in turns allows optimizing the front-end to maximize the SLAM accuracy. One could think about this as a *bilevel optimization* problem, *i.e.*, an upper-level optimization process subject

to a lower-level optimization — in particular, a neural network based-optimization to train the front-end, subject to a geometry-based optimization that computes the SLAM solution for a given front-end output.

The ability to compute gradients end-to-end through an optimization is the core of solving a bilevel optimization problem, which allows neural models to take advantage of geometric priors captured by the optimization. The flexibility of such a scheme has led to promising state-of-the-art results in a wide range of applications such as structure from motion [1084], motion planning [88, 1212], SLAM [519, 1086], bundle adjustment [1072, 1262], state estimation [1239, 188], and image alignment [724].

In this chapter, we illustrate the basics of how to differentiate through nonlinear least squares problems, such as the ones arising in SLAM. Specifically, Section 4.1.1 restates the non-linear least square (NLS) problem. Section 4.2 describes how to differentiate through the NLS problem. Section 4.3 shows how to differentiate problems defined on manifold. Section 4.4 discusses numerical challenges of the above differentiation and introduces related machine learning libraries. Finally, Section 4.5 provides examples of differentiable optimization in contemporary SLAM systems.

4.1.1 Recap on Nonlinear Least Squares

Non-linear least squares (NLS) estimate the parameters of a model by minimizing the sum of the squares of the mismatch between observed values and those predicted by the model. Unlike linear least squares, NLS involves a model that is non-linear in the parameters. Beyond our factors graphs in Chapter 1, this approach is widely used in many fields such as statistics, physics, and engineering, where it is useful for fitting complex models to data when the relationship between variables is not straightforward, enabling more accurate and robust predictions.

Specifically, NLS aim to find variables $\mathbf{x} \in \mathbb{R}^n$ by solving:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}) = \arg \min_{\mathbf{x}} = \frac{1}{2} \sum_i \|\underbrace{w_i \mathbf{c}_i(\mathbf{x}_i)}_{\mathbf{r}_i(\mathbf{x}_i)}\|^2, \quad (4.1)$$

where the objective $\mathcal{L}(\mathbf{x})$ is a sum of squared vector-valued residual terms \mathbf{r}_i , each a function of $\mathbf{x}_i \subset \mathbf{x}$ that are (non-disjoint) subsets of the optimization variables $\mathbf{x} = \{\mathbf{x}_i\}$. While for now we assume \mathbf{x}_i to be vectors, later in the chapter we generalize the discussion to the case where the variables belong to a manifold. For flexibility, here we represent a residual $\mathbf{r}_i(\mathbf{x}_i) = w_i \mathbf{c}_i(\mathbf{x}_i)$ as a product of a weight w_i and vector cost \mathbf{c}_i .

As explained in Chapter 1, a NLS is normally solved by iteratively linearizing the nonlinear objective around the current variables to get the linear system $(\sum_i \mathbf{J}_i^\top \mathbf{J}_i) \delta \mathbf{x} = (\sum_i \mathbf{J}_i^\top \mathbf{r}_i)$, then solving the linear system to find the update $\delta \mathbf{x}$, and finally updating the variables $\mathbf{x} \leftarrow \mathbf{x} + \delta \mathbf{x}$, until convergence. We have also

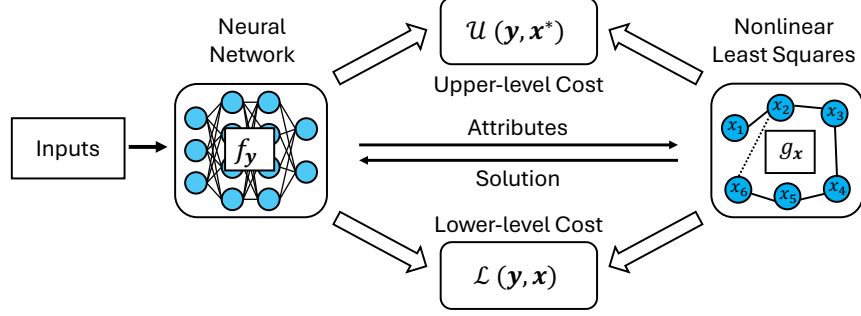


Figure 4.1 A modern SLAM system often involves both neural networks and nonlinear least squares. To eliminate compound errors introduced by optimizing the two modules separately, we can optimize the system in an end-to-end manner by formulating the entire system as a bilevel optimization, which involves an upper-level cost and a lower-level cost.

commented in Chapter 2 that the addition in the update step is more generally a retraction mapping for variables that belong to a manifold. In the linear system, $\mathbf{J}_i = [\partial \mathbf{r}_i / \partial \mathbf{x}_i]$ are the Jacobians of residuals with respect to the variables. This iterative method above, called Gauss-Newton (GN), is a nonlinear optimizer that is (approximately) second-order, since $\sum_i \mathbf{J}_i^\top \mathbf{J}_i$ is an approximation of the Hessian. To improve robustness and convergence, variations like Levenberg-Marquardt (LM) dampen the linear system, while others adjust the step size for the update with line search, *e.g.*, Dogleg introduced in Chapter 1.

4.2 Differentiation Through Nonlinear Least Squares

To seamlessly merge deep learning with nonlinear least squares, differentiable nonlinear least squares (DNLS) are often required to solve the optimization problem illustrated in Figure 4.1. This necessitates gradients of the solution \mathbf{x}^* with respect to any upper-level neural model parameters \mathbf{y} that parameterize the objective $\mathcal{U}(\mathbf{x}; \mathbf{y})$ and, in turn, any costs $c_i(\mathbf{x}_i; \mathbf{y})$ or initialization for variables $\mathbf{x}_{\text{init}}(\mathbf{y})$. The goal is to learn these parameters \mathbf{y} end-to-end with a lower-level learning objective \mathcal{L} defined as a function of \mathbf{x} . This results in a bilevel optimization (BLO), which can be written as:

$$\mathbf{y}^* = \arg \min_{\mathbf{y} \in \Theta} \mathcal{U}(\mathbf{y}, \mathbf{x}^*), \quad (4.2a)$$

$$\text{s. t. } \mathbf{x}^* = \arg \min_{\mathbf{x} \in \Psi} \mathcal{L}(\mathbf{y}, \mathbf{x}), \quad (4.2b)$$

where $\mathcal{L} : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a lower-level (LL) cost, $\mathcal{U} : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a upper-level (UL) cost, $\mathbf{x} \in \Psi$ and $\mathbf{y} \in \Theta$ are the feasible sets.

In practice, the variables \mathbf{x} are often parameters with explicit physical meanings

such as camera poses, while \mathbf{y} are parameters without physical meanings such as weights in a neural network. We next present two examples to explain this.

