

Lecture 7 - Structure from motion II

- Offline SfM →
 - ① No realtime requirements, all images are available at once.
 - ② Eg. as basis for dense 3D model reconstruction.
- Realtime SfM →
 - ① Realtime requirements, images are available one by one, output required at each time step.
 - ② Eg. for mobile augmented reality in unknown environments.
- Comparison.

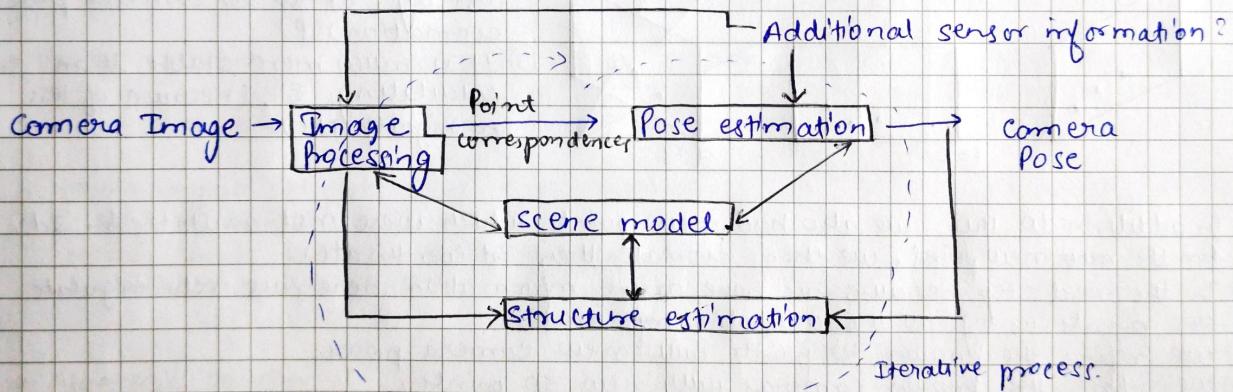
Offline SfM

- ① uses many features and computationally intense image processing methods
- ② forward and backward feature matching
- ③ global bundle adjustment

Realtime SfM

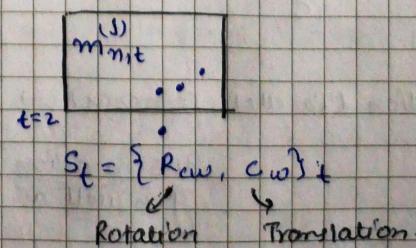
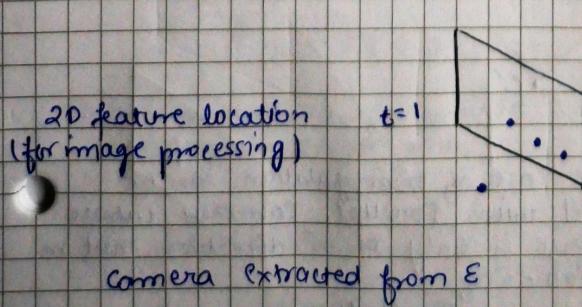
- ① Reduced no. of features and lightweight image processing methods.
- ② sequential / recursive processing scheme.
- ③ local bundle adjustment (keyframe/window based).

- Structure and motion - Basic Pipeline.



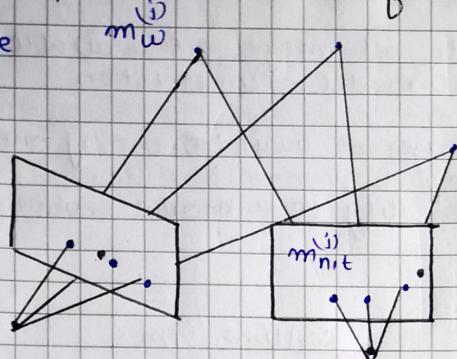
- Realtime SfM → Alternating estimation of Camera Poses and 3D feature locations (triangulation) and from a (continuous) image sequence.
- We assume that the cameras are calibrated. (we know K, intrinsic camera parameters) and we know how to extract the camera position and orientation from "E". (which defines the geometry b/w 2 images).

- ① How do we initialize the first camera poses?



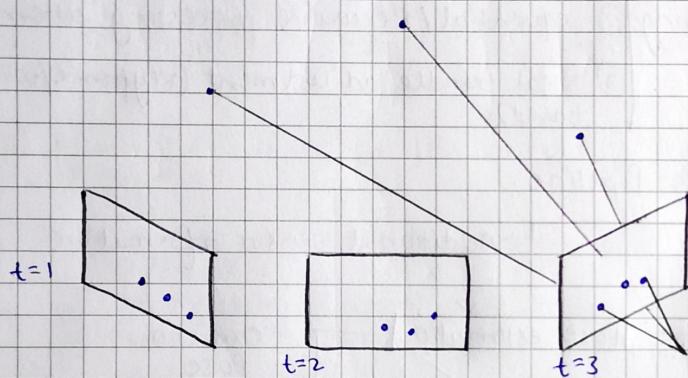
(2) Now, we have the two first camera poses, how do we go on?

