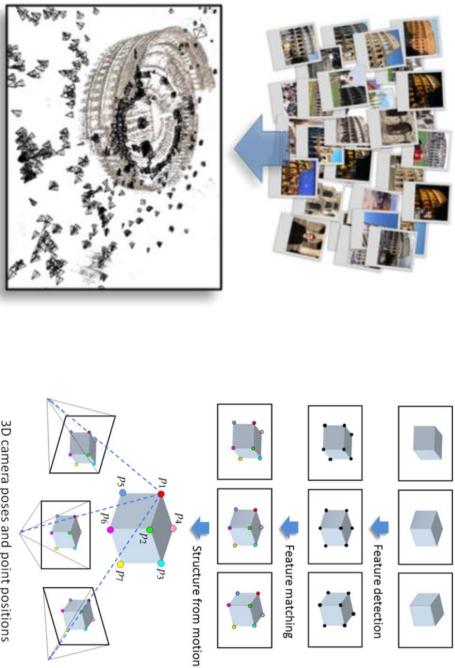


modern SFM

Topic 04:

Structure from Motion (SFM)

- big picture
- Euclidean & affine structure
- the factorization algorithm
- projective structure
- projective factorization
- SFM in Photo Tourism
- projective matrices & Euclidean upgrade from F
- incremental SFM



Snavey et al, "Scene Reconstruction and Visualization from community photo collections", Proc IEEE, 2010

modern SFM

3

example results: colliseum

4

Photo Tourism

Exploring photo collections in 3D

Noah Snavely Steven M. Seitz Richard Szeliski
University of Washington Microsoft Research

SIGGRAPH 2006

Agarwal et al, Building Rome in a Day, ICCV09

Ultimate objective: bundle adjustment

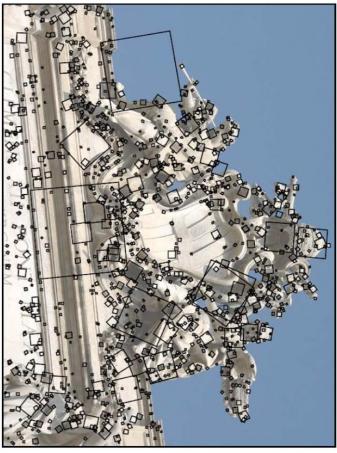
6

We wish to solve for the point positions $\mathcal{P} \equiv \{\vec{P}_n\}_{n=1}^N$ and the camera matrices $\mathcal{M} \equiv \{M_j\}_{j=1}^J$ by minimizing

$$\mathcal{O}(\mathcal{M}, \mathcal{P}) \equiv \sum_{j,n} \left\| \begin{pmatrix} x_{j,n} \\ y_{j,n} \\ z_{j,n} \end{pmatrix} - \frac{1}{\vec{e}_3^T M_j \vec{P}_n} (I_2, \vec{0}) M_j \vec{P}_n \right\|^2, \quad (4)$$

where the camera matrices M_j must be of the form $M_{in,j} M_{ex,j}$, as above in (2) and (3).

This nonlinear LS optimization problem is called **bundle adjustment**. For Gaussian IID measurement noise it is a ML estimator.



basic process

keypoint detection & matching, image selection

8

basic process

keypoint detection & matching, image selection

7

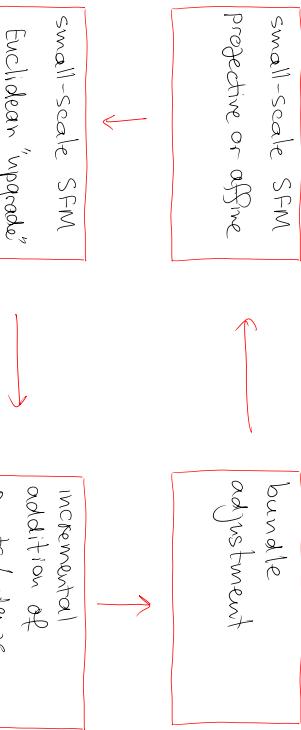
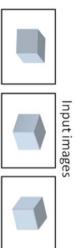
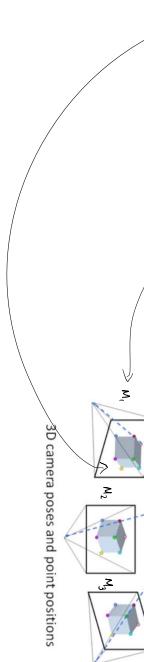


Fig. 4. Example set of detected SIFT features. Each detected SIFT feature is displayed as a black box centered on the detected feature location. SIFT detects a canonical scale and orientation for each feature, depicted by scaling and rotating each box.

The image points, $\vec{p}_{j,n}$, and the 3D scene points, \vec{P}_n , are related by perspective projection,

$$\vec{p}_{j,n} = \frac{1}{z_{j,n}} M_j \vec{P}_n. \quad (1)$$

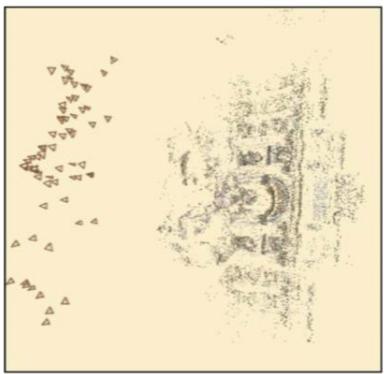
- $\vec{p}_{j,n} = (x_{j,n}, y_{j,n}, 1)^T$ is given in homogeneous pixel coordinates;
- $\vec{P}_n = (P_{n,1}, P_{n,2}, P_{n,3}, 1)^T$ is also in homogeneous coordinates;
- $M_j = M_{in,j} M_{ex,j}$ is the 3×4 camera matrix formed from the product of the intrinsic and extrinsic calibration matrices;
- $z_{j,n}$ is the projective depth, $z_{j,n} = \vec{e}_3^T M_j \vec{P}_n$, where $\vec{e}_3^T = (0, 0, 1)$ (i.e., \vec{e}_3 is the third standard unit vector).



5

basic process

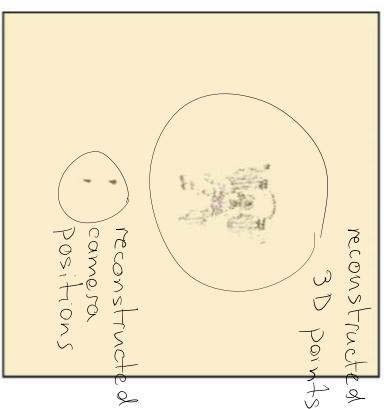
keypoint detection & matching, image selection



bundle
adjustment

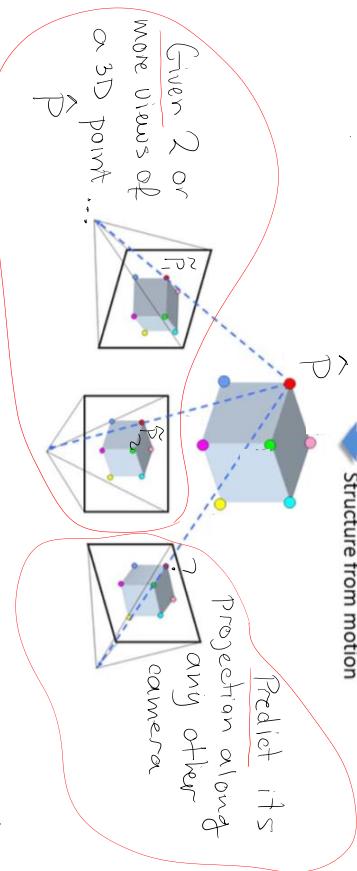
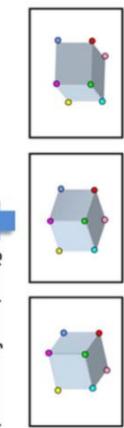
incremental
addition of
points / views

small-scale SFM
projective or affine
Euclidean "upgrade"



SFM: intuitive/geometric definition

12



Assumptions : ① we can measure the distance between two points in the same image
② projection model (affine/perspective..)

basic process

keypoint detection & matching, image selection

Topic 04:

11

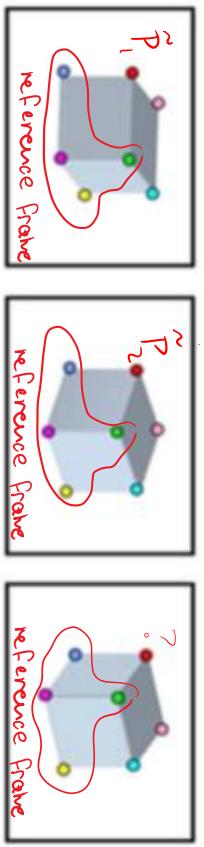
Structure from Motion (SFM)

- big picture
- Euclidean & affine structure: definitions
- computing affine structure from image projections
- the factorization algorithm
- projective structure
- projective factorization
- SFM in Photo Tourism
- projective matrices & Euclidean upgrade from F
- incremental SFM

