

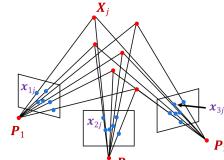
---



# Lec 18. structure from motion

## ① problem formulation:

- Given  $m$  images of  $n$  fixed 3D points such that (ignoring visibility)
$$x_{ij} \cong P_i X_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$
- Problem: estimate  $m$  projection matrices  $P_i$  and  $n$  3D points  $X_j$  from the  $mn$  correspondences  $x_{ij}$



## — ambiguity:

- If we scale the entire scene by some factor  $k$  and, at the same time, scale the camera matrices by the factor of  $1/k$ , the projections of the scene points remain exactly the same:

$$x \cong PX = \left(\frac{1}{k} P\right) (kX)$$

- Without a reference measurement, it is impossible to recover the absolute scale of the scene!
- In general, if we transform the scene using a transformation  $Q$  and apply the inverse transformation to the camera matrices, then the image observations do not change:

$$x \cong PX = (PQ^{-1})(QX)$$

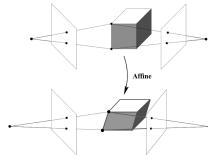
→ we can reconstruct it into a projective ambiguity Q: a general full-rank  $4 \times 4$  matrix

## ② affine ambiguity

- If we impose parallelism constraints, we can get a reconstruction up to an affine ambiguity:

$$x \cong PX = (PQ_A^{-1})(Q_AX)$$

$$Q_A = \begin{bmatrix} A \\ 0^T \\ 1 \end{bmatrix}$$

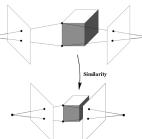


## 2) Similarity ambiguity

- A reconstruction that obeys orthogonality constraints on camera parameters and/or scene

$$x \cong PX = (PQ_S^{-1})(Q_S X)$$

$$Q_S = \begin{bmatrix} S^R \\ 0^T \\ 1 \end{bmatrix}$$



## ③ Affine structure from motion

### i) general affine projection:

$$P = \begin{bmatrix} a_{11} & a_{12} & a_{13} & t_1 \\ a_{21} & a_{22} & a_{23} & t_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} A & T \\ 0^T & 1 \end{bmatrix}$$

in nonhomogeneous cord.:  $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} = AX + t$

Given  $m$  images of  $n$  fixed 3D points such that:

$$x_{ij} = A_i \cdot X_j + t_i, \quad i = 1 \dots m, \quad j = 1 \dots n$$

problem use  $m \cdot n$  correspondences  $x_{ij}$  to estimate  $m$  projection matrix  $A_i$

and translation vector  $t_i$ , and  $n$  points  $X_j$

$$x_{ij} = \begin{bmatrix} A & T \\ 0^T & 1 \end{bmatrix} \cdot \begin{bmatrix} x_i \\ 1 \end{bmatrix} = \begin{bmatrix} A & T \\ 0^T & 1 \end{bmatrix} Q^{-1} Q \begin{bmatrix} x_j \\ 1 \end{bmatrix}$$

# Affine structure from motion pipeline:

## ① normalize for each view (image)

- First, center the data by subtracting the centroid of the image points in each view:

$$\begin{aligned}\hat{x}_{ij} &= x_{ij} - \frac{1}{n} \sum_{k=1}^n x_{ik} \\ &= A_i X_j + t_i - \frac{1}{n} \sum_{k=1}^n (A_i X_k + t_i) \\ &= A_i \left( X_j - \frac{1}{n} \sum_{k=1}^n X_k \right) \\ &= A_i \hat{X}_j\end{aligned}$$

codes:

so there is no translation vector here

## ② $D = MS \rightarrow X = A \cdot X$

- Let's create a  $2m \times n$  data (measurement) matrix:

$$D = \begin{bmatrix} \hat{x}_{11} & \hat{x}_{12} & \cdots & \hat{x}_{1n} \\ \hat{x}_{21} & \hat{x}_{22} & \cdots & \hat{x}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{x}_{m1} & \hat{x}_{m2} & \cdots & \hat{x}_{mn} \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_m \end{bmatrix} [X_1 \ X_2 \ \cdots \ X_n] \quad \begin{matrix} X \\ Y \\ Z \end{matrix} \quad (3 \times n)$$

points ( $3 \times n$ )

$M$  cameras ( $2m \times 3$ )

- What must be the rank of the measurement matrix  $D = MS$ ?