Example 4.1 Imagine a SLAM system that leverages a neural network (parameterized by \mathbf{y}) for feature extraction/matching, while utilizing bundle adjustment (BA) for pose estimation (parameterized by \mathbf{x}), which take the feature matching as an input. In this example, the UL cost (4.2a) can be feature matching error for optimizing the network, while the LL cost L (4.2b) can be the reprojection error for BA. Intuitively, the optimal solution \mathbf{x}^* for the camera poses and landmark positions plays the role of a supervisory signal in the neural network training. Therefore, optimizing the BLO (4.2) allows us to further reduce the matching error via back-propagating the BA reprojection errors [1262].

Example 4.2 Imagine a full SLAM system that uses a neural network for front-end pose estimation, while leverages pose-graph optimization (PGO) as the back-end to eliminate odometry drifts. In this example, both UL and LL costs can be the pose-graph error. The difference is the UL cost optimizes the network parameterized by \mathbf{y} , while the LL cost optimizes the camera poses parameterized by \mathbf{x} . As a result, the front-end network can leverage global geometric knowledge obtained through pose-graph optimization by back-propagating the pose residuals from the back-end PGO [354].

BLO is a long-standing and well-researched problem [1148, 524, 686]. Solving a BLO often relies on gradient-descent techniques. Specifically, the UL optimization performs updates in the form $\mathbf{y} \leftarrow \mathbf{y} + \delta\mathbf{y}$, where $\delta\mathbf{y}$ is a step in the direction of the negative gradient. Therefore, we need compute the gradient of \mathcal{U} with respect to the UL variable \mathbf{y} , which can be written as

$$\nabla_{\mathbf{y}}\mathcal{U} = \frac{\partial\mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial\mathbf{y}} + \frac{\partial\mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial\mathbf{x}^*} \frac{\partial\mathbf{x}^*(\mathbf{y})}{\partial\mathbf{y}}, \quad (4.3)$$

where the term $\frac{\partial\mathbf{x}^*(\mathbf{y})}{\partial\mathbf{y}}$ involves indirect gradient computation. Since other direct gradient terms in (4.3) are easy to obtain, the challenge of solving a BLO (4.2) is to compute the term $\frac{\partial\mathbf{x}^*(\mathbf{y})}{\partial\mathbf{y}}$. For this purpose, a series of techniques have been developed from either explicit or implicit perspectives. This involves recurrent differentiation through dynamical systems or implicit differentiation theory, which are often referred to as **unrolled differentiation** and **implicit differentiation**, respectively. These algorithms have been summarized in [686, 1148] and here we list a generic framework incorporating both methods in Algorithm 1. We next explain the unrolled differentiation and implicit differentiation, respectively.

4.2.1 Unrolled Differentiation

Unrolled Differentiation needs automatic differentiation (AutoDiff) through the LL optimization to solve a BLO problem. Specifically, given an initialization $\mathbf{x}_0 =$

Algorithm 1 Solving BLO by *Unrolled Differentiation* or *Implicit Differentiation*.

- 1: **Initialization:** $\mathbf{y}_0, \mathbf{x}_0$.
- 2: **while** Not Convergent ($\|\mathbf{y}_{k+1} - \mathbf{y}_k\|$ is large enough) **do**
- 3: Obtain \mathbf{x}_T by solving (4.2b) by a generic optimizer \mathcal{O} with T steps.
- 4: Efficient estimation of upper-level gradients in (4.3) via
- 5: **Unrolled Differentiation:** $\hat{\nabla}_{\mathbf{y}_k} U = \frac{\partial U(\mathbf{y}_k, \mathbf{x}_T)}{\partial \mathbf{y}_k}$ via AutoDiff in (4.7).
- 6: **Implicit Differentiation (Algorithm 2):** Compute

$$\hat{\nabla}_{\mathbf{y}_k} \mathcal{U} = \frac{\partial \mathcal{U}}{\partial \mathbf{y}_k} \Big|_{\mathbf{x}_T} + \frac{\partial \mathcal{U}}{\partial \mathbf{x}^*} \frac{\partial \mathbf{x}^*}{\partial \mathbf{y}_k} \Big|_{\mathbf{x}_T},$$

where the implicit derivatives $\frac{\partial \mathbf{x}^*}{\partial \mathbf{y}_k}$ can be obtained by solving an equation derived via lower-level optimality conditions (surveyed in following sections).

- 7:
 - 8: Compute \mathbf{y}_{k+1} via gradients using $\hat{\nabla}_{\mathbf{y}_k} U$.
 - 9: **end while**
-

$\Phi_0(\mathbf{y})$ at step $t = 0$, the iterative process of unrolled LL optimization is

$$\mathbf{x}_t = \Phi_t(\mathbf{x}_{t-1}; \mathbf{y}), \quad t = 1, \dots, T, \quad (4.4)$$

where Φ_t denotes an updating scheme based on the LL problem at the t -th step and T is the number of iterations. One updating scheme is the gradient descent:

$$\Phi_t(\mathbf{x}_{t-1}; \mathbf{y}) = \mathbf{x}_{t-1} - \eta_t \cdot \frac{\partial \mathcal{L}(\mathbf{x}_{t-1}, \mathbf{y})}{\partial \mathbf{x}_{t-1}}, \quad (4.5)$$

where η_t is a learning rate and the term $\frac{\partial \mathcal{L}(\mathbf{x}_{t-1}, \mathbf{y})}{\partial \mathbf{x}_{t-1}}$ can be computed from AutoDiff.¹ Therefore, we can compute the $\nabla_{\mathbf{y}} \mathcal{U}(\mathbf{y})$ by substituting \mathbf{x}_T approximately for \mathbf{x}^* and the full unrolled system can be defined as

$$\mathbf{x}^* \approx \mathbf{x}_T = \Phi(\mathbf{y}) = (\Phi_T \circ \dots \circ \Phi_1 \circ \Phi_0)(\mathbf{y}), \quad (4.6)$$

where the symbol \circ denotes the function composition. As a result, we only need to consider the following problem instead of a bilevel optimization in (4.2):

$$\min_{\mathbf{y} \in \Theta} \mathcal{U}(\mathbf{y}, \Phi(\mathbf{y})), \quad (4.7)$$

which needs to compute $\frac{\partial \Phi(\mathbf{y})}{\partial \mathbf{y}}$ via AutoDiff instead of calculating (4.3). It is worth noting that there exist two approaches for computing the recurrent gradients, one of which corresponds to backward propagation in a reverse-mode way [864], and the other corresponds to the forward-mode way [934]. We omit the details of the two approaches of AutoDiff and refer the readers to the AutoDiff libraries such as

¹ While here we mention gradient descent, the same ideas can be extended to other iterative optimization methods, such as Gauss-Newton.

PyTorch [856] for deep learning and PyPose [1147] and Theseus [880] for SLAM. A review of these approaches can also be found in Liu et al. [686].