9

SFM: intuitive/geometric definition

14



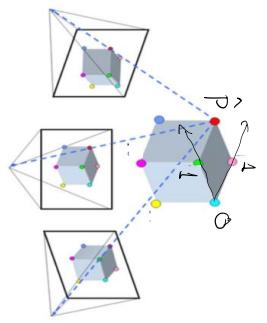
- Geometric intuition of the SFM process
- ① identify a set of image points that will serve as a reference frame
 - ② express point \hat{P} of the reference frame (i.e. assign it coordinates)
 - ③ use those coordinates and the projection of the reference points to predict \hat{P} 's projection

15

Euclidean reference frames in 2D

Conventions:

- $\hat{P} = (P_1, P_2)$
- \hat{P} : 3D point
- P : image projection
- P_i : \hat{P} 's coordinates

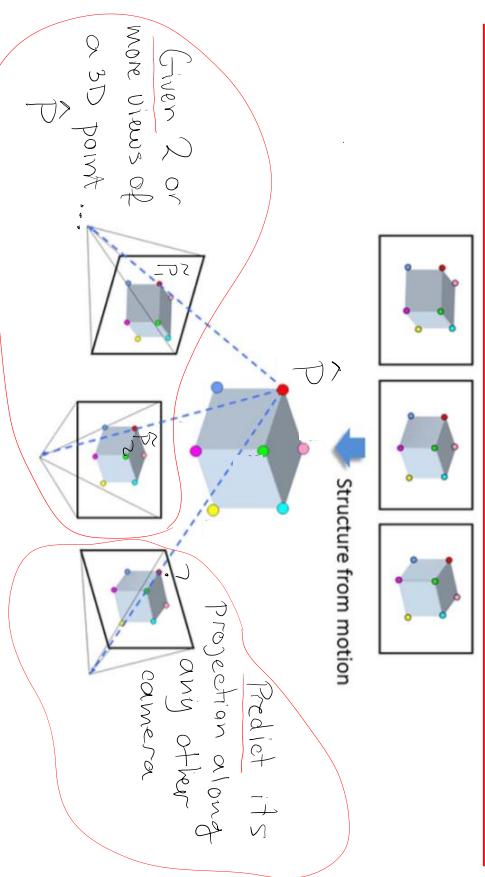


SFM: intuitive/geometric definition

In most cases the reference points are not chosen explicitly — they are implicit in the way SFM is formulated in various algorithms

- Geometric intuition of the SFM process
- ① identify a set of image points that will serve as a reference frame
 - ② express point \hat{P} in terms of the reference frame (i.e. assign it coordinates)
 - ③ use those coordinates and the projection of the reference points to predict \hat{P} 's projection

15



SFM: intuitive/geometric definition

13

Euclidean reference frames in 2D

18

Euclidean

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

st.

$$P_1 + P_2 + P_3 = 1$$

} called an affine sum
of 3 points

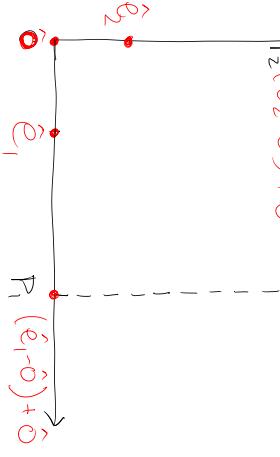
To be Euclidean, the

$$P_2 (\hat{e}_2 - \hat{o}) + \hat{o} = \hat{P}$$

representation of \hat{P} must

① be an affine sum
of 3 basis points

② the lengths of
line segments $\hat{o}\hat{e}_1$
and $\hat{o}\hat{e}_2$ must be 1
③ angle between line
segments $\hat{o}\hat{e}_1$ and $\hat{o}\hat{e}_2$ is $\frac{\pi}{2}$



Euclidean reference frames in 2D

Euclidean

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

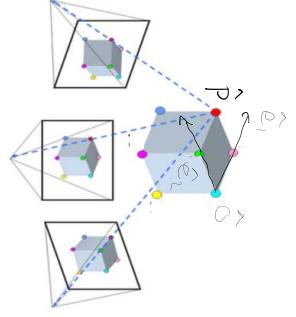
st.

$$P_1 + P_2 + P_3 = 1$$

To be Euclidean, the
representation of \hat{P} must

① be an affine sum
of 3 basis points

② the lengths of
line segments $\hat{o}\hat{e}_1$
and $\hat{o}\hat{e}_2$ must be 1
③ angle between line
segments $\hat{o}\hat{e}_1$ and $\hat{o}\hat{e}_2$ is $\frac{\pi}{2}$



Euclidean reference frames in 2D

Euclidean

$$\hat{P} = \hat{o} + P_1 (\hat{e}_1 - \hat{o}) + P_2 (\hat{e}_2 - \hat{o})$$

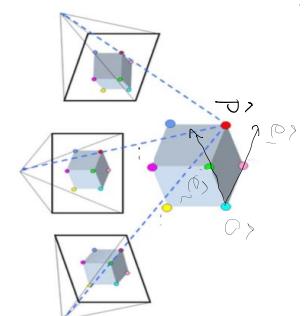
st.

$$P_1 + P_2 = 1$$

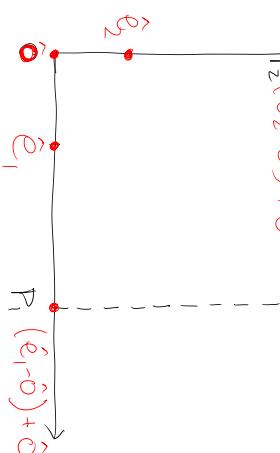
To be Euclidean, the
representation of \hat{P} must

① be an affine sum
of 3 basis points

② the lengths of
line segments $\hat{o}\hat{e}_1$
and $\hat{o}\hat{e}_2$ must be 1
③ angle between line
segments $\hat{o}\hat{e}_1$ and $\hat{o}\hat{e}_2$ is $\frac{\pi}{2}$



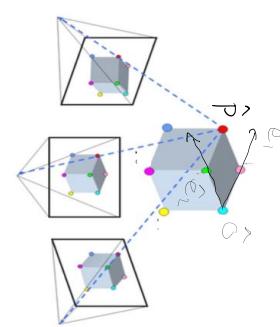
} called an affine sum
of 3 points



Euclidean

$\hat{P} = \hat{o} + P_1 (\hat{e}_1 - \hat{o}) + P_2 (\hat{e}_2 - \hat{o})$
=
 $(1 - P_1 - P_2) \hat{o} + P_1 \hat{e}_1 + P_2 \hat{e}_2$

origin
unit along $\hat{o}\hat{e}_1$,
unit along $\hat{o}\hat{e}_2$



19

Euclidean reference frames in 2D

17

Euclidean reference frames in 2D

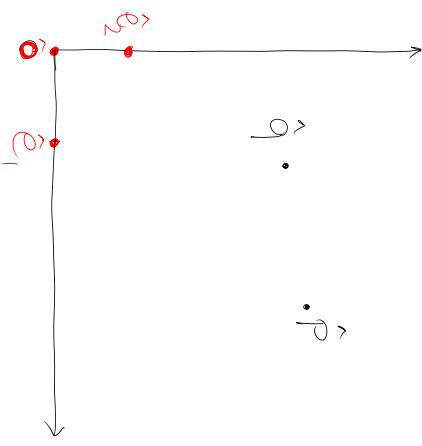
Euclidean

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

$$\hat{q} = q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}$$

$\hat{o}, \hat{e}_1, \hat{e}_2$ define
a Euclidean frame
of reference

e.g. suppose we rigidly
rotate the cube



24

Euclidean reference frames in 2D

Euclidean

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

$$\hat{q} = q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}$$

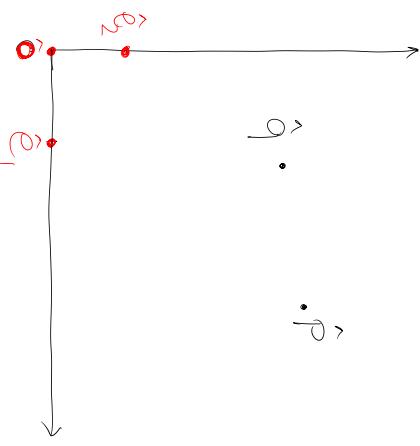
$\hat{o}, \hat{e}_1, \hat{e}_2$ define
a Euclidean frame
of reference

Euclidean structure:

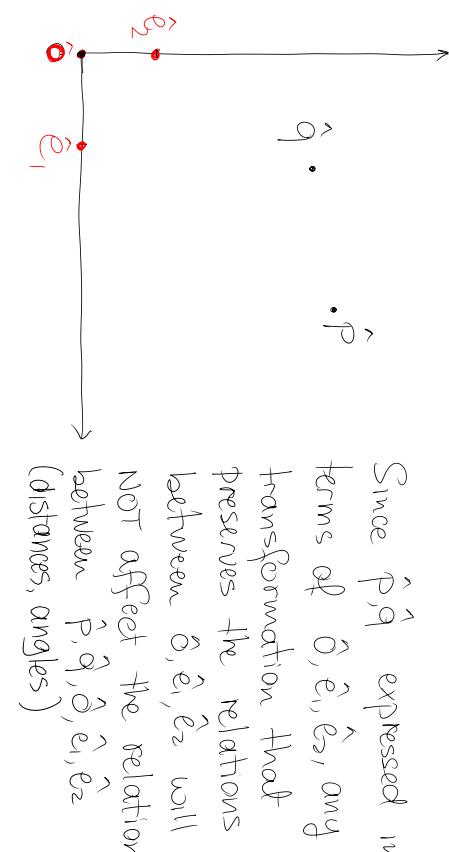
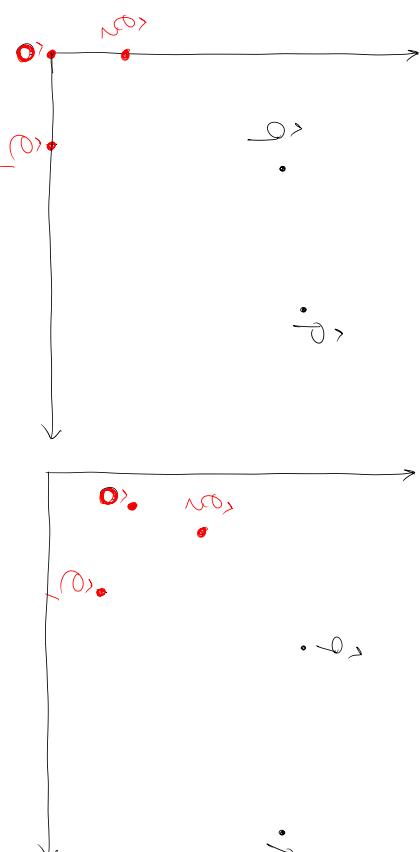
An assignment of
coordinates relative to
a Euclidean frame of
reference

represented by

$$(P_1, P_2, P_3) \quad (P_1+P_2+P_3=1)$$



q_1, q_2, q_3 ($q_1+q_2+q_3=1$)



23

Euclidean reference frames in 2D

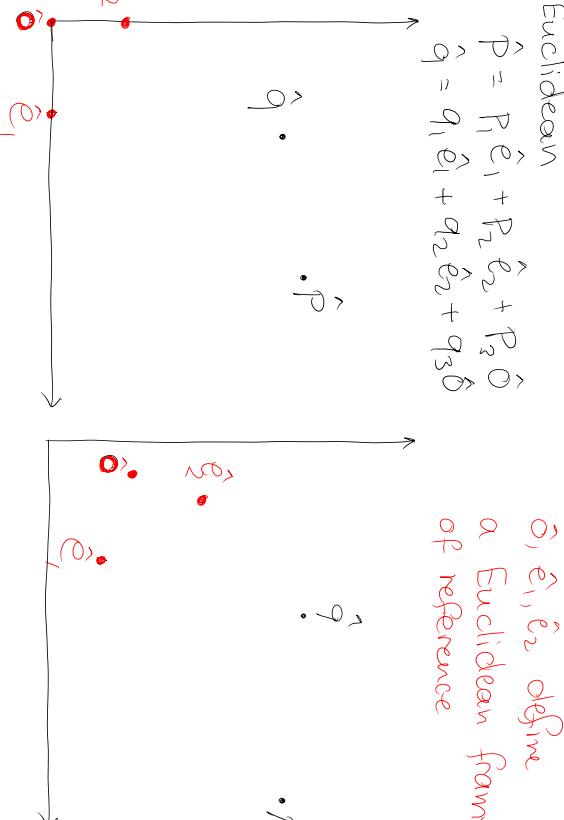
Euclidean

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

$$\hat{q} = q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}$$

$\hat{o}, \hat{e}_1, \hat{e}_2$ define
a Euclidean frame
of reference

Since \hat{p}, \hat{q} expressed in
terms of $\hat{o}, \hat{e}_1, \hat{e}_2$, any
transformation that
preserves the relations
between $\hat{o}, \hat{e}_1, \hat{e}_2$ will
not affect the relation
(distances, angles)



q_1, q_2, q_3 ($q_1+q_2+q_3=1$)

21

affine reference frame in 2D

26

affine reference frame in 2D

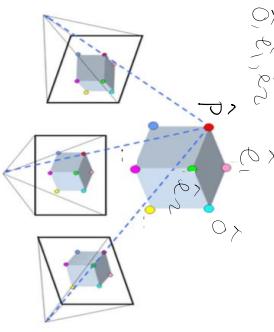
25

AFFine

$$\begin{aligned}\hat{P} &= P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o} \\ \hat{q} &= q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}\end{aligned}$$

e.g. if we don't know the 3D relation between

An affine frame of reference is defined by any non-collinear points $\hat{o}, \hat{e}_1, \hat{e}_2$ on the plane



The Euclidean (aka metric) upgrade

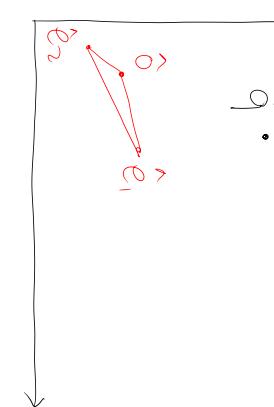
AFFine

$$\begin{aligned}\hat{P} &= P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o} \\ \hat{q} &= q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}\end{aligned}\quad (*)$$

e.g. if we don't know the 3D relation between

An assignment of coordinates relative to an affine frame of reference

An assignment of coordinates relative to an affine frame of reference



affine structure in 2D

AFFine

$$\begin{aligned}\hat{P} &= P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o} \\ \hat{q} &= q_1 \hat{e}_1 + q_2 \hat{e}_2 + q_3 \hat{o}\end{aligned}$$

An affine frame of reference is defined by any non-collinear points $\hat{o}, \hat{e}_1, \hat{e}_2$ on the plane

An affine frame of reference is defined by any non-collinear points $\hat{o}, \hat{e}_1, \hat{e}_2$ on the plane

An affine structure is represented by P_1, P_2, P_3 ($P_1+P_2+P_3=1$)

Once these coordinates are known, we can compute the Euclidean coords of \hat{P}, \hat{q} using Eq(*)

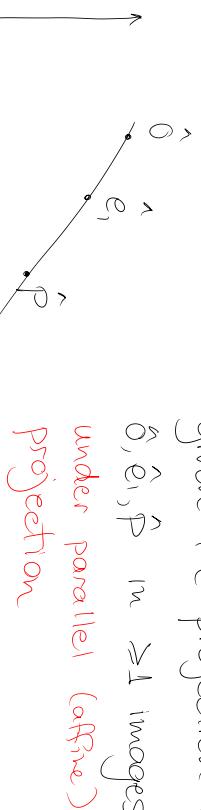
Whether structure is affine or Euclidean depends on the frame of reference, not the coords of \hat{P}, \hat{q}

Affine Structure

computing 1D affine structure from 1D images

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$$

Goal: Compute (P_1, P_2)
given the projection of
 $\hat{o}, \hat{e}_1, \hat{P}$ in ≥ 1 images



computing 1D affine structure from 1D images

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$$

Goal: Compute (P_1, P_2)
given the projection of
 $\hat{o}, \hat{e}_1, \hat{P}$ in ≥ 1 images
under parallel (affine)
projection

Question: is one
known image sufficient to
compute P_1 ?

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$$

Goal: Compute (P_1, P_2)
given the projection of
 $\hat{o}, \hat{e}_1, \hat{P}$ in ≥ 1 images
under parallel (affine)
projection

$$\begin{aligned} \hat{P} &= P_1 \hat{e}_1 + (1-P_1) \hat{o} \\ &= P_1 \hat{e}_1 + \hat{o} - P_1 \hat{o} \\ &= P_1 (\hat{e}_1 - \hat{o}) + \hat{o} \end{aligned}$$

$\Rightarrow P_1 = \frac{\text{length}(\hat{o} - \hat{P})}{\text{length}(\hat{o} - \hat{e}_1)}$

Topic 04: Structure from Motion (SfM)