$m$  camera  
 $n$  points

$$\hat{x}_{11} = \begin{pmatrix} \hat{x}_{11} \\ \hat{y}_{11} \end{pmatrix}$$

$2m \times 3 \rightarrow P = \begin{bmatrix} a_{11} & a_{12} & a_{13} & t_1 \\ a_{21} & a_{22} & a_{23} & t_2 \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & a_{m3} & t_m \end{bmatrix} \cdot \begin{bmatrix} A & T \\ 0 & I \end{bmatrix}$

分解  
Factorizing the measurement matrix

given  $D$ , how to figure out  $M, S$ ?

method: SVD

$$D_{2m \times n} = U_{2m \times 2m} \times \Sigma_{2m \times n} \times V^T_{n \times n}$$

second step

$$D_{2m \times n} = M_{2m \times 3} \times Q_{3 \times 3} \times Q^{-1}_{3 \times 3} \times S_{3 \times n}$$

We can estimate  $Q$  to give the camera matrices in  $M$  desirable properties, like orthographic projection

$$D_{2m \times n} = U_{2m \times 3} \times \Sigma_{3 \times n} \times V^T_{n \times n}$$

- Keep top 3 singular values:
- This is the closest approximation of  $D$  with a rank-3 matrix in terms of Frobenius norm

$$\Sigma_3_{3 \times 3} \times V^T_{3 \times n}$$

- What to do about  $\Sigma_3$ ?
- One solution:  $M = U_3 \Sigma_3^{1/2}, S = \Sigma_3^{1/2} V^T_3$

This is a good way to eliminate noise  
"rank Theorem for Noisy Measurement"

- One possible solution:

$$D_{2m \times n} = M_{2m \times 3} \times S_{3 \times n}$$

$M = U_3 \Sigma_3^{1/2}$

- Are there other ...

$$D = \begin{bmatrix} A_1 Q \\ A_2 Q \\ \vdots \\ A_m Q \end{bmatrix} [Q^{-1} X_1 \ Q^{-1} X_2 \ \cdots \ Q^{-1} X_n]_{3 \times m}$$

that each camera matrix  $A_i Q$  represent orthographic projection

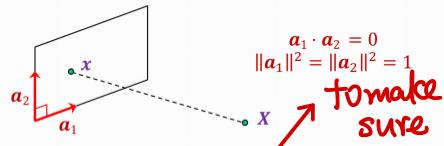
To be specific:

- Let  $a_1$  and  $a_2$  be the rows of a  $2 \times 3$  orthographic projection matrix. Then



- This translates into  $3m$  constraints on the 9 entries of  $Q$ :  $(A_i Q)(A_i Q)^T = A_i(QQ^T)A_i^T = I_{2 \times 2}, \quad i = 1, \dots, m$
- Are the constraints linear?
- First, solve for  $L = QQ^T$
- Recover  $Q$  from  $L$  by Cholesky decomposition
- Update  $M$  to  $MQ$ ,  $S$  to  $Q^{-1}S$

- Let  $\mathbf{a}_1$  and  $\mathbf{a}_2$  be the rows of a  $2 \times 3$  orthographic projection matrix. Then



- This translates into  $3m$  constraints on the 9 entries of  $\mathbf{Q}$ :
- $$(A_i \mathbf{Q})(A_i \mathbf{Q})^T = A_i (\mathbf{Q} \mathbf{Q}^T) A_i^T = I_{2 \times 2}, \quad i = 1, \dots, m$$
- Are the constraints linear?
  - First, solve for  $L = \mathbf{Q} \mathbf{Q}^T$
  - Recover  $\mathbf{Q}$  from  $L$  by Cholesky decomposition
  - Update  $M$  to  $M\mathbf{Q}$ ,  $S$  to  $\mathbf{Q}^{-1}\mathbf{S}$

algorithm :

- Compute the singular-value decomposition  $\tilde{\mathbf{W}} = \mathbf{O}_1 \Sigma \mathbf{O}_2$ .
- Define  $\hat{\mathbf{R}} = \mathbf{O}_1' (\Sigma)^{1/2}$  and  $\hat{\mathbf{S}} = (\Sigma)^{1/2} \mathbf{O}_2'$ , where the primes refer to the block partitioning defined in (13).
- Compute the matrix  $\mathbf{Q}$  in equations (15) by imposing the metric constraints (equations (16)).
- Compute the rotation matrix  $\mathbf{R}$  and the shape matrix  $\mathbf{S}$  as  $\mathbf{R} = \hat{\mathbf{R}} \mathbf{Q}$  and  $\mathbf{S} = \mathbf{Q}^{-1} \hat{\mathbf{S}}$ .
- If desired, align the first camera reference system with the world reference system by forming the products  $\mathbf{R}\mathbf{R}_0$  and  $\mathbf{R}_0^T \mathbf{S}$ , where the orthonormal matrix  $\mathbf{R}_0 = [\mathbf{i}_1 \ \mathbf{j}_1 \ \mathbf{k}_1]$  rotates the first camera reference system into the identity matrix.

