



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



INFORMATION ENGINEERING DEPARTMENT
COMPUTER ENGINEERING - AI & ROBOTICS

3D Data Processing 2023/2024
Francesco Crisci 2076739

Contents

1	Intoduction concepts	2
1.1	Sensors	2
1.1.1	Proximity sensors	2
1.1.2	Sonar	2
1.1.3	LiDAR	2
1.1.4	Radar	2
1.1.5	Time-of-flight Cameras	2
1.1.6	Stereo Cameras	2
1.1.7	RGB-D Cameras	3
1.2	Linear Algebra and Rigid-Body Transofrmationas	3
1.2.1	Rigid Objects	3
1.2.2	Rotation Matrix	3
1.3	Homographies	4
1.3.1	Back-Projections	4
1.3.2	Scene to Image Homography	4
1.3.3	Normalization	5
1.3.4	Pure Rotation Homography	5
1.3.5	Image to Image Homography	5
1.3.6	Robust estimation: RaNSAC Algorithm	5
1.4	Camera Calibration	6
1.4.1	Pseudoinverse	6
1.4.2	Nonlinear Least Squares Problems	6
1.4.3	Gauss-Newton	7
1.4.4	Leverberg-Marquardt	7
1.5	Pinhole Camera Calibration	7
1.5.1	Data	7
1.5.2	Initial Guess for K	7
1.5.3	Initial Guesses for R_i, t_i	8
1.5.4	Refined Calibration	8
1.5.5	PnP: Pwerspective-n-Point Problem	8

Chapter 1

Introduction concepts

1.1 Sensors

A point cloud is a data structure used to represent a collection of **multi-dimensional points** and is commonly used to represent three-dimensional data. The points usually represent the X,Y and Z geometric coordinates of a sampled surface. Each point can hold additional information: RGB colors, intensity values and so on.

1.1.1 Proximity sensors

An emitter transmits a light wavelength. Light eventually reflected back from a nearby object, e.g. perceived by a photoresistor.

1.1.2 Sonar

Senses the range of an obstacle from the round-trip time of an ultrasound pulse. The narrower the cone, the higher the angular resolution of the sensor.

1.1.3 LiDAR

It stands for Light Detection And Ranging. We have two types: **Continuous wave**, and **Pulse-based**: it measures directly the time-of-flight, i.e. the round-trip-time of a pulse of light. Needs very short laser pulses and high temporal accuracy.

Solid-state LiDAR uses non-mechanical elements to deflect the laser beams. Usually an optical phase array (OPA) of closely spaced optical antennas.

1.1.4 Radar

Radio Detection and Ranging are similar to LiDARs, but emit radio waves. **FMCW**(Frequency-Modulated Continuous-Wave) radars emit a signal where the frequency varies linearly over time. The difference between the frequency of the signal sent and that received is linearly related to the distance from the object that generated the reflected signal.

Lower angular and range resolution/accuracy w.r.t. LiDARs, but more robust to environmental conditions such as snow, fog, ... and able to compute the relative velocity between the sensor and the objects. In **imaging radar**, multiple FMCW transceivers are cascaded together to generate a dense 3D point cloud of the surrounding environment.

1.1.5 Time-of-flight Cameras

Meeting point between LiDARs and digital cameras. IR light emitters illuminate the scene, sensed by a 2D imaging sensor sensitive to the projected light.

1.1.6 Stereo Cameras

Two or more usually identical, rigidly mounted cameras, framing the same scene from different point of view. Knowing the rigid body transformation between the two cameras and the pixel onto which a 3D point is projected in both cameras, it is possible to estimate the depth of the 3D point by using triangulation techniques.

Active stereo cameras: To create visual saliency also in homogeneous surfaces. Couple the digital cameras with a dense pattern projector that projects into the scene some visible textured pattern.

Structured Light Cameras: SL cameras employ the same operating principle of stereo cameras, but one of the two cameras is replaced with a light projector that illuminates the scene with a textured visual pattern.

