



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



INFORMATION ENGINEERING DEPARTMENT
COMPUTER ENGINEERING - AI & ROBOTICS

3D Data Processing 2023/2024
Francesco Crisci 2076739

Contents

1	Introduction concepts	4
1.1	Sensors	4
1.1.1	Proximity sensors	4
1.1.2	Sonar	4
1.1.3	LiDAR	4
1.1.4	Radar	4
1.1.5	Time-of-flight Cameras	4
1.1.6	Stereo Cameras	4
1.1.7	RGB-D Cameras	5
1.2	Linear Algebra and Rigid-Body Transformations	5
1.2.1	Rigid Objects	5
1.2.2	Rotation Matrix	5
1.3	Homographies	6
1.3.1	Back-Projections	6
1.3.2	Scene to Image Homography	6
1.3.3	Normalization	7
1.3.4	Pure Rotation Homography	7
1.3.5	Image to Image Homography	7
1.3.6	Robust estimation: RaNSAC Algorithm	7
1.4	Camera Calibration	8
1.4.1	Pseudoinverse	8
1.4.2	Nonlinear Least Squares Problems	8
1.4.3	Gauss-Newton	9
1.4.4	Leverberg-Marquardt	9
1.5	Pinhole Camera Calibration	9
1.5.1	Data	9
1.5.2	Initial Guess for K	9
1.5.3	Initial Guesses for R_i, t_i	10
1.5.4	Refined Calibration	10
1.5.5	PnP: Pwerspective-n-Point Problem	10
2	Epipolar Geometry	11
2.1	Epipoles and Epipolar Lines	12
2.1.1	Essential Test	12
2.2	Pose Recovery from E	12
2.3	Fundamental Matrix	12
2.4	Robust Fundamental Matrix Estimation	13
2.5	Alternative Solutions	14
2.6	Essential Matrix Estimation	14
2.7	Stereo Rectification	14
2.7.1	Trinagulation in a Rectilinear Rig	15
3	Stereo Matching	16
3.1	Matching Model	16
3.2	Depth uncertainty	16
3.3	Block Matching	17
3.3.1	Matching cost	17
3.4	Global Methods	17
3.4.1	Belief-propagation matching	18
3.4.2	Semi-Global Matching	18
3.5	PatchMatch	19

3.5.1	PatchMatch Propagation	19
3.6	PatchMatch Stereo	19
3.6.1	Left-Right Consistency	20
3.6.2	Sub-pixel Accuracy	20
3.7	Beyond Classical Stereo Matching	20
3.7.1	Active Stereo	20
3.7.2	Projectors	21
3.7.3	Structured-Light Sensors	21
3.7.4	Temporal-Encoding	21
3.8	Stereo Matching with Deep Learning	21
3.8.1	Matching Cost with a CNN	21
3.8.2	MC-CNN Dataset	21
3.8.3	MC-CNN Loss Function	22
3.8.4	MC-CNN Implementation Details	22
3.9	Content-CNN	22
3.9.1	Content-CNN Loss Function	22
3.9.2	Content-CNN Implementation Details	23
3.9.3	3D CNN on Volumetric Data	23
3.10	GC-Net	23
3.10.1	GC-Net: Estimate Disparity	23
3.10.2	GC-Net Extensions	23
3.11	RAFT-STEREO	24
4	Structure From Motion	25
4.1	Correspondence Search	25
4.2	Geometric Verification	26
4.3	Scene Graph	26
4.4	Incremental Mapping	26
4.5	Initialization	26
4.6	Initial Traingulation	26
4.7	PnP-based Image Registration	27
4.8	Next Best View Selection	27
4.9	Bundle Adjustment	27
4.9.1	BA Complexity	28
4.10	Schur Complement Trick	28
4.11	Robust Cost Function	29
5	Visual SLAM	30
5.1	Direct Vs Indirect Methods	30
5.2	Matching Vs Tracking	30
5.3	Dense Vs Sparse	30
5.4	Visual SLAM Pipeline Modules	30
5.5	Visual SLAM Front-end Approaches	30
5.5.1	Epipolar Geometry	31
5.5.2	Probabilistic Filtering	31
5.5.3	Local Bundle Adjustment	31
5.5.4	ORB SLAM	32
5.5.5	BoW-based Place Recognition	33
5.5.6	Loop Correction	33
5.5.7	Pose Graph Optimization	33
5.5.8	Direct Approaches	33
5.6	DSO: Direc Sparse Odometry	34
5.6.1	DSO Photometric Error	34
6	3D Data Representation	35
6.1	Depth maps and Range Images	35
6.2	Point Clouds	35
6.2.1	Organized Point Clouds	35
6.3	Meshes	35
6.3.1	Manifolds	36
6.3.2	Triangular Mesh	36
6.3.3	Implicit Surfaces	36
6.4	Voxels	36

6.5	Octrees	36
6.6	Volumetric signed Distance Function	36
6.7	Surface Normals	36
6.7.1	Surface Normals from Point Clouds	36
6.7.2	PCA-Recup	37
6.8	K-D trees	37
6.8.1	Example of K-D Tree Construction	37
6.8.2	Example: K-D Tree NN Query	37

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:

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

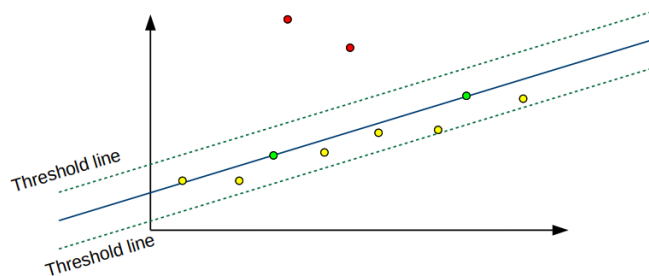
1.3.5 Image to Image Homography

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:

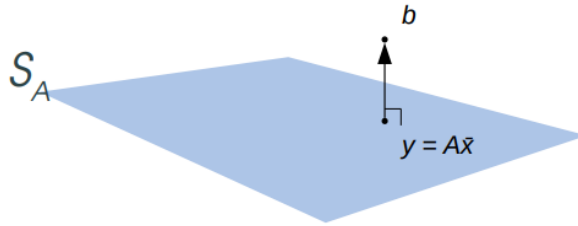
$$N = \frac{\log(1 - p)}{\log(1 - (1 - \epsilon)^s)}$$

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:

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

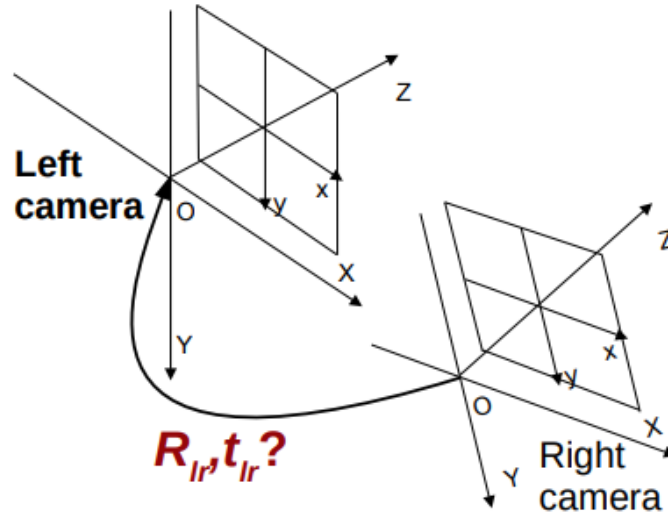
The ideas for the minimization: we use the LM algorithm to solve the iterative problem; we parametrize the rotations with **axis-angle representation**

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.

Chapter 2

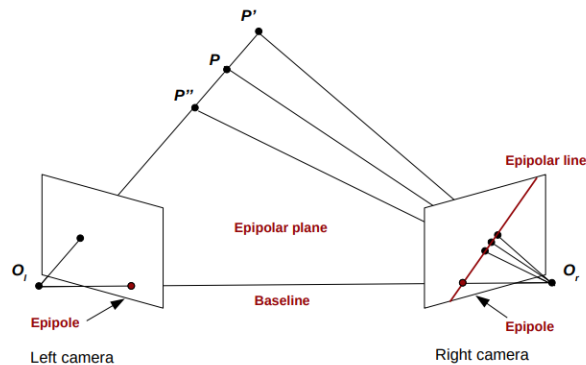
Epipolar Geometry



Stereo pairs or stereo rig is composed by two calibrated cameras mounted on a rigid support. Common practice is to choose the optical center of one camera to be the origin of the stereo rig coordinate system. To **calibrate** it we use a checkerboard and, for each i -th view, normalize points and extract R, t for both cameras. For each view find an estimate of R_{lr}, t_{lr} . From all views, compute an average R'_{lr}, t'_{lr} and set it as initial guess. Solve a non linear least square problem to compute a refined calibration $T_{lr} = R_{lr}, t_{lr}$ with **reprojection errors** computed in both cameras:

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

3D reconstruction from stereo. **Correspondence problem**: For a point p in the left image, find the corresponding point p' in the right image. **Reconstruction problem**: Given two corresponding points p and p' , find the 3D coordinates of corresponding scene point P . Since p_l, p_r and t_{lr} coplanar, we have that $\tilde{p}_l(t_{lr} \times R_{lr} \tilde{p}_r) = 0$.



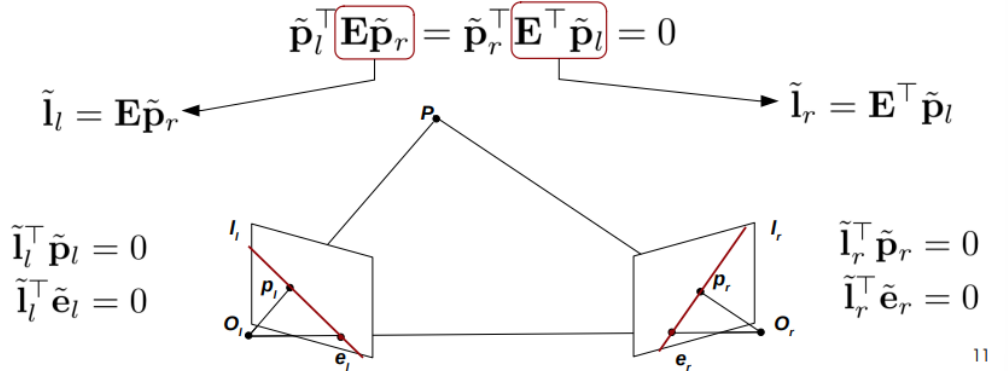
From there we can obtain the following constraint: $\tilde{p}_l^T \hat{t}_{lr} R_{lr} \tilde{p}_r = 0$, obtaining the essential constraint:

$$\tilde{p}_l^T E \tilde{p}_r = \tilde{p}_r^T E^T \tilde{p}_l = 0 \quad E = \hat{t}_{lr} R_{lr}$$

The essential matrix E encapsulates the rotation and translation associated with the relative pose of the two cameras.

