
Notes on Visual SLAM 14 lectures – Chapter 9

1 State estimation from probabilistic perspective

SLAM process: motion and observation equations:

$$\begin{cases} x_k = f(x_{k-1}, u_k) + w_k \\ z_{k,j} = h(y_j, x_k) + v_{k,j}, \quad k = 1, \dots, N, j = 1, \dots, M \end{cases} \quad (1)$$

A. in the observation equation, only when x_k sees y_j will generate a real observation equation, but due to a large number of visual SLAM feature points, the number of observation equations will be much larger

B. we may not have a device to measure the motion

- (a) assume there is no motion
- (b) assume the camera does not move
- (c) assume the camera is moving at a constant speed

in the absence of motion equations, this is similar to SfM (structure from motion) denote x_k as all the unknowns at the moment of k (camera pose and all landmarks):

$$x_k \triangleq \{x_k, y_1, \dots, y_m\} \quad (2)$$

and we would have the re-written equations:

$$\begin{cases} x_k = f(x_{k-1}, u_k) + w_k \\ z_{k,j} = h(x_k) + v_{k,j}, \quad k = 1, \dots, N \end{cases} \quad (3)$$

we hope to use the data from 0 to k to estimate the current state distribution:

$$p(x_k | x_0, u_{1:k}, z_{1:k}) \quad (4)$$

according to Bayes' rule, and the fact that $P(A|BC) = \frac{P(A|C)P(B|C)}{P(A)}$, denote $A = x_k, B = z_k, C = x_0, u_{1:k}, z_{1:k-1}$:

$$\begin{aligned} P(x_k | x_0, u_{1:k}, z_{1:k}) &= P(x_k | x_0, u_{1:k}, z_k, z_{1:k-1}) \\ &= \frac{P(z_k | x_k, x_0, u_{1:k}, z_{1:k-1}) P(x_k, x_0, u_{1:k}, z_{1:k-1})}{P(x_0, u_{1:k}, z_k, z_{1:k-1})} \\ &= \frac{P(z_k | x_k) P(x_k | x_0, u_{1:k}, z_{1:k-1}) P(x_0, u_{1:k}, z_{1:k-1})}{P(x_0, u_{1:k}, z_k, z_{1:k-1})} \\ &= \frac{P(z_k | x_k) P(x_k | x_0, u_{1:k}, z_{1:k-1})}{\frac{P(x_0, u_{1:k}, z_k, z_{1:k-1})}{P(x_0, u_{1:k}, z_{1:k-1})}} \\ &= \frac{P(z_k | x_k) P(x_k | x_0, u_{1:k}, z_{1:k-1})}{P(z_k | x_0, u_{1:k}, z_{1:k-1})} \\ &= \frac{P(z_k | x_k) P(x_k | x_0, u_{1:k}, z_{1:k-1})}{P(z_k)} \quad \text{since } z_k \text{ only depends on } x_k \\ &\propto P(z_k | x_k, x_0, u_{1:k}, z_{1:k-1}) P(x_k | x_0, u_{1:k}, z_{1:k-1}) \end{aligned}$$

the first term is likelihood, and the second term is prior. To estimate prior, we can expand it:

$$P(x_k|x_0, u_{1:k}, z_{1:k-1}) = \int P(x_k|x_{k-1}, x_0, u_{1:k}, z_{1:k-1})P(x_{k-1}|x_0, u_{1:k}, z_{1:k-1})dx_{k-1} \quad (5)$$

1.1 Linear systems and Kalman filter

The first part of equation 5 can be simplified as:

$$P(x_k|x_{k-1}, x_0, u_{1:k}, z_{1:k-1}) = P(x_k|x_{k-1}, u_k) \quad (6)$$

where we omit the states earlier than $k-1$ since they are not related to the k -th state. The second part can be simplified as:

$$P(x_{k-1}|x_0, u_{1:k}, z_{1:k-1}) = P(x_{k-1}|x_0, u_{1:k-1}, z_{1:k-1}) \quad (7)$$