The pattern is known and each pattern patch represents a unique binary matrix.

Time-Multiplexing coding: Projecting a series of light patterns so that every pixel is encoded with the codeword identified by the pattern sequence. Most common structure of the patterns is a sequence of stripes in(de)creasing their width by the time.

1.1.7 RGB-D Cameras

RGB-D cameras are sensor ensembles composed by a ToF or a SL camera rigidly coupled in the same chassis with a color camera. Each pixel of the RGB image is associated with a corresponding pixel in depth map.

1.2 Linear Algebra and Rigid-Body Transformations

When dealing with 3D data, we often need to change the coordinate frame, move sensors, move objects in front of sensors, etc...

Many 3D sensors are based on projective geometry.

Euclidean Space: E^N , the set of n-dimensional points that satisfy the axioms of Euclidean geometry. Relationships can be expressed in terms of angles and distances. No special point, no additions between points, only differences (**vectors**).

Linear Vector Space: A set V of objects over the field F is a vector space if its elements are closed under **scalar multiplication** and **vector summation**. Given two vectors $u, v \in V$ and $\alpha, \beta \in \mathbb{F}$, we have $\alpha u + \beta v \in V$.

Euclidean Vector Space: a linear space denoted by \mathbb{R}^N is called **real vector space**, each element defined by an ordered tuple. Any point in E^N can be identified with a vector in \mathbb{R}^N .

Dot Product: Given two vectors u, v in \mathbb{R}^N , we define the canonical inner product, or dot product, as:

$$\langle u, v \rangle = u^T v = u_1 v_1 + \dots + u_n v_n$$

Distances and Angles: dot product is used to measure distances: Euclidean norm or L2-norm is just the square root of the dot product. Dot product is also used to measure angles.

Subspaces: a subset $w \subseteq V$ of a vector space V is called a subspace if the zero vector belongs to W . given a set of vectors S , the spanned subspace is :

$$\text{span}(S) = \left\{ \sum_{i=1}^n \alpha_i v_i \mid v_i \in S, \alpha_i \in \mathbb{R} \right\}$$

S is linearly independent if $\sum_{i=1}^n \alpha_i b_i = 0 \Rightarrow \alpha_i = 0 \forall i = 1, \dots, m$.

Basis: A set of vectors S is said to be a basis for V if it is linearly independent and it spans the whole space V . There are other notions from linear algebra, but it is the fifth year seeing this stuff, so i'm gonna skip them. Feel free to look at the slides.

1.2.1 Rigid Objects

A non-deformable object O in 3D can be modeled by a set of points in a Euclidean vector space. We assume that the coordinate frame of this space represents an inertial frame we call the **world** frame. If the object starts to move, its points move accordingly, i.e., they change their coordinates $p(t)$ w.r.t. the world frame over time. Distances between any of these two points $u(t)$, $q(t)$ must be preserved over time as the object moves:

$$\|u(t_1) - q(t_1)\| = \|u(t_2) - q(t_2)\|, \forall t_1, t_2 \in \mathbb{R}$$

A transformation that preserves distances is called a Euclidean transformation. Not enough to describe a rigid body movement. Consider reflection: it does not preserve, generally, the cross product.

Individual objects points cannot move relative to each other. Attach a coordinate frame to the object. Rigidbody motion w.r.t. the world frame can be described by the motion of the origin of the object frame and the rotation of the object frame.

Object pose is defined by translational and rotational part. The translational part is a vector t between the origin of the world frame and that of the object frame. The rotational part is a relative orientation R of the object frame axes w.r.t. fixed world frame axes. In practice, we apply R first, then add the translation t .

Suppose the object is performing a pure rotation around its origin o . We may always assume that the origin of the world frame is the center of rotation o .

