

An Implementation of Visual-inertial Simultaneous localization and mapping (SLAM) Based on Extended Kalman Filter and Stereo Camera

Wenran, Tian

In this paper, our main object is to build a visual-inertial simultaneous localization and mapping (SLAM) on a moving robot, and based on Extended Kalman Filter (EKF). There is a IMU on the robot to obtain the movement information, and a stereo camera to extract the feature and then obtain the localization information. And after implementing EKF on visual-inertial SLAM, we get a reasonable result.

Index Terms—EKF, SLAM, Stereo Camera

I. INTRODUCTION

IN nowadays, self-driving car is a dramatically hot field. With more and more sophisticated algorithms were implemented, we can make the self-driving car much more robust than last century. Sensing and detection of the environment is an essential link in self-driving car control because with an accurate location, we can give the control parameters more specifically and subtly. Our main object is to build a visual-inertial simultaneous localization and mapping (SLAM) on a moving robot, and based on Extended Kalman Filter (EKF). There is a IMU on the robot to obtain the movement information, and a stereo camera to extract the feature and then obtain the localization information.

II. PROBLEM FORMULATION

A. Kalman Filter

The Kalman Filter is a kind of special case of Bayes filter. Specifically, there are several assumptions in Kalman Filter:

1. The prior pdf $p_{t|t}$ is Gaussian
2. The motion model is linear in the state \mathbf{x}_t with Gaussian noise \mathbf{w}_t
3. The observation model is linear in the state \mathbf{x}_t with Gaussian noise \mathbf{v}_t
4. The motion noise \mathbf{w}_t and observation noise \mathbf{v}_t are independent of each other, of the state \mathbf{x}_t , and across time.

However, with those assumptions, it's not easy to give an analytical solution of a Kalman filter model. Because if do that, we need to solve integrals of some equations multiple with Gaussian distribution, which can be quite complicated. On

another hand, motion model and observation model are not often linear models. For these reasons, we will choose the Extended Kalman Filter in this paper.

Extended Kalman Filter has the following formats:

$$\mathbf{x}_t | \mathbf{z}_{0:t}, \mathbf{u}_{0:t-1} \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$$

Motion model:

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t), \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, W)$$

$$F_t = \frac{df(\boldsymbol{\mu}_{t|t}, \mathbf{u}_t, \mathbf{0})}{d\mathbf{x}}, \quad Q_t = \frac{df(\boldsymbol{\mu}_{t|t}, \mathbf{u}_t, \mathbf{0})}{d\mathbf{w}}$$

Prediction:

$$\boldsymbol{\mu}_{t+1|t} = f(\boldsymbol{\mu}_{t|t}, \mathbf{u}_t, \mathbf{0})$$

$$\boldsymbol{\Sigma}_{t+1|t} = F_t \boldsymbol{\Sigma}_{t|t} F_t^T + Q_t W Q_t^T$$

Obs. model:

$$\mathbf{z}_t = h(\mathbf{x}_t, \mathbf{v}_t) \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, V)$$

$$H_{t+1} = \frac{dh(\boldsymbol{\mu}_{t+1|t}, \mathbf{0})}{d\mathbf{x}}, \quad R_{t+1} = \frac{dh(\boldsymbol{\mu}_{t+1|t}, \mathbf{0})}{d\mathbf{v}}$$

Kalman Gain:

$$K_{t+1|t} = \boldsymbol{\Sigma}_{t+1|t} H_{t+1}^T (H_{t+1} \boldsymbol{\Sigma}_{t+1|t} H_{t+1}^T + R_{t+1} V R_{t+1}^T)^{-1}$$

Update:

$$\boldsymbol{\mu}_{t+1|t+1} = \boldsymbol{\mu}_{t+1|t} + K_{t+1|t} (z_{t+1} - h(\boldsymbol{\mu}_{t+1|t}, \mathbf{0}))$$

$$\boldsymbol{\Sigma}_{t+1|t+1} = (I - K_{t+1|t} H_{t+1}) \boldsymbol{\Sigma}_{t+1|t}$$

B. Target Problems

IMU-based Localization via EKF Prediction: Implement the EKF prediction step, based on the SE(3) kinematics and the linear and angular velocity measurements to estimate the pose $\mathbf{T}_t \in \text{SE}(3)$ of the IMU over time t .