4.2.2 Truncated Unrolled Differentiation

The reverse and forward modes are two precise recurrent gradient calculation methods but are time-consuming with the full iterative propagation. This is due to the complicated long-term dependencies of the UL problem on \mathbf{x}_t , where $t = 0, 1, \dots, T$. This difficulty is further aggravated when both \mathbf{y} and \mathbf{x} are high-dimensional vectors. To overcome this challenge, the truncated unrolled differentiation has been investigated as a way to compute high-quality approximate gradients with significantly less computation time and memory. Specifically, by ignoring the long-term dependencies and approximating the gradient of (4.5) with partial history, i.e., storing only the last M iterations ($t = T, T-1, \dots, T-M$), we can significantly reduce the time and space complexity. It has been proved by Shaban et al. [997] that using fewer backward steps to compute the gradients could perform comparably to optimization with the exact one, while requiring much less memory and computation.

In case of more stringent computational and memory constraints, truncated unrolled differentiation is still often a bottleneck in modern robotic applications. Therefore, researchers have also tried to further simplify the truncated differentiation by only performing a one-step iteration in (4.4) to remove the recursive structure [683], i.e.,

$$\nabla_{\mathbf{y}} \mathcal{U} = \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}_1(\mathbf{y}))}{\partial \mathbf{y}} + \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}_1(\mathbf{y}))}{\partial \mathbf{x}_1} \frac{\partial \mathbf{x}_1(\mathbf{y})}{\partial \mathbf{y}}, \quad (4.8)$$

where the term $\frac{\partial \mathbf{x}_1(\mathbf{y})}{\partial \mathbf{y}}$ is a Hessian that can be calculated from (4.5) as

$$\frac{\partial \mathbf{x}_1(\mathbf{y})}{\partial \mathbf{y}} = -\frac{\partial^2 \mathcal{L}(\mathbf{x}_0, \mathbf{y})}{\partial \mathbf{x}_0 \partial \mathbf{y}}. \quad (4.9)$$

Since calculating a Hessian is time-consuming in some applications, we can resort to numerical solutions that apply small perturbations to the variables \mathbf{x} and calculate an approximation of the second term in (4.8) as a whole:

$$\frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}_1(\mathbf{y}))}{\partial \mathbf{x}_1} \frac{\partial \mathbf{x}_1(\mathbf{y})}{\partial \mathbf{y}} \approx \frac{\frac{\partial \mathcal{L}(\mathbf{x}_0^+, \mathbf{y})}{\partial \mathbf{y}} - \frac{\mathcal{L}(\mathbf{x}_0^-, \mathbf{y})}{\partial \mathbf{y}}}{2\epsilon}, \quad (4.10)$$

where ϵ is a small scalar and $\mathbf{x}_0^\pm = \mathbf{x}_0 \pm \epsilon \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}_1(\mathbf{y}))}{\partial \mathbf{x}_1}$ is a small perturbation. This bypasses an explicit calculation of the Jacobian $\frac{\partial \mathbf{x}_1(\mathbf{y})}{\partial \mathbf{y}}$. Nevertheless, we need to pay attention to the perturbation model if non-Euclidean variables are involved, e.g., variables belonging to Lie Groups. Fortunately, the AutoDiff of Lie Group for Hessian-vector and Jacobian-vector multiplications are supported in modern libraries, such as PyPose [1147], which will be introduced in Section 4.4.

4.2.3 Implicit Differentiation

It is intuitive that the term $\frac{\partial \mathbf{x}^*(\mathbf{y})}{\partial \mathbf{y}}$ in (4.3) is dependent on the LL cost (4.2b), thus *implicit differentiation* can be used to derive a solution to the gradient.

Example 4.3 In calculus, *implicit differentiation* refers to the method makes use of the chain rule to differentiate implicit function. To differentiate an implicit function $y(x)$, defined by an equation $R(x, y) = 0$, it is not generally possible to solve it explicitly for y and then differentiate. Instead, one can totally differentiate $R(x, y) = 0$ with respect to x and then solve the resulting linear equation for $\frac{dy}{dx}$ to explicitly get the derivative in terms of x and y . For instance, consider an implicit function $x + y + 5 = 0$, differentiating it with respect to x on its both sides gives $\frac{dy}{dx} + \frac{dx}{dx} + \frac{d}{dx}(5) = 0 \Rightarrow \frac{dy}{dx} + 1 + 0 = 0$. Solving for $\frac{dy}{dx}$ gives $\frac{dy}{dx} = -1$.

Assume the LL cost \mathcal{L} is at least twice differentiable w.r.t. both \mathbf{y} and \mathbf{x} , then we have $\frac{\partial \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y})} = 0$ due to the optimality condition where \mathbf{x}^* is a stationary point. Derive the equation $\frac{\partial \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y})} = 0$ on both sides w.r.t. \mathbf{y} giving us

$$\frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{y}} + \frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{x}^*(\mathbf{y})} \cdot \frac{\partial \mathbf{x}^*(\mathbf{y})}{\partial \mathbf{y}} = 0. \quad (4.11)$$

This leads to the indirect gradient $\frac{\partial \mathbf{x}^*(\mathbf{y})}{\partial \mathbf{y}}$ as

$$\frac{\partial \mathbf{x}^*(\mathbf{y})}{\partial \mathbf{y}} = - \left(\frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{x}^*(\mathbf{y})} \right)^{-1} \frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{y}}, \quad (4.12)$$

The strength of (4.12) is that we convert the indirect gradient among the variables \mathbf{y} and \mathbf{x} to direct gradients of \mathcal{L} at the cost of an inversion of Hessian matrix. However, the weakness is that a Hessian is often too large to compute, thus it is common to solve a linear system leveraging the fast Hessian-vector product.

Example 4.4 Assume both UL and LU costs have a network with merely 1 million (10^6) parameters (32-bit float numbers), thus each network only needs a space of $10^6 \times 4\text{Byte} = 4\text{MB}$ to store, while their Hessian matrix needs a space of $(10^6)^2 \times 4\text{Byte} = 4\text{TByte}$ to store. This indicates that a Hessian matrix cannot even be explicitly stored in the memory of a low-power computer, thus directly calculating its inversion is more impractical.

