

• Large-scale 3D Reconstruction

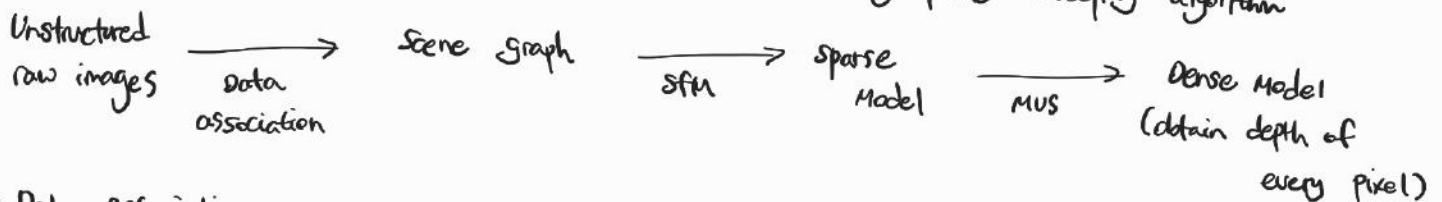
Given set of images $\{I_1, \dots, I_n\}$, find $\left\{ \begin{array}{l} \text{motion of cameras w.r.t. } P_w \{P_1, \dots, P_n\} \\ \text{3D scene points } \{x_1, \dots, x_n\} \end{array} \right.$

• 3-step pipeline

1) Data association: check two images are related using image correspondences and robust two-view geometry
keypoint/descriptor matching
RANSAC, pose estimation (EF)

2) Structure-from-Motion (SfM): initial 3D reconstruction by triangulation, then refine w/
bundle adjustment
PnP, EPnP
Gauss-Newton, Levenberg-Marquardt

3) Multi-view Stereo (MVS): get dense 3D model using plane sweeping algorithm



• Data association

Connect images w/ overlapping views \rightarrow Connected components represent a single 3D scene
 \downarrow steps

1) Extract image keypoints/descriptors e.g. ORB, SIFT, SURF

2) Establish keypoint correspondences

3) Take geometric verification (find outliers) \because keypoints matching is purely based on appearance
 \hookrightarrow Given $x_i \leftrightarrow x_j$ correspondences, take RANSAC-based two-view geometry (EF w/o outliers)
if inlier counts $>$ threshold, image pairs added to scene graph

• Problem: these steps are computationally expensive!

for N images and K keypoints/image, querying 1 image $= O(NK^2)$

\downarrow Solution: use bag-of-words image retrieval

• Place Recognition / Image Retrieval

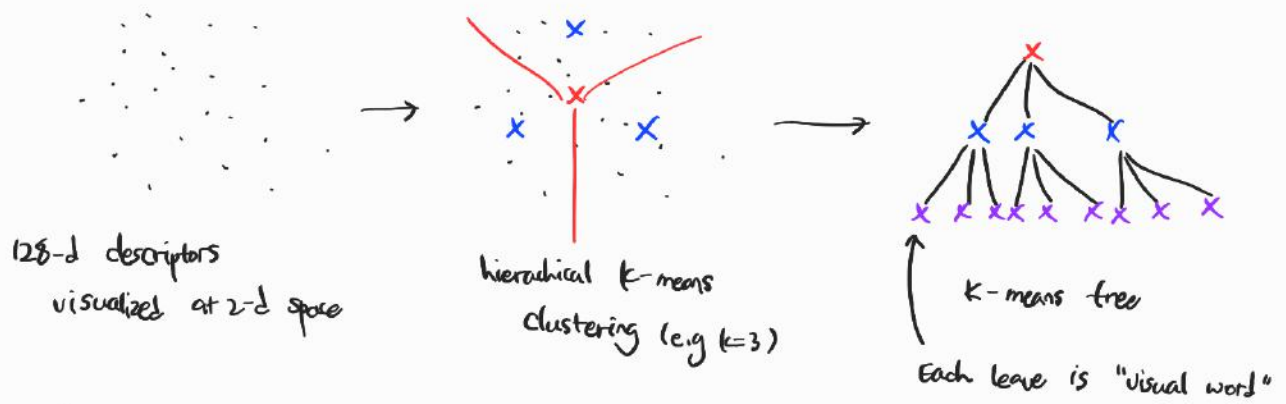
efficient tree-based search s.t. query image features \rightarrow find associated image from database

