

OVERVIEW

We propose a novel framework for unknown object pose estimation in a Sliding Window Filter(SWF) manner with the following object properties :

1. No prior information of shape and size are available
2. No motion assumption is required
3. No sensors are equipped on the object

The structure and pose of object on SE(3) are estimated simultaneously. Gauss-Newton (GN) method is implemented for each window with an initial guess generated by OPnP[1] algorithm.

PROBLEM SETUP

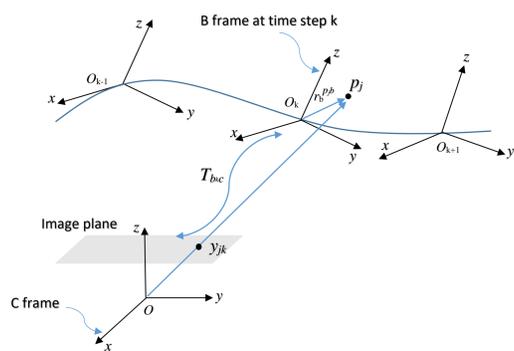


Figure 1: Problem Setup

As shown in Fig.1, the estimated states are pose T_k and structure of the feature point p_j :

$$T_k = T_{c,b_k} = \begin{bmatrix} R_{c,b_k} & t_{c,b_k} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad p_j = \begin{bmatrix} r_b^{p_j} \\ 1 \end{bmatrix} \quad (1)$$

where $k = 1, \dots, K$, K is the window size of SWF; $j = 1, \dots, M$, M is the total number of visible feature points of one window. The measurements y_{jk} , which are corrupted by zero mean Gaussian noise n_{jk} , are pixel coordinates (u_{jk}, v_{jk}) and depth d_{jk} from RGB-D camera. The observation model is

$$y_{jk} = \pi(x_{jk}) + n_{jk}, n_{jk} \sim N(\mathbf{0}, Q_{jk}) \quad (2)$$

with the following shorthand: $x = \{T_1, \dots, T_K, p_1, \dots, p_M\}$, as well as $x_{jk} = \{T_k, p_j\}$ for k th pose and j th feature point.

POSE ESTIMATION ON SE(3)

The measurement inconsistency error of j th feature point from k th pose of a sliding window is defined as:

$$e_{y,jk}(x) = y_{jk} - \pi(x_{jk}) \quad (3)$$

where $z(\pi_{jk})$ transforms feature coordinates from object frame to camera frame, and then projects the coordinates from camera frame into image plane. The objective function is to minimize the sum of squared inconsistency error given states x :

$$J(x) = \frac{1}{2} \sum_{j,k} e_{y,jk}(x)^T Q_{jk}^{-1} e_{y,jk}(x) \quad (4)$$

The states at operating points are perturbed by:

$$T_k = \exp(\epsilon_k^\wedge) T_{op,k} \approx (1 + \epsilon_k^\wedge) T_{op,k} \quad (5)$$

$$p_j = p_{op,j} + D\zeta_j$$

where operator \wedge is defined as

$$\epsilon_{6 \times 1}^\wedge = \begin{bmatrix} \rho_{3 \times 1} \\ \phi_{3 \times 1} \end{bmatrix}^\wedge = \begin{bmatrix} \phi^\times & \rho_{3 \times 1} \\ 0_{1 \times 3} & 0 \end{bmatrix}_{4 \times 4} \quad (6)$$

\times is infinitesimal rotation and \exp is exponential forward mapping on SE(3). To implement GN method, the observation model is linearized as:

$$\pi(x_{jk}) \approx \pi(x_{op,jk}) + \Pi_{jk} \delta x_{jk} \quad (7)$$

where Π_{jk} is the Jacobian of function π with respect to perturbed states. Thus, the objective function can then be rearranged as:

$$J(x) \approx J(x_{op}) - b^T \delta x + \frac{1}{2} \delta x^T A \delta x \quad (8)$$

$$b = H^T Q^{-1} e(x_{op}) \quad (9)$$

$$A = H^T Q^{-1} H$$

CONCLUSION

The experiment shows that the proposed SWF-based framework can estimate the pose of a unknown object with different shape and arbitrary trajectory accurately. Besides, the framework is robust to the number of feature point within each window compared with OPnP algorithm.