Landmark Mapping via EKF Update: assume that the predicted IMU trajectory from above is correct and focus on estimating the landmark positions. In detail, we will implement an EKF with the unknown landmark positions $\mathbf{m} \in R^{3 \times M}$ as a state and perform EKF update steps after every visual observation \mathbf{z}_t in order to keep track of the mean and covariance of \mathbf{m} . Besides, we are assuming that the landmarks are static so it is not necessary to implement a prediction step. Moreover, since the sensor does not move sufficiently along the

z-axis, the estimation of the z coordinate of the landmarks will not be very good and we will focus only on estimating their xy coordinates.

III. TECHNICAL APPROACH:

A. IMU and Trejectory

The IMU will provide the instantaneous velocity $\mathbf{v}_t \in R^3$ and angular velocity $\boldsymbol{\omega}_t \in R^3$ in body frame. Then, we use $\boldsymbol{\zeta}_t$ to denote the generalized velocity $\boldsymbol{\zeta}_t = [\mathbf{v}_t \ \boldsymbol{\omega}_t]^T$. We can calculate the instantaneous transformation to implement the prediction step in Kalman filter. Specifically, we have:

$$T_{t+1} = T_t \exp(\tau \hat{\boldsymbol{\zeta}}_t)$$

$$\hat{\boldsymbol{\zeta}}_t = \begin{bmatrix} \hat{\boldsymbol{\omega}}_t & \mathbf{v}_t \\ \mathbf{0} & 0 \end{bmatrix}$$

In the EKF prediction step, in time t we have:

$$T_t = \boldsymbol{\mu}_{t|t} \exp(\delta \hat{\boldsymbol{\mu}}_{t|t}), \quad \delta \boldsymbol{\mu}_{t|t} \sim \mathcal{N}(\mathbf{0}, \Sigma_{t|t})$$

$$\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t} \exp(\tau \hat{\boldsymbol{u}}_t)$$

$$\delta \boldsymbol{\mu}_{t+1|t} = \exp(-\tau \hat{\boldsymbol{u}}_t) \delta \boldsymbol{\mu}_{t|t} + \mathbf{w}_t$$

$$\hat{\boldsymbol{u}}_t = \begin{bmatrix} \hat{\boldsymbol{\omega}}_t & \hat{\mathbf{v}}_t \\ \mathbf{0} & \hat{\boldsymbol{\omega}}_t \end{bmatrix} \in R^{6 \times 6}$$

After implement this prediction step on the given data and tried the dead-reckoning method, we can have the result in Figure 1.

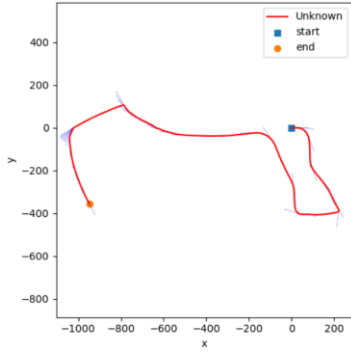


Fig. 1.

B. Visual Mapping

Firstly, we use \mathbf{m}_j to denote the j-th feature in world frame, for the relationship between observation $\mathbf{z}_{t,i}$ we have:

$$\mathbf{z}_{t,i} = h(T_t, \mathbf{m}_j) = K\pi(T_{OI}T_t^{-1}\mathbf{m}_j) + \mathbf{v}_{t,i}$$

In the visual mapping step, because we assume the landmarks will not move during sensor collecting data. So, we will ignore the prediction step, and implement our EKF visual mapping algorithm in the following several steps:

- Initialize the prior parameters $\boldsymbol{\mu}_0 \in R^{3M}$, $\Sigma_0 \in R^{3M \times 3M}$ as all zeros. We have the calibration matrix K, extrinsic $T_{OI} \in SE(3)$.
- Predicted observations based on $\boldsymbol{\mu}_t$ and known correspondence Δ_{t+1} :

$$\tilde{\mathbf{z}}_{t+1,i} = K\pi(T_{OI}T_t^{-1}\boldsymbol{\mu}_{t,j}) \in R^4, \text{ for } i = 1, 2, \dots, N_{t+1}$$