Recollect that our goal is to compute the gradient in (4.3), substituting (4.12) into (4.3) gives us:

$$\begin{aligned} \nabla_{\mathbf{y}} \mathcal{U} &= \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial \mathbf{y}} - \underbrace{\frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial \mathbf{x}^*}}_{\mathbf{v}^\top} \underbrace{\left(\frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{x}^*(\mathbf{y})} \right)^{-1}}_{(\mathbf{H}^\top)^{-1}} \frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{y}} \\ &= \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial \mathbf{y}} - \mathbf{q}^\top \cdot \frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{x}^*(\mathbf{y}) \partial \mathbf{y}} \end{aligned} \quad (4.13)$$

Then we can solve the linear system $\mathbf{H}\mathbf{q} = \mathbf{v}$ for \mathbf{q}^\top by optimizing

$$\mathbf{q}^* = \min \arg_{\mathbf{q}} Q(\mathbf{q}) = \min \arg_{\mathbf{q}} \frac{1}{2} \mathbf{q}^\top \mathbf{H} \mathbf{q} - \mathbf{q}^\top \mathbf{v}, \quad (4.14)$$

using efficient linear solvers such as a simple gradient descent or conjugate gradient method [467]. For gradient descent, we need to compute the gradient of Q as $\frac{\partial Q(\mathbf{q})}{\partial \mathbf{q}} = \mathbf{H}\mathbf{q} - \mathbf{v}$, where $\mathbf{H}\mathbf{q}$ can be computed using the fast Hessian-vector product, i.e., a Hessian-vector product is the gradient of a gradient-vector product:

$$\mathbf{H}\mathbf{q} = \frac{\partial^2 \mathcal{L}}{\partial \mathbf{x} \partial \mathbf{x}} \cdot \mathbf{q} = \frac{\partial \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \cdot \mathbf{q} \right)}{\partial \mathbf{x}}, \quad (4.15)$$

where $\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \cdot \mathbf{q}$ is a scalar. This means that the Hessian matrix \mathbf{H} is not explicitly computed or stored. We summarize the computation of implicit differentiation with linear systems in Algorithm 2. The algorithm using a conjugate gradient is similar.

Algorithm 2 Computing *Implicit Differentiation* via Linear System.

- 1: **Input:** The current UL variable \mathbf{y} and the optimal LL variable \mathbf{x}^* .
- 2: **Initialization:** $k = 1$, learning rate η .
- 3: **while** Not Convergent ($\|\mathbf{q}_k - \mathbf{q}_{k-1}\|$ is large enough) **do**
- 4: Perform gradient descent:

$$\mathbf{q}_k = \mathbf{q}_{k-1} - \eta (\mathbf{H}\mathbf{q}_{k-1} - \mathbf{v}), \quad (4.16)$$

where $\mathbf{H}\mathbf{q}_{k-1}$ is computed via the fast Hessian-vector product.

- 5: **end while**
- 6: Assign $\mathbf{q} = \mathbf{q}_k$
- 7: Compute $\nabla_{\mathbf{y}} \mathcal{U}$ in (4.3) as:

$$\nabla_{\mathbf{y}} \mathcal{U} = \frac{\partial \mathcal{U}(\mathbf{y}, \mathbf{x}^*)}{\partial \mathbf{y}} - \underbrace{\left(\frac{\partial^2 \mathcal{L}(\mathbf{x}^*(\mathbf{y}), \mathbf{y})}{\partial \mathbf{y} \partial \mathbf{x}^*(\mathbf{y})} \cdot \mathbf{q} \right)^\top}_{(\mathbf{H}_{\mathbf{y}\mathbf{x}} \cdot \mathbf{q})^\top}, \quad (4.17)$$

where $\mathbf{H}_{\mathbf{y}\mathbf{x}} \cdot \mathbf{q}$ can also be computed efficiently using the Hessian-vector product.

Approximations. Implicit differentiation is complicated to implement but there is one approximation, which is to ignore the implicit components and only use the direct part $\hat{\nabla}_{\mathbf{y}} U \approx \frac{\partial U}{\partial \mathbf{y}} \Big|_{\mathbf{x}_T}$. This is equivalent to taking the solution \mathbf{x}_T from the LL optimization as *constants* in the UL problem. Such an approximation is more efficient but introduces an error term

$$\epsilon \sim \left| \frac{\partial U}{\partial \mathbf{x}^*} \frac{\partial \mathbf{x}^*}{\partial \mathbf{y}} \right|. \quad (4.18)$$

Nevertheless, it is useful when the implicit gradients contain products of *small* second-order derivatives, which depends on the specific NLS problems.

4.3 Differentiation on Manifold

Given that the state of a system within SLAM is bound to evolve on specific manifolds, optimization on manifolds plays a crucial role in solving back-end SLAM problems. We next derive the Jacobians required to differentiate with respect to variables belonging to Lie groups, which is an essential step for differentiation on the manifold.

4.3.1 Derivatives on the Lie Group

Since we introduced the basic concepts of Lie group, Lie algebra, and their basic operations (*e.g.*, exponential and logarithmic maps) in Chapter 2, this section will briefly recap those concepts but mainly focus on the definition of their derivatives, which is essential for solving a differentiable optimization problem.

Consider a Lie group's manifold \mathcal{M} , each point χ on this smooth manifold possesses a unique tangent space, denoted by $T_{\chi}\mathcal{M}$, where the fundamental principles of calculus are valid. The Lie algebra, represented as \mathfrak{m} , is a vector space that can be locally defined to the point χ as $\mathfrak{m} = T_{\chi}\mathcal{M}$. The exponential map $\exp : \mathfrak{m} \rightarrow \mathcal{M}$ projects elements from the Lie algebra to the Lie group, while the logarithmic map $\log : \mathcal{M} \rightarrow \mathfrak{m}$ serves as its inverse, establishing a bi-directional relationship:

$$\chi = \exp(\tau^{\wedge}) \Leftrightarrow \tau^{\wedge} = \log(\chi), \quad (4.19)$$

where $\hat{\cdot}$ is a linear invertible map, and $\tau^{\wedge} \in \mathfrak{m}$. By representing the coordinates within the Lie algebra as vectors τ in \mathbb{R}^n , we can define mappings between vector τ and the Lie group χ :

$$\chi = \text{Exp}(\tau) \Leftrightarrow \tau = \text{Log}(\chi), \quad (4.20)$$

where we redefined the exponential and logarithm maps to directly use a vector as input and output, respectively.

To calculate derivatives on Lie groups, it is crucial to first understand the relative change between two manifold elements, say χ_1 and χ_2 . These changes are quantified by first defining the \oplus and \ominus operators, which capture the concept of displacement on the manifold, as described in (4.21) and (4.22) below:

$$\begin{aligned} \chi_2 &= \chi_1 \oplus \tau \triangleq \chi_1 \circ \text{Exp}(\tau), \\ \tau &= \chi_2 \ominus \chi_1 \triangleq \text{Log}(\chi_1^{-1} \circ \chi_2). \end{aligned} \quad (4.21)$$

The placement of τ on the right-hand side in (4.21) signifies that it is expressed in the local frame at χ_1 . Conversely, the left operators in (4.22) reflect a global frame perspective:

$$\begin{aligned} \chi_2 &= \varepsilon \oplus \chi_1 \triangleq \text{Exp}(\varepsilon) \circ \chi_1, \\ \varepsilon &= \chi_2 \ominus \chi_1 \triangleq \text{Log}(\chi_2 \circ \chi_1^{-1}), \end{aligned} \quad (4.22)$$

where ε is expressed in the global frame. Both τ and ε can be viewed as incremental perturbations to the manifold elements. By using corresponding composition operators \oplus and \ominus , the variations are expressed as vectors in the tangent space.