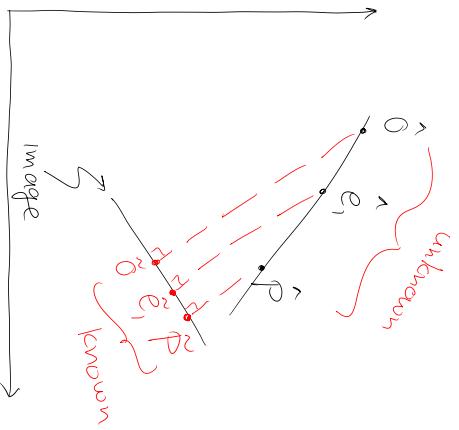
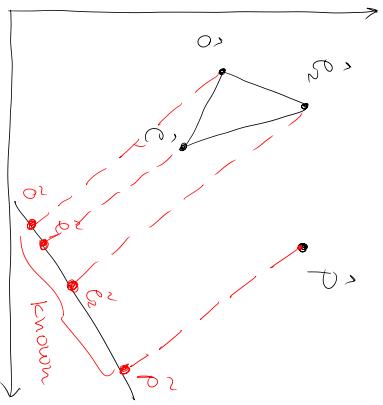
- big picture
- Euclidean & affine structure: definitions
- the factorization algorithm
- projective structure
- projective factorization
- SfM in Photo Tourism
- projective matrices & Euclidean upgrade from F
- incremental SfM

computing 1D affine structure from 1 image

34

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$$

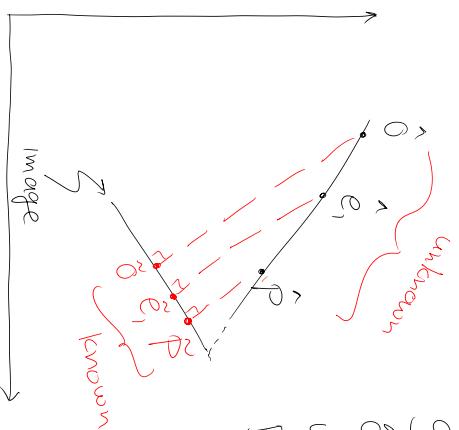
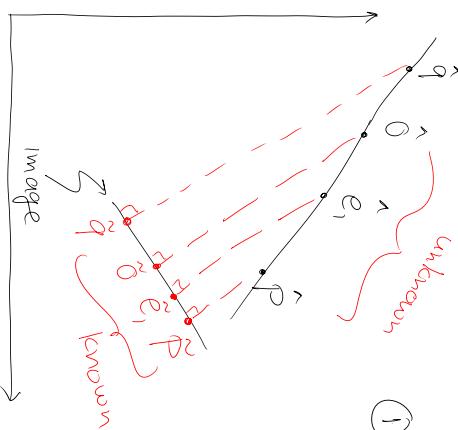
Goal: Compute (P_1, P_2)



computing 2D affine structure from 1D images

35

$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$ Goal: Compute (P_1, P_2, P_3)
given the projection of
 $\hat{o}, \hat{e}_1, \hat{e}_2$ in ≥ 1 images
under parallel (affine)
projection



computing 1D affine structure from 1 image

35

$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$ Goal: Compute (P_1, P_2)
given the projection of
 $\hat{o}, \hat{e}_1, \hat{e}_2$ in ≥ 1 images
under parallel (affine)
projection

① Yes, 1 orthographic view of
≥2 collinear points is
sufficient to compute
affine 1D structure

② distances between
any points \hat{p}, \hat{q}
known up to a scale
factor determined by
 $\text{len}(\hat{o}\hat{e}_1)$

Given the projection of
 $\hat{o}, \hat{e}_1, \hat{p}$ in ≥ 1 images
under parallel (affine)
projection

$$\frac{\text{len}(\hat{o}\hat{p})}{\text{len}(\hat{o}\hat{e}_1)} = \frac{\text{len}(\hat{o}\hat{p})}{\text{len}(\hat{o}\hat{e}_1)} = P_1$$

We also know that
 $P_1 + P_2 = 1$

computing 1D affine structure from 1D images

33

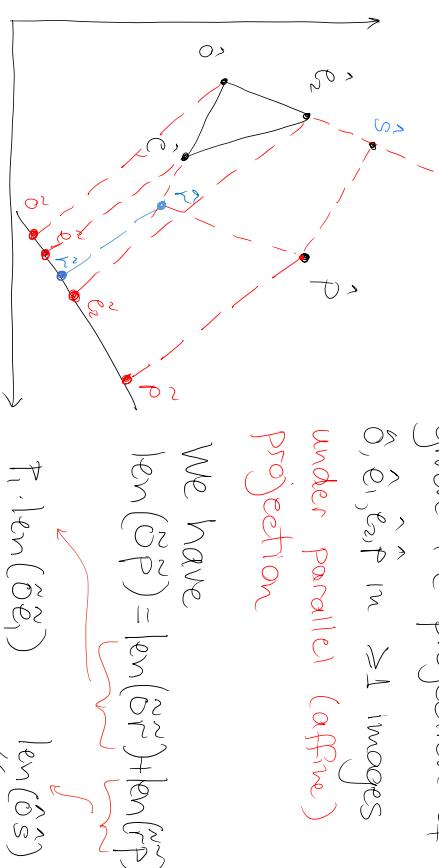
computing 2D affine structure from 1D images

38

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

Goal: Compute (P_1, P_2, P_3)

given the projection of
 $\hat{o}, \hat{e}_1, \hat{e}_2, \hat{p}$ in ≥ 1 images
 under parallel (affine) projection



computing 2D affine structure from 1D images

40

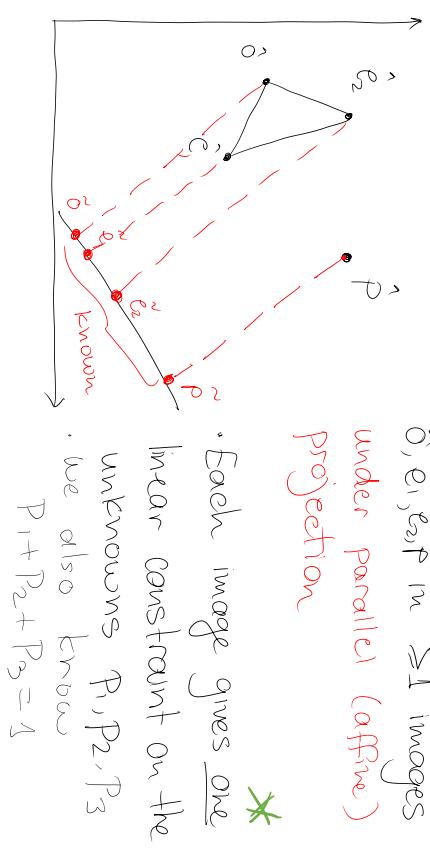
$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

Goal: Compute (P_1, P_2, P_3)

given the projection of
 $\hat{o}, \hat{e}_1, \hat{e}_2, \hat{p}$ in ≥ 1 images
 under parallel (affine) projection

★

- Each image gives one linear constraint on the unknowns P_1, P_2, P_3
- we also know $P_1 + P_2 + P_3 = 1$



computing 2D affine structure from 1D images

39

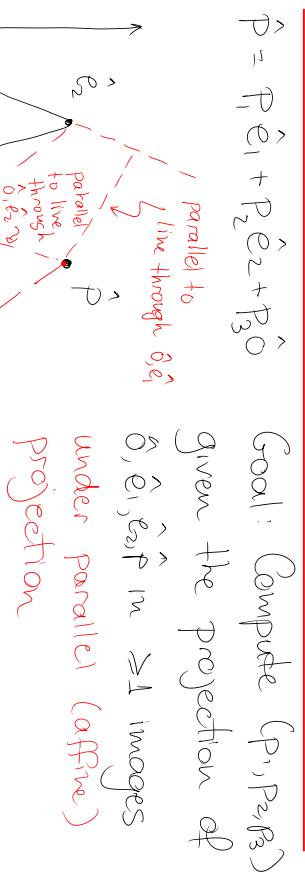
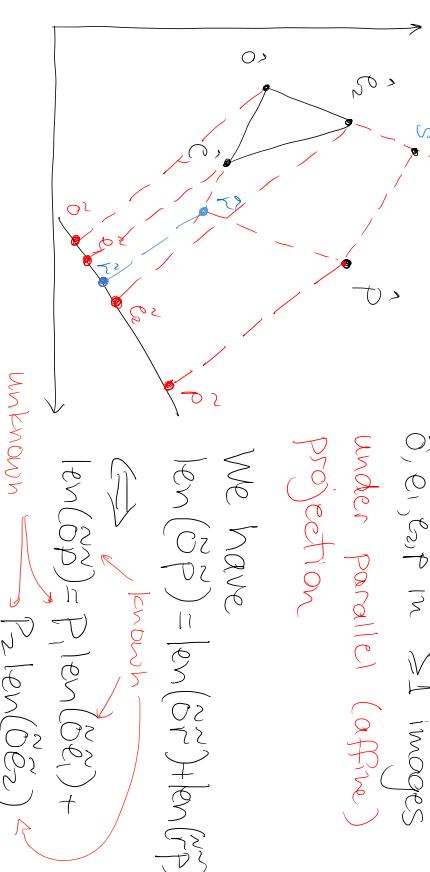
$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

Goal: Compute (P_1, P_2, P_3)

given the projection of
 $\hat{o}, \hat{e}_1, \hat{e}_2, \hat{p}$ in ≥ 1 images
 under parallel (affine) projection