• Steps

1) Extract image keypoints and descriptors from training data

\hookrightarrow can use crawled images for training

2) Perform K-means hierarchical clustering on descriptors



3) Build inverted file index

For each image, extract image features and (tree-) search for closest leaf nodes
Build frequency count for each image

e.g. $\begin{matrix} \times & \times & \times & & \times & \times & \times & & \times & \times & \times \\ 1 & 2 & 0 & & 0 & 0 & 1 & & 3 & 1 & 5 \end{matrix}$ Image #1 counts

4) Given query image, retrieve visually similar image

e.g. $\begin{matrix} \times & \times & \times & & \times & \times & \times & & \times & \times & \times \\ 1 & 2 & 0 & & 0 & 0 & 1 & & 3 & 1 & 5 \end{matrix}$ Image #1 counts

$\begin{matrix} \updownarrow & \updownarrow & & & \updownarrow \\ 1 & 1 & 1 & 2 & 0 & 1 & 6 & 0 & 0 \end{matrix}$ Query image counts

Match = 3

* frequency counts can be normalized i.e. high frequency node ↓ weight
∴ non-discriminative node

• Time complexity

For K image features per image, N # of images, b/L width/depth of K -means tree,

Naive data association : $O(NK^2)$
Image retrieval : $O(KbL)$ → Much lower!

∴ Data association $\left\{ \begin{array}{l} \text{Take image retrieval to eliminate infeasible edges in scene graph} \\ \text{Take geometric verification for remaining edges} \end{array} \right.$

• SfM paradigms

1) Incremental 2) Global 3) Hierarchical

• Incremental sfm

• Initialization

1) choose two non-panoramic views ($\|u\| \neq 0$) w/ highest inlier keypoint matching → Seed pair

∴ difficult to derive H.Ef from pure rotation \hookrightarrow ∴ ↑ accuracy for least-square problem

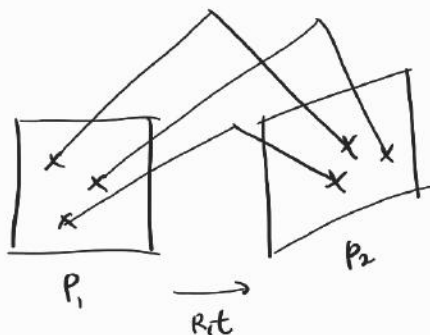
2) Take 4-point algorithm (planar scene, H) or 8-point algorithm (non-planar scene, E/F)

→ Decompose to R, t or P, P'

3) Set scale of translation to 1 ($\|t\|=1$ or $\|P'\|=1$)

4) Triangulate initial 2D-2D correspondences to 3D points

5) Apply bundle adjustment to refine camera matrices & 3D points



Initialization of
first two views (P_1, P_2)

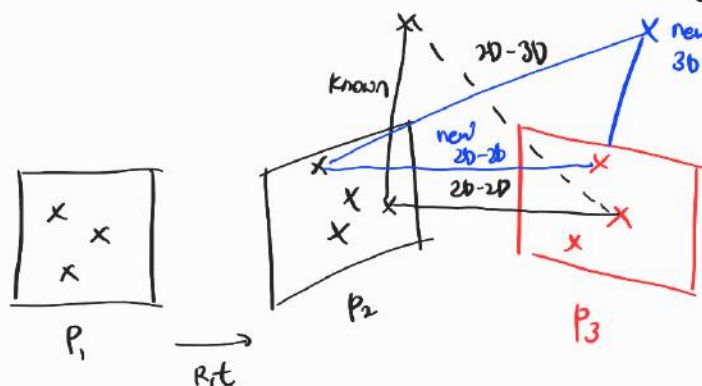
• Subsequent views

1) Find 2D-2D correspondence of existing-new views

2) Since 3D point of existing view is known, find 2D(of new view)-3D correspondences