With the right \oplus and \ominus operators in place, we use the Jacobian matrix \mathbf{J} to describe perturbations on manifolds. The Jacobian captures the essence of infinitesimal perturbations τ within the tangent space \mathfrak{m} :

$$\begin{aligned} \frac{\partial f(\chi)}{\partial \chi} &\triangleq \lim_{\tau \rightarrow 0} \frac{f(\chi \oplus \tau) \ominus f(\chi)}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{f(\chi \circ \text{Exp}(\tau)) \ominus f(\chi)}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{\text{Log}(f(\chi)^{-1} \circ f(\chi \circ \text{Exp}(\tau)))}{\tau}. \end{aligned} \quad (4.23)$$

Let $g(\tau) = \text{Log}(f(\chi)^{-1} \circ f(\chi \circ \text{Exp}(\tau)))$, then the right Jacobian \mathbf{J}_R can be expressed as the derivative of $g(\tau)$ at $\tau = 0$:

$$\frac{\partial f(\chi)}{\partial \chi} = \mathbf{J}_R = \left. \frac{\partial g(\tau)}{\partial \tau} \right|_{\tau=0}. \quad (4.24)$$

In this way, the derivatives of $f(\chi)$ with respect to χ in the manifold are represented by the Jacobian matrix $\mathbf{J}_R \in \mathbb{R}^{m \times n}$, where m and n are the dimensions of the Lie groups \mathcal{M} and \mathcal{N} , respectively. The right Jacobian matrix performs a linear mapping from the tangent space \mathfrak{m} to the tangent space $\mathfrak{n} = T_{f(\chi)}\mathcal{N}$.

Similarly, consider an infinitesimal perturbation $\varepsilon \in T_g\mathcal{M}$ applied to the Lie group element χ , the left Jacobian \mathbf{J}_L can be defined with the left plus and minus operators:

$$\begin{aligned} \frac{\partial f(\chi)}{\partial \chi} &\triangleq \lim_{\varepsilon \rightarrow 0} \frac{f(\varepsilon \oplus \chi) \ominus f(\chi)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{f(\text{Exp}(\varepsilon) \circ \chi) \ominus f(\chi)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\text{Log}(f(\text{Exp}(\varepsilon) \circ \chi) \circ f(\chi)^{-1})}{\varepsilon} \\ &= \left. \frac{\partial \text{Log}(f(\text{Exp}(\varepsilon) \circ \chi) \circ f(\chi)^{-1})}{\partial \varepsilon} \right|_{\varepsilon=0}. \end{aligned} \quad (4.25)$$

The resulting left Jacobian $\mathbf{J}_L \in \mathbb{R}^{n \times m}$ is also a linear mapping, but in the global tangent space from $T_g\mathcal{M}$ to $T_g\mathcal{N}$.

To delve into the local perturbations around a point χ_1 , we consider perturbations τ as $\tau = \chi \ominus \chi_1$, with χ being a perturbed version of χ_1 . The covariance matrices Σ_χ defined on the tangent space are derived using the expectation operator \mathbb{E} , enabling the representation of uncertainties and their propagation:

$$\Sigma_\chi \triangleq \mathbb{E}[\tau \tau^\top] = \mathbb{E}[(\chi \ominus \chi_1)(\chi \ominus \chi_1)^\top]. \quad (4.26)$$

These covariance matrices facilitate the establishment of Gaussian distributions on the manifold, expressed as $\chi \sim \mathcal{N}(\chi_1, \Sigma_\chi)$. It is important to note that the covariance matrices Σ_χ are defined on the tangent space $T_{\chi_1}\mathcal{M}$, which allows the uncertainty in the manifold to be represented by a vector and be propagated in the form of covariance matrices.

Example 4.5 Consider a robot equipped with an inertial measurement unit (IMU) and a camera. Given noisy observations \mathbf{R}_{IMU} and \mathbf{R}_{Cam} from both sensors, the orientation of the robot can be estimated by minimizing the discrepancy between the measurements, which can be formulated as a nonlinear least squares problem on the manifold $SO(3)$:

$$\hat{\mathbf{R}} = \arg \min_{\mathbf{R} \in SO(3)} f(\mathbf{R}, \mathbf{R}_{\text{IMU}}, \mathbf{R}_{\text{Cam}}), \quad (4.27)$$

where $f(\cdot)$ is the cost function that quantifies the differences between the estimated orientation \mathbf{R} and the sensor measurements \mathbf{R}_{IMU} and \mathbf{R}_{Cam} . With the Jacobian matrices in place, the optimization on the manifold $SO(3)$ for the pose estimation can be effectively managed. The cost function $f(\cdot)$ can be detailed as:

$$f(\mathbf{R}) = \|\text{Log}(\mathbf{R}_{\text{IMU}}^{-1}\mathbf{R})\|^2 + \|\text{Log}(\mathbf{R}_{\text{Cam}}^{-1}\mathbf{R})\|^2. \quad (4.28)$$

To minimize $f(\mathbf{R})$, we need to compute its gradient with respect to \mathbf{R} on the manifold $SO(3)$. The gradient can be derived using the right Jacobian \mathbf{J}_R as:

$$\begin{aligned} \nabla f(\mathbf{R}) = & 2 \left(\frac{\partial \text{Log}(\mathbf{R}_{\text{IMU}}^{-1}\mathbf{R})}{\partial \mathbf{R}} \right)^\top \text{Log}(\mathbf{R}_{\text{IMU}}^{-1}\mathbf{R}) \\ & + 2 \left(\frac{\partial \text{Log}(\mathbf{R}_{\text{Cam}}^{-1}\mathbf{R})}{\partial \mathbf{R}} \right)^\top \text{Log}(\mathbf{R}_{\text{Cam}}^{-1}\mathbf{R}). \end{aligned} \quad (4.29)$$

The gradient $\nabla f(\mathbf{R})$ can be used in conjunction with optimization algorithms like gradient descent which moves along the tangent space and reprojecting back to the manifold to update the pose \mathbf{R} iteratively:

$$\mathbf{R}_{k+1} = \mathbf{R}_k \text{Exp}(-\alpha \nabla f(\mathbf{R})), \quad (4.30)$$

where α is the step size. This iterative process continues until the cost function $f(\mathbf{R})$ converges to a minimum, providing an optimal pose estimate $\hat{\mathbf{R}}$ that aligns with the sensor measurements.

4.3.2 Differentiation Operations on Manifold

For typical manifold operations, we can derive closed-form expressions for the Jacobians associated with inversion, composition, and group actions. These expressions

facilitate a comprehensive approach to optimization in SLAM, by enabling the computation of function derivatives on manifolds with the chain rule:

$$\frac{\partial \mathcal{Z}}{\partial \chi} = \frac{\partial \mathcal{Z}}{\partial \mathcal{Y}} \frac{\partial \mathcal{Y}}{\partial \chi}, \quad (4.31)$$

where $\mathcal{Z} = g(\mathcal{Y})$, and $\mathcal{Y} = f(\chi)$.