★

- Each image gives one linear constraint on the unknowns P_1, P_2, P_3
- we also know $P_1 + P_2 + P_3 = 1$



computing 2D affine structure from 1D images

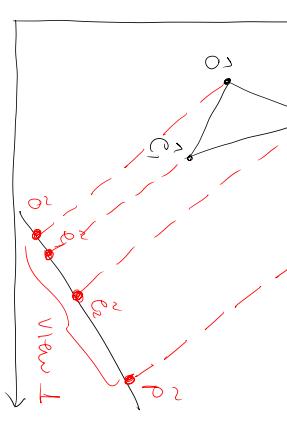
37

computing 2D affine structure: 2 views, 4 points

41

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

Note that the second view can be acquired under any arbitrary rearrangement of the reference



computing 3D affine structure

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{e}_3 + P_4 \hat{o} \quad \sum_{i=1}^4 P_i = 1$$

$\nabla \hat{P} = P_1 \nabla \hat{e}_1 + P_2 \nabla \hat{e}_2 + P_3 \nabla \hat{e}_3 + P_4 \nabla \hat{o}$

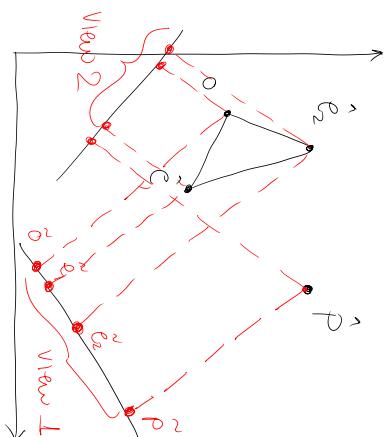
Image plane

42

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

Goal: Compute (P_1, P_2, P_3) given the projection of $\hat{o}, \hat{e}_1, \hat{e}_2, \hat{p}$ in ≥ 1 images

Structure recovery is possible for any 2 views of an affinely distorted 2D pointset under parallel (affine) projection



computing 3D affine structure

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{e}_3 + P_4 \hat{o} \quad \sum_{i=1}^4 P_i = 1$$

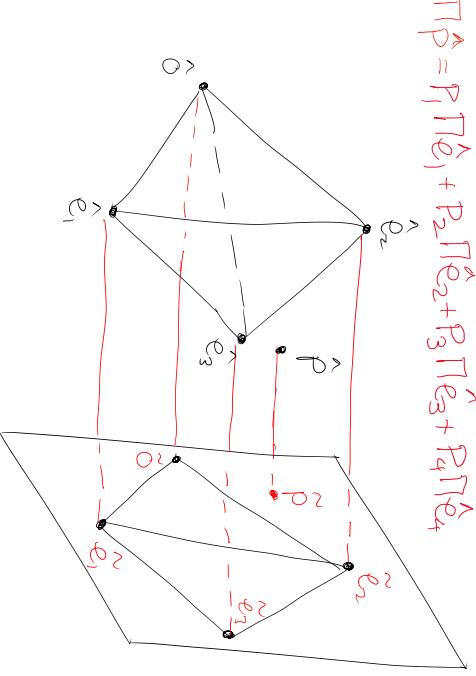
$\nabla \hat{P} = P_1 \nabla \hat{e}_1 + P_2 \nabla \hat{e}_2 + P_3 \nabla \hat{e}_3 + P_4 \nabla \hat{o}$

Image plane

43

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{o}$$

\Rightarrow 2 views of $N \geq 4$ points necessary & sufficient to recover affine 2D structure



44

computing 3D affine structure

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{e}_3 + P_4 \hat{o} \quad \sum_{i=1}^4 P_i = 1$$

$\nabla \hat{P} = P_1 \nabla \hat{e}_1 + P_2 \nabla \hat{e}_2 + P_3 \nabla \hat{e}_3 + P_4 \nabla \hat{o}$

Image plane

45

3D affine structure: constraints from 1 view

46

computing 3D affine structure

$$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2 + P_3 \hat{e}_3 + P_4 \hat{o} \quad \sum_{i=1}^4 P_i = 1$$

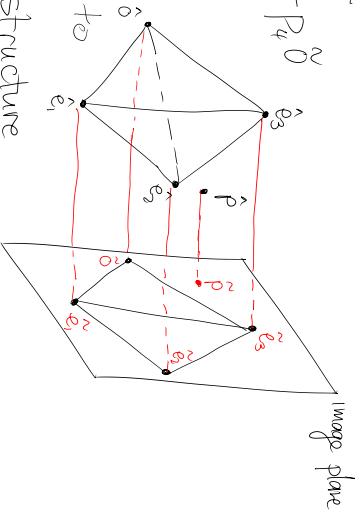
Goal: compute P_1, \dots, P_4 from $N \geq 1$ views

One view gives us:

$$\tilde{P} = \tilde{P}_1 \tilde{e}_1 + \tilde{P}_2 \tilde{e}_2 + \tilde{P}_3 \tilde{e}_3 + \tilde{P}_4 \tilde{o}$$

$$\Rightarrow \sum_{i=1}^4 \tilde{P}_i = 1$$

\Rightarrow 2 constraints
in 4 unknowns
 \Rightarrow need $N \geq 2$ views to
get 3D affine structure



48

3D affine structure: to read further

47

J. J. Koenderink and A. J. van Doorn,
"Affine Structure from Motion,"
J. Opt. Soc. Am., Vol. 8, no. 2, 1991

Structure from Motion (SfM)

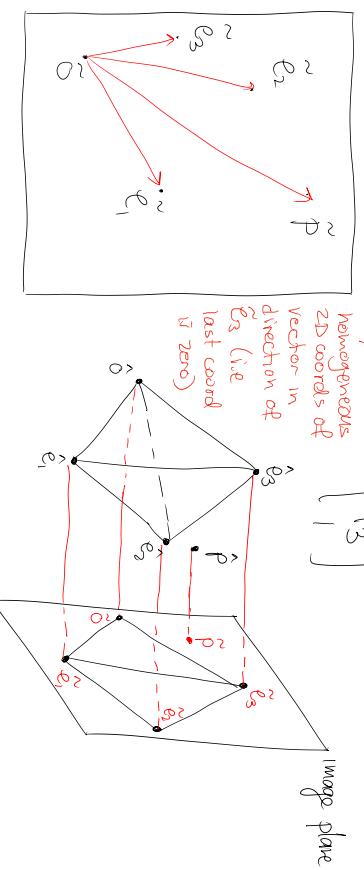
- big picture
- Euclidean & affine structure
- the factorization algorithm
- projective structure
- projective factorization
- SfM in Photo Tourism
 - projective matrices & Euclidean upgrade from F
 - incremental SfM

computing affine 3D structure from N views

50

$$\hat{\vec{P}} = \begin{bmatrix} \text{3x4 projection matrix} \\ \vec{e}_1 - \vec{o} & \vec{e}_2 - \vec{o} & \vec{e}_3 - \vec{o} & \vec{0} \end{bmatrix} \begin{bmatrix} \text{homogeneous 2D coords of origin} \\ P_1 \\ P_2 \\ P_3 \end{bmatrix}$$

non-homogeneous 3D coords
cf. \vec{P}



computing affine 3D structure from N views

52

- For the i-th input image

$$\vec{P}_i = \begin{bmatrix} \vec{e}_1 - \vec{o} & \vec{e}_2 - \vec{o} & \vec{e}_3 - \vec{o} & \vec{0} \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ 1 \end{bmatrix}$$

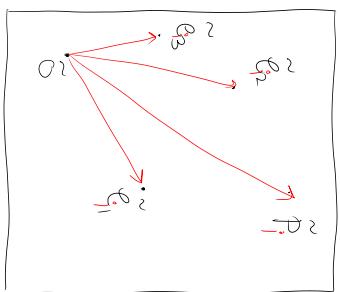
last row is always [0 0 0 1] so redundant

- Equation becomes (in non-homogeneous coords):

$$\begin{bmatrix} \vec{P}_i^x \\ \vec{P}_i^y \end{bmatrix} - \vec{o} = \begin{bmatrix} \vec{e}_1 - \vec{o} & \vec{e}_{i2} - \vec{o} & \vec{e}_{i3} - \vec{o} \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix}$$

2D vectors

\vec{P}_i^z
2D vector
 \vec{P}_i
3D vector
of affine coords



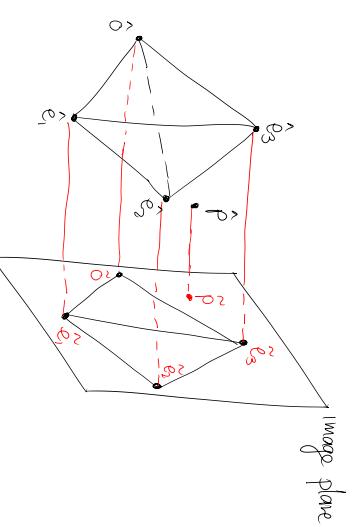
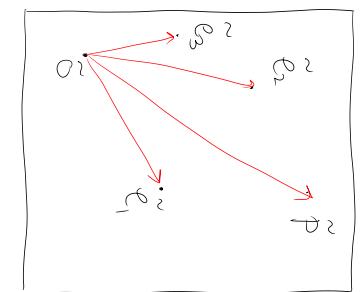
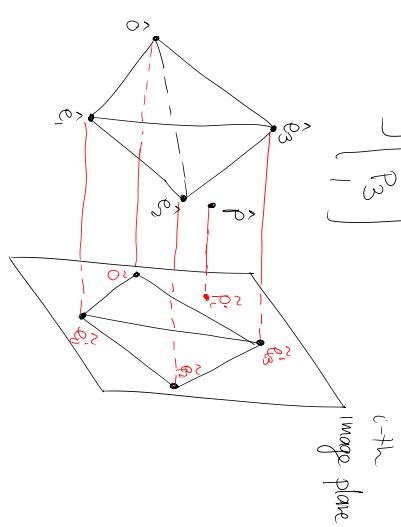
computing affine 3D structure from N views