Linear Gaussian system:

$$\begin{cases} x_k = A_k x_{k-1} + u_k + w_k \\ z_k = C_k x_k + v_k, \quad k = 1, \dots, N \end{cases} \quad (8)$$

and we assume that the states and noises are all Gaussian, so $w_k \sim N(0, R)$ and $v_k \sim N(0, Q)$ Using Markov property, suppose we know the posterior state estimation at time $k-1$ \hat{x}_{k-1} and its covariance \hat{P}_{k-1} , now we want to estimate the posterior distribution of x_k based on input and observation data at time k . We use \tilde{x}_k to denote the prior distribution and \hat{x}_k to denote its posterior distribution. Prior of x_k through the equation of motion:

$$P(x_k|x_0, u_{1:k}, z_{1:k-1}) = N(A_k \hat{x}_{k-1} + u_k, A_k \hat{P}_{k-1} A_k^T + R) \quad (9)$$

The second part can be quickly proved:

$$\begin{aligned} \Sigma_k &= E((x_k - E(x_k))(x_k - E(x_k))^T) \\ &= E((A_k x_{k-1} + u_k - A_k E(x_{k-1}) - u_k)(A_k x_{k-1} + u_k - A_k E(x_{k-1}) - u_k)^T) \\ &= E((A_k(x_{k-1} - E(x_{k-1}))) (A_k(x_{k-1} - E(x_{k-1}))))^T) \\ &= E(A_k(x_{k-1} - E(x_{k-1}))(x_{k-1} - E(x_{k-1}))^T A_k^T) \\ &= A_k E((x_{k-1} - E(x_{k-1}))(x_{k-1} - E(x_{k-1}))^T) A_k^T \\ &= A_k \Sigma_{k-1} A_k^T \end{aligned}$$

this step is called prediction:

$$\tilde{x}_k = A_k \tilde{x}_{k-1} + u_k, \quad \tilde{P}_k = A_k \hat{P}_{k-1} A_k^T + R \quad (10)$$

from observation equation, we can calculate what kind of observation data should be generated in a certain state:

$$P(z_k|x_k) = N(C_k x_k, Q) \quad (11)$$

we want to get the posterior $P(x_k|z_k)$, we have the result of $x_k \sim N(\hat{x}_k, \hat{P}_k)$:

$$N(\hat{x}_k, \hat{P}_k) = \eta N(C_k x_k, Q) \cdot N(x_k, P_k) \quad (12)$$

in which η is the normalization factor that makes the integral of the distribution equal to one. expand the exponential part as:

$$(x_k - \hat{x}_k)^T \hat{P}_k^{-1} (x_k - \hat{x}_k) = (z_k - C_k \hat{x}_k)^T Q^{-1} (z_k - C_k \hat{x}_k) + (x_k - \tilde{x}_k)^T \tilde{P}_k^{-1} (x_k - \tilde{x}_k) \quad (13)$$

In order to compute $\hat{x}_k \hat{P}_k$ on the left side, we expand the quadratics and compare their first-order and second-order coefficients of x_k , we have:

$$\hat{P}_k^{-1} = C_k^T Q^{-1} C_k + \tilde{P}_k^{-1} \quad (14)$$

which gives the relationship of the covariance matrix. Define an intermediate variable for convenience in the following derivation:

$$K = \hat{P}_k C_k^T Q^{-1} \quad (15)$$

multiply \hat{P}_k on both sides of equation 14:

$$I = \hat{P}_k C_k^T C_k + \hat{P}_k P^{-1} = K C_k + \hat{P}_k P^{-1} \quad (16)$$

then we have:

$$\hat{P}_k = (I - K C_k) \check{P}_k \quad (17)$$

compared the first-order coefficients in equation 13:

$$-2\hat{x}_k^T \hat{P}_k^{-1} x_k = -2z_k^T Q^{-1} x_k - 2\check{x}_k^T \check{P}_k^{-1} x_k \quad (18)$$

take the coefficients and transpose them:

$$\hat{P}_k^{-1} \hat{x}_k = C_k^T Q^{-1} z_k + \check{P}_k^{-1} \check{x}_k \quad (19)$$

multiply \hat{P}_k on both sides:

$$\begin{aligned} \hat{x}_k &= \hat{P}_k C_k^T Q^{-1} z_k + \hat{P}_k \check{P}_k^{-1} \check{x}_k \\ &= K z_k + (I - K C_k) \check{x}_k \\ &= \check{x}_k + K(z_k - C_k \check{x}_k) \end{aligned}$$

the general steps are:

A. Predict:

$$\check{x}_k = A_k \hat{x}_{k-1} + u_k, \quad \check{P}_k = A_k \hat{P}_{k-1} A_k^T + R \quad (20)$$

B. Update: Calculate K which is the Kalman gain. Following the definition of K , we have:

$$\begin{aligned} K Q_k &= \hat{P}_k C_k^T \\ &= (\check{P}_k - K C_k \check{P}_k) C_k^T \\ K(Q_k + C_k \check{P}_k C_k^T) &= \check{P}_k C_k^T \\ K &= \check{P}_k C_k^T (C_k \check{P}_k C_k^T + Q_k)^{-1} \end{aligned}$$

and calculate the posterior:

$$\hat{P}_k = (I - K C_k) \check{P}_k$$

1.2 Nonlinear systems and EKF

First-order Taylor expansion of the motion equation and the observation equation near a working point. Let the mean and covariance matrix at time $k-1$ be \hat{x}_{k-1} and \hat{P}_{k-1} . At the moment k , we do the linearization:

$$x_k \approx f(\hat{x}_{k-1}, u_k) + \left. \frac{\partial f}{\partial x_{k-1}} \right|_{\hat{x}_{k-1}} (x_{k-1} - \hat{x}_{k-1}) + w_k \quad (21)$$

in which $F = \left. \frac{\partial f}{\partial x_{k-1}} \right|_{\hat{x}_{k-1}}$ and for the observation model:

$$z_k \approx h(x_k) + \left. \frac{\partial h}{\partial x_k} \right|_{\hat{x}_{k-1}} (x_k - \hat{x}_k) + n_k \quad (22)$$

in which $H = \frac{\partial h}{\partial x_k} \bigg|_{\hat{x}_{k-1}}$ Then the prediction part becomes:

$$P(x_k|x_0, u_{1:k}, z_{1:k-1}) = N(f(\hat{x}_{k-1}, u_k), F\hat{P}_{k-1}F^T + R_k) \quad (23)$$

in which the prior mean and covariance are:

$$x_k = f(\hat{x}_{k-1}, u_k), \quad P_k = F\hat{P}_{k-1}F^T + R_k \quad (24)$$

then for the observation part we have:

$$P(z_k|x_k) = N(h(x_k) + H(x_k - \hat{x}_k), Q_k) \quad (25)$$

Define a Kalman gain K_k :

$$K_k = \check{P}_k H^T (H \check{P}_k H^T + Q_k)^{-1} \quad (26)$$

and the posterior can be written as:

$$\hat{x}_k = x_k + K_k(z_k - h(x_k)), \quad \hat{P}_k = (I - K_k H) \check{P}_k \quad (27)$$

In SLAM, it gives the maximum a posteriori estimate (MAP) under a single linearization step.

2 Bundle Adjustment and Graph Optimization

2.1 Projection model and cost function

A. Transform the world coordinates of point p into the camera frame using extrinsics:

$$P' = Rp + t = [X', Y', Z'] \quad (28)$$

B. Then project P' into the normalized plane and get the normalized coordinates:

$$P_c = [u_c, v_c, 1]^T = [X'/Z', Y'/Z', 1]^T \quad (29)$$

C. Apply distortion model (only radial distortion here):

$$\begin{cases} u'_c = u_c(1 + k_1 r_c^2 + k_2 r_c^4) \\ v'_c = v_c(1 + k_1 r_c^2 + k_2 r_c^4) \end{cases} \quad (30)$$

D. compute the pixel coordinates using intrinsics:

$$\begin{cases} u_s = f_x u'_c + c_x \\ v_s = f_y v'_c + c_y \end{cases} \quad (31)$$

we denote this entire process as the observation equation:

$$z = h(x, y) \quad (32)$$

and the observation data is the pixel coordinate $z \triangleq [u_s, v_s]^T$, then the error of this observation becomes:

$$e = z - h(T, p) \quad (33)$$

denote z_{ij} as the data generated by observing landmark P_j at the pose T_i , then the overall cost function:

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \|e_{ij}\|^2 = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \|z_{ij} - h(T_i, p_j)\|^2 \quad (34)$$

which is equivalent to adjusting the pose and road signs at the same time, which is the so-called BA.