Jacobians of inversion can be derived through the application of the function $f(\chi) = \chi^{-1}$ with (4.23) for the right Jacobian \mathbf{J}_R , which leads to:

$$\begin{aligned} \frac{\partial \chi^{-1}}{\partial \chi} &\triangleq \lim_{\tau \rightarrow 0} \frac{\text{Log}((\chi^{-1})^{-1}(\chi \text{Exp}(\tau))^{-1})}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{\text{Log}(\chi \text{Exp}(\tau)^{-1} \chi^{-1})}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{(\chi(-\tau)^\wedge \chi^{-1})^\vee}{\tau}. \end{aligned} \quad (4.32)$$

Jacobians of composition can be derived through the application of the function $f(\chi) = \chi \circ \chi_1$ with the Equation (4.23). The derivative of the composition operator $\chi \circ \chi_1$ with respect to χ is:

$$\begin{aligned} \frac{\partial(\chi \circ \chi_1)}{\partial \chi} &\triangleq \lim_{\tau \rightarrow 0} \frac{\text{Log}((\chi \chi_1)^{-1}(\chi \text{Exp}(\tau) \chi_1))}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{\text{Log}(\chi_1^{-1} \text{Exp}(\tau) \chi_1)}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{(\chi_1^{-1} \tau^\wedge \chi_1)^\vee}{\tau}. \end{aligned} \quad (4.33)$$

The derivative of the composition operator $\chi \circ \chi_1$ with respect to χ_1 is:

$$\begin{aligned} \frac{\partial(\chi \circ \chi_1)}{\partial \chi_1} &\triangleq \lim_{\tau \rightarrow 0} \frac{\text{Log}((\chi \chi_1)^{-1}(\chi \chi_1 \text{Exp}(\tau)))}{\tau} \\ &= \lim_{\tau \rightarrow 0} \frac{\text{Log}(\text{Exp}(\tau))}{\tau} \\ &= \mathbf{I}. \end{aligned} \quad (4.34)$$

Jacobians of the manifold \mathcal{M} are characterized by the right Jacobian of χ which is derived from the exponential map of $\tau \in \mathbb{R}^m$. This is expressed as:

$$\mathbf{J}_r(\tau) \triangleq \frac{\tau \partial \text{Exp}(\tau)}{\partial \tau}. \quad (4.35)$$

The right Jacobian conveys minor changes in τ to modifications in the local tangent space at $\text{Exp}(\tau)$. Similarly, the left Jacobian of χ maps changes of τ to variations within the global tangent space of the manifold. This is expressed as:

$$\mathbf{J}_l(\tau) \triangleq \frac{\epsilon \partial \text{Exp}(\tau)}{\partial \tau}. \quad (4.36)$$

Jacobians of group action depends on the specific group action set $v \in \mathcal{V}$. The group action is defined as:

$$\begin{aligned} \mathbf{J}_{\chi}^{\chi \cdot v} &\triangleq \frac{{}^x D\chi \cdot v}{D\chi}, \\ \mathbf{J}_v^{\chi \cdot v} &\triangleq \frac{{}^v D\chi \cdot v}{Dv}, \end{aligned} \quad (4.37)$$

where $\chi \in \mathcal{M}$ and $v \in \mathcal{V}$.

Example 4.6 Consider a robotic arm with two joints, \mathbf{R}_1 and \mathbf{R}_2 , each represented by an element in $SO(3)$. The final orientation of the robot's end-effector is determined by the composition of the joint rotations:

$$\mathbf{R} = \mathbf{R}_1 \circ \mathbf{R}_2. \quad (4.38)$$

To evaluate the impact of small perturbations $\boldsymbol{\tau}$ in \mathbf{R}_1 and \mathbf{R}_2 on the end-effector orientation \mathbf{R} . It can be quantified using the Jacobians of composition:

$$\begin{aligned} \frac{\partial(\mathbf{R}_1 \circ \mathbf{R}_2)}{\partial \mathbf{R}_1} &= \lim_{\boldsymbol{\tau} \rightarrow 0} \frac{(\mathbf{R}_2^{-1} \boldsymbol{\tau} \wedge \mathbf{R}_2)^\vee}{\boldsymbol{\tau}}, \\ \frac{\partial(\mathbf{R}_1 \circ \mathbf{R}_2)}{\partial \mathbf{R}_2} &= \mathbf{I}. \end{aligned} \quad (4.39)$$

This example implies that adjustments to the first joint \mathbf{R}_1 affect the final orientation \mathbf{R} through a transformation influenced by the current state of the second joint \mathbf{R}_2 . However, changes in the second joint \mathbf{R}_2 directly impact \mathbf{R} without being influenced by the first joint \mathbf{R}_1 .

4.4 Modern Libraries

4.4.1 Numerical Challenges of Automatic Differentiation

Automatic Differentiation (AutoDiff) is a cornerstone technique for computing derivatives accurately and efficiently in various optimization contexts, including differentiation on manifolds. Differentiation on manifolds poses unique challenges due to the complex geometrical properties inherent in manifold structures, which can affect the performance and applicability of AutoDiff. In differential optimization, these challenges become pronounced as AutoDiff interacts with the curved space of manifolds, potentially introducing numerical instability and inaccuracies.

This section delves into the specific numerical issues that arise when using automatic differentiation for manifold-based optimization tasks. Particular attention will be paid to the complexities involved in maintaining numerical stability and precision in the presence of manifold constraints, such as those found in constrained optimization and in systems defined by differential equations on manifolds. For simplicity, we will take the PyPose library [1147] as an example, which defines a

general data structure, **LieTensor** for Lie Group and Lie Algebra. Specifically, we will show its numerical challenges and how PyPose tackle this challenge.

Analytical Foundations of Exponential Mapping to Quaternions. The exponential map is a fundamental concept in the theory of Lie groups and is particularly critical when transitioning between Lie algebras and Lie groups represented by quaternions. This mapping enables the translation of angular velocities from the algebraic structure in \mathbb{R}^3 to rotational orientations in the group of unit quaternions \mathbb{S}^3 . Analytically, the exponential map for quaternions is derived from the Rodrigues' rotation formula, which relates a vector in \mathbb{R}^3 to the corresponding rotation. Given a vector \mathbf{x} in \mathbb{R}^3 , representing the axis of rotation scaled by the rotation angle, the quaternion representation of the rotation is given by:

$$\text{Exp}(\boldsymbol{\nu}) = \left[\sin\left(\frac{\|\boldsymbol{\nu}\|}{2}\right) \frac{\boldsymbol{\nu}^\top}{\|\boldsymbol{\nu}\|}, \cos\left(\frac{\|\boldsymbol{\nu}\|}{2}\right) \right]^\top \quad (4.40)$$

where $\|\boldsymbol{\nu}\|$ represents the magnitude of $\boldsymbol{\nu}$, corresponding to the angle of rotation, and $\frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|}$ is the unit vector in the direction of $\boldsymbol{\nu}$.