1.2.2 Rotation Matrix

All objects coordinates p_o are related to \mathbb{R}^3 standard basis in the object frame:

$$p_o = [x_o \ y_o \ z_o]^T = x_o e_{o,1} + y_o e_{o,2} + z_o e_{o,3}$$

To represent the position of any point w.r.t. world, we just need to present $e_{o,1}, e_{o,2}, e_{o,3} : T(e_{o,1} = r_1), T(e_{o,2} = r_2), T(e_{o,3} = r_3)$ More math on the slides, including the slides "FROM OBJECTS TO CAMERA".



1.3 Homographies

Some recap about distortion models, and how to correct the distortion. To fit the image, the input k matrix may be changed while performing image undistortion.

1.3.1 Back-Projections

Given a depth map how can we project back the rays into 3D points? For each pixel (u,v) : we look at the pixel dept $d=\text{depth}(u,v)$, we normalize the values:

$$u_n = \frac{u - u_c}{\alpha_u}, \quad v_n = \frac{v - v_c}{\alpha_v}$$

and we calculate back-projection computed as:

$$\begin{bmatrix} Z \\ Y \\ X \end{bmatrix} = d \begin{bmatrix} u_n \\ v_n \\ 1 \end{bmatrix}$$

1.3.2 Scene to Image Homography

A **homography** is a transformation in $GL(n)/\mathbb{R}$ that maps points and lines to lines. Very useful model when dealing with planar scenes.

Let Q and q be a 3D point on a planar board and its 2D projection into the image plane, respectively:

$$\tilde{q} = AT\tilde{Q} = K[R|t]\tilde{Q}$$

Without loss of generality, we can choose to define the object plane so that $Z=0$:

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K[r_1 \ r_2 \ r_3 \ t] \begin{pmatrix} X \\ Y \\ 0 \\ 1 \end{pmatrix} = K[r_1 \ r_2 \ t] \begin{pmatrix} X \\ Y \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = H \begin{pmatrix} X \\ Y \\ 1 \end{pmatrix}$$

Due to the scale equivalence of homogeneous coordinates, H is defined up to a scalar facto (8 DOFs)

$$\lambda \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = H \begin{pmatrix} X \\ Y \\ 1 \end{pmatrix} \rightarrow \lambda H = H \forall \lambda \in \mathbb{R}$$

If we normalize the image coordinates, we have $H = [r_1 \ r_2 \ t]$

If the rigid body transformation is unknown, but the camera is calibrated and we are able to estimate H from data:

$$H = [r_1 \ r_2 \ t]$$

we can extract R and t as:

- Scale H with $\lambda = 1/|h_1|$ to try to guarantee the condition of orthonormality
- t is just the third column of H
- R can be computed in closed form:

$$r_3 = r_1 \times r_2$$

from

$$\lambda \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

In homogeneous coordinates we have:

$$x = \frac{h_{11}X + h_{12}Y + h_{13}}{h_{31}X + h_{32}Y + h_{33}} \Rightarrow x(h_{31}X + h_{32}Y + h_{33}) - h_{11}X - h_{12}Y - h_{13} = 0$$

$$y = \frac{h_{21}X + h_{22}Y + h_{23}}{h_{31}X + h_{32}Y + h_{33}} \Rightarrow y(h_{31}X + h_{32}Y + h_{33}) - h_{21}X - h_{22}Y - h_{23} = 0$$

From a set of four correspondences we obtain 8 equations with unknowns the elements of H . Stacking H as a 9×1 column vector h , we have to solve the following underdetermined system $Ah=0$.

