

UKRAINIAN CATHOLIC UNIVERSITY

FACULTY OF APPLIED SCIENCES

COMPUTER SCIENCE PROGRAMME

Simultaneous Localization and Mapping

Linear Algebra final project report

Authors:

Bohdan HLOVATSKYI, Mykhailo PASICHNYK, Stefan-Yuriy
MALYK

May 2022



APPLIED
SCIENCES
FACULTY ●

Abstract

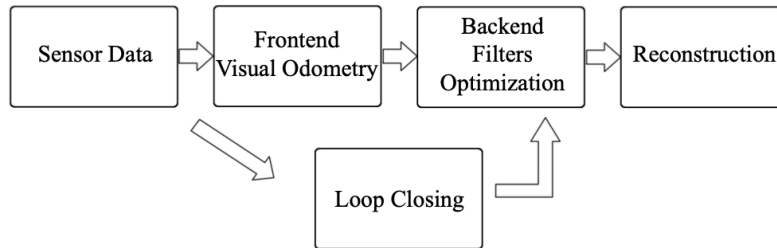
SLAM or Simultaneous Localization and Mapping - problem of constructing a map of unknown environment simultaneously tracking the object's location. Such algorithms utilize the data obtained from sensors, such as video camera, lidar etc to construct a 3-D model of the nearby objects. They are used in variety of fields: autonomous driving and robots, augmented reality and much more. The algorithms that rely on video data are referred to as Visual SLAM. They obtain the valuable information by means of Computer Vision and different techniques of Signal Processing. We implemented such algorithm by means of Python and pyopencv as our Linear Algebra course work. Source code: [1].

1 Introduction

Visual SLAM is a family of algorithms that solve the problem of Simultaneous Localization and Mapping by means of Computer Vision. Depending on the type of input video stream there are different types of visual SLAM:

- Monocular SLAM. Such SLAM takes as an input video stream from a single RGB camera. This is the simplest type of SLAM, though it has a lot of limitations, the most important of which is the fact that we can construct environment map only to some scale factor. In other words, we cannot obtain the exact depth of featured points in an environment from single image source.
- Stereo SLAM. Input video stream comes from two cameras. This allows to get the exact depth of the featured points.
- RGBD. Specific kind of SLAM that gets the input data from RGBD camera - the one that included depth of objects within an image. There are different technologies for such cameras though they are not particularly important for the algorithm itself.

The pipeline of SLAM algorithms will usually look in the following way:



In this work we are mainly interested in Frontend Visual Odometry, goal of which is to track the object that produces the source data and build non-optimized 3-D cloud map of environment. Backend optimization and loop closing are powerful techniques to optimize the SLAM and get rid of the noise that brings a lot of inaccuracy in the SLAM algorithm. Those algorithm, such as Kalman and Particle filters, are beyond the scope of Computer vision and are related mainly to signal processing, thus won't be focused on here.

2 Problem setting

2.1 Problem formulation

With the given setting, we are interested mainly in two parts:

1. Modeling the object's motion. In an abstract setting this can be modeled as $x_k = f(x_{k-1}, u_k, w_k)$, where x_k denotes a current position of an object, x_{k-1} - its previous position, u_k - change of its state and w_k - noise.
2. Collecting the landmarks of a nearby environment. The observed data $z_{k,j}$ will then be described by this equation: $z_{k,j} = f(y_j, x_k, v_{k,j})$. Here y_j - landmark point, x_k - its coordinates, $v_{k,j}$ - noise of the observation.

These two equations summarize the SLAM. For the Euclidean space the motion model can be expressed as the so-called special Euclidean Group:

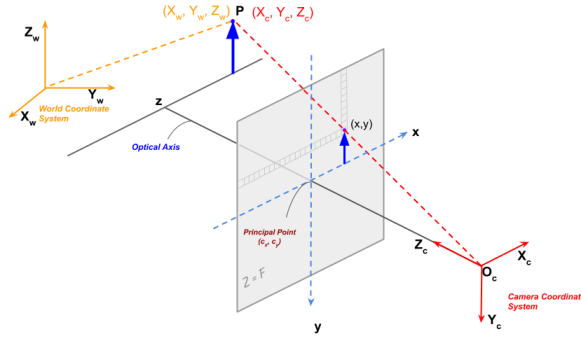
$$SE(3) = \{T = \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} \in \mathbf{R}^{4 \times 4}\},$$

where R - rotation matrix, t - translation vector. It allows to express state a_1 as

$$a_1 = Ra_0 + t.$$

Important note here is that we associate the state of an object with the location of camera that produces some input data.