2.1 Epipoles and Epipolar Lines

Remembering that in homogeneous coordinates line-point incidence can be verified as: $I^T p = p^T I = 0$ **The**



essential matrix is a projective mapping that maps a point onto a line.

2.1.1 Essential Test

To check if a match $p_l \leftrightarrow p_r$ is correct in a stereo rig compute:

$$\tilde{I}_l = E \tilde{p}_r$$

and verify if the distance between I_l and p_l is less than a certain threshold.

Essential matrices belong to a special set of matrices called Essential Space:

$$\epsilon = \{\hat{t}R | R \in S(3), t \in \mathbb{R}^3\} \subset \mathbb{R}^{3 \times 3}$$

A nonzero 3×3 matrix E is an essential matrix if and only if E has the following singular value decomposition:

$$E = U \Sigma V^T \quad \text{with} \quad \Sigma = \text{diag}\{\sigma, \sigma, 0\}$$

for some $\sigma \in \mathbb{R}_+$ and U, V in $SO(3)$. E is rank 2, i.e., **is not invertible**.

2.2 Pose Recovery from E

From SVD it is possible to compute up to four possible decompositions of E . By defining:

$$w = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

the solution is given by:

$$R = UWV^T \quad \text{or} \quad UW^T V^T \quad t = \pm u_3$$

To select a feasible solution, test with a single 3D point the cheirality constraint. When recovering position from the Essential matrix we get 4 possible solutions: use cheirality constraint to select one.

2.3 Fundamental Matrix

Map to a non distorted real camera:

$$\tilde{x}_l = K_l \tilde{p}_l \quad \tilde{x}_r = K_r \tilde{p}_r$$

from: $\tilde{p}_l^T E \tilde{p}_r = 0$, we have: $\tilde{x}_l K_l^{-T} E K_r^{-1} \tilde{x}_r = 0$

$$\tilde{x}_l F \tilde{x}_r = 0 \quad F = K_l^{-T} E K_r^{-1}$$

The fundamental matrix F encapsulates the rotation, translation and the intrinsic parameters. As the essential matrix, F maps points (pixels) to lines:

$$I' = Fx \quad I = F^T x'$$

Viceversa is not possible since F it has rank 2, i.e. is not invertible.

F has 7 DoF: like the homography is defined up to a scale factor but, being rank 2, is not invertible. The SVD is given as:

$$F = U \text{diag}(\sigma_1, \sigma_2, 0) V^T$$

$$x'^T F x = 0, \quad x = [x, y, 1]^T \quad \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \quad x' = [x', y', 1]^T$$

We obtain:

$$x'x f_{11} + x'y f_{12} + x' f_{13} + y'x f_{21} + y'y f_{22} + y' f_{23} + x f_{31} + y f_{32} + f_{33} = 0$$

One equation for each correspondence. Get n correspondences and arrange the following linear system:

$$\begin{bmatrix} x'_1 x_1 & x'_1 y_1 & x'_1 & y'_1 x_1 & y'_1 y_1 & y'_1 & x_1 & y_1 & 1 \\ x'_n x_n & x'_n y_n & x'_n & y'_n x_n & y'_n y_n & y'_n & x_n & y_n & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0 \Rightarrow Af = 0$$

As for H , we have to solve a homogeneous linear system, with F determined up to a scalar factor:

- $n=8$: minimal system, solution is the non trivial nullspace of A (**Eight-Point Algorithm**)
- $n > 8$: overdetermined system, solution with linear least squares.

In both case, solve with SVD taking the last column of V corresponding to the smallest singular value: again, **coordinates normalization is essential**.

In general the solution for F will not have zero determinant (i.e., rank 2). solution:

- Take the estimated F' , and compute its SVD, $F' = EDV^T$
- Defined D_n as being equal to D but with the smallest singular value set to 0
- Recompute $F = ED_n V^T$

2.4 Robust Fundamental Matrix Estimation

- Normalize data
- Repeat the following steps N times:
 - Select a set of eight correspondences at random
 - Estimate F solving the undetermined system $af=0$
 - Find the number of inliers that support F
 - * For each corresponding $\{p, p'\}$, map p to the corresponding epipolar line, i.e.: $I' = Hp$
 - * If the distance between p' and the line I' is less than a threshold, add the correspondence to the inliers associated to F
- Find the F with the highest number of inliers among the N trials
- Re-compute a least square estimate of F using all m inliers solving an overdetermined system (A $m \times 9$ matrix)
- Eventually enforce Rank 2 to F and denormalize F



2.5 Alternative Solutions

Seven-Point algorithm: Minimal solution: use seven correspondences and impose $\det(F)=0$.

- Use SVD and take the last two columns of V that define a 2-dimensional null-space of solutions: $F = \alpha F_1 + (1 - \alpha)F_2$
- The additional constraint $\det(\alpha F_1 + (1 - \alpha)F_2) = 0$ leads to a cubic polynomial equation in α which has 1 or 3 solutions: solve this equation, e.g. by exploiting some polynomial root-finding method

Non-linear approaches: often used for refinement. given an initial guess of F , setup a non-linear least squares problem where the squared distances between image points and epipolar lines are minimized.

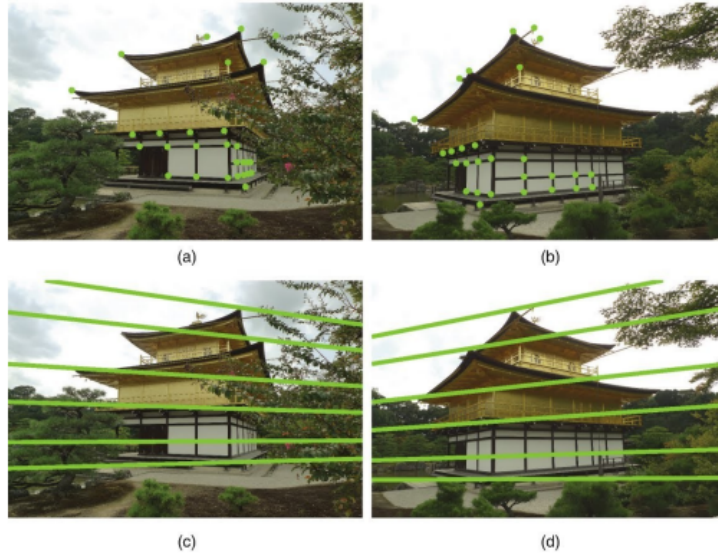
2.6 Essential Matrix Estimation

The **Eight-Point Algorithm** can be used also for computation of E , as long as:

- We use normalized coordinates in the correspondences;
- We enforce the affinity with the essential space:
 - Take the estimated E' , and compute its SVD, $E' = UDV^T$, with $D = \text{diag}\{\sigma_1, \sigma_2, \sigma_3\}$
 - Define $D_n = \{\sigma, \sigma, 0\}$, with $\sigma = (\sigma_1 + \sigma_2)/2$
 - Recompute $E = UD_nV^T$

The minimal solver is represented by the **Five-Point Algorithm**, that uses only five correspondences.

2.7 Stereo Rectification



We assume a calibrated, non distorted rig with some vergence. We wish to map the image points onto a pair of image planes that are parallel to the baseline and in the same plane. we aim to perform only **pure rotations** around the cameras: each pixel in the rotated version of the camera is mapped to an original camera pixel by means of a homography.

- Take a plane Π perpendicular to the baseline.
- project the unit vectors z_i and z_j of both optical axes into Π and take the bisector z_{ij}
- compute y_{ij} from x_{ij} and z_{ij} and compose:

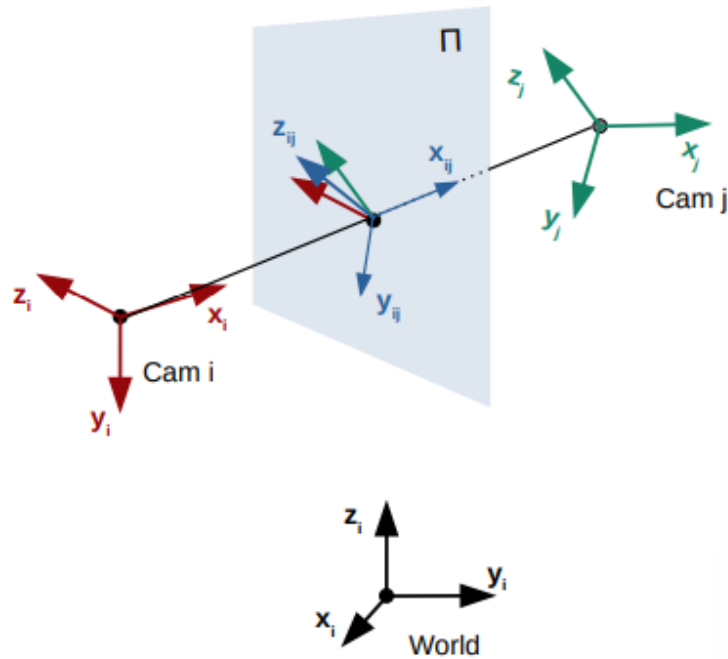
$$R_{ij} = (x_{ij} y_{ij} z_{ij})$$

- Rotate both cameras into their new viewing direction are as follows:

$$R_i^* = R_{ij} R_i^T \quad \text{and} \quad R_j^* = R_{ij} R_j^T$$

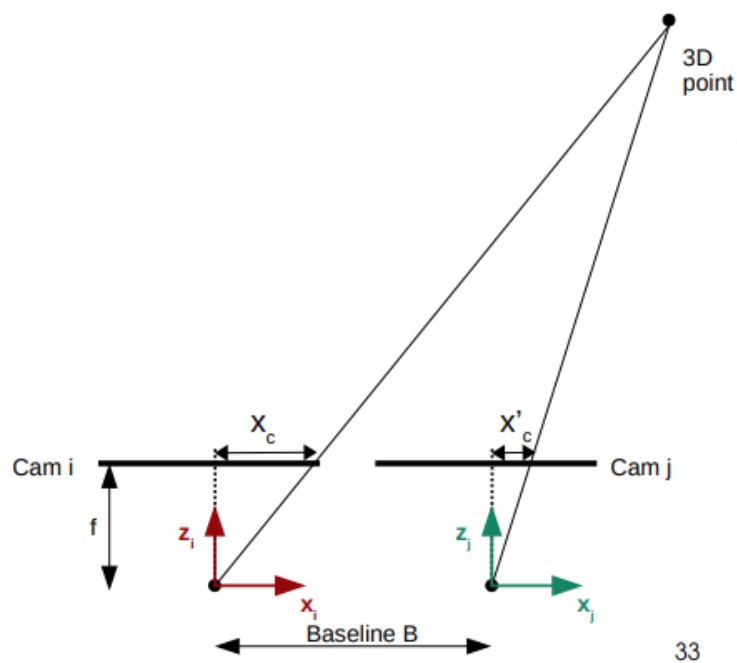
- Map to the original pixel are given by:

$$H = KRK^{-1}$$



2.7.1 Trinagulation in a Rectilinear Rig

The **Goal** is to get Z from stereo correspondences. For a rectilinear stereo rig, corresponding left and right image points lie in the same scanline. Two rays from the camera centers through the left and right image points can be back-projected in 3D. The two rays and the baseline lie on the epipolar plane.