The singular Value decomposition (SVD) of a $m \times n$ real matrix is a factorization of the form:

$$M = U\Sigma V^T$$

To solve the system $Ax=0$ we can use SVD, obtaining $U\Sigma V^T x = 0$ and set as solution the **last column of V** . The minimum number of correspondences required to estimate a homography is four. Better to use more than four correspondences. A number of matches n greater than the minimum will not match exactly the same homography due to image noise: **use a least square solution for a over-determined system of $2n$ linear equations**. We can written the problem as:

$$\arg \min_x ((Ax)^A x), \text{ subject to } \|x\| = 1$$

Let $y = V^T x$, we can rewrite the problem after SVD as:

$$\arg \min_y (y^T \Sigma^T \Sigma y), \text{ subject to } \|y\| = 1$$

The straightforward solution is $\hat{y} = (0, \dots, 0, 1)^T$. From $\hat{x} = V\hat{y}$ the solution is the **last column of V** .

1.3.3 Normalization

SVD-based homography estimation performs badly in presence of noise: **use data normalization prior to the application of SVD**. Translate and scale both image and world coordinates to avoid orders of magnitude difference between the columns of the data matrix:

$$x_n = N_i x, \quad X_n = N_s X$$

Shift the coordinates so that the average coordinate is $(0,0)$; Scale the coordinates so that the average distance from the origin is 1. The homography relative to the original coordinates can be recovered as:

$$H = N_i^{-1} \tilde{H} N_s$$

Quick summary

To estimate H from a set of correspondences $Ah=0$, normalize the data, compute the SVD and take as solution H' the last column of V . Denormalize H' to H .

1.3.4 Pure Rotation Homography

If two views are separated by a pure rotation around the camera center:

$$x = K[I|0]X, \quad x' = K[R|0]X$$

Mapping between corresponding pixels in the two images are defined by a homography:

$$x' = K R K^{-1} x = H x$$

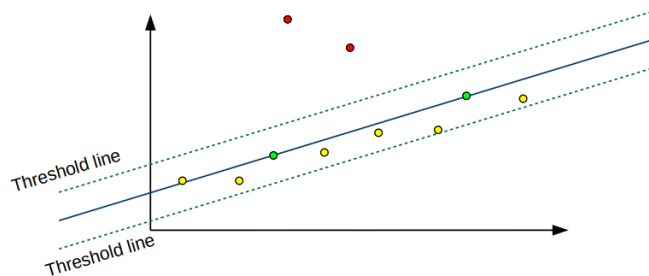
1.3.5 Image to Image Homography

If we have two view of a planar scene:

$$\lambda_x x = H_x x, \quad \lambda_{x'} x' = H_{x'} X, \quad H = H_{x'} H_x^{-1} \Rightarrow \lambda x' = H x$$

From 2D-2D correspondences, it is possible to decompose H into the rigid body transformation R, t that relates the two cameras. Specifically, up to 4 solutions for such decomposition, with at most two physically possible.

1.3.6 Robust estimation: RANSAC Algorithm



Random Sample consensus: repeat the following steps N times: select two points at random (**putative points**), fit a line **l**, find the number of inliers that support **l**. The line with most inliers wins. Fit a new line with all the selected inliers.

For Homography Estimation, we need to repeat the following steps N times:

- Select a set of four correspondences at Random
- Estimate H solving the **undetermined** system $Ah=0$
- Find the number of inlier that support H :
 - For each correspondence $\{P, p\}$, project 3D point P with H , i.e.: $p'=HP$
 - If the distance between p and p' is less than a threshold, add the correspondence to the inliers set associated to H .
- Find the H with the highest number of inliers among the N trials
- Re-compute a least squares estimate of H using all m inliers solving an **overdetermined** system.

How many putative correspondences N we need to get at least one good sample with high probability? Probability p of getting at least a good sample after N trials:

$$p = 1 - (1 - (1 - \epsilon)^s)^N$$

where ϵ is the fraction of outliers and s is the number of required correspondences to compute a solution. We have:

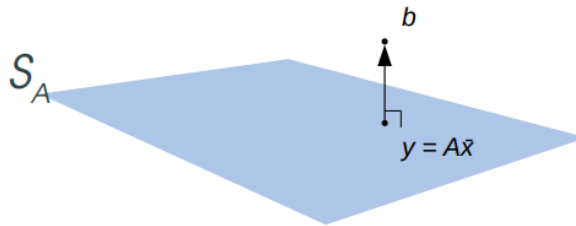
$$N = \frac{\log(1 - p)}{\log(1 - (1 - \epsilon)^s)}$$