3d feature location



Triangulate the 3D points.
(we have the structure of the scene)

(3) When we have a new frame i.e. estimate the next camera pose.



(1) The first step is to propagate the info. from the 2nd frame to the 3rd / new frame (by tracking of 2D points).

(2) Once we have these points in the new frame, we will do the camera pose estimation directly (3d to 2d camera pose estimation) (P)

(3) It is much more stable than calculating " \mathcal{E} ", (because of less ambiguity, etc)

(4) In addition to this, we also had new points which were not available in 3d. For the new new point, we then estimate their 3d coordinate.

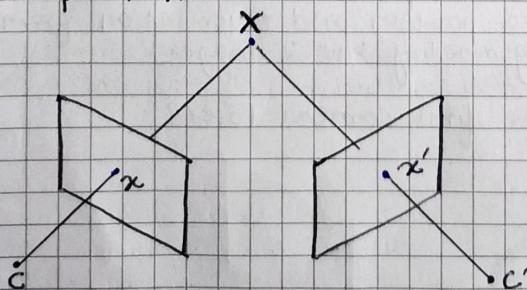
(5) In the next step, because we have more information here, we triangulate our points using all the available information.

(6) We refine the known 3D points with new camera poses.

(7) We refine the known cameras with new 3D points.

Triangulation

objective \rightarrow Given at least 2 known cameras (C and C') and 2 corresponding feature points (x and x') (i.e. the 2 camera views), estimate the 3D point X .



\Rightarrow Triangulation (Parallel cameras) \rightarrow A simple case of triangulation can be understood with parallel cameras where the motion of camera is only in x direction not in y as well as z direction - \downarrow translation.

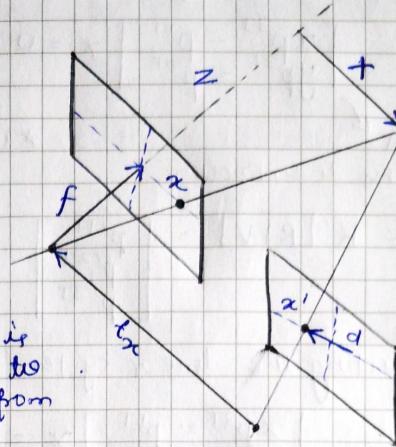
$$K = K' = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad R = I \quad t = \begin{pmatrix} tx \\ 0 \\ 0 \end{pmatrix}$$

Then $y = y'$ (epipolar lines) and depth z can be computed from disparity:

$$\frac{z}{fx} = \frac{tx}{d} \quad (d = x' - x)$$

$$z = \frac{fx}{d}$$

we have a depth which is inversely proportional to "d" (shift of the point from image to the other)
disparity of image point



* derivation via equal triangles $\frac{x}{f} = \frac{x}{z}$ $\frac{x'}{f} = \frac{x + tx}{z}$

$$\frac{x'}{f} = \frac{x}{f} + \frac{tx}{z}$$

* feature displacement (disparity) is inversely proportional to depth.
as $d \rightarrow 0$, $z \rightarrow \infty$

⇒ Vector solution → Let's say instead of stereo cameras, we have two separate cameras. Due to noise, the two rays (which are projected from the image point) might not meet.

→ we compute the mid point of the shortest line between the two rays.

→ This can be further solved algebraically by least squares method.

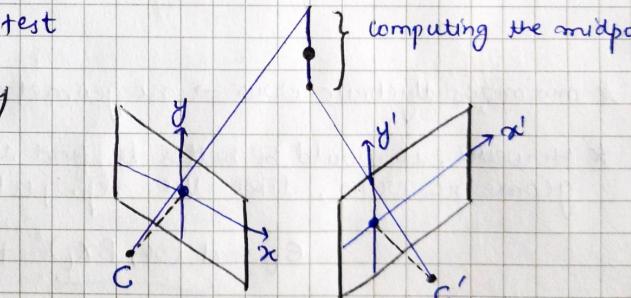
⇒ Algebraic solution.