$$x'_c = f \frac{X}{Z}, \quad x_c = f \frac{X+B}{Z}, \quad d = x_c - x'_c = \frac{fB}{Z}, \quad Z = \frac{fB}{d}$$

Chapter 3

Stereo Matching

Search in a pair of stereo images for corresponding pixels that are projections of the same 3D point. Epipolar constraints reduce the search to one dimension. Epipolar constraints reduce the search to one dimension.

Sparse Matching: extract salient points and compute descriptors. Match according to the minimum distance between descriptors. It is suitable for wide baseline stereo rigs.

Dense matching: Find correspondences for all points that are projected in both images.

The **assumptions** are that most scene points are visible from both viewpoints. Corresponding image regions are similar. We consider rectified rectilinear stereo rigs: the epipolar lines are the row of the images (**scanlines**).

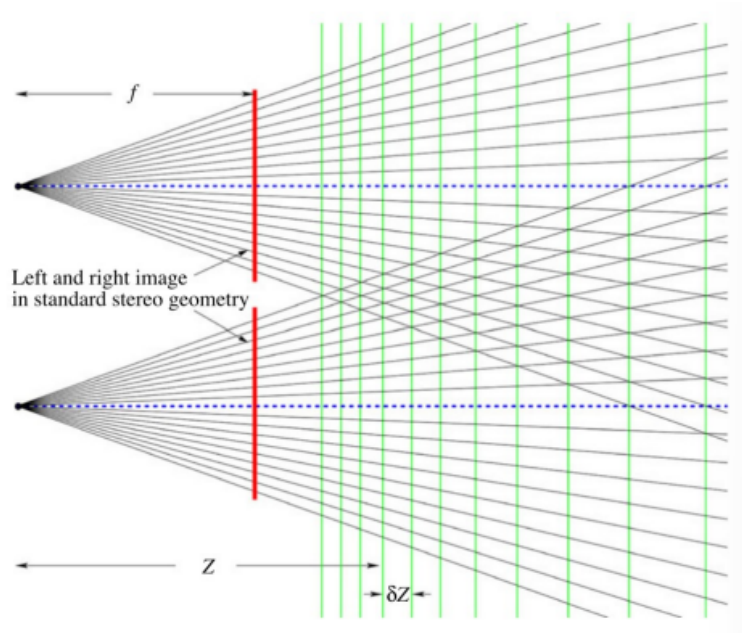
Challenges: Occlusion: some 3D points are projected in one image but not in the other; **Non-lambertian surfaces:** the appearance of the surfaces may depend on the point of view; **Perspective transformations:** a shape may appear different from a different point of view.

3.1 Matching Model

For each pixel (x,y) in the **base image B**, we search for a corresponding pixel in the **match image M**. We set a maximum disparity d and add a special value. Stereo matching as a labeling problem:

- Left base image: $L = \{+1, 0, -1, \dots, -d_{max}\}$ $0 \leq -d \leq \min\{d_{max}, x - 1\}$
- Right base image: $L = \{-1, 0, +1, \dots, d_{max}\}$ $0 \leq d \leq \min\{d_{max}, N_{cols} - x\}$

3.2 Depth uncertainty



A disparity of 0 between corresponding pixels $x_L - x_R = 0$ means 3D point at infinity. The distances δZ between subsequent depth layers increases nonlinearly.

3.3 Block Matching



For each pixel, slide a **support window** along the match scanline for the whole search interval and compare the pixels inside the window with corresponding pixels in the window of the base image.

3.3.1 Matching cost

- Sum of squared differences:

$$E_{SSD}(p, d) = \sum_{i=-l}^l \sum_{j=-k}^k [B(x+i, y+j) - M(x+d+i, y+j)]^2$$

- Sum of absolute differences:

$$E_{SAD}(p, d) = \sum_{i=-l}^l \sum_{j=-k}^k |B(x+i, y+j) - M(x+d+i, y+j)|$$

- Zero-mean normalized cross correlation:

$$E_{ZNCC}(x, d) = 1 - \frac{\sum_{i=-l}^l \sum_{j=-k}^k [B_{x+i, y+j} - \bar{B}_x][M_{x+d+i, y+j} - \bar{M}_{x+d}]}{\sqrt{\sigma_{B,x}^2 \sigma_{M,x+d}^2}}$$

$$\sigma_{B,x}^2 = \sum_{i=-l}^l \sum_{j=-k}^k [B_{x+i, y+j} - \bar{B}_x]^2 \quad \sigma_{M,x+d}^2 = \sum_{i=-l}^l \sum_{j=-k}^k [M_{x+d+i, y+j} - \bar{M}_{x+d}]^2$$

- Census transform:
 - Binary number: one bit for each sample in the Census window
 - Bit value: comparison against the central value
 - Compare the descriptor using Hamming distance

According to the window size we have different result: smaller window brings high accuracy but more noise, larger window bring higher reliability but less details.

Local cost calculation is generally ambiguous because of textureless surfaces, occlusions, repetitions and non-lambertian surfaces.

To improve the block matching we can use adaptive windows: start with a standard window size and then adapt the window according to the current disparity estimate and iterate. We could also use hierarchical stereo matching, or exploit non-local correspondence constraints (**Global Methods**).

3.4 Global Methods

Enforce smoothness constraint, i.e. disparity is **piecewise smooth**. Enforce ordering constraints and greater computational complexity.

3.4.1 Belief-propagation matching

$$E(f) = \sum_{p \in \Omega} \left[E_{data}(p, f_p) + \sum_{q \in A(p)} E_{smooth}(f_p - f_q) \right]$$

Corresponding pixels should have a similar intensity. Neighboring pixels in the image should have a similar disparity space, e.g.:

$$E_{smooth}(l - h) = E_{smooth}(a) = \begin{cases} 0 & \text{if } a = 0 \\ c & \text{otherwise} \end{cases}$$

$$E_{smooth}(l - h) = E_{smooth}(a) = b|l - h| = b|a|$$

$$E_{smooth}(l - h) = E_{smooth}(a) = \min\{b|l - c|, c\} = \min\{b|a|, c\}$$

Usually 4-pixels adjacency set. the set of all possible labelings has a cardinality $(d_{MAX})^{w \times h}$: intractable exhaustive solution. Solve by using efficient message-passing algorithms. It is accurate but very slow.

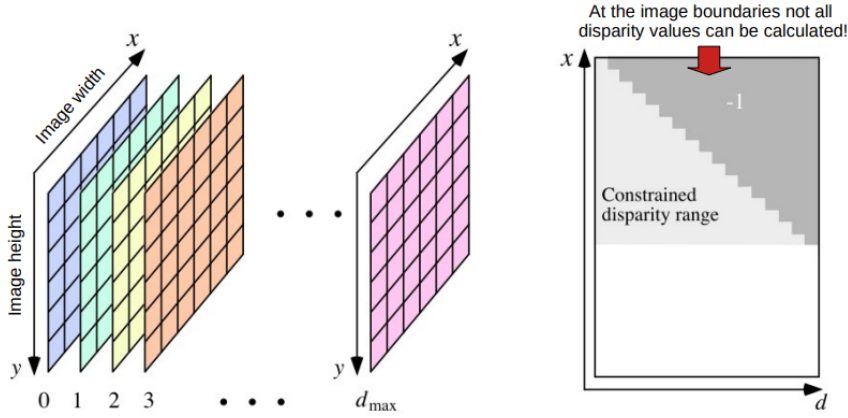


Figure 3.1: Cost volume

Area of influence

Set used by the stereo matcher for each pixel to decide about its label via the smoothness constraint: in global methods we use the whole image. the idea is to reduce the area of influence, e.g. by considering only some directions.

3.4.2 Semi-Global Matching

Consider the following error function:

$$E(f) = \sum_{p \in \Omega} [E_{data}(p, f_p) + \sum_{q \in A(p)} E_{smooth}(f_p, f_q)]$$

with:

$$E_{smooth}(f_p, f_q) = \chi_1 f(f_p, f_q) + \chi_2(f_p, f_q)$$

with

$$\chi_1(f_p, f_q) = \begin{cases} c_1 & \text{if } |f_p - f_q| = 1 \\ 0 & \text{otherwise} \end{cases} \quad \chi_2(f_p, f_q) = \begin{cases} c_2(p, q) & \text{if } |f_p - f_q| > 1 \\ 0 & \text{otherwise} \end{cases}$$

This method approximate global methods by aggregating costs for a number of directions. the minimization along individual image rows can be performed efficiently in polynomial time using Dynamic programming.