- Jacobian of $\tilde{\mathbf{z}}_{t+1,i}$ with respect to \mathbf{m}_j evaluated at $\boldsymbol{\mu}_{t,j}$:

$$H_{t+1,i,j} = \begin{cases} K \frac{d\pi}{dq}(T_{OI}T_t^{-1}\boldsymbol{\mu}_{t,j}) T_{OI}T_t^{-1}P^T, & \text{if } \Delta_t(j) = i \\ \mathbf{0}, & \in R^{4 \times 3}, \text{ otherwise} \end{cases}$$

- EKF update:

$$K_{t+1} = \Sigma_t H_{t+1}^T (H_{t+1} \Sigma_t H_{t+1}^T + I \otimes V)^{-1}$$

$$\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + K_{t+1|t}(\mathbf{z}_{t+1} - \tilde{\mathbf{z}}_{t+1})$$

$$\Sigma_{t+1} = (I - K_{t+1} H_{t+1}) \Sigma_t$$

It needs to be emphasized that during the computation, we may not able to update the entire Σ_t matrix each time. If we have 10000 landmarks, the size of Σ_t will be $(3 \times 10000) \times (3 \times 10000) = 9 \times 10^8$. If we use float32 to store Σ_t , the size of cache to store Σ_t will be about 3.35 GiB, which will be quite computational costly if we take it as operand during the iteration. On the other hand, in EKF update step, we even need to compute an inverse of a matrix, which we require even more computational resource. So, to simplify the computation, in our code we will just find those landmarks will be affected in an iteration, and crop the corresponding variance in Σ_t to compute. In addition, the H_t and K_t will also be implemented within a certain range in an iteration.

For the initial state, there is no value of $\boldsymbol{\mu}_0 \in R^{3M}$ because we did not do any update step. Our strategy is, once we get a point from observation that is not initialized in $\boldsymbol{\mu}$, we will transform the point in stereo camera coordinate to world frame, then use the value in world frame to initialize it. So this means we will initial the location of each landmark separately.

And this time, we will add some noise to the IMU data and solve the trajectory of the robot. After we implement the aforementioned steps, we can get the point cloud image in Figure 2:

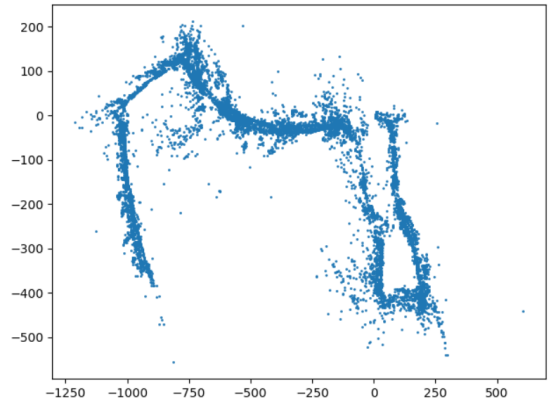


Fig. 2.

C. Visual-Initial SLAM

Here we will combine the methods in two previous section and get a complete visual-initial SLAM. In each time of iteration, the execution steps will be:

- Predict the localization
- Update/Initialize the mapping
- Update the localization

The update processes for landmarks are the same as we described in visual mapping section, so following EKF processes are all about the localization.

- Prediction step:

$$\boldsymbol{\mu} \in SE(3), \quad \Sigma \in R^{6 \times 6}$$

$$\boldsymbol{\mu}_{t+1|t} = \boldsymbol{\mu}_{t|t} \exp(\tau \hat{\boldsymbol{u}}_t)$$

$$\delta \boldsymbol{\mu}_{t+1|t} = \exp(-\tau \hat{\boldsymbol{u}}_t) \delta \boldsymbol{\mu}_{t|t} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, W)$$

$$\Sigma_{t+1|t} = E(\delta \boldsymbol{\mu}_{t+1|t} \delta \boldsymbol{\mu}_{t+1|t}^T)$$

$$= \exp(-\tau \hat{\boldsymbol{u}}_t) \Sigma_{t|t} \exp(-\tau \hat{\boldsymbol{u}}_t)^T + W$$