2.2 Solve bundle adjustment

Optimize all variables together:

$$x = [T_1, \dots, T_m, p_1, \dots, p_n]^T \quad (35)$$

when we give an increment to the optimization variable, the objective function becomes:

$$\frac{1}{2} \|f(x + \Delta x)\|^2 \approx \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \|e_{ij} + F_{ij} \Delta \xi_i + E_{ij} \Delta p_j\|^2 \quad (36)$$

in which F_{ij} is the partial derivative of the entire cost function to the i-th pose, and E_{ij} is the partial derivative of function to the j-th landmark. Put the camera pose variable together:

$$x_c = [\xi_1, \xi_2, \dots, \xi_m]^T \in \mathbb{R}^{6m} \quad (37)$$

and also the landmarks together:

$$x_p = [p_1, p_2, \dots, p_n]^T \in \mathbb{R}^{3n} \quad (38)$$

then equation 36 becomes:

$$\frac{1}{2} \|f(x + \Delta x)\|^2 = \frac{1}{2} \|e + F \Delta x_c + E \Delta x_p\|^2 \quad (39)$$

we will face the incremental equation:

$$H \Delta x = g \quad (40)$$

H is either $J^T J$ or $J^T J + \lambda I$ depending on whether it's Gauss Newton method or Levenberg-Marquardt method. The Jacobian matrix can be divided into two parts:

$$J = [F \quad E] \quad (41)$$

and we have

$$H = J^T J = \begin{bmatrix} F^T F & F^T E \\ E^T F & E^T E \end{bmatrix} \quad (42)$$

the inversion of H has $O(n^3)$ complexity

2.3 Sparsity

Consider one of the error term e_{ij} , note that this error term only describes the residual about p_j in T_i , and only involves the i-th camera pose and the j-th landmark. The derivatives of all the remaining variables are 0. The Jacobian matrix corresponding to the error term has the following form:

$$J_{ij}(x) = (0_{2 \times 6}, \dots, 0_{2 \times 6}, \frac{\partial e_{ij}}{\partial T_i}, 0_{2 \times 6}, \dots, 0_{2 \times 3}, \dots, 0_{2 \times 3}, \frac{\partial e_{ij}}{\partial p_j}, 0_{2 \times 3}, \dots, 0_{2 \times 3}) \quad (43)$$

which causes the sparsity. In the above image, the non-zero blocks are at (i,i), (i,j), (j,i), (j,j):

The diagram shows the Jacobian matrix J and the Hessian matrix H as block matrices. The Jacobian matrix J is a row of blocks: $\square \square \square \square \square \square \square$. The third and sixth blocks are highlighted in blue and labeled i and j respectively with arrows. The Hessian matrix H is a square matrix with four quadrants of non-zero blocks, labeled i and j . The top-left quadrant is i , the top-right is j , the bottom-left is j , and the bottom-right is i . The non-zero blocks are highlighted in blue.

$$H = \sum_{i,j} J_{ij}^T J_{ij} \quad (44)$$

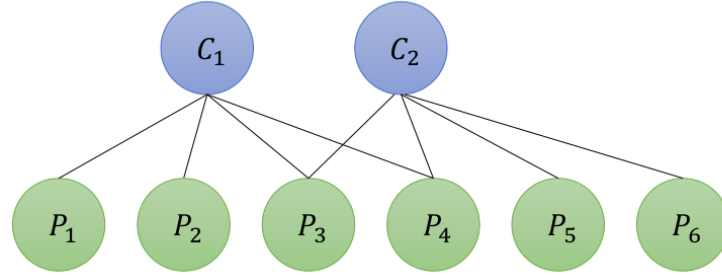
we can divide H into blocks:

$$H = \begin{bmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{bmatrix} \quad (45)$$

in which H_{11} is only related to camera pose and H_{22} is only related to landmarks. when we iterate over the i, j index, the following properties hold:

- A. No matter how i, j changes, H_{11} is always a block-diagonal matrix, with only non-zero blocks at $H_{i,i}$
- B. same for H_{22}
- C. H_{12} or H_{21} maybe sparse or dense depending on the specific observation data

suppose there are two camera poses (C_1, C_2) and 6 landmarks ($P_1, P_2, P_3, P_4, P_5, P_6$) in the scene. The variables corresponding to these cameras and point clouds are T_i and p_j . Suppose camera C_1 observes landmarks P_1, P_2, P_3, P_4 and camera C_2 observes landmarks P_3, P_4, P_5, P_6 . the cost function the becomes:



$$\frac{1}{2}(\|e_{11}\|^2 + \|e_{12}\|^2 + \|e_{13}\|^2 + \|e_{14}\|^2 + \|e_{23}\|^2 + \|e_{24}\|^2 + \|e_{25}\|^2 + \|e_{26}\|^2) \quad (46)$$