- For each direction, start from one end point and go toward p
- For each pixel along the direction, update the following dynamic programming equation:

$$E(p_i, d) = E_{data}(p_i, d) + E_{smooth}(p_i, p_{i-1}) - \min_{0 \leq \Delta \leq d_{max}} E(p_{i-1}, \Delta)$$

where

$$E_{smooth}(p, q) = \min \begin{cases} E(q, f_q) & \text{if } f_p = f_q \\ E(e, f_q) + c_1 & \text{if } |f_p - f_q| = 1 \\ \min_{0 \leq \Delta \leq d_{max}} E(q, \Delta) + c_2(p, q) & \text{if } |f_p - f_q| > 1 \end{cases}$$

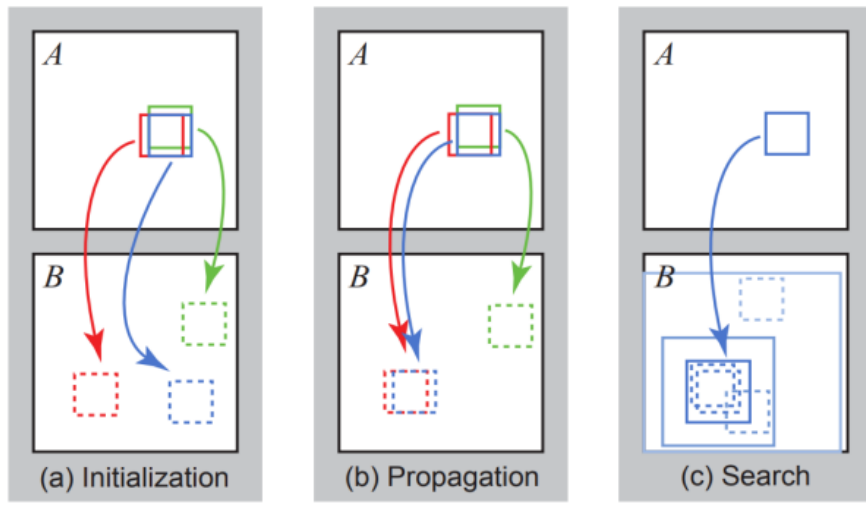
For each direction i (start with $E(p_0, d) = E_{data}(p_0, d)$). Add the last column of each integration matrix associated to a specific direction and **get the minimum**.

This method is really efficient, suitable for real-time applications. Basic algorithm performs badly in large untextured regions.

3.5 PatchMatch

Problem definition: given images A and b, find for every patch in A the nearest neighbor in B under some patch distance metric: Nearest-Neighbor Field (NNF): holds dx, dy **offsets**.

PatchMatch: a randomized algorithm for rapidly finding correspondences between image patches exploiting **spatial coherence** between images.



3.5.1 PatchMatch Propagation

Scan each patch in row major order, from top-left pixel:

- We attempt to improve $f(x, y) = (dx, dy)$ using the known offset of $f(x-1, y)$ and $f(x, y-1)$, assuming that the patch offsets are likely to be the same
- Let $D(p, f)$ denote the patch distance between the patch at $p = (x, y)$ in A and patch $p + f = (x, y) + (dx, dy)$ in B. We take the new value for $f(x, y)$ to be:

$$\arg \min_f \{D(p, f(x, y)), D(p, f(x-1, y)), D(p, f(x, y-1))\}$$

Repeat the scan in inverse order, from bottom-right pixel

We can summarize the PatchMatch algorithm as: repeat until convergence: initialization, forward and backward propagation, then random search.

3.6 PatchMatch Stereo

A support window is typically centered on a pixel of the base image, but:

- The support window constrains pixels that lie on a different surface than the center pixel
- The window captures a surface that is slanted

PatchMatch Stereo tries to address the second problem, by estimating a 3D local plane that approximates the slanted surface at each image point: Hence estimating a better support window in the match image.

For each pixel p of both frames, we search a plane f_p . Disparity of p is computed as:

$$d_p = a_{f_p} p_x + b_{f_p} p_y + c_{f_p}$$

where a_f, b_f, c_f are the plane parameters being estimated, p_x, p_y are the pixel coordinates. The cost for matching pixel p with pixel q in the other image according to plane f is computed as:

$$m(p, f) = \sum_{q \in W_p} w(p, q) \rho(q, q - (a_f q_x + b_f q_y + c_f))$$

where:

- w_p is the square support window centered on pixel p
- $\rho(q, q')$ is the function that computes the pixel dissimilarity
- $w(p, q)$ is the weighting function

The algorithm is the following:

- RRandom initialization
- Repeat until convergence:
 - **spatial propagation**: we evaluate whether assigning to p the plane f_q of spatial neighbor pixel q decrease the matching cost. If $m(p, f_q) < m(p, f_p)$, accept the new plane
 - **View propagation**: we check all pixels p' of the second view that have our current pixel p as a matching point according to their current plane. if $m(p, f_p) < m(p', f_p)$ accept the new plane.
 - **Plane refinement**: sample a new plane f_{rand} around the current one. If $m(p, f_{rand}) < m(p, f_p)$ accept the new plane.

It is a good trade-off between accuracy and computational complexity.

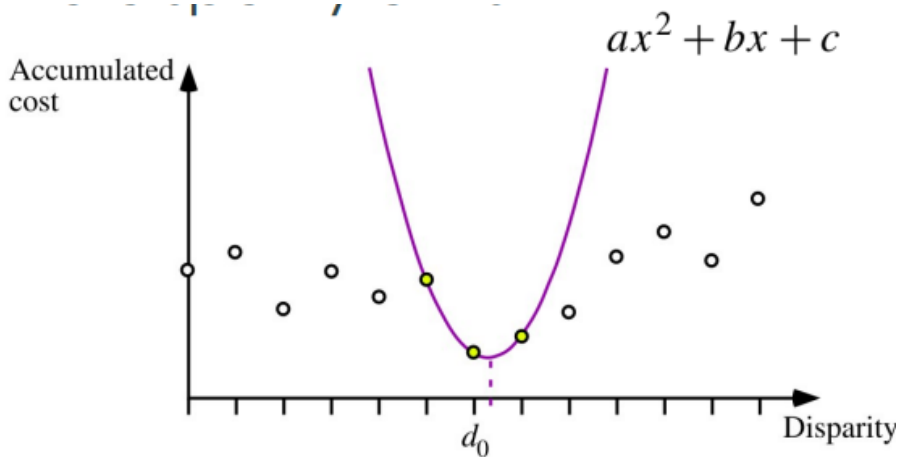
3.6.1 Left-Right Consistency

Left-to-right matching and then right-to-left matching, and compare disparities. Used also as a confidence measure:

$$\Gamma_1(p) = \frac{1}{|f_p^{(R)} - f_p^{(L)}| + 1}$$

3.6.2 Sub-pixel Accuracy

Fit a parabola to accumulate costs in a neighbourhood of the disparity d_0 , and take the disparity at its minimum.



3.7 Beyond Classical Stereo Matching

3.7.1 Active Stereo

Add a projector that produces an highly textured light pattern, enhancing the texture of the scene. Projector-camera extrinsic parameters are not required. Usually a Near Infrared (NIR) projector is used: measurement cameras are sensitive to NIR and a "texture" camera equipped with IR-cut filter can be used to acquire RGB information.



3.7.2 Projectors

Can be of different types:

- Coherent laser light + diffractive optical elements
- Light projector + mask
- Digital Light Processing (DLP) projectors

3.7.3 Structured-Light Sensors

Conceptually a passive stereo system where one of the measurement cameras is replaced by a **calibrated projector**. Projector-camera extrinsic parameters are required and the knowledge of the spatial and temporal structure of the pattern is required.

3.7.4 Temporal-Encoding

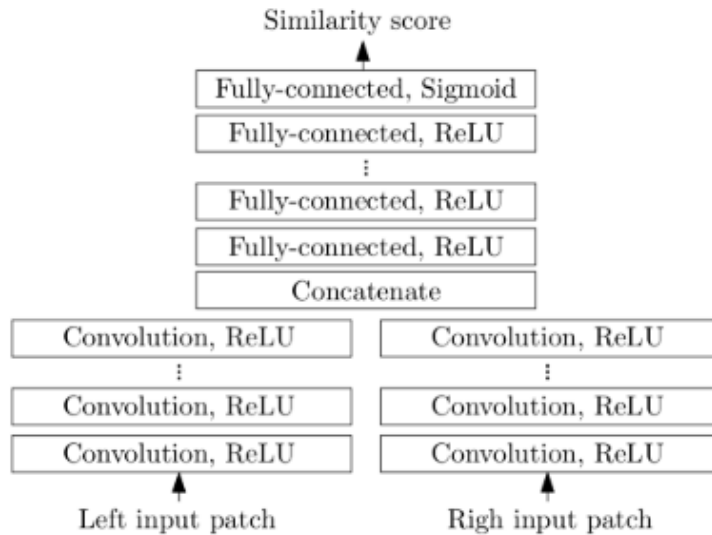
For reasons of error reduction, it is recommended that the projected slides follow the **Gray code** rather than the usual binary code: consecutive integers are represented by binary numbers that differ in one digit only.

<i>Integers</i>	0	1	2	3	4	5	6	7
<i>Usual code</i>	000	001	010	011	100	101	110	111
<i>Gray code</i>	000	001	011	010	110	100	101	111

3.8 Stereo Matching with Deep Learning

3.8.1 Matching Cost with a CNN

MC-CNN: given two small patches, compute the matching cost with a Convolutional Neural Network. A set of examples is required.



3.8.2 MC-CNN Dataset

A set of positive and negative examples:

- Each example is a pair of **pair of patches**, one from the left and one from the right image:

$$\langle \mathcal{P}_{n \times n}^L, \mathcal{P}_{n \times n(\mathbf{q})}^R \rangle$$

- A positive example is obtained by setting the right patch close to the correct match (o_{pos} "small")

$$\mathbf{q} = (x - d + o_{pos}, y)$$

- A negative example is obtained by randomly setting the right patch far from the correct match (o_{neg} "big")

$$\mathbf{q} = (x - d + o_{neg}, y)$$

3.8.3 MC-CNN Loss Function

We use the binary cross-entropy:

$$t \log(s) + (1 - t) \log(1 - s)$$

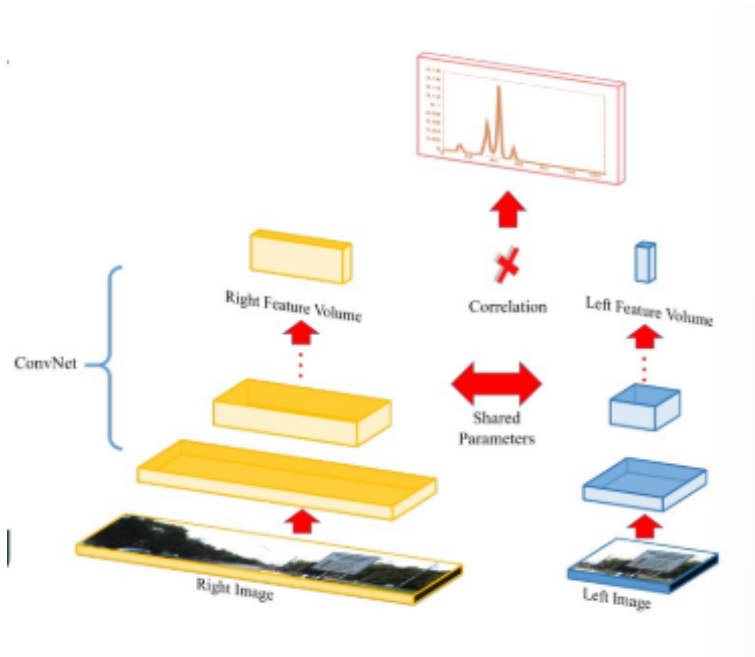
where s is the output of the network for one training example and t is the example label (1 if it is positive class, 0 otherwise).