3) Solve PnP to find P_3

4) Take 2D-2D correspondences (that x has matching at P_1) to triangulate new 3D points



...
repeat w/
many other views

• Refinement stage

$$\min_{P, x} \|x - \pi(P, x)\|$$

Multi-view non-linear optimization using bundle-adjustment
⇒ Minimize reprojection error

• Global SfM

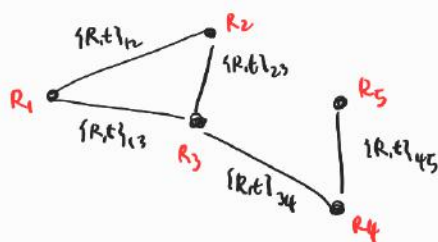
1) Compute relative poses $\{R_i, t_{ij}\}$ of all edges in scene graph (by decomposing calculated $E(F)$)

2) Estimate global rotations

$$\min_R \|R_{ij} - R_i R_j^T\|$$

$$R \in SO(3) \xrightarrow{R = \exp([\omega]_x)} \omega \in \mathbb{R}^3 \quad ([\omega]_x \in so(3))$$

$$\xleftarrow{[\omega]_x = \log(R)}$$



Scene graph

$\{R_i, t_{ij}\}$: known relative

R_i, t_i : unknown global

Then, $R_{ij} = R_i R_j^T \rightarrow w_{ij} = w_j - w_i = [0 \dots \underset{i\text{th}}{-1} \dots \underset{j\text{th}}{1} \dots 0] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} = A w_{\text{global}}$

\therefore Least-squares problem: $A w_{\text{global}} = w_{\text{rel}}$ where stacked $w_{\text{rel}} \in \mathbb{R}^{3 \times (\# \text{ of relative motion})}$

3) Estimate global translations

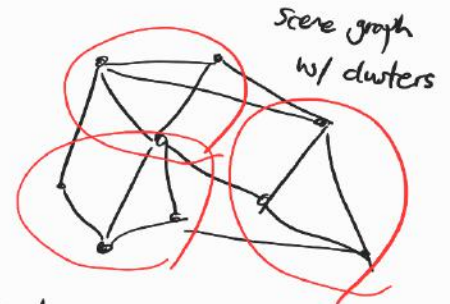
$$\min_t \left\| t_{ij} - \frac{t_i - t_j}{\|t_i - t_j\|} \right\| \rightarrow Ax = b \text{ problem}$$

\uparrow
Relative t_{ij} normalized s.t. $\|t_{ij}\| = 1$ \therefore up-to-scale

4) Triangulate to 3D using R_i, t_i & refine using bundle adjustment

• Hierarchical SfM

- 1) Hierarchical clustering of scene graph
- 2) Reconstruct clusters independently (Incremental / Global SfM)
- 3) Merge clusters using similarity transformations

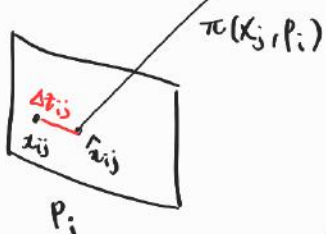


\hookrightarrow absolute orientation (R, t) + find scale using point correspondences

• Bundle Adjustment

\hookrightarrow Refine step to jointly optimize 3D points X_j and camera poses P_i

$X_j \rightarrow$ 3D point was triangulated during initialization of SfM



Cost function

$$\argmin_{P, X} \sum_i \sum_j \frac{\|x_{ij} - \pi(P_i, X_j)\|^2}{w_{ij}}$$

where w_{ij} : Measurement error covariance

$$\pi(P_i, X_j) = \frac{P_i X_j}{\hat{x}_3} \text{ to make homogenous form } \begin{pmatrix} \hat{x}_1 / \hat{x}_3 \\ \hat{x}_2 / \hat{x}_3 \\ 1 \end{pmatrix}$$

another form

$$\argmin_P \sum_i \sum_j \Delta z_{ij}^T w_{ij} \Delta z_{ij}$$

where P is $(12N \times 3M)$ vector w/ $12N$ $P_1 \dots P_N$ parameters and $3M$ $X_1 \dots X_j$ 3D points