$$\mathbf{Q}^{-1} \cdot \mathbf{S}$$

$$(3 \times 3) \quad (3 \times 2 \times 5)$$

how to deal with  $\mathbf{Q}$  ?

$$\text{Now } \hat{\mathbf{w}} = \mathbf{O}_1' \tilde{\Sigma}' \mathbf{O}_2'$$

$$\text{init, } \hat{\mathbf{R}} = \mathbf{O}_1' \cdot [\tilde{\Sigma}]^{\frac{1}{2}} \quad \hat{\mathbf{S}} = [\tilde{\Sigma}]^{\frac{1}{2}} \cdot \mathbf{O}_2'$$

$$\hat{\mathbf{w}} = \hat{\mathbf{R}} \cdot \hat{\mathbf{S}} = \underbrace{\hat{\mathbf{R}} \cdot \mathbf{Q}}_{\text{rotation matrix}} \cdot \mathbf{Q}^T \cdot \hat{\mathbf{S}} \rightarrow \text{shape matrix}$$

$$A_i \quad (2 \times 3) \quad A_i^T \quad (3 \times 2)$$

$$\mathbf{Q} \quad (3 \times 3) \quad \mathbf{Q}^T \quad (3 \times 3)$$

$$A \quad (2m \times 3) \quad \underbrace{\mathbf{Q} \cdot \mathbf{Q}^T}_{3 \times 3} \cdot A^T \quad (3 \times 2m) = I \quad (2m \times 2m)$$

$$\therefore L = A^{-1} \cdot I \cdot A^{-T}$$

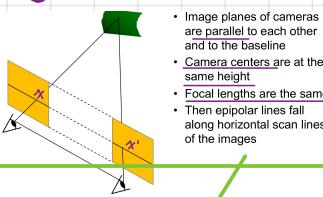
# Lec 19 stereo vision

— problem formulation: given stereo pair (assumed calibrated)  
figure out dense depth map

## ① Basic stereo matching algorithm



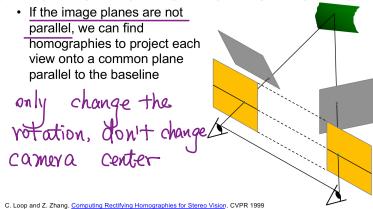
- For each pixel in the first image
  - Find corresponding epipolar line in the right image ①
  - Examine all pixels on the epipolar line and pick the best match ②
  - Triangulate the matches to get depth information ③
- Simplest case: epipolar lines are corresponding scanlines
  - When does this happen?



$$x^T \cdot E \cdot x = 0, E = [t \cdot I]R \text{ and } R = I, t = (t, 0, 0)$$

$$\therefore E = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -t \\ 0 & t & 0 \end{bmatrix} \quad \therefore (u', v', 1) \begin{pmatrix} 0 \\ -t \\ tv \end{pmatrix} = 0 \quad u = u', v = v' \quad y\text{-coord. is same}$$

$$(u', v', 1) \cdot \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -t \\ 0 & t & 0 \end{bmatrix} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = 0 \quad -tv + tv = 0$$



C. Loop and Z. Zhang, Computing Rectifying Homographies for Stereo Vision, CVPR 1999

Triangulation: not using  $X = P \cdot X$  (since we don't know  $P$ )

Depth from disparity

$$\frac{x}{f} = \frac{B_1}{z}, \quad \frac{x'}{f} = \frac{B_2}{z}$$

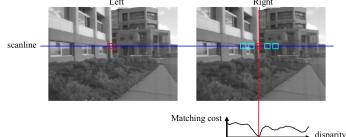
$$\frac{x - x'}{f} = \frac{B_1 - B_2}{z}$$

$$x - x' = \frac{f(B_1 - B_2)}{z}$$

$$z = \frac{fB}{x - x'}$$

## matching:

### 1) first method: Local stereo matching algorithm



- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

- Smaller window:
  - More detail
  - More noise
- Larger window:
  - Smoother disparity maps
  - Less detail

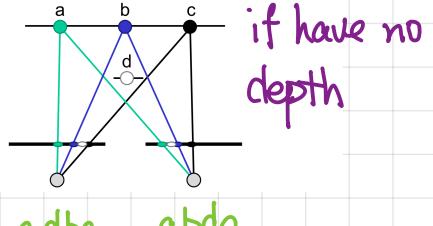
### 2) Beyond local stereo matching Non-local constraint:

#### ① Uniqueness

Each point in one image should match at most one point in the other image

#### ② ordering

Corresponding points should appear in the same order  
Is ordering always preserved in real life?



#### ③ smoothness