51

- For the i-th input image

$$\vec{P}_i = \begin{bmatrix} \vec{e}_1 - \vec{o} & \vec{e}_2 - \vec{o} & \vec{e}_3 - \vec{o} & \vec{0} \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ 1 \end{bmatrix}$$

i-th



computing affine 3D structure from N views

49

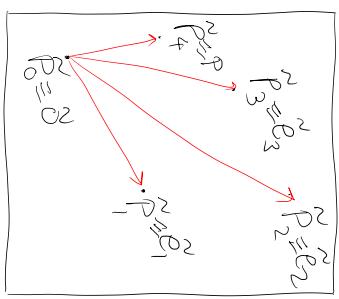
$$\hat{\vec{P}} = \vec{P}_1(\hat{\vec{e}}_1 - \hat{\vec{o}}) + \vec{P}_2(\hat{\vec{e}}_2 - \hat{\vec{o}}) + \vec{P}_3(\hat{\vec{e}}_3 - \hat{\vec{o}}) + \hat{\vec{o}}$$

$$\hat{\vec{P}} = \vec{P}_1(\vec{e}_1 - \vec{o}) + \vec{P}_2(\vec{e}_2 - \vec{o}) + \vec{P}_3(\vec{e}_3 - \vec{o}) + \vec{o}$$

computing affine 3D structure for M+1 points

54

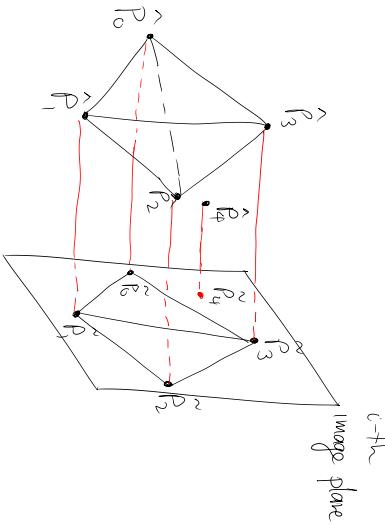
For $M+1$ points in a frame, 4 act as reference frame and the rest play the role of point \hat{P}



affine 3D structure for N views, M+1 points

55

$$\begin{aligned}
 & \text{For } N \text{ frames and } M+1 \text{ points} \\
 & \left[\begin{array}{c|c|c|c}
 \bar{P}_1^x & \bar{P}_1^y & \bar{P}_1^z & \bar{P}_M^x \\
 \bar{P}_2^x & \bar{P}_2^y & \bar{P}_2^z & \bar{P}_M^y \\
 \bar{P}_3^x & \bar{P}_3^y & \bar{P}_3^z & \bar{P}_M^z \\
 \hline
 \bar{P}_N^x & \bar{P}_N^y & \bar{P}_N^z &
 \end{array} \right] = \left[\begin{array}{c|c|c|c}
 \bar{P}_1^x & \bar{P}_1^y & \bar{P}_1^z & \bar{P}_M^x \\
 \bar{P}_2^x & \bar{P}_2^y & \bar{P}_2^z & \bar{P}_M^y \\
 \bar{P}_3^x & \bar{P}_3^y & \bar{P}_3^z & \bar{P}_M^z \\
 \hline
 \bar{P}_N^x & \bar{P}_N^y & \bar{P}_N^z &
 \end{array} \right] \\
 & \text{Projection of } \hat{e}_1 \text{ in first image in } N\text{-th image} \\
 & \text{if the first 3 points are chosen as reference points } \hat{e}_1, \hat{e}_2, \hat{e}_3 \text{ their affine coords are as shown} \\
 & \text{all } N \text{ views centered on a single point}
 \end{aligned}$$



affine 3D structure for N views, M+1 points

55

$$\begin{aligned}
 & \text{For } N \text{ frames and } M+1 \text{ points} \\
 & \left[\begin{array}{c|c|c|c}
 \bar{P}_1^x & \bar{P}_1^y & \bar{P}_1^z & \bar{P}_M^x \\
 \bar{P}_2^x & \bar{P}_2^y & \bar{P}_2^z & \bar{P}_M^y \\
 \bar{P}_3^x & \bar{P}_3^y & \bar{P}_3^z & \bar{P}_M^z \\
 \hline
 \bar{P}_N^x & \bar{P}_N^y & \bar{P}_N^z &
 \end{array} \right] = \left[\begin{array}{c|c|c|c}
 \bar{e}_{11}^x & \bar{e}_{12}^x & \bar{e}_{13}^x & \bar{P}_1^x \\
 \bar{e}_{11}^y & \bar{e}_{12}^y & \bar{e}_{13}^y & \bar{P}_2^x \\
 \bar{e}_{11}^z & \bar{e}_{12}^z & \bar{e}_{13}^z & \bar{P}_3^x \\
 \hline
 \bar{P}_3^z & \bar{P}_2^z & \bar{P}_1^z &
 \end{array} \right]
 \end{aligned}$$

where $\bar{P}_i = \hat{P}_i - \hat{O}$

$$\bar{e}_{ij} = \hat{e}_j - \hat{O}$$

are inhomogeneous coordinates relative to projection δ of origin \hat{O}

$$\begin{aligned}
 & \text{a single view of all points} \\
 & \text{projection of } \hat{e}_3 \text{ in } N\text{-th image} \\
 & \text{if the first 3 points are chosen as reference points } \hat{e}_1, \hat{e}_2, \hat{e}_3 \text{ their affine coords are as shown} \\
 & \text{all } N \text{ views centered on a single point}
 \end{aligned}$$

3D affine coordinates of M -th point

factorization algorithm: the rank-3 constraint

58

✓ Data matrix C

$$\begin{bmatrix} P_1^x & \dots & P_M^x \\ P_1^y & \dots & P_M^y \\ P_1^z & \dots & P_M^z \end{bmatrix} = \begin{bmatrix} P_{11}^x & P_{11}^y & P_{11}^z \\ P_{21}^x & P_{21}^y & P_{21}^z \\ \vdots & \vdots & \vdots \\ P_{N1}^x & P_{N1}^y & P_{N1}^z \end{bmatrix} \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix}$$

$$\begin{bmatrix} P_1^x & \dots & P_M^x \\ P_1^y & \dots & P_M^y \\ P_1^z & \dots & P_M^z \end{bmatrix} = \begin{bmatrix} P_{1N}^x & P_{1N}^y & P_{1N}^z \\ P_{2N}^x & P_{2N}^y & P_{2N}^z \\ \vdots & \vdots & \vdots \\ P_{NN}^x & P_{NN}^y & P_{NN}^z \end{bmatrix} \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix}$$

* Matrix has rank 3 for $M+1 \geq 4$ centered points

shape matrix
D
← camera motion matrix M

factorization algorithm: affine structure from SVD⁶⁰

We seek a decomposition of C

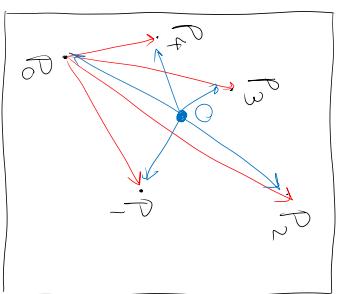
$$C = \begin{bmatrix} \text{Data matrix} \\ \text{Camera matrix} \\ \text{Shape matrix} \\ D \end{bmatrix}$$

② Define the basis via SVD

Idea: No need to commit to a specific set of points for reference frame

factorization algorithm: basic steps

59



① Under orthography,

the projection of the centroid of $M+1$ points is the centroid of their projection
⇒ choose centroid for origin