3.8.4 MC-CNN Implementation Details

MC-CNN don't use the raw matching cost, but calculates the average of the costs over a fixed window of neighboring pixels. Use Semiglobal Matching to enforce smoothness constraints. Perform sub-pixel enhancement and smooth the disparity map by using bilateral filter: Smooth without blurring the edges.

3.9 Content-CNN

Inspired by MC-CNN but without concatenation: just a dot product. Computes a 64-dimensional feature representation for every pixel of both images. For each possible disparity, match each pixel by calculating the dot product between corresponding **features**.



3.9.1 Content-CNN Loss Function

We train the network with cross entropy: for all possible patches (i) and all possible disparity values (y_i)

$$\min_w \sum_{i, y_i} p_{gt}(y_i) \log p_i(y_i, \mathbf{w})$$

$$p_{gt}(y_i) = \begin{cases} \lambda_1 & \text{if } y_i = y_i^{GT} \\ \lambda_2 & \text{if } |y_i - y_i^{GT}| = 1 \\ \lambda_3 & \text{if } |y_i - y_i^{GT}| = 2 \\ 0 & \text{otherwise} \end{cases}$$

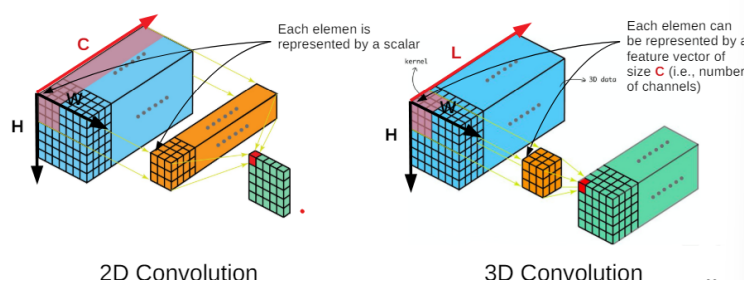
where $\lambda_1 > \lambda_2 > \lambda_3$

3.9.2 Content-CNN Implementation Details

The cost aggregation: average pooling over a 5x5 window. We use Semiglobal Matching to enforce smoothness constraints. Left-right consistency check and subpixel enhancement.

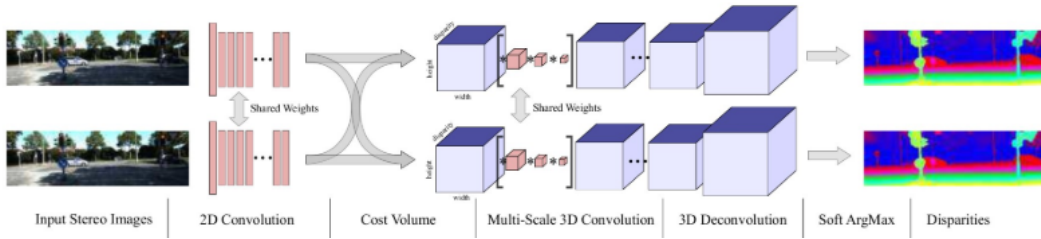
3.9.3 3D CNN on Volumetric Data

3D convolution uses possibly 4D kernels.



3.10 GC-Net

GC-Net: **end-to-end** mapping from an image pair to disparity. It uses 3D convolutions to learn the cost volume. Regress sub-pixel disparity values from the disparity cost volume. Assemble a feature volume by



concatenating corresponding left and right features on each disparity.

We aim to compute a $W \times H \times D$ cost volume: we use 3-D convolutions to incorporate context and we extend the encoder-decoder idea to 3D, in order to increase feature's receptive field while reducing computation.

3.10.1 GC-Net: Estimate Disparity

Classical Argmin operator over the cost volumes has two problem: no sub-pixel disparities and it is not differentiable. GC-Net propose to use a soft argmin:

$$\text{soft arg min} := \sum_{d=0}^{D_{max}} d \times \delta(-c_d)$$

where

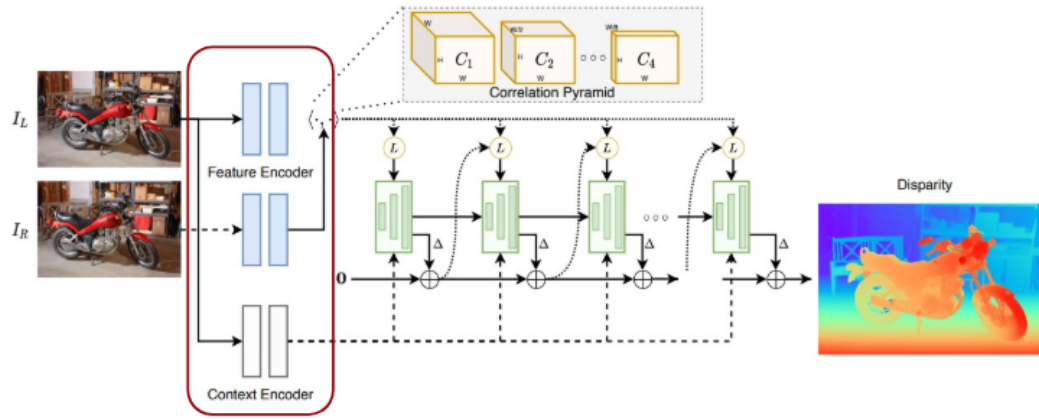
$$\delta(x_i) = \frac{\exp^{x_j}}{\sum_i \exp^{x_i}}$$

that is a weighted average of all possible softmax costs. The loss function is:

$$Loss = \frac{1}{N} \sum_{n=1}^N \|d_n - \hat{d}_n\|_1$$

3.10.2 GC-Net Extensions

SsSMNet (Self-Supervised Stereo Matching) minimize the warping error without using ground truth disparities. It is online and has self-improving capabilities.



3.11 RAFT-STEREO

RAFT-Stereo: Multilevel Recurrent Field Transform for Stereo Matching uses only 2D convolutions and, also it is trained on synthetic data, and provides better generalization and accuracy. It is composed by a feature extractor, a feature correlation pyramid and a Gated Recurrent Unit (GRU) applied several times. Set of residual blocks and downsampling layers producing a feature per pixel. The cost volumes obtained by calculating the scalar products between the pairs of pixels corresponding to all disparities, "downscaled" in several pyramid levels.

The current estimate of disparity are used to index the correlation volume, producing a set of correlation features. These features, the current disparity map, and context features are then concatenated and injected into the GRU, that iteratively refine the disparity map.

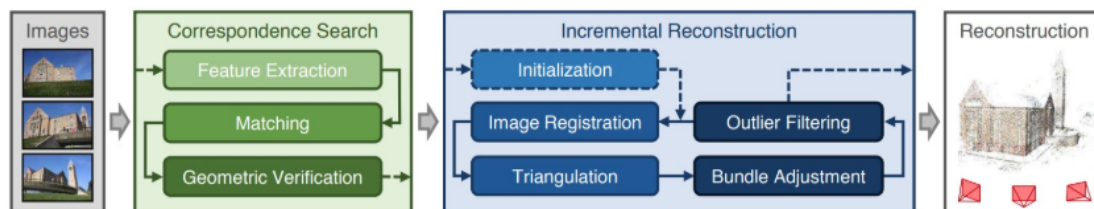
Chapter 4

Structure From Motion

Structure form motion is the process of reconstructing the 3D structure of a scene from its projections into a series of images taken from different viewpoints.

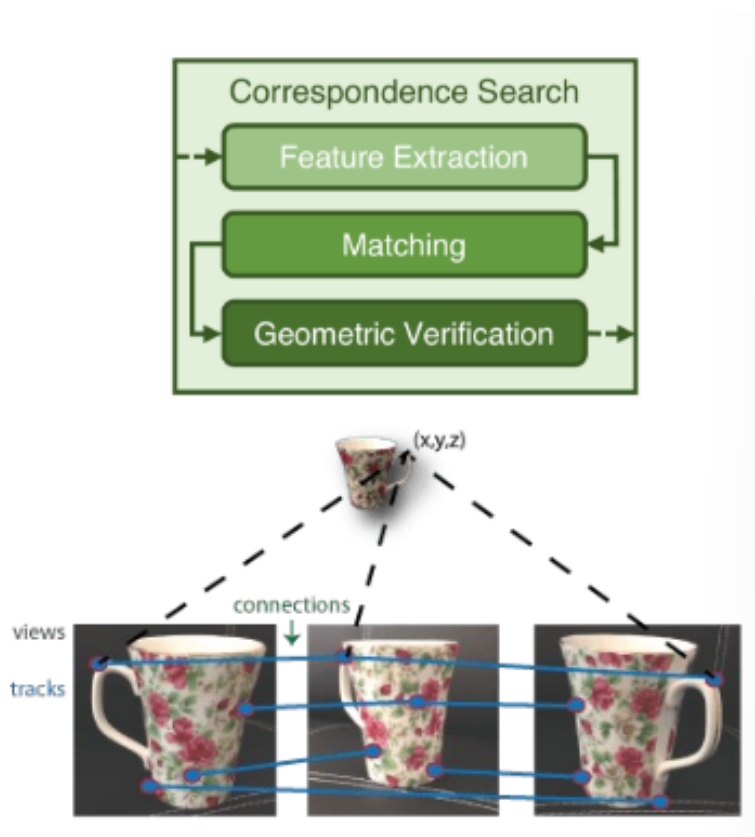
Input: a set of images that frame the same scene from different points of view. Each image should have some field of view overlap with other images: images can be totally unordered and taken with different cameras. The camera calibration parameters may not be available.

Output: camera positions, camera calibrations, and sparse 3D scene structure. Incremental SfM is a sequential



processing pipeline with an iterative reconstruction component.

4.1 Correspondence Search





For each image extract invariant features with associated descriptors. Match features between every image pair: this is just a naive approach since an exhaustive matching has computational complexity $O(N_I^2 N_{F_i}^2)$. Many efficient approaches have been introduced to improve matching.

4.2 Geometric Verification

Try all models:

- Essential matrix E
- Fundamental matrix F
- Homography matrix H

For each model, use RANSAC to estimate the best estimate: Sample N minimal sets and find the maximum number of inliers that support the model. Select the best model E, F or H based on certain criteria: its inliers define the inlier matches between an image pair.

4.3 Scene Graph

After geometric verification, we can build a scene graph with images as nodes and verified pairs of images as edges. It provides correspondences between a pair of images, number of correspondences per image, number of observations in an image and the correspondence of an observation with all other images.

4.4 Incremental Mapping

Iterative try and error process. Reconstruction may never recover from a bad initialization and also from a bad image registration.

4.5 Initialization

Find a good initial pair, possible basic algorithm:

- Sort images by decreasing number of correspondences, i.e., images with more correspondences appear in the front of the list.
- Select the first image in the list as the first seed image.
- Complete the seed pair with an image that is connected to the first seed image and has not been selected before.

Try to recover a rigid body transformation between the seed pair images, from the best model selected during geometric verification: in the uncalibrated case, this most likely leads to a ill-defined reconstruction. If failed, restart the initialization process. If the motion is a pure rotation estimated from H, restart the initialization process: No triangulation can be obtained from a pure rotation.

4.6 Initial Triangulation

Using point correspondences and the initial R,t triangulate point to obtain their 3D coordinates.

For a non rectilinear stereo rigs, no single solution due to noise. We have R,t and the calibrations parameters: given a 3D point, we know how to project it in both views! To use our points coordinates x, we need to normalize the homogeneous points:

$$x = \alpha P \tilde{X}$$

Cross product of two vectors of same direction is zero:

$$x \times PX = 0$$

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \alpha \begin{bmatrix} p_1 & p_2 & p_3 & p_4 \\ p_5 & p_6 & p_7 & p_8 \\ p_9 & p_{10} & p_{11} & p_{12} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \alpha \begin{bmatrix} p_1^T X \\ p_2^T X \\ p_3^T X \end{bmatrix} \Rightarrow x \times PX = 0$$

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \times \begin{bmatrix} p_1^T X \\ p_2^T X \\ p_3^T X \end{bmatrix} = \begin{bmatrix} yp_3^T X - p_2^T X \\ p_1^T X - xp_3^T X \\ xp_2^T X - yp_1^T X \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

The third line is a linear combination of the first and second lines.

$$\begin{bmatrix} yp_3^T - p_2^T \\ p_1^T - xp_3^T \end{bmatrix} X = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Use both points:

$$\begin{bmatrix} yp_3^T - p_2^T \\ p_1^T - xp_3^T \\ y'p_3^T - p_2^T \\ p_1^T - x'p_3^T \end{bmatrix} X = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow AX = 0$$

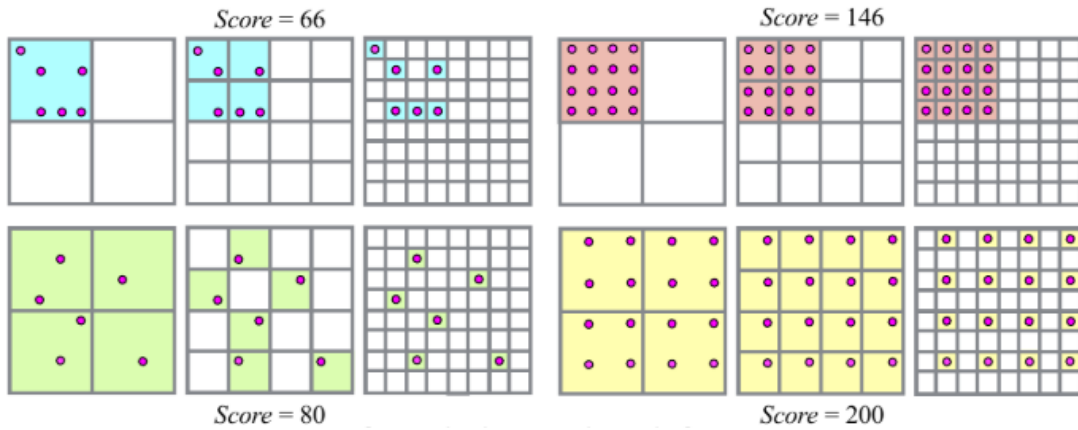
Is a homogeneous linear system. Use SVD and take the last column of V.

4.7 PnP-based Image Registration

A newly registered image must observe existing scene points. Use Perspective-n-Point using feature correspondences to triangulated points in already registered iamges. Estiamte both R', t' and intrinsic parameters. Use a modified version of PnP that exploits RANSAC and a minimal pose solver.

The newly registered image may also increase scene coverage by extending the set of points X trhough triangulation.

4.8 Next Best View Selection



Several criteria to select a still not registered image, among others:

- Sees the maximum of already registered points, or;
- Use visivility score.

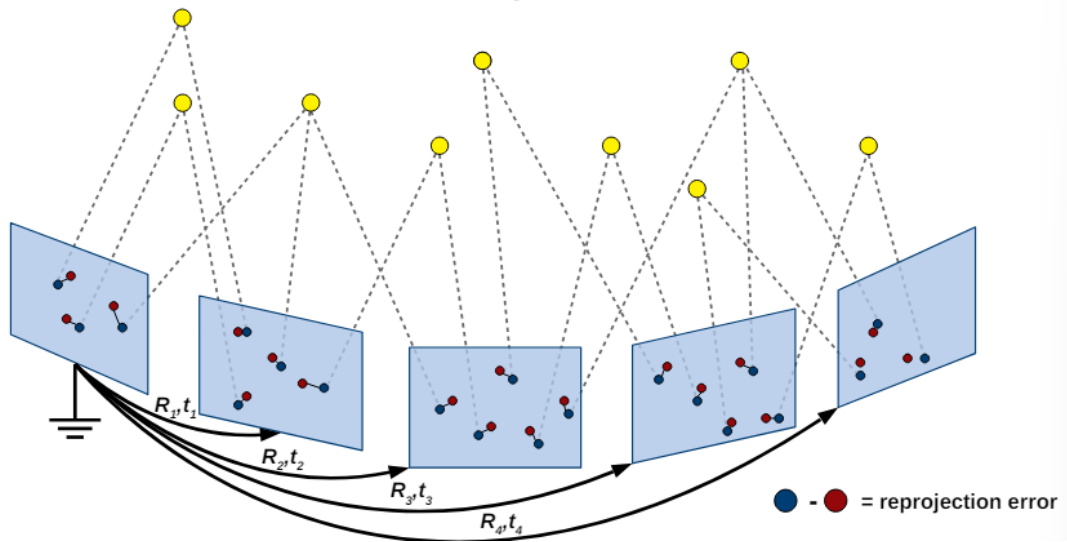
4.9 Bundle Adjustment

We may iterate with PnP and triangulation steps to register new images and new points but they are very correlated procedures: uncertainties in the camera pose propagate to triangulated points and vice-versa. **Bundle Adjustment (BA): joint non-linear refinement of camera parameters and point parameters.** BA goal: find a suitable parameters set (transformations $T_i = R_i, t_i$, calibration parameters θ_i , 3D points positions $P_j, i = 0, \dots, N_c - 1, j = 0, \dots, N_p - 1$) that minimizes the sum of squares of reprojection errors (one for each observation pf each camera):

$$\arg \min_{\theta_0, \dots, \theta_{N_c}, T_1, \dots, T_{N_c-1}, P_0, \dots, P_{N_p-1}} \sum_{i=0, j=0}^{N_c-1, N_0(i)-1} \rho(\|K_i f_i(A_n T_i \tilde{P}_{Obs(i,j)}) - \tilde{p}_{i,j}\|^2)$$

Camera positions ($\# = 6(N_c - 1)$):

- Set the coordinate system of one of the two seed cameras as reference **world frame**



- Use a suitable rotation representation

Calibration parameters (if not given) ($\# = 3N_c$):

- Fix the principal point at the image center (principal point calibration is an ill-posed problem)
- Optimize for each camera the focal length in pixel and some distortion parameters

Scene points ($\# = 3N_p$): simply parametrized as 3D points.

4.9.1 BA Complexity

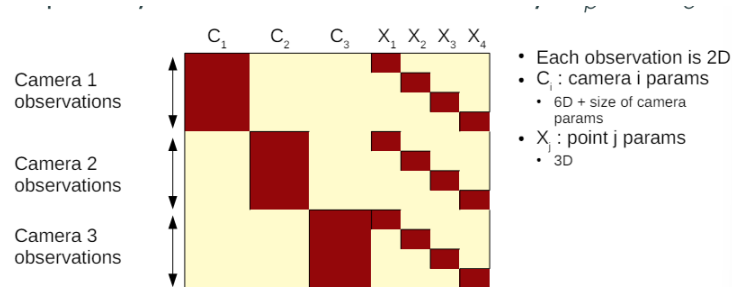
LM update:

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

with

$$J = \begin{bmatrix} \frac{\partial e_1}{\partial a_1} & \cdots & \frac{\partial e_1}{\partial a_m} \\ \vdots & \vdots & \vdots \\ \frac{\partial e_n}{\partial a_1} & \cdots & \frac{\partial e_n}{\partial a_m} \end{bmatrix} \quad e_k = K_i f_i(A_n T_i \tilde{P}_{Obs(i,j)}) - \tilde{p}_{i,j}$$

4.10 Schur Complement Trick



The idea is to efficiently solve BA: exploit J sparsity and the fact that usually $N_p \gg N_c$. We can rewrite the LM step as:

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

and let define $H = J^T J + \lambda I \Rightarrow H \delta a^{LM} = -g$.