1.4 Camera Calibration

Let the following overdetermined linear system $Ax=b$. we aim to find a solution for the following least square problem:

$$\arg \min_x \|Ax - b\|^2$$

The n columns of A span a subspace $S_A \subseteq R^M$ with dimension up to n . Let $\bar{x} \in R^N$ a solution, with $y = A\bar{x} \in S_A$: y is the point in S_A closest to b . $y-b$ should be perpendicular to each of the column of A :



$$A^T(Ax - b) = 0 \Rightarrow A^T Ax = A^T b$$

1.4.1 Pseudoinverse

New linear system, with $A^T A$ $n \times n$ symmetric matrix with same solution as $Ax=b$. It defines the so-called normal equations of a least square problem: if $A^T A$ is invertible, the solution is just:

$$x = (A^T A)^{-1} A^T b$$

defined as the **pseudoinverse of A**. Otherwise, solve the normal equations.

1.4.2 Nonlinear Least Squares Problems

Let f a nonlinear model function:

$$f : x \rightarrow f(x, a)$$

f depends on a set of parameters: $a = [a_1, \dots, a_m]^T$. Dataset composed by n pairs $\{X_i, Y_i\}$ ("data points"), where x_i is an independent variable (input) and y_i an observation given x_i (output). "Goodness" of the model parameters with respect to the i -th data point is defined by:

$$\text{residual} \rightarrow e_i(a) = \|y_i - f(x_i, a)\|$$

Nonlinear least-square methods find the parameters vector a that minimizes the sum of squared residuals:

$$E(a) = \frac{1}{2} \sum_{i=1}^n \|y_i - f(x_i, a)\|^2 = \frac{1}{2} \sum_{i=1}^n e_i(a)^2$$



We assume to have a "good" initial estimate a_0 of the parameters vector. the goal of an iterative least square method is to choose at each iteration an update δa of the current state a_k that reduces the value of the cost function.

$$a_{k+1} = a_k + \delta a$$

We might stop after the number of iterations, step size or something else.

1.4.3 Gauss-Newton

Gauss-Newton is an efficient iterative method used to minimize a sum of squared function values:

$$E(a + \delta a) \approx E(a) + dE_a \delta a + \frac{1}{2!} (\delta a^T d^2 E_a \delta a)$$

We want to find an update δa that minimizes the error function E . The Gauss-Newton update:

$$\delta a^{GN} = -(J^T J)^{-1} J^T e(a)$$

The Pros are that it generally ensures fast convergence also for nonlinear model function but the convergence is not guaranteed: an update may increase the cost function.

We have notions of gradient descent here, but i'm gonna omit them since we saw it literally everywhere.

1.4.4 Levenberg-Marquardt

Levenberg-Marquardt algorithm combines Gauss Newton and Gradient Descent, trying to exploit the strength of both these approaches:

$$\delta a^{LM} = -(J^T J + \lambda I)^{-1} J^T e(a)$$

It tunes the damping factor λ at each iteration to ensure rapid progress even where Gauss-Newton fails. If λ increases we tend to converge toward the Gradient Descent, if it decreases we go toward the Gauss Newton algorithm.

1.5 Pinhole Camera Calibration

Camera calibration estimates intrinsic and extrinsic parameters of a given camera. **Intrinsic** parameters are the focal length, coordinates of the principal point, distance parameters, etc... **Extrinsic** parameters identify the pose of the camera, e.g. in a world coordinate system.