① Equation for 1 camera view

$$\text{image point } x \propto Px \xrightarrow{\text{upto scale}} \lambda x = Px \xrightarrow{\text{scale factor}} x \propto \begin{bmatrix} p_1^T \\ p_2^T \\ p_3^T \end{bmatrix} x = 0$$

Registered matrix $\xrightarrow{\text{3D point (3DOF)}}$ we know λ, x, P . We have to calculate the i^{th} row of P .

(3×4 camera matrix) $\xrightarrow{\text{3D point } x}$ doing cross prod on both the sides



② multiple camera views

$$x \propto Px = 0 \quad \text{and} \quad x' \propto P'x = 0$$

③ One camera view can provide 3 equations, only 2 linearly independent i.e. dropping the 3rd row

$$P = K[R|t] = \begin{bmatrix} p_1^T \\ p_2^T \\ p_3^T \\ p \end{bmatrix}$$

$$\begin{aligned} x(p_3^T x) - (p_1^T x) &= 0 \\ y(p_3^T x) - (p_2^T x) &= 0 \\ px(p_3^T x) - y(p_1^T x) &= 0 \end{aligned}$$

④ Construction of a linear system.

$$\begin{bmatrix} xp^{3T} - p'^T \\ yp^{3T} - p'^{2T} \end{bmatrix} x = 0$$

2×4 matrix A
i.e. 2 equations 4 unknowns

* when we increase the camera view, i.e.

$$P' = K' [R|t] = \begin{bmatrix} p'^{1T} \\ p'^{2T} \\ p'^{3T} \end{bmatrix}$$

$$x(p'^{3T}x) - (p'^{1T}x) = 0$$

$$y(p'^{3T}x) - (p'^{2T}x) = 0$$

again constructing the linear system.

4 equations,
4 unknowns.

$$\begin{cases} \begin{bmatrix} xp^{3T} - p'^T \\ yp^{3T} - p'^{2T} \end{bmatrix} x = 0 & A_1 \\ \vdots & \vdots \\ \begin{bmatrix} xp'^{3T} - p'^{1T} \\ yp'^{3T} - p'^{2T} \end{bmatrix} x = 0 & A_2 \\ \vdots & \vdots \end{cases}$$

⑤ Stacking equations for all camera views:

$$\begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_N \end{bmatrix} x = 0$$

(i) Again, the linear system ($AX=0$) can be solved by SVD.

(ii) Recalling →

① minimize $\|Ax\|_2^2$, subject to
 $\|x\|_2^2 \leq 1$

② then normalize the 3D point x by dividing the 4th entry of x .

* minimizes algebraic error → no geometrical interpretation.

* However, it would be better to find the 3D point that minimizes a meaningful geometric error, like the reprojection error.

$$E_i = \underbrace{d(x, Px_i)^2}_{\text{camera 1}} + \underbrace{d(x', P'x_i)^2}_{\text{camera 2}}$$

should be 0 or as small as possible.

⇒ minimizing the geometric error

objective → Estimate the 3D point \hat{x} , which exactly satisfies the supplied camera geometry Pond P' , so that it projects as

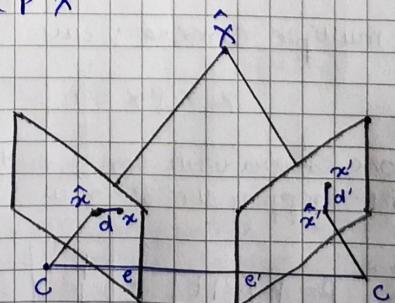
$$\hat{x} \propto Px$$

noise

$$\hat{x}' \propto P'x$$

① If the measurement noise in the image points is Gaussian with mean equal to zero, minimizing the reprojection error gives the maximum likelihood estimation of x .

② where \hat{x} and \hat{x}' are closest to actual image measurements.



$$\min d(x, \hat{x})^2 + d(x', \hat{x}')^2 \text{ subject to } \hat{x} \propto P \hat{x} \text{ and } \hat{x}' \propto P' \hat{x}'$$

- ③ this assumes perfect camera poses meaning P is error free.
 ④ Non-linear problem \rightarrow can be solved by Levenberg - Marquardt method.

\Rightarrow Triangulation - Properties.

- ① The smaller the angle θ_{par} (small baseline, big distance), the bigger the reconstruction uncertainty.

②

\Rightarrow Bundle Adjustment.

- ① The whole idea of BA is to optimize all the above computations that we just studied.