One of the challenges of implementing a differentiable **LieTensor** is that one often need to calculate numerically problematic terms such as $\frac{\sin \boldsymbol{\nu}}{\boldsymbol{\nu}}$ in (4.40) for the Exp and Log mapping [1087]. The direct computation of sine and cosine functions for very small angles can lead to precision issues due to the finite representation of floating-point numbers in computer systems. To manage these issues and maintain numerical stability, PyPose takes the Taylor expansion to avoid calculating the division by zero.

$$\text{Exp}(\boldsymbol{\nu}) = \begin{cases} \left[\boldsymbol{\nu}^\top \gamma_e, \cos\left(\frac{\|\boldsymbol{\nu}\|}{2}\right) \right]^\top & \text{if } \|\boldsymbol{\nu}\| > \text{eps} \\ \left[\boldsymbol{\nu}^\top \gamma_o, 1 - \frac{\|\boldsymbol{\nu}\|^2}{8} + \frac{\|\boldsymbol{\nu}\|^4}{384} \right]^\top & \text{otherwise,} \end{cases} \quad (4.41)$$

where $\gamma_e = \frac{\sin(\frac{\|\boldsymbol{\nu}\|}{2})}{\|\boldsymbol{\nu}\|}$ when $\|\boldsymbol{\nu}\|$ is significant, and $\gamma_o = \frac{1}{2} - \frac{\|\boldsymbol{\nu}\|^2}{48} + \frac{\|\boldsymbol{\nu}\|^4}{3840}$ for small $\|\boldsymbol{\nu}\|$, ensuring precise calculations across all ranges of rotation magnitudes. Here, eps is the smallest machine number where $1 + \text{eps} \neq 1$. This analytical-to-numerical progression demonstrates the importance of accurate and stable methods for computing exponential maps in applications that require high fidelity in rotation representation, such as in 3D graphics, robotics, and aerospace engineering.

LieTensor is different from the existing libraries in several aspects: **(1)** PyPose supports auto-diff for any order gradient and is compatible with most popular devices, such as CPU, GPU, TPU, and Apple silicon GPU, while other libraries like LieTorch [1087] implement customized CUDA kernels and only support 1st-order gradient. **(2)** **LieTensor** supports parallel computing of gradient with the **vmap** operator, which allows it to compute Jacobian matrices much faster. **(3)** Libraries such as LieTorch, JaxLie [1239], and Theseus only support Lie groups, while Py-

Pose supports both Lie groups and Lie algebras. As a result, one can directly call the `Exp` and `Log` maps from a `LieTensor` instance, which is more flexible and user-friendly. Moreover, the gradient with respect to both types can be automatically calculated and back-propagated. The readers may find a list of supported `LieTensor` operations in [2] and the tutorial of PyPose is available in [4]. The usages of a `LieTensor` and its automatic differentiation can be found at <https://github.com/pypose/slambook-snippets/blob/main/lietensor.ipynb>.

4.4.2 Implementation of Differentiable Optimization

To enable end-to-end learning with bilevel optimization, one need to integrate general optimizers beyond the gradient-based methods such as SGD [936] and Adam [587] required by neural methods, since many problems in SLAM such as bundle adjustment and factor graph optimization require other optimizations algorithms such as constrained or 2nd-order optimization [58]. Moreover, practical problems have outliers, hence one needs to robustify the loss as described in Chapter 3. Next we consider an Iteratively Reweighted Least Squares (IRLS) approach to SLAM as introduced in Section 3.3, and present the intuition behind the optimization-oriented interfaces of PyPose, including `solver`, `kernel`, `corrector`, and `strategy` for using the 2nd-order Levenberg-Marquardt (LM) optimizer.

Let us start by considering a weighted least square problem:

$$\min_{\mathbf{y}} \sum_i (\mathbf{h}_i(\mathbf{x}_i) - \mathbf{z}_i)^T \boldsymbol{\Sigma}_i (\mathbf{h}_i(\mathbf{x}_i) - \mathbf{z}_i), \quad (4.42)$$

where $\mathbf{h}(\cdot)$ is a regression model (`Module`), $\mathbf{x} \in \mathbb{R}^n$ is the parameters to be optimized, \mathbf{h}_i deontes prediction for the i -th input sample, $\boldsymbol{\Sigma}_i \in \mathbb{R}^{d \times d}$ is a square information matrix. The solution to (4.42) of an LM algorithm is computed by iteratively updating an estimate \mathbf{x}_t via $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + \boldsymbol{\delta}_t$, where the update step $\boldsymbol{\delta}_t$ is computed as:

$$\sum_i (\boldsymbol{\Lambda}_i + \lambda \cdot \text{diag}(\boldsymbol{\Lambda}_i)) \boldsymbol{\delta}_t = - \sum_i \mathbf{J}_i^T \boldsymbol{\Sigma}_i \mathbf{r}_i, \quad (4.43)$$

where $\mathbf{r}_i = \mathbf{h}_i(\mathbf{x}_i) - \mathbf{z}_i$ is the i -th residual error, \mathbf{J}_i is the Jacobian of \mathbf{h} computed at \mathbf{x}_{t-1} , $\boldsymbol{\Lambda}_i$ is an approximated Hessian matrix computed as $\boldsymbol{\Lambda}_i = \mathbf{J}_i^T \boldsymbol{\Sigma}_i \mathbf{J}_i$, and λ is a damping factor. To find step $\boldsymbol{\delta}_t$, one needs a linear `solver`:

$$\mathbf{A} \cdot \boldsymbol{\delta}_t = \boldsymbol{\beta}, \quad (4.44)$$

where $\mathbf{A} = \sum_i (\boldsymbol{\Lambda}_i + \lambda \cdot \text{diag}(\boldsymbol{\Lambda}_i))$, $\boldsymbol{\beta} = - \sum_i \mathbf{J}_i^T \boldsymbol{\Sigma}_i \mathbf{r}_i$. In practice, the square matrix \mathbf{A} is often positive-definite, so we could leverage standard linear solvers such as Cholesky. If the Jacobian \mathbf{J}_i is large and sparse, we may also use sparse solvers such as sparse Cholesky [167] or preconditioned conjugate gradient (PCG) [467] solver. In practice, one often introduces robust `kernel` functions $\rho : \mathbb{R} \mapsto \mathbb{R}$ into (4.42) to

reduce the effect of outliers:

$$\min_{\mathbf{y}} \sum_i \rho(\mathbf{r}_i^T \boldsymbol{\Sigma}_i \mathbf{r}_i), \quad (4.45)$$

where ρ is designed to down-weight measurements with large residuals \mathbf{r}_i . In this case, we need to adjust (4.43) to account for the presence of the robust kernel. A popular way is to use an IRLS method, Triggs' correction [1111], which is also adopted by the Ceres [22] library. However, it needs 2nd-order derivative of the kernel function ρ , which is always negative. This can lead 2nd-order optimizers including LM to be unstable [1111]. Alternatively, PyPose introduces an IRLS method, **FastTriggs**, which is faster yet more stable than **Triggs** by only involving the 1st-order derivative:

$$\mathbf{r}_i^\rho = \sqrt{\rho'(c_i)} \mathbf{r}_i, \quad \mathbf{J}_i^\rho = \sqrt{\rho'(c_i)} \mathbf{J}_i, \quad (4.46)$$

where $c_i = \mathbf{r}_i^T \boldsymbol{\Sigma}_i \mathbf{r}_i$, \mathbf{r}_i^ρ and \mathbf{J}_i^ρ are the corrected model residual and Jacobian due to the introduction of kernel functions, respectively. More details about **FastTriggs** and its proof can be found in [1], while IRLS was introduced in Section 3.3.