To calibrate the camera we do the following steps:

- Use a known planar pattern
- Collect a sequence of images of the pattern in several positions, and extract all corners from images
- Exploit the scene-to-image homography constraints to extract an initial guess of the intrinsic parameters matrix K (we assume no distortion for now)
- Normalize all pixel projections with K and for each image i extract the homography and recover an initial guess of the intrinsic parameters (R_i, y_i)
- Recover the final intrinsic and extrinsic parameters solving an iterative least square problem using K and (R_i, y_i) as initial guess. (We optimize also the distortion parameters here)

1.5.1 Data

Collect k images, each one framing a checkerboard with m internal corners from different point of views. For each image i , extract the m 2D corners projections: $p_{i,0}, \dots, p_{i,m-1}$

1.5.2 Initial Guess for K

We need to estimate an initial guess for K , i.e.:

$$\theta_K = \alpha_u, \alpha_v, u_c, v_c, s$$

For each image i , use the m 3D-2D matches to estimate a scene-to-image homography H_i

$$H = [h_1 \ h_2 \ h_3]$$

Remembering that

$$\lambda_H h_1 = K r_1, \quad \lambda_H h_2 = K r_2$$



From orthonormality of rotation matrix R :

$$\begin{aligned} r_1^T r_2 &= h_1^T K^{-T} K^{-1} h_2 = 0 \\ r_1^T r_1 &= r_2^T r_2 \Rightarrow h_1^T K^{-T} K^{-1} h_1 = h_2^T K^{-T} K^{-1} h_2 \end{aligned}$$

Exactly two constraints for each homography. We define:

$$B = K^{-T} K^{-1} = \begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \\ B_{31} & B_{32} & B_{33} \end{bmatrix}$$

Rearrange each constraint as:

$$h_i^T B h_j = v_{ij}^T b$$

where

$$\begin{aligned} v_{ij} &= [h_{1i}h_{1j}, h_{1i}h_{2j} + h_{2i}h_{1j}, h_{2i}h_{2j}, h_{3i}h_{1j} + h_{1i}h_{3j}, h_{3i}h_{2j} + h_{2i}h_{3j}, h_{3i}h_{3j}]^T \\ b &= [B_{11}, B_{12}, B_{22}, B_{13}, B_{23}, B_{33}]^T \end{aligned}$$

Two constraints can be rewritten as:

$$\begin{bmatrix} V_{12}^T \\ (v_{11} - v_{22})^T \end{bmatrix} b = 0$$

that is a 2×6 system of equations.

Given k images arrange a $2k \times 6$ homogeneous linear system: $Vb=0$. Solve it by using SVD and get B . We applied the **Direct Linear Transformation (DLT) algorithm**, which solves a set of variables from a set of similarity relations. Extract the five intrinsic camera parameters plus one unknown scale factor from the 6 equations defined by:

$$B = \lambda_B K^{-T} K^{-1}$$

1.5.3 Initial Guesses for R_i, t_i

For each checkerboard position, we normalize all the image points using the initial guess camera intrinsics computed in the previous step. Remembering that $H = [r_1 \ r_2 \ t]$ we can extract R, t .

1.5.4 Refined Calibration

At each optimization iteration, given the current estimate of the transformations T_i and the calibration parameters, we can project into each image the m 3D corners and compare them with the corresponding pixels, and optimize again looking for better T_i and calibration parameters that minimize the cost function:

$$\arg \min_{\theta_K, \theta_d, T_0, \dots, T_{k-1}} \sum_{i=0, j=0}^{k-1, m-1} \|K f_d(A_n T_i \tilde{P}_j) - \tilde{p}_{i,j}\|^2$$

The ideas for the minimization: we use the LM algorithm to solve the iterative problem; we parametrize the rotations with **axis-angle representation**

1.5.5 PnP: Pwerspective-n-Point Problem

Given a rigid object and a set of 3D-2D correspondences, how to estimate the rigid body transformation R, t with a **calibrated camera**? Find an initial guess for R, t , then solve with LM a nonlinear least-square problem with parameters R, t using as starting point the initial guess.