Let J_{11} be the Jacobian matrix corresponding to e_{11} , and it is not difficult to see that the partial derivatives of ξ_2 and landmarks:

$$J_{11} = \frac{\partial e_{11}}{\partial x} = \left(\frac{\partial e_{11}}{\partial \xi_1}, 0_{2 \times 6}, \frac{\partial e_{11}}{\partial p_1}, 0_{2 \times 3}, 0_{2 \times 3}, 0_{2 \times 3}, 0_{2 \times 3}, 0_{2 \times 3} \right) \quad (47)$$

consider there are m cameras and n landmarks, since there are far more landmarks than cameras, we have

$$J = \begin{bmatrix} J_{11} \\ J_{12} \\ J_{13} \\ J_{14} \\ J_{23} \\ J_{24} \\ J_{25} \\ J_{26} \end{bmatrix} = \begin{bmatrix} \begin{matrix} C_1 & C_2 & P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \end{matrix} \\ \begin{matrix} \text{[blocks]} \end{matrix} \end{bmatrix} \quad H = J^T J = \begin{bmatrix} \text{[blocks]} \end{bmatrix}$$

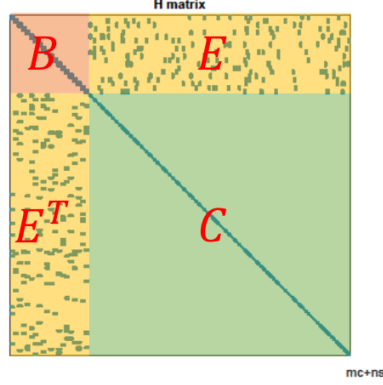
$n \gg m$. The actual H matrix will be arrow-like.

2.4 Schur trick

The linear equation $Hx = g$ can be rewritten as:

$$\begin{bmatrix} B & E \\ E^T & C \end{bmatrix} \begin{bmatrix} \Delta x_c \\ \Delta x_p \end{bmatrix} = \begin{bmatrix} v \\ w \end{bmatrix} \quad (48)$$

in which B is a block-diagonal matrix, the dimension of each diagonal block is the same as the dimension of the camera pose. The number of diagonal blocks is the number of camera variables. C is often much larger than B , with each block being a 3×3 matrix. Perform Gaussian elimination on the linear equation, we multiply a coefficient matrix on the left side:



$$\begin{bmatrix} I & -EC^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} B & E \\ E^T & C \end{bmatrix} \begin{bmatrix} \Delta x_c \\ \Delta x_p \end{bmatrix} = \begin{bmatrix} I & -EC^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} \quad (49)$$

rearrange it:

$$\begin{bmatrix} B - EC^{-1}E^T & 0 \\ E^T & C \end{bmatrix} \begin{bmatrix} \Delta x_c \\ \Delta x_p \end{bmatrix} = \begin{bmatrix} v - EC^{-1}w \\ w \end{bmatrix} \quad (50)$$

take the first row of the equation and get the incremental equation for the pose part:

$$[B - EC^{-1}E^T] \Delta x_c = v - EC^{-1}w \quad (51)$$

thus we can solve Δx_c first, then plug into the original equation and solve Δx_p later ($\Delta x_p = C^{-1}(w - E^T \Delta x_c)$). This is called Schur elimination. Another way of marginalization is by Cholesky decomposition. Eliminating camera variables is also commonly used in SLAM.

2.5 Robust kernels

Anomaly data might introduce a large gradient, thus eliminating the influence of other correct edges.

Kernel functions: makes sure error of each edge will not be big enough to cover other edges, replace \mathcal{L}_2 norm. For example, the Huber kernel:

$$H(e) = \begin{cases} \frac{1}{2}e^2 & \text{when } |e| \leq \delta \\ \delta(|e| - \frac{1}{2}\delta) & \text{otherwise} \end{cases} \quad (52)$$