2.2 Camera model



Any point in space can be represented by three coordinate systems:

- The World Coordinate System - point $\mathbf{P} = (X_w, Y_w, Z_w)$. For example if our object is in a room we can treat floor as $x - y$ axis and a wall as z axis
- The Camera Coordinate System - (X_c, Y_c, Z_c) the same object in 3D space but from camera, where X_c points to the right, Y_c down and Z_c towards the camera
- The point projection to the plane (image itself)

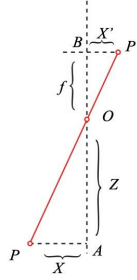
The world and camera coordinate systems are related by the rotation matrix \mathbf{R} and translation vector \mathbf{t} . Thus

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = \mathbf{R} \begin{bmatrix} X_w \\ Y_w \\ Z_w \end{bmatrix} + \mathbf{t}$$

We can rewrite statement above to make it more compact and call such $[\mathbf{R}|\mathbf{t}]$ matrix the **Extrinsic matrix**

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = [\mathbf{R}|\mathbf{t}] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

Now that we have our point \mathbf{P} in the camera coordinate system we can project it onto the plane.



Our projected point (x, y) can be found using following equations:

$$x = f \frac{X_c}{Z_c}, y = f \frac{Y_c}{Z_c}$$

where f is the focal length of the camera. One important thing to mention is that the image \mathbf{P} going through camera's optical center would be inverted. However, since usage of such camera would be inconvenient camera's software performs some transformation. Thus, this step could be omitted and we can stick to the equations above. Next we construct such a matrix:

$$\begin{bmatrix} x' \\ y' \\ z' \end{bmatrix} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & f \end{bmatrix} \begin{bmatrix} Y_c \\ X_c \\ Z_c \end{bmatrix}$$

It would be referred as the **Intrinsic Matrix**. Another important observation to make is that the optical center of the camera could be shifted from the origin of the camera coordinate system. Thus we need to modify the Intrinsic matrix as such:

$$\mathbf{K} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & f \end{bmatrix}$$

Combining all of the above we can get coordinates of the object on our image (u, v) :

$$\begin{bmatrix} u' \\ v' \\ w' \end{bmatrix} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & f \end{bmatrix} \begin{bmatrix} Y_c \\ X_c \\ Z_c \end{bmatrix}$$

where

$$u = \frac{u'}{w'}, v = \frac{v'}{w'}$$

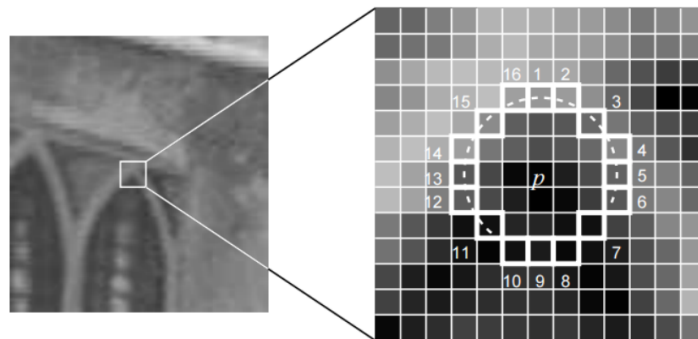
2.3 Feature detection

Firstly let define what a feature is. It is a some pattern in the image which has several proprieties:

1. **Repeatability**: The same feature can be found in different images.
2. **Distinctiveness**: Different features have different expressions.
3. **Efficiency**: In the same image, the number of feature points should be far smaller than the number of pixels.
4. **Locality**: The feature is only related to a small image area.