A simple LM optimizer may not converge to the global optimum if the initial guess is too far from the optimum. For this reason, we often need other strategies such as adaptive damping, dogleg, and trust region methods [704] to restrict each step, preventing it from stepping “too far”. To adopt those strategies, one may simply pass a **strategy** instance, e.g., **TrustRegion**, to an optimizer. In summary, PyPose supports easy extensions for the aforementioned algorithms by simply passing **optimizer** arguments to their constructor, including **solver**, **strategy**, **kernel**, and **corrector**. A list of available algorithms and examples can be found in [3]. The usages of a 2nd-order optimizer can be found at <https://github.com/pypose/slambook-snippets/blob/main/optimization.ipynb>.

4.4.3 Related Open-source Libraries

Open-source libraries related to differentiable optimization can be divided into three groups: (1) linear algebra, (2) machine learning libraries, and (3) specialized optimization libraries.

Linear Algebra Libraries are essential to machine learning and robotics research. NumPy [829], a linear algebra library for Python, offers comprehensive operations on vectors and matrices while enjoying higher running speed due to its underlying well-optimized C code. Eigen [420], a high performance C++ linear algebra library, has been used in many projects such as TensorFlow [9], Ceres [22], GTSAM [260], and g²o [410]. ArrayFire [734], a GPU acceleration library for C, C++, Fortran, and Python, contains simple APIs and provides GPU-tuned functions.

Machine Learning Libraries focus more on operations on tensors (i.e., high-dimensional matrices) and automatic differentiation. Early machine learning frameworks, such as Torch [236], OpenNN [836], and MATLAB [745], provide primitive tools for researchers to develop neural networks. However, they only support CPU computation and lack concise APIs, which plague engineers using them in applications. A few years later, deep learning frameworks such as Chainer [1105], Theano [31], and Caffe [525] arose to handle the increasing size and complexity of neural networks while supporting multi-GPU training with convenient APIs for users to build and train their neural networks. Furthermore, the recent frameworks, such as TensorFlow [9], PyTorch [856], and MXNet [190], provide a comprehensive and flexible ecosystem (e.g., APIs for multiple programming languages, distributed data parallel training, and facilitating tools for benchmark and deployment). Gvnn [437] introduced differentiable transformation layers into Torch-based framework, leading to end-to-end geometric learning. JAX [116] can automatically differentiate native Python and NumPy functions and is an extensible system for composable function transformations. In many ways, the existence of these frameworks facilitated and promoted the growth of deep learning. Recently, more efforts have been taken to combine standard optimization tools with deep learning. Recent work like Theseus [880] and CvxpyLayer [25] showed how to embed differentiable optimization within deep neural networks. PyPose [1147] incorporates 2nd-order optimizers such as Gaussian-Newton and Levenberg-Marquardt and can compute any order gradients of Lie groups and Lie algebras, which are essential to robotics.

Other Specialized Optimization Libraries have been developed and leveraged in robotics. To mention a few, Ceres [22] is an open-source C++ library for large-scale nonlinear least squares optimization problems and has been widely used in SLAM. Pyomo [441] and JuMP [300] are optimization frameworks that have been widely used due to their flexibility in supporting a diverse set of tools for constructing, solving, and analyzing optimization models. CasADi [39] has been used to solve many real-world control problems in robotics due to its fast and effective implementations of different numerical methods for optimal control. Pose and factor-graph optimization also play an important role in robotics. For example, g²o [410] and GTSAM [260] are open-source C++ frameworks for graph-based nonlinear optimization, which provide concise APIs for constructing new problems and have been leveraged to solve several optimization problems in SLAM.

Optimization libraries have also been widely used in robotic control problems. To name a few, IPOPT [1143] is an open-source C++ solver for nonlinear programming problems based on interior-point methods and is widely used in robotics and control. Similarly, OpenOCL [598] supports a large class of optimization problems such as continuous time, discrete time, constrained, unconstrained, multi-phase, and trajectory optimization problems for real-time model-predictive control. Another library for large-scale optimal control and estimation problems is CT [391], which provides standard interfaces for different optimal control solvers and can be

extended to a broad class of dynamical systems in robotic applications. Drake [1083] has solvers for common control problems and that can be directly integrated with its simulation tool boxes. Its system completeness made it favorable to researchers.

4.5 Final Considerations & Recent Trends

Deep learning methods have witnessed significant development in recent years [1281]. As data-driven approaches, they are believed to perform better on visual tracking than traditional handcrafted features. Most studies on the subject employed end-to-end structures, including both supervised methods such as DeepVO [1161] and TartanVO [1165] and unsupervised methods such as UnDeepVO [663] and Unsupervised VIO [1171]. It is generally observed that the supervised approaches achieve higher performance compared to their unsupervised counterparts since they can learn from a diverse range of ground truths such as pose, flow, and depth. Nevertheless, obtaining such ground truths in the real world is a labor-consuming process [1164].

Recently, hybrid methods have received increasing attention as they integrate the strengths of both geometry-based and deep-learning approaches. Several studies have explored the potential of integrating Bundle Adjustment (BA) with deep learning methods to impose topological consistency between frames, such as attaching a BA layer to a learning network such as BA-Net [1072] and DROID-SLAM [1086]. Additionally, some works focused on compressing image features into codes (embedded features) and optimizing the pose-code graph during inference such as DeepFactors [244]. Furthermore, DiffPoseNet [848] is proposed to predict poses and normal flows using networks and fine-tune the coarse predictions through a Cheirality layer. However, in these works, the learning-based methods and geometry-based optimization are decoupled and separately used in different sub-modules. The lack of integration between the front-end and back-end may result in sub-optimal performance. Besides, they only back-propagate the pose error “through” bundle adjustment, thus the supervision is from the ground truth poses. In this case, BA is just a special layer of the network. Recently, iSLAM [354] connects the front-end and back-end bidirectionally and enforces the learning model to learn from geometric optimization through a bilevel optimization framework, which achieves performance improvement without external supervision. Some other tasks can also be formulated as bilevel optimization, *e.g.*, reinforcement learning [1033], local planning [1212], global planning [193], feature matching [1262], and multi-robot routing [426].