Idea: No need to commit to a specific set of points for reference frame

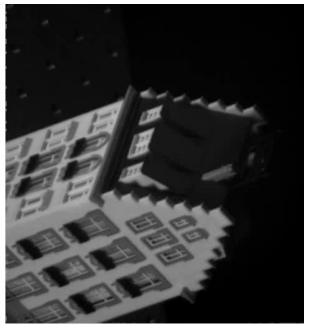
factorization algorithm: data centering

57

⇒ choose centroid for origin

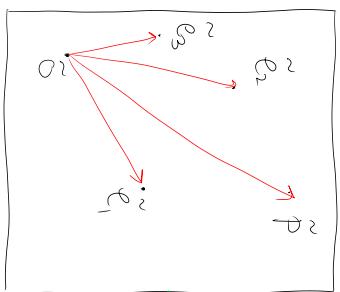
factorization algorithm: basic steps *

66

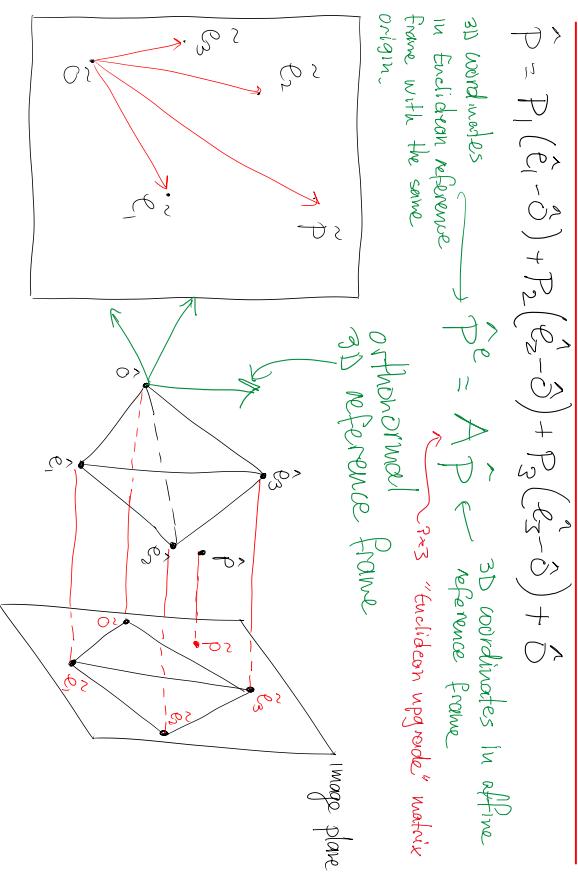


factorization algorithm: basic steps

- 6) Optionally:
- ① Subtract centroid from all point projections \hat{p}_j
 - ② Form matrix C
 - ③ Compute SVD of C
 - ④ Compute SVD of C
 - ⑤ Return $M = W \Sigma V^\top$
- $C = W \Sigma V^\top$
- "Metric/Euclidean Upgrade"



computing affine 3D structure from N views



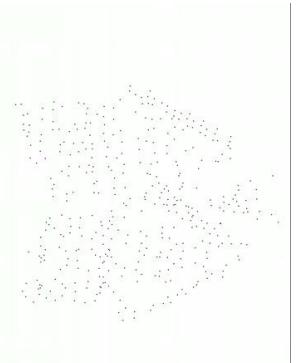
67

- ① Subtract centroid from all point projections \hat{p}_j
- ② Form matrix C
- ③ Compute SVD of C

$$C = W \Sigma V^\top$$

factorization algorithm: basic steps *

65



factorization algorithm: Euclidean upgrade

74

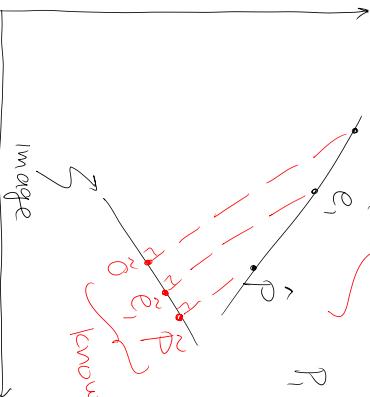
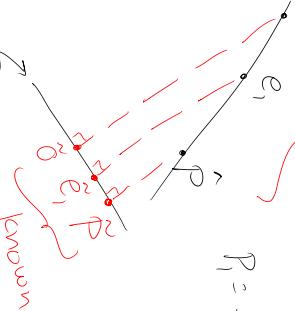
Putting it all together:

$$\begin{aligned}
 & \xrightarrow{\text{affine camera matrices}} M = W U_q \Lambda_q^{1/2} V^T \quad \xrightarrow{\text{Euclidean upgrade}} \text{unknown global scale} \\
 D = & \boxed{R^T \frac{1}{r} \left[\begin{array}{c|c} \Lambda_q^{1/2} & U_q V^T \\ \hline U_q & \sum \sqrt{ }^T \end{array} \right]} \quad \xrightarrow{\text{rotation}} \text{unknown global} \\
 & \xrightarrow{\text{unknown rotation \& scale}} \xrightarrow{\text{upgrade inverse}} \text{affine shape matrix}
 \end{aligned}$$

affine-invariant representations

$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{e}_2$ The affine coordinate P_1 is an affine invariant because it is not affected by affine transformations & projections

$$P_1 = \frac{\text{len}(\hat{P} \hat{o})}{\text{len}(\hat{e}_1 \hat{o})} = \frac{\text{len}(\hat{P} \tilde{o})}{\text{len}(\hat{e}_1 \tilde{o})}$$



75

Topic 04:

Structure from Motion (SfM)

- big picture

- Euclidean & affine structure

- the factorization algorithm

- **projective structure**

- projective factorization

- SfM in Photo Tourism

- projective matrices & Euclidean upgrade from F

- incremental SfM

- ④ We now perform an eigenvalue decomposition of Q to get $Q = U_q \Lambda_q U_q^T$ (assuming non-negative eigenvalues) $\Rightarrow A = U_q \Lambda_q^{1/2}$
- ③ Define $Q = AA^T$ (= symmetric positive definite so 6 unknown elements) For $N \geq 3$ we can solve for Q up to scale
- ② Then $W_i A A^T W_i^T = S_i \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} R R^T \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} S_i = \begin{bmatrix} S_i^{-2} & 0 \\ 0 & S_i^{-2} \end{bmatrix}$ \Rightarrow each W_i gives 2 linear constraints on AA^T homogeneous

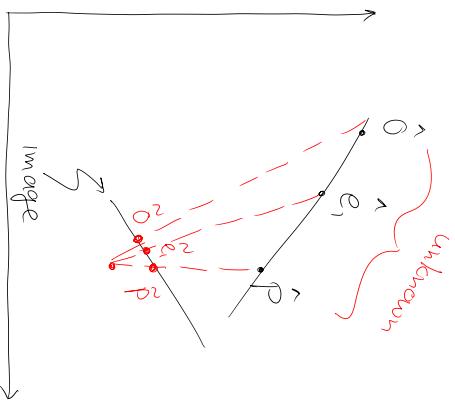
76

73

projectively-invariant representations

78

Goal: Represent \hat{P} in a way that is NOT affected by perspective projection



the cross-ratio is a projective invariant

① Need a 3rd point, \hat{u} , to define "unit" on reference frame

Goal: Represent \hat{P} in a way that is NOT affected by perspective projection

② Then the cross-ratio is invariant:

$$\text{cross}(\hat{o}, \hat{e}_1, \hat{p}, \hat{u}) = \frac{\text{len}(\hat{p}\hat{o})}{\text{len}(\hat{p}\hat{e}_1)} \cdot \frac{\text{len}(\hat{u}\hat{o})}{\text{len}(\hat{u}\hat{e}_1)}$$

Note: The cross-ratio depends on angles between rays, not point projections.

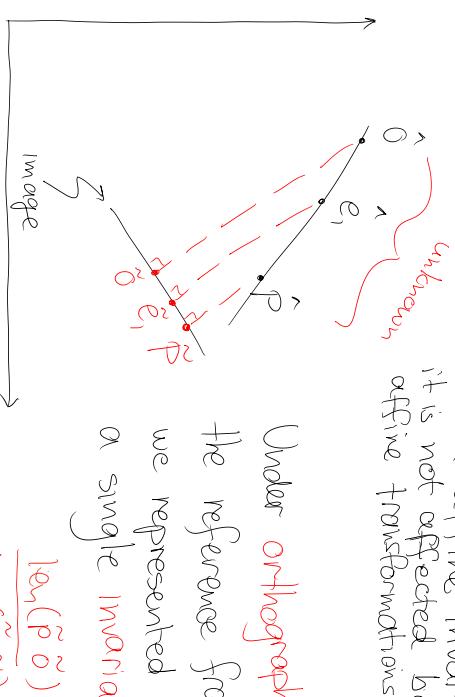
80

projectively-invariant representations

Goal: Represent \hat{P} in a way that is NOT affected by perspective projection

But ratios of distances are NOT preserved under perspective projection!