Historically there were 3 most popular approaches in algorithms of feature detection:

1. **SIFT** (Scale-Invariant Feature Transform) : is one of the most classic. It is fully consider the changes in illumination, scale, and rotation during the image transformation, but is so calculating expensive, so rarely used in SLAM.
2. **FAST** : is extremely fast to compute. Has following algorithm:
 - (a) Select pixel p in the image, assuming its brightness as I_p
 - (b) Set a threshold T (for example, 20% of I_p).
 - (c) Take the pixel p as the center, and select the 16 pixels on a circle with a radius of 3.
 - (d) If there are consecutive N points on the selected circle whose brightness is greater than $I_p + T$ or less than $I_p - T$, then the central pixel p can be considered a feature point. N usually takes 12, which is FAST-12. Also N could be 9 / 11 / etc, with respective names of FAST.
 - (e) Iterate through the above four steps on each pixel.

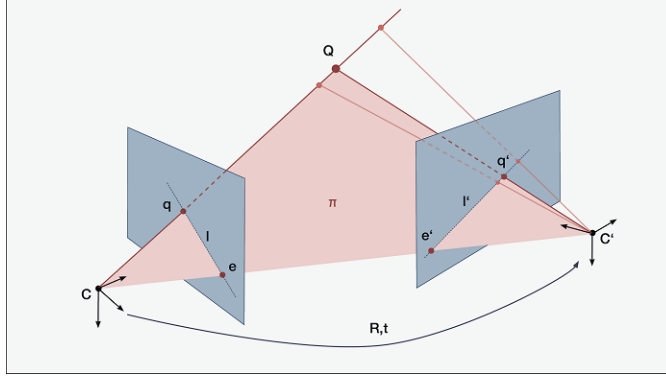


[Fig 1] : Illustration of of FAST algorithm

It has problem that the detector does not have descriptors (fixed pattern, which try to find). Also, as it has fixed the radius of the circle = 3, there is a scaling problem. This two problems fix ORB algorithm.

2.4 Epipolar geometry

The object motion can be found from the correspondence between two consecutive video frames.



According to the camera model, each of the landmarks Q will be projected onto 2-D plane of pixels as q, q' . Motion, described by the rotation matrix R and transformation vector t can be reproduced than via Epipolar constraint.

Pixel points q, q' are described by the following equations:

$$\begin{cases} s_1 p_1 = K P \\ s_2 p_2 = K(RP + t) \end{cases} \quad (1)$$

As was discussed in the setting, due to the limitations of the approach, we can represent the points only to some constant factor. Moreover, it is worth mentioning that p_1, p_2 will be the points matched by the feature matcher.

To find the coordinates on the normalized plane of two pixels, we can simply take the inverse of the matrix representing a camera: $x_{1,2} = K^{-1}p_{1,2}$. Then the equation above can be rewritten as $x_2 = Rx_1 + t$. This can be further transformed into

$$p_2^T K^{-T} t \times R K^{-1} p_1 = 0.$$

This equation is referred to as epipolar constraint. For the sake of simplicity, the notions of Fundamental and Essential matrices are introduced:

$$\begin{cases} E = t \times R \\ F = K^{-T} t \times R K^{-1} \\ x_2^T E x_1 = 0 \\ p_2^T F p_1 = 0 \end{cases} \quad (2)$$

It is clear that usage of Essential matrix will make the calculations more accurate, as it uses precomputed calibration parameters for the camera. The process of camera pose estimation can then be summarized as:

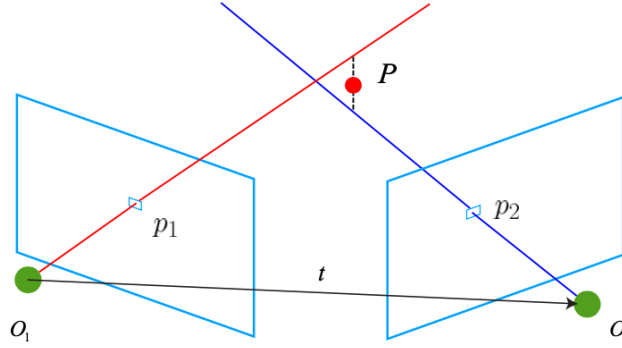
1. Find E or F from the known pixel positions.
2. Find R, t based on those matrices.