- Update step:

$$\tilde{\mathbf{z}}_{t+1,i} = K\pi(T_{OI}\mu_{t+1|t}^{-1}\mathbf{m}_j) \in R^4, \text{ for } i = 1, 2, \dots, N_{t+1}$$

$$H_{t+1,i} = -K \frac{d\pi}{dq}(T_{OI}\mu_{t+1|t}^{-1}\mathbf{m}_j)T_{OI}(\mu_{t+1|t}^{-1}\mathbf{m}_j)^{\odot} \in R^{4 \times 6}$$

For the " \odot " operation:

$$\mathbf{p} \in R^3, \begin{bmatrix} \mathbf{p} \\ 1 \end{bmatrix}^{\odot} = \begin{bmatrix} I & -\hat{\mathbf{p}} \\ \mathbf{0} & 0 \end{bmatrix} \in R^{4 \times 6}$$

Then we stack all the $\tilde{\mathbf{z}}_{t+1,i}$ to \mathbf{z}_{t+1} , all the $H_{t+1,i}$ to \mathbf{H} . We can have:

$$\mathbf{z}_{t+1} \in R^{4N_t}, \quad \mathbf{H}_{t+1} \in R^{4N_t \times 6}$$

$$K_{t+1} = \Sigma_{t+1|t} H_{t+1}^T (H_{t+1} \Sigma_{t+1} H_{t+1}^T + I \otimes V)^{-1} \in R^{6 \times 4N_t}$$

$$\boldsymbol{\zeta}_t = K_{t+1}(\mathbf{z}_{t+1} - \tilde{\mathbf{z}}_{t+1}) \in R^6$$

$$\boldsymbol{\mu}_{t+1|t+1} = \boldsymbol{\mu}_{t+1|t} \exp(\boldsymbol{\zeta}_t)$$

$$\Sigma_{t+1|t+1} = (I - K_{t+1} H_{t+1}) \Sigma_{t+1|t}$$

And again, in our code, we will still use the same tactics as we did in Visual Mapping step. In each iteration, we will compute the part of Σ_t that will be influenced, instead of the entire Σ_t . In addition, we will treat the sigma of position Σ_t^P and sigma of landmarks Σ_t^M separately instead of using a combined Sigma matrix $\Sigma_t \in R^{(3M+6) \times (3M+6)}$.

IV. RESULTS

After we execute the Visual-Initial SLAM, we can have the point cloud of SLAM and trajectory that in Figure 3:

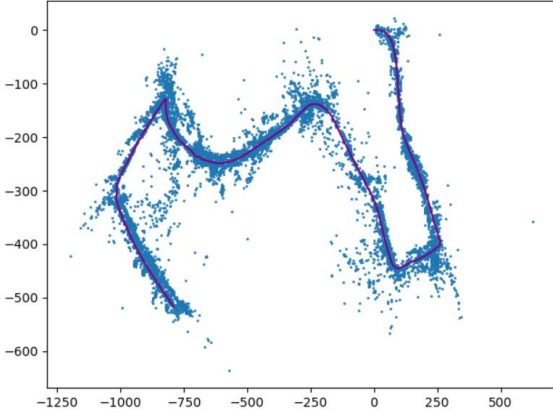


Fig. 3. The purple line stands for trajectory, the blue points stands for features in stereo camera.

Although we get a reasonable result, we find that the program takes about 50mins to run. After analysis, we hold the opinion that is because it contains two steps of finding the inverse matrix in each iteration:

$$K_{t+1} = \Sigma_{t+1|t} H_{t+1}^T (H_{t+1} \Sigma_{t+1} H_{t+1}^T + I \otimes V)^{-1}$$

As we know, it is quite computational costly to inverse a large matrix. In the code, even if we crop the matrix to compute in each iteration, we still need to inverse a matrix with size about $(50 \times 4, 50 \times 4)$. If the computational complexity of `np.linalg.inv` is $O(n^3)$, it will require $O(200^3)$ times of iteration to run. We think it is worth to discussion furtherly about how to replace the inverse item in the computation of K_{t+1} to speed up the algorithm in the future.