$$H = [H_{10}^T, \dots, H_{M0}^T, H_{11}^T, \dots, H_{MK}^T]^T$$

$$Q = \text{diag}\{Q_{10}, \dots, Q_{M0}, Q_{11}, \dots, Q_{MK}\} \quad (10)$$

$$e(x_{op}) = [e_{y,10}(x_{op})^T, \dots, e_{y,MK}(x_{op})^T]^T$$

$$H_{jk} = \Pi_{jk} P_{jk}$$

P_{jk} is a projection matrix to select the visible j th measurement at time step k of the overall perturbed state. The minimum value of objective function is then calculated by:

$$A \delta x^* = b \quad (11)$$

Then the operating points are updated by:

$$T_{op,k} \leftarrow \exp(\epsilon_k^\wedge) T_{op,k} \quad (12)$$

$$p_{op,j} \leftarrow p_{op,j} + D\zeta_j^*$$

When SWF with window size K slides along the time line, the above algorithm runs iteratively for each window until convergence. In the first window, the overall estimated states will be recorded with initial state T_{b_0c} and p_{j_0} generated by Schmidt Orthogonalization. For the rest of windows, the initial state of window l is the first estimated state in window $l - 1$ with initial guess generated by OPnP[1] based algorithm, as summarized follows:

Algorithm 1 Initial Guess Generation

- 1: Match the feature points of time step $k - 1$ and time step k . Set the 3D position $p_{j(k-1)}$ of matched feature points as the initial guess \check{p}_{jk1} ;
- 2: Use the initial guess \check{p}_{jk} and measurement y_{jk} to solve the initial pose \check{T}_{b_kc} by OPnP; when the measurement is insufficient $\check{T}_{b_kc} = T_{b_{(k-1)c}}$;
- 3: Use observation model and initial pose \check{T}_{b_kc} to compute 3D position \check{p}_{jk2} of unmatched feature points. Combine the \check{p}_{jk1} and \check{p}_{jk2} together to form initial points \check{p}_{jk} ;

REFERENCES

- [1] Yinqiang Zheng, Yubin Kuang, Shigeki Sugimoto, Kalle Astrom, and Masatoshi Okutomi. Revisiting the pnp problem: A fast, general and optimal solution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2344–2351, 2013.

RESULTS

The pose error is defined as translational error and rotational error with estimated value x^* :

$$\delta r_k = [\delta r_{x,k} \quad \delta r_{y,k} \quad \delta r_{z,k}]^T := \bar{r}_c^{b_k c^*} - r_c^{b_k c}$$

$$\delta \theta_k^\times = \begin{bmatrix} \delta \theta_{x,k} \\ \delta \theta_{y,k} \\ \delta \theta_{z,k} \end{bmatrix}^\times := \mathbf{1} - \bar{R}_{b_k c}^* R_{b_k c}^T \quad (13)$$

The objects, estimated errors and estimated trajectories are shown in Fig.2.

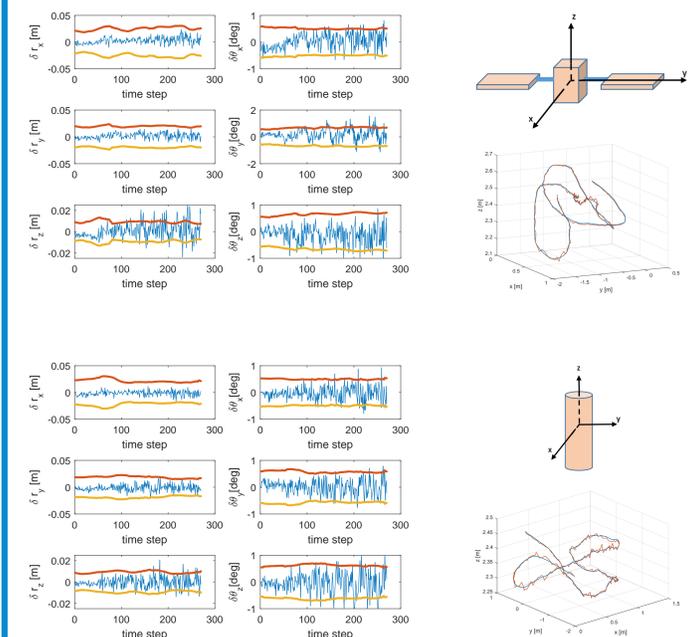


Figure 2: Results of two objects

To compare with OPnP in cases of different feature point, the following errors are defined as follows and Fig.3 shows the comparison result.

$$e_{rk} = \|\delta r_k\|_2, \quad e_{\theta k} = \|\delta \theta_k\|_2, \quad (14)$$

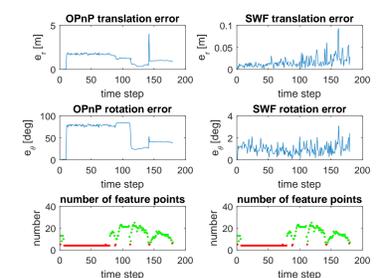


Figure 3: Comparison of OPnP