In practice, the Essential matrix is found via Eight-Point algorithm and then R, t are found via SVD and triangulation to get the result that would satisfy the constraints given by the known projections of landmarks.

$$\begin{cases} R = UR_{Z(\pm\pi/2)}V^T \\ t = UR_{Z(\pm\pi/2)}\Sigma U^T \end{cases} \quad (3)$$

2.5 Triangulation

The second goal of SLAM is to estimate the relative positions of the feature points. Knowing the relative position of the camera and the relative positions of those feature points, we can reconstruct their global coordinates and construct the feature map of the environment.



From the epipolar constraint: $s_2x_2 = s_1Rx_1 + t$. Here $x_{1,2}$ - normalised coordinates (coordinates in the pixels' plane) of two feature points. Knowing R, t , we can now calculate $s_{1,2}$ that will represent the 'depth' of a pixel.

$$s_2x_2 \times x_2 = 0 = s_1x_2 \times Rx_1 + x_2 \times t, \quad (4)$$

which can be solved for s_1 . Geometrically, we seek to find a 3D point on the ray O_1p_1 to make its projection close to p_2 .

3 Overview of approaches

As was mentioned earlier, there exists different types of SLAM algorithms that depend on different type of input video stream. Though as we are focusing particularly on visual SLAM in this work, we would consider different approached to SLAM problem from the perspective of Computer Vision:

- The most popular implementation of visual SLAM is OpenVSLAM available open-source. It utilises notions of gloval and local key frame maps. For the local map it finds points via triangulation and for the global one uses pose-graph optimization. It also contain a global optimization module, which optimizes mathcing against global map. [7]
- Interesting to mention is the approach proposed by the ORB-SLAM. The authors use beg-of-words representation of encoded descriptors for each of the feature points

in order to optimise the process of feature matching against the stored feature map. It has also a logical continuation as ORB-SLAM 2, which introduced ability to work with stereo input data. [3]

- Speaking of the camera calibration, the classical approach to it is using a chess board' photos from different perspective in order to estimate the intrinsic parameters. [4]. There are also some research conducted in the field of camera self-calibration. [5].
- It is worth mentioning that the research in the SLAM field is being actively maintained. The good example of it is creation of general frameworks such as expressed in the cited papeer. [8].

4 Implementation

Data: stream of video frames

initialization;

calibrate camera;

while true do

 get frame;

 get the features;

 match the features between current and previous frame;

 find essential matrix;

 find rotation matrix and transformation;

 store the movement and point map;

 update previous frame;

end

Algorithm 1: General algorithm of SLAM

The algorithm was implemented in Python via pyopencv. The visualizer of the SLAM was written via OpegGL. For the feature extraction the ORB was used due to its computational efficiency. Feature matching was done in brute-force manner because limitation on the number of features was only several hundreds. Fundamental matrix was then constructed as we lack the camera parameters for the test video.



5 Experiments

6 Limitations

There are two main limitations of such approach:

1. Triangulation is only possible if the enough amount of translation is present. If the amount of translation is small, pixel's change would result in big depth uncertainty. On the other hand, if translation is too big, the algorithm may fail to match the features. This problem will be especially seen in cases of so-called pure rotation, when there is no translation. Though in not so extreme cases, there are means to make this problem less visible - delayed triangulation etc.
2. Assumption about static sight of view brings a lot of limitations as well. For instance, the algorithm cannot be applied to autopilots as the moving object would make it think that the camera has moved.

7 Conclusions

References

- [1] <https://github.com/bohdanhlovatskyi/OhISee>
- [2] http://www.cvlibs.net/datasets/kitti/eval_odometry.php
- [3] <http://webdiis.unizar.es/~raulmur/MurMontielTardosTR015.pdf>
- [4] <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tr98-71.pdf>
- [5] <https://users.aalto.fi/~kannalj1/publications/wacv2016.pdf>
- [6] <https://github.com/gaoxiang12/slambook-en>
- [7] <https://arxiv.org/pdf/1910.01122.pdf>
- [8] <https://arxiv.org/pdf/1902.07995.pdf>