• Iterative Estimation Methods

For non-linear function $f: x = f(p)$ where $x \in \mathbb{R}^N$ is measurement vector

\downarrow objective $p \in \mathbb{R}^M$ is Euclidean parameter vector

Find \hat{p} most nearly satisfies x closest to true value \bar{x}

e.g. Bundle Adjustment

- Measured value x : observed image coordinate x_{ij}
- Parameter \hat{p} : p, x
- Redefine objective : minimize error ϵ

Find \hat{p} s.t. $x = f(\hat{p}) - \epsilon$ or $\arg\min_{\hat{p}} g(\hat{p}) = \frac{1}{2} \|\epsilon\|^2 = \frac{1}{2} \epsilon^T \epsilon$

• Iterative estimation methods

1) Newton's Method

initial estimated value $p_0 \rightarrow$ error $g(p_0)$

Taylor series expansion : $g(p_0 + \Delta) = g(p_0) + g_p \Delta + \Delta^T g_{pp} \Delta / 2 + \dots$

where $g_p = \frac{\partial g}{\partial p} \Big|_{p=p_0}$, $g_{pp} = \frac{\partial^2 g}{\partial p^2} \Big|_{p=p_0}$

* Jacobian and Hessian Matrices

For $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Jacobian J : 1st order partial derivatives
- Hessian H : 2nd order

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad \text{or} \quad J_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n} \quad \text{or} \quad H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

Then, for $p_1 = p_0 + \Delta$:

$$g(p_1) = g(p_0) + \underbrace{g_p}_{J} \Delta + \underbrace{\Delta^T g_{pp} \Delta / 2}_H \quad \because \text{Taylor series expansion}$$

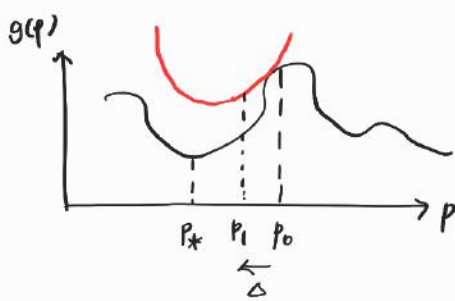
By differentiating w.r.t. Δ :

$$0 = g_p + g_{pp} \Delta \rightarrow g_{pp} \Delta = -g_p \quad (\text{inhomogeneous linear equation } Ax=b \text{ form})$$

Solve for Δ by g_{pp}^+ or least-square approach

\therefore Iterate $p_{i+1} = p_i + \Delta$ until convergence

• Remarks



Assumption: approx. quadratic cost function near minimum

∴ If p_0 far from p_* , assumption fails → slow/fail convergence

⊖ ↑ computational cost calculating H

2) Gauss - Newton method

$$g(p) = \frac{1}{2} \|\varepsilon(p)\|^2 = \varepsilon(p)^T \varepsilon(p) / 2$$

$$\frac{dg}{dp} = g_p = \varepsilon_p^T \varepsilon \quad \text{where} \quad \frac{d\varepsilon}{dp} = \varepsilon_p = \frac{d(X - f(p))}{dp} = \frac{df(p)}{dp} = f_p = J$$

$$\therefore g_p = J^T \varepsilon$$

$$\frac{d^2g}{dp^2} = g_{pp} = \varepsilon_p^T \varepsilon_p + \varepsilon_p^T \varepsilon \simeq \varepsilon_p^T \varepsilon_p = J^T J$$

↑ Hessian matrix that is computationally expensive
since second-order ε_{pp} is negligible, ignore it!

$$\Rightarrow g_{pp} \Delta = -g_p \quad \text{becomes} \quad J^T J \Delta = -J^T \varepsilon \quad \text{"normal equation"}$$