Are there projectively invariant quantities?



$$\frac{\text{len}(\hat{p}\hat{o})}{\text{len}(\hat{e}_1\hat{o})}$$

$\hat{P} = P_1 \hat{e}_1 + P_2 \hat{o}$ The affine coordinate P_1 is an affine invariant because it is not affected by affine transformations & projections

Under orthography, given the reference frame $\hat{O}\hat{e}_1$ we represent \hat{P} using a single invariant ratio

79

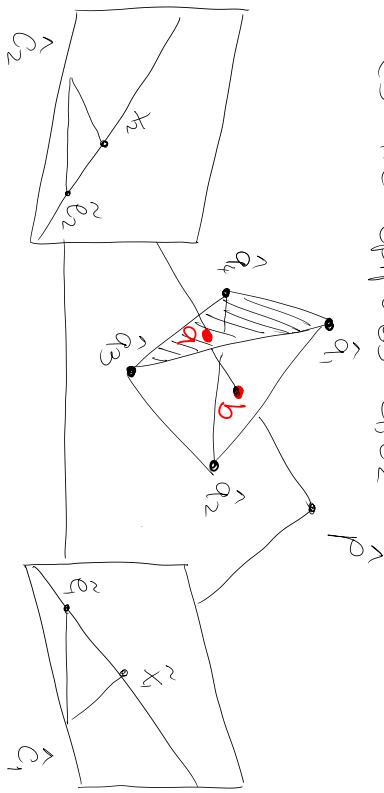
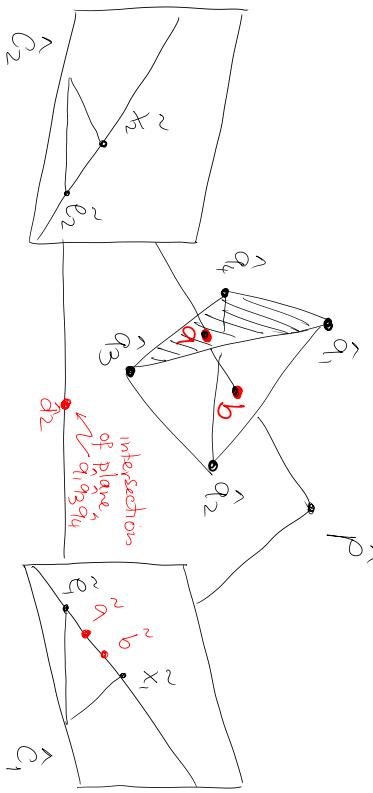
affine-invariant representations

77

Let us reconstruct a point \hat{p} given

- (a) the projection of 4 non-coplanar points in 2 views : $\hat{q}_1, \hat{q}_2, \hat{q}_3, \hat{q}_4$

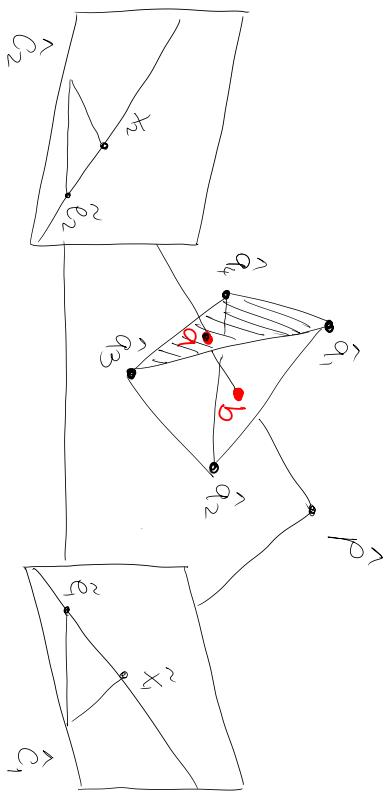
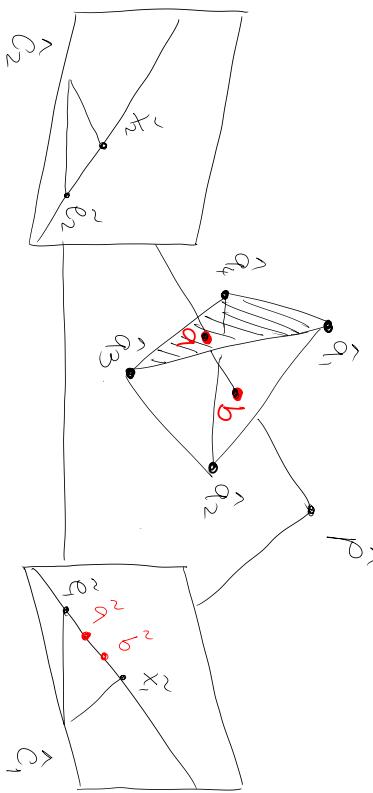
- (b) the epipoles \hat{e}_1, \hat{e}_2



feasibility of computing projective structure from 2 views: a geometric proof

Localizing \hat{a} :

- a lies on the plane containing $\hat{q}_1, \hat{q}_3, \hat{q}_4, \hat{d}_2$
- these points project to $\hat{q}_1, \hat{q}_3, \hat{q}_4, e_1$ in image 1 and $\hat{q}_1, \hat{q}_3, \hat{q}_4, \hat{e}_2$ in image 2



Let us reconstruct a point \hat{p} given

- (a) the projection of 4 non-coplanar points in 2 views : $\hat{q}_1, \hat{q}_2, \hat{q}_3, \hat{q}_4$

- (b) cross(\hat{C}_2, a, b, \hat{p}) = cross($\hat{C}_1, \hat{a}, \hat{b}, \hat{p}$)

We will do this by computing α cross-ratio that localizes \hat{p} uniquely: can compute this as long as we can locate \hat{a}, \hat{b}

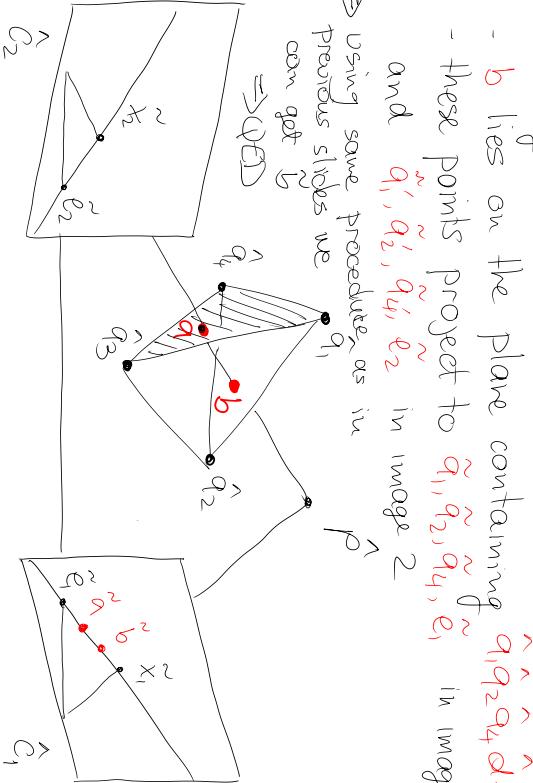
feasibility of computing projective structure from 2 views: a geometric proof

Localizing \hat{b}

- \hat{b} lies on the plane containing $\hat{q}_1, \hat{q}_2, \hat{q}_4, \hat{d}_1$
- these points project to $\hat{q}_1, \hat{q}_2, \hat{q}_4, \hat{e}_1$ in image 1 and $\hat{q}'_1, \hat{q}'_2, \hat{q}'_4, \hat{e}'_1$ in image 2

\Rightarrow using same procedure as in previous slides we can get \hat{b}

$\Rightarrow QED$

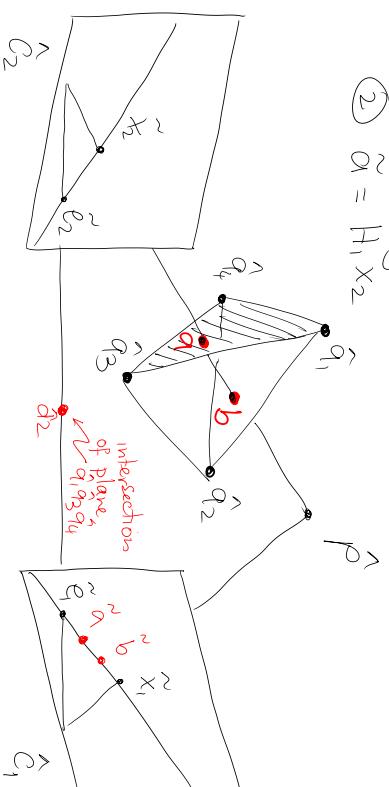


88

Localizing \hat{a} :

- \hat{a} lies on the plane containing $\hat{q}_1, \hat{q}_3, \hat{q}_4, \hat{d}_2$
- \Rightarrow ① compute the homography H_1 mapping image 2 to image 1

$$\textcircled{2} \quad \hat{a} = H_1 \hat{x}_2$$



87

Topic 04:

projective structure from 2 views: to read further

For more details on this construction:

A. Shashua, "Projective Depth: A geometric invariant for 3D reconstruction from two perspective/orthographic views and for visual recognition," Proc. ICCV 1993

- big picture
- Euclidean & affine structure
- the factorization algorithm
- projective structure
- projective factorization
- SFM in Photo Tourism
 - projective matrices & Euclidean upgrade from F
 - incremental SFM

Topic 04:

projective factorization (see notes)

Suppose we knew the projective depths $z_{j,n}$. We could then build a new data matrix C in terms of the 3D points $z_{j,n}\vec{p}_{j,n}$. Accordingly, we can form $C = [z_{j,n}\vec{p}_{j,n}]$, by stacking the depth-scaled image points just as we did in the orthographic case. Now C is $3J \times N$.

Structure from Motion (SfM)

- big picture
- Euclidean & affine structure
- the factorization algorithm
- projective structure
- projective factorization
- **SfM in Photo Tourism**
 - projective matrices & Euclidean upgrade from F
 - incremental SfM

Of course, this assumes