Partitioning the step update between cameras and structure:

$$\begin{bmatrix} C & E \\ E^T & P \end{bmatrix} \begin{bmatrix} \delta a_c \\ \delta a_p \end{bmatrix} = \begin{bmatrix} v \\ w \end{bmatrix}$$

Multiply both sides by:

$$\begin{bmatrix} I & -EP^{-1} \\ 0 & I \end{bmatrix} \Rightarrow \underbrace{C - EP^{-1}E^T}_{\text{Schur complement of P}} \delta a_c = v - EP^{-1}w$$

P block diagonal matrix: calculating the inverse of P by inverting each of its 3×3 blocks is cheap. Solve for δa_c then backsubstituting it to obtain the value of δa_p . It has complexity of $O(N_c^3)$.



4.11 Robust Cost Function

Standard least square assumes gaussian errors: high influence from high errors. To account for potential outliers, it is better to choose a proper loss function ρ .

Chapter 5

Visual SLAM

Visual Simultaneous Localization and Mapping has the goal of estimating the camera trajectory while reconstructing the environment in real-time. Visual odometry focus on incremental estimation/local consistency, while Visual SLAM focus on globally consistent estimation.

SfM is more general than Visual SLAM and VO and tackles the problem of 3D reconstruction of both the structure and camera poses from possibly **unordered images** sets taken from possibly different and/or **un-calibrated cameras**. Visual SLAM is a particular case of SfM:

- it focuses on the **sequential, real-time estimation** of the 3D motion and the 3D structures as a new frame arrives, by using a typically calibrated camera.
- Used for **robot navigation**.

5.1 Direct Vs Indirect Methods

Indirect methods:

- Images are pre-processed to generate an intermediate representation by means of point correspondences
- Correspondences are interpreted as noisy measurements inside an estimation problem.

Direct methods: Skips the pre-processing step and directly use the actual sensor values, i.e. light received from a certain direction over a certain time period.

5.2 Matching Vs Tracking

Matching: extract features for every image, compute a descriptor for each feature, and match them between images looking for similar descriptors.

Tracking: extract features in the first image; for each feature, look for a similar patch in a neighborhood on the next image; update the feature location and repeat from the previous step.

5.3 Dense Vs Sparse

Sparse methods use and reconstruct only a selected set of independent, salient points. Dense methods attempt to use and reconstruct all pixels in the 2D image domain. Semi-dense methods use and reconstruct largely connected and well-constructed subset of pixels.

5.4 Visual SLAM Pipeline Modules

The visual SLAM front-end is composed by position tracking and Local Mapping, while Visual SLAM back-end is composed by Loop closure detection and Global Optimization.

5.5 Visual SLAM Front-end Approaches

We can use Epipolar geometry, filtering, Local bundle adjustment and direct approaches.



5.5.1 Epipolar Geometry

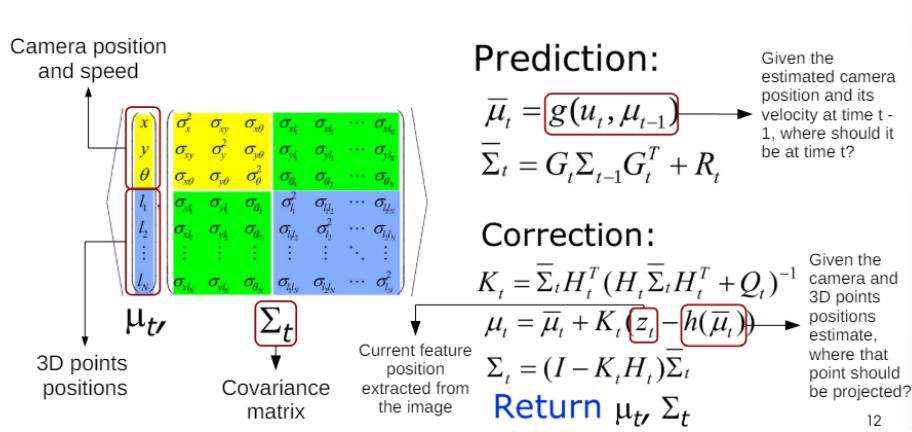
Extract and match a set of salient features. It infers ego-motion using the epipolar geometry (essential, fundamental or homography matrix) inside a sample consensus framework. We estimate 3D features positions using triangulation. We track consecutive motions with PnP and/or by exploiting epipolar geometry constraints + scale registration.

5.5.2 Probabilistic Filtering

Define a state that includes camera position and features locations, along with other components used in the kinematic equations. For each image:

- Track a set of features;
- Update the state and **its uncertainty** using the kinematic equations;
- Predict the position of the tracked features after this motion update and compute the re-projection errors.
- "Correct" the state, e.g. using probabilistic tools.

Visual SLAM with Kalman Filter



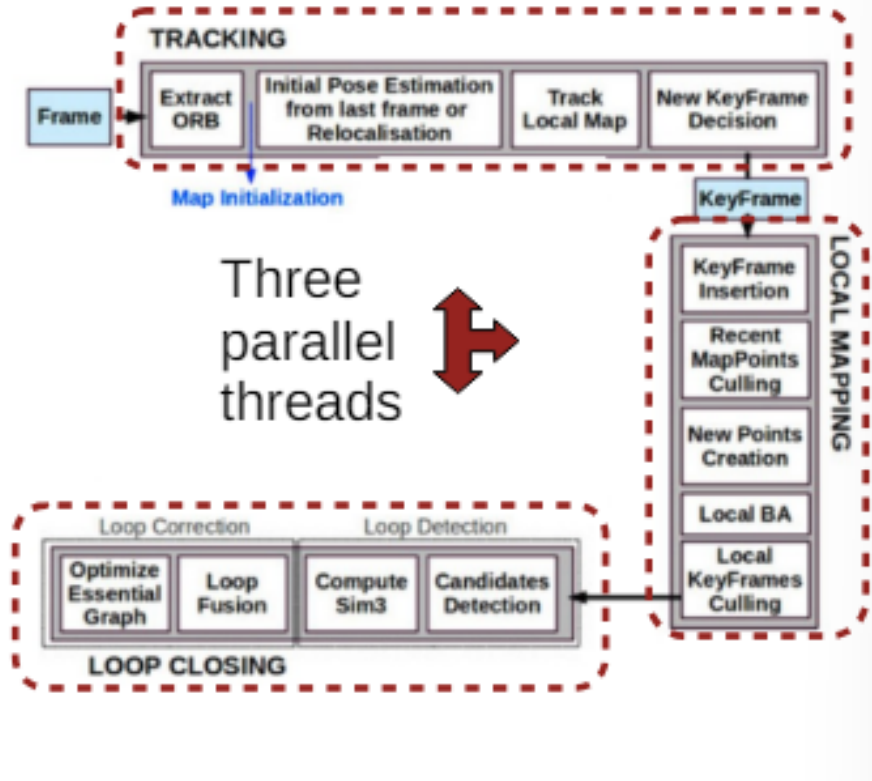
- Initialize first motion and first 3D points;
- Repeat:
 - **Predict** the camera position at time t given the position and velocity estimates at time $t-1$
 - **Correct** the estimation by integrating features tracks;
 - Newly observed points should be integrated in the state, points that are no longer visible should be removed.

The challenges of visual SLAM with Kalman Filters are that incorrect tracks lead the filter to divergence. We can track only up to a few hundred points. With a single camera the estimated motions are defined up to a scale factor, in addition we may have cumulative scale drift: a possible solution is to use stereo vision or to add an IMU.

5.5.3 Local Bundle Adjustment

Initialize first motion and first 3D points. For each image:

- **Matching:** estimate feature positions and descriptors and match them with a previous **keyframe** (i.e., an image whose descriptors are stored for mapping and closure tasks)
- **Motion tracing:** estimate the current camera position by solving a simplified bundle adjustment problem that only optimizes the camera position (i.e., 3D points positions are fixed)
- **Local Mapping:** when required, select a new keyframe and estimate the camera position along the 3D points positions by solving a bundle adjustment problem over a fixed sequence of k newest keyframes.



5.5.4 ORB SLAM

ORB SLAM uses ORB features for position tracking, local mapping and loop-closure detection. ORB features for each **keyframe** are kept in memory for mapping and loop-closure. Motion tracking is run at each frame and local mapping is run when a new keyframe has been selected. Loop closure detection is tested at each keyframe. We can define two problems: Motion-Only Bundle Adjustment:

$$\{R, t\} = \arg \min_{R, t} \sum_{i \in \mathcal{X}} \rho \left(\left\| x_{(\cdot)}^i - \pi_{(\cdot)}(RX^i + t) \right\|_{\Sigma}^2 \right)$$

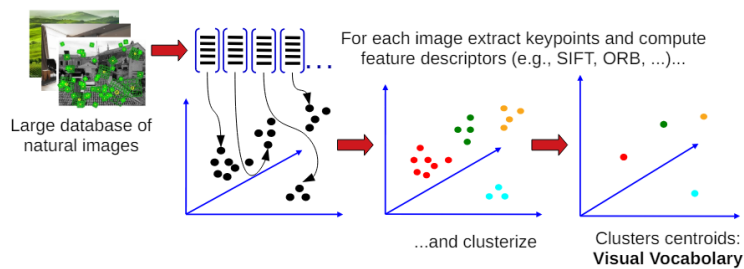
or Local Bundle Adjustment:

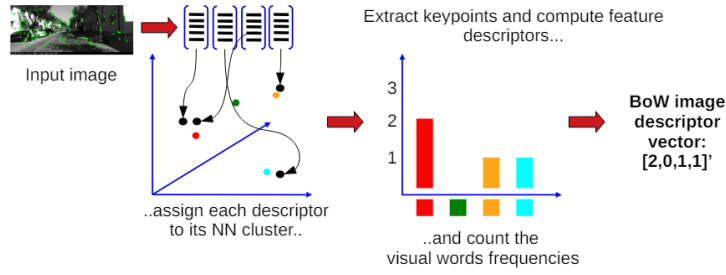
$$\{X^i, R_l, t_l | i \in \mathcal{P}_L, l \in \mathcal{K}_L\} = \arg \min_{X^i, R_l, t_l} \sum_{k \in \mathcal{K}_L \cup \mathcal{K}_F} \sum_{j \in \mathcal{X}_k} \rho(E_{kj})$$

$$E_{kj} = \left(\left\| x_{(\cdot)}^j - \pi_{(\cdot)}(R_k X^j + t_k) \right\|_{\Sigma}^2 \right)$$

Loop-Closure Detection

A robust recognition is an essential requirement for long-term Visual SLAM. Loop closure detections provide correct data association to obtain consistent maps. ORB-SLAM, as several other visual SLAM approaches, performs loop closure detection with **Bag of Words (BoW)** techniques. BoW summarizes the content of an image by the visual words it contains. visual words correspond to a discretization of the descriptor space, known as the visual vocabulary.





5.5.5 BoW-based Place Recognition

Online place recognition: compare the bag of words vector of the current image against the bag of words vectors of the images in the database (DB holds the BoWs of the **keyframes** collected so far). The image with the maximum individual score is considered as a loop candidate for the current image. A loop candidate is accepted if it is consistent with the previous k queries and it passes a geometrical check. ORB SLAM uses a vocabulary tree to speed-up images search.

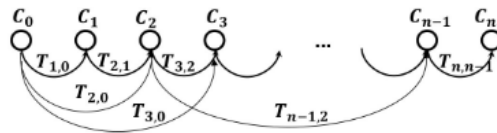
5.5.6 Loop Correction

Once a loop closure has been detected, the map should be corrected accordingly:

- Match the descriptors belonging to the current image I_{cur} with the descriptors off the matched keyframe I_{loop} representing the loop closure.
- Recover from the map the 3D points associated with the I_{loop} descriptors.
- Recover the rigid body transformation $\mathbf{T}=(\mathbf{R},\mathbf{T})$ between I_{cur} and I_{loop} e.g., by registering I_{curr} in the map with PnP.
- Use this transformation to correct the full path: ORB SLAM like many other approaches achieves this by optimizing a **pose graph**.

5.5.7 Pose Graph Optimization

A pose graph optimizes the camera positions C_0, C_1, \dots associated to each collected keyframe. It only con-



siders the estimates of the rigid body transformations T between the nodes: minimize the sum off squared transformation errors in place of SfM reprojection errors:

$$\sum_{e_{ij}} \|C_i - T_{e_{ij}} C_j\|^2$$

5.5.8 Direct Approaches

Camera motion and the scene structure are computed directly from **image intensity diiscrepancies using optimization techniques**. For each involved pixel:

- take note of its previous intensity y_i
- Un-project it back to 3D using the current depth guess d_i
- Re-project onto the next image using the current transformation T guess
- Take note of its current intensity y'_i
- Compute the per-point residual $y_i - y'_i$
- Iteratively optimize T, d_1, \dots, d_n using least-square over such residuals



5.6 DSO: Direc Sparse Odometry

Direct method that jointly optimizes camera poses, sparse point positions, **camera intrinsics**, **photometric calibration** (such lens attenuation, gamma correction and exposure times). Optimization is performed in a sliding window of keyframes. Fast real-time CPU implementation.

5.6.1 DSO Photometric Error

Photometric error of a point p in a reference image I_i observed in a target image I_j as weighted SSD over small neighborhood of pixel N_p :

$$E_{p_j} := \sum_{p \in N_p} \left\| I_j[p'] - \frac{t_j}{t_i} I_i[p] \right\|$$

$$p' = \Pi_c(R\Pi^{-1}(p, d_p) + t) \quad \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} := T_j T_i^{-1}$$

The full photometric error over all considered frames and points is given by:

$$E_{photo} := \sum_{i \in \mathcal{F}} \sum_{p \in \mathcal{P}_i} \sum_{j \in obs(p)} E_{pj}$$

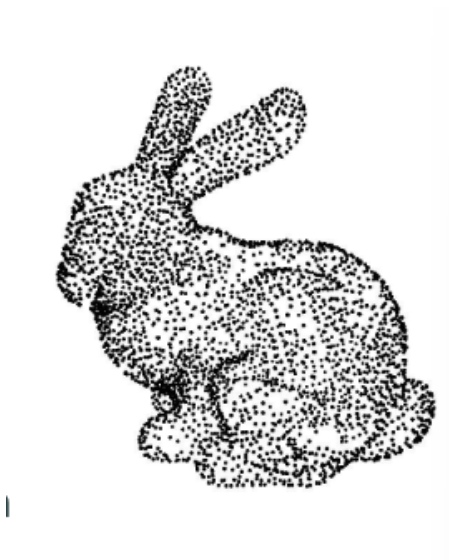
Chapter 6

3D Data Representation

6.1 Depth maps and Range Images

A depth map contains the distances of points in the scene to a plane. a range image contains the distance of points in the scene to another point.

6.2 Point Clouds



Set of 3D vertices with coordinates (x,y,z) . Common raw output of several 3D scanners. Generally unstructured, arbitrary points order. we have no information about vertex connectivity. No clue about the underlying surface from which the points are sampled. Any permutation of the points yields the same point cloud. The connectivity of a given point can be approximated with its k -nearest neighbors. Surface normals can be estimated by assuming piecewise planar surfaces and fitting a plane at one point and some of its neighbors.

6.2.1 Organized Point Clouds

Also called structured point clouds are in 2-D matrix of 3D points: $m \times n \times 3$ array and, for each $i=1,\dots,m$ and $j=1,\dots,n$, there is a corresponding 3D vertex v_{ij} . Points are organized according to some spatial relationships between them. Very efficient memory layout. We can use an additional 2D arrays of size $m \times n$: presence or absence of 3D data, confidence values, point colors and so on.

There is an implicit mesh connectivity between neighbouring points, and single non-boundary points are always degree 6.

6.3 Meshes

Represent or approximate a 2D manifold in 3D space. A set of vertices connected to each other in order to form some polygons. **triangula mesh**: very common data structure usually employed by 3D renders: each triangle has a single surface normal. for each surface point it is possible to get a normal vector e.g., from the vertices normals or directly from the face normal. Normals can then be used by rendering algorithms.



6.3.1 Manifolds

In mathematics, a manifold is a geometrical space not necessarily equipped with a metric to measure distances that locally resembles a Euclidean space. Each point of n-dimensional manifold has a neighborhood that is homeomorphic to an open subset of n-dimensional Euclidean space. one-dimensional manifolds includes lines and circles. Two-dimensional manifolds also called surfaces, i.e. planes, spheres and the torus.

6.3.2 Triangular Mesh

Closed: the surface completely encloses a volume. **Open:** the mesh contains boundary edges which are used by only one triangle. **Genus** of a surface: number of cuts we can make by means of closed, non-intersecting curves without disconnecting the manifold.

It is possible to associate a color to vertices of a mesh. The points that are inside a triangle can be painted with a color obtained by interpolating the colors of the vertices of the triangle.

Often, meshes are augmented with texture coordinates used to map from the surface to a 2D planar image. Texture coordinates range from 0 to 1 in the x and y axis.

6.3.3 Implicit Surfaces

The set of all points which satisfy the function $f(x,y,z)=0$. Typically, function values greater than zero indicate a point which is outside the object, while negative values indicate a point which is inside; this is known as a **Signed Distance Function (SDF)**. In general, obtaining a function $f()$ which exactly describes a general surface is difficult.

6.4 Voxels

Voxels are the 3D analogue of pixels. At each voxel, a boolean or a scalar is stored which indicates the presence or absence of the volume at that location. Voxels grids are used to represent properties at discrete spatial locations over a 3D volume. Precision depends on resolution. Data storage requirements are high. Can be easily used within deep learning models.

6.5 Octrees

A space-efficient representation for voxel grids. It uses adaptive resolution depending on the complexity of the surface throughout the volume. Idea; subdivide voxels which are not fully occupied by the volume or fully empty. Represented by trees: the root represents the complete volume. Nodes have eight children. Leaves are fully occupied or empty voxels.

6.6 Volumetric signed Distance Function

The idea is to represent a 3D environment as a voxel grid in which each voxel stores the signed distance to the nearest surface. The surface can be recovered by interpolation. In order to represent non-convex or multiple, disconnected surfaces, the SDF can be truncated between a minimum and maximum value (TSDF).

6.7 Surface Normals

Surface normals are important properties of a geometric surface used in several problems: 3D shape recognition; 3D ego-motion estimation/reconstruction from point-clouds; Object pose estimation; Surface classification; In computer graphics to apply the correct light sources that generate shadings and other visual effects.

If the surface is truly piecewise planar associate a face normal n_f to each triangular facet.

$$n_f = \frac{(v_2 - v_1) \times (v_3 - v_1)}{\|(v_2 - v_1) \times (v_3 - v_1)\|}$$

If the surface is assumed smooth the surface normal is approximated at each vertex by means of weighted averages of the adjacent facet normals.

6.7.1 Surface Normals from Point Clouds

No information about vertex connectivity. The idea is that for each point p_i we select k-nearest neighbors and fit a plane. Given a point cloud $P = \{p_1, \dots, p_n\}, p_i \in \mathbb{R}^3$, for each p_i we would like to estimate a normal vector



n_i from a set of k points in its neighborhood plus the point p_i itself:

$$Q_i^+ = [p_1, q_{i1}, \dots, q_{ik}]^T$$

Idea 1: fit a local plane to Q_i^+ :

$$S_i = n_{ix}x + n_{iy}y + n_{iz}z + d$$

i.e, solve:

$$\min_{b_i} \left\| [Q_i^+ \ 1_{k+1}] b_i \right\|_2$$

with

$$b_i = [n_i^T, d]^T$$

Solved by computing SVD for Q_i^+ and taking the last column of V .

Idea 2: **minimize the variance** by removing the empirical mean from the data matrix Q_i^+ :

$$\min_{n_i} \left\| [Q_i^+ \ \bar{Q}_i^+] n_i \right\|_2$$

we may use SVD to solve this system.

6.7.2 PCA-Recap

Orthogonal projection of data onto a lower dimension linear space that maximizes variance of projected data and minimizes mean squared distance between original data points and their projection.

6.8 K-D trees

Brute force k-nearest neighbor search has $= (Kn^2)$: **K-dimensional trees** can be used for organising points to speed up this search to $O(kn \log n)$. A K-D tree is a binary tree in which every leaf node is represented by one or more K-dimensional point. K for 3D point clouds is just 3. A K-D tree is a space partitioning data structure for organizing points. Partitioning is done by setting planes orthogonal to one of the axes.

This simplified implementation provides an **approximated NN**, since the actual NN may not be included in the reached leaf. The algorithm should be modified by adding a backtracking step to find the exact nearest neighbour. The tree should be as balanced as possible. Select the better dimension where operate the split.

6.8.1 Example of K-D Tree Construction

Construction for 1-NN: built the K-D tree from all 3D points pick a random dimension, find the median, split the data, repeat.

6.8.2 Example: K-D Tree NN Query

For each query point, find the NN by navigating the whole tree up to a leaf.