If weighted cost function i.e. $g(p) = \frac{1}{2} \|\varepsilon(p)\|_{\Sigma}^2 = \varepsilon(p)^T \Sigma \varepsilon(p) / 2$:

$$\underline{J^T \Sigma^{-1} J} \Delta = -J^T \Sigma^{-1} \varepsilon \rightarrow Ax = b \text{ form}$$

Symmetric, positive definite matrix

3) Gradient Descent

Recall $g_p = \varepsilon_p^T \varepsilon \rightarrow -g_p = -\varepsilon_p^T \varepsilon$ is most rapid decrease of cost function

∴ Iteratively decrease to negative gradient direction

$$\hookrightarrow \lambda \Delta = -g_p \quad \text{where } \lambda \text{ controls step size}$$

* Newton's method w/ $H \simeq \lambda I \equiv$ Gradient Descent

⊖ slow convergence ⊕ ↓ computational cost

4) Levenberg - Marquardt method $\simeq 2) + 3)$

Recall normal equations $J^T J \Delta = -J^T \varepsilon$

$$\rightarrow \text{Augmented normal equations: } (J^T J + \lambda I) \Delta = -J^T \varepsilon$$

• Updating λ

$$1. \text{ Initialization: } \lambda = 10^{-3} \cdot \text{avg}(\text{trace}(J^T J))$$

$$2. (J^T J + \lambda I) \Delta = -J^T \varepsilon \rightarrow \text{calculate } \Delta \rightarrow \text{update } p \ (p' \leftarrow p + \Delta)$$

Compare $\varepsilon(p)$ and $\varepsilon(p')$

if $\varepsilon(p) > \varepsilon(p')$: error \downarrow , accept p' , $\lambda' \leftarrow \lambda/10$
else : error \uparrow , reject p' , $\lambda' \leftarrow 10\lambda$

• Justification of LM method

(i) Small λ

$$(\cancel{J^T J} + \lambda I) \Delta = -J^T \varepsilon \longrightarrow \text{Normal equation } J^T J \Delta = -J^T \varepsilon$$

\therefore Small λ = close to min = Gauss-Newton method fits well

(ii) Large λ

$$(\cancel{J^T J} + \lambda I) \Delta = -J^T \varepsilon \longrightarrow \lambda \Delta = -J^T \varepsilon = -g_p$$

\therefore Large λ = far from min = Gradient Descent works well

• Sparse Levenberg-Marquardt algorithm

Problem of L-M method : intractable for large # of parameters to optimize e.g. BA

\therefore Solving normal equation $J^T J \Delta = -J^T \varepsilon$ takes $O(N^3)$

Solution:
sparse L-M

sparse matrix

\rightarrow find better method to calculate inverse!

Partition $p \in \mathbb{R}^M$ into a, b s.t. $p = (a^T, b^T)^T \rightarrow$ e.g. BA : camera/3D point parameter partition

Then, $J = [\partial \hat{x} / \partial p] = [A | B]$ where $A = [\partial \hat{x} / \partial a]$, $B = [\partial \hat{x} / \partial b]$

Normal equation $J^T J \delta = -J^T \varepsilon$ into block form

$$\begin{bmatrix} U^* & W \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta a \\ \delta b \end{pmatrix} = \begin{pmatrix} \varepsilon_A \\ \varepsilon_B \end{pmatrix}$$

Multiplying $\begin{bmatrix} I & -WU^{*-1} \\ 0 & I \end{bmatrix}$ on both sides

$$\begin{bmatrix} U^* - WU^{*-1}W^T & 0 \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta a \\ \delta b \end{pmatrix} = \begin{pmatrix} \varepsilon_A - WU^{*-1}\varepsilon_B \\ \varepsilon_B \end{pmatrix}$$

By making it to 0, first equation only depends to δa

$\therefore (U^* - WU^{*-1}W^T) \delta a = \varepsilon_A - WU^{*-1}\varepsilon_B \rightarrow$ solve δa

* $(U^* - WU^{*-1}W^T)$ is smaller & sparser than $J^T J$ from original normal equation