- ② Valid poses estimation $[R_1 | t_1]$, $[R_2 | t_2]$, $[R_3 | t_3]$... and 3D points $x^1, x^2, x^3 \dots$ must let the re-projection close to the observations, i.e. to minimize the reprojection error.

$$\underset{\mathbf{x}^T, R_i, t_i}{\operatorname{argmin}} \sum_i \| \mathbf{x}_i^T - K[R_i | t_i] \mathbf{x}^T \|_2^2 \rightarrow \text{cost function.}$$

- ③ Global Bundle Adjustment (BA):

\rightarrow Refine the visual reconstruction to produce jointly optimal camera poses and 3D points.

estimate. $\underset{\mathbf{x}^T, R_i, t_i}{\operatorname{argmin}} \sum_i \sum_j \mathbf{r}_i^T$

minimize over parameter vector containing all camera poses and 3D points

over all frames and 3D points

Residual / Reprojection error.

- ④ For each camera, we have 6 parameters, + for every 3D point, we have 3 parameters.

\rightarrow 6R+3L parameters must be estimated. (for optimization)
 \rightarrow matrices are sparse (having lot of zeros), and therefore can be highly optimized.

- ⑤ minimizing the cost function:

$$\underset{\mathbf{x}^T, R_i, t_i}{\operatorname{argmin}} \sum_i \sum_j \| \mathbf{x}_i^T - K[R_i | t_i] \mathbf{x}^T \|_2^2 \quad \text{3 non-linear optimization.}$$

\Rightarrow Drift (A big problem)

- ① Both off-line SfM and SLAM will drift (diverge) very quickly.
 ② measurement has errors.
 ③ Uncertainties in camera poses propagate to the triangulated 3D points and vice-versa.

\Rightarrow Because of the above factors we have to perform drift correction.

\Rightarrow Dnft - (Initializing with 3 images) (instead of two)

- ① The length of the translation vector is set to 1 ($\|t\|=1$) for 2 views.

② this defines the scaling (an "arbitrary" scaling of the world).
 → for the 3rd view, the scaling must be computed and coherent with the defined scale.

③ R_{12}, R_{13}, R_{23} and t_{12}, t_{13}, t_{23} are the rotation and translation vectors among the three views 1, 2 and 3.

$$P_1 = K_1(I, o)$$

$$P_2 = K_2(R_{12}, t_{12})$$

$$P_3 = K_3(R_{13}, \lambda_{13} t_{13}) \rightarrow \text{scaling factor b/w view 1 and 3.}$$

$$\lambda_{13} = \frac{(t_{23} \times t_{13})^T (t_{23} \times R_{23} t_{12})}{\|t_{23} \times t_{13}\|^2}$$

⇒ Drift reduction - methods

① offline Sfm: If we have an offline system, the reprojection error is minimized through global bundle adjustment, where the estimated 3D points and camera poses will be refined globally.

② Online Sfm (Realtime)

① filter based → estimated 3D points and camera poses could be refined recursively in Extended Kalman Filter (EKF)

② keyframe based → for improving realtime efficiency, local bundle adjustment is used only on keyframes (selected good frames), combined with global bundle adjustment, when it is necessary.

③ loop closure → loop closure is based on the concept of recognizing the place where the camera visited before, close the loop and eliminate drift.

④ feature level → ① Reduce drifts in feature tracking / matching
 ② Extend feature tracks
 ③ Reacquire lost features. (when a person enters a scene, some features gets lost,).

⑤ geometric level → (careful, precise, robust) 3D point initialization.
 → How?

→ Triangulate over a whole set of camera views (> 2)

→ Enforce a minimal angle θ between view rays.

→ Use RANSAC to eliminate outliers.

→ Use only well reconstructed points for further estimation.

summary → Incorporate uncertainties, e.g. simple stochastic model and WLS (weighted least square) estimation: All entities modelled as Gaussian Random Variables. (since, the noise is gaussian)

⇒ Initial triangulation using RANSAC:

→ Trial and error.

→ Algorithm

- ① if $\theta \geq \theta_0$: store the new camera view
- ② if size of inliers $> n$

(a) until no more valid results or more options found

(a.1) Triangulate the pairs of camera views

(a.2) Validate the results.

(3) If a valid result is found:

(a) Compute a least squares (LS) estimate from all inliers

(b) Compute a weighted least squares (wLS) estimate from all inliers.

⇒ 3D point refinement.

→ Incorporate new camera view, each time the feature is observed in an image.

→ methods:

→ Repeated triangulation

→ Recursive filtering (E.g. using Extended Kalman filter)

⇒ Loop closure → Detection and recognition of already visited and mapped area after the camera exploring a long time. It is done by understanding the scene.

⇒ Scene Understanding.

① The content of an image can be inferred from the frequency of words.

→ visual words = independent features

→ Represent the images based on a histogram of word occurrences (bag)

→ Each detected feature is assigned to the closest entry in the codebook.