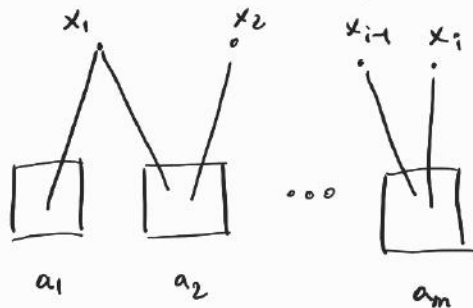
$W^T \delta a + V^* \delta b = \varepsilon_B \rightarrow$ solve δb

↓ update $P = (a + \delta a)^T, (b + \delta b)^T)^T$

Again like L-M method $\begin{cases} \text{if } \varepsilon(p) > \varepsilon(p') : \text{error} \downarrow, \text{accept } p', \lambda' \leftarrow \lambda/10 \\ \text{else : error} \uparrow, \text{reject } p', \lambda' \leftarrow 10\lambda \end{cases}$

• Example : Bundle Adjustment

↳ sparse L-M is advantageous due to lack of parameter interaction



$X = \{x_1, x_2, \dots, x_i\}$ where $x_i = (x_{i1}^T, x_{i2}^T, \dots, x_{im}^T)^T$

x_{ij} : image of x_i to j -th camera

$a = (a_1^T, a_2^T, \dots, a_m^T)^T$ where $a_j \in \mathbb{R}^{12 \times 1}$ projective matrix

↓ sparsity

$\frac{\partial \hat{x}_{ij}}{\partial a_k} = 0$ unless $j = k$ ∴ image projection \hat{x}_{ij} is only dependent to j -th camera

$\frac{\partial \hat{x}_{ij}}{\partial x_k} = 0$ unless $i = k$ ∴ image projection \hat{x}_{ij} is only dependent to i -th 3D point

∴ J and $H (= J^T J)$ are both sparse

$J =$

	p_1	p_2	p_3	x_1	x_2	x_3
x_{i1}	///	0	0	///	0	0
x_{i2}	0	///	0	0	///	0
x_{i3}	0	0	///	0	0	///

eg 3 cameras, 3 3D points

All non-diagonal entries of sub-matrix is zero

$H = J^T J =$

///	0	0
0	///	0
0	0	///

Hessian matrix is also the adjacency matrix of a graph

↓

$\begin{bmatrix} U^* - WU^{*-1}W^T & 0 \\ W^T & U^* \end{bmatrix} \begin{pmatrix} \delta a \\ \delta b \end{pmatrix} = \begin{pmatrix} \varepsilon_A - WU^{*-1}\varepsilon_B \\ \varepsilon_B \end{pmatrix}$

Schur's complement Camera parameters 3D point params

Since normally $|\delta a| < |\delta b|$, finding inverse of $U^* - WU^{*-1}W^T$ is computationally small

* Schur's complement : Block diagonal, sparse, symmetric positive Matrix

$$(U^* - WU^{*-1}W^T) \delta a = \epsilon_A - WU^{*-1}\epsilon_B \rightarrow \text{solve } x (= \delta a) \text{ of } Ax = b$$

Sparse matrix factorization

(i) $A = LU$: $L(Ux) = b \rightarrow Ly = b, y = Ux$

(ii) $A = QR$: $Q(Rx) = b \rightarrow y = Q^T b, y = Rx$

(iii) $A = LL^T$ (cholesky) : $L(L^T x) = b$

Iterative methods

(i) Conjugate gradient

(ii) Gauss-Seidel

• Problem of Fill-in

After matrix factorization, factorized matrix becomes dense

e.g. sparse $A \xrightarrow{A=LL^T} \text{Dense } L$

↓ Solution: Reorder sparse matrix s.t. single non-zero element per row & column

$$Ax = b \xrightarrow[\text{Permutation matrix}]{P} (P^T A P)(P^T x) = P^T b$$

But this step is *up-complete!*

↓ Approximate solutions

(i) Minimum degree

(ii) Column approximate minimum degree permutation

(iii) Reverse Cuthill-McKee

(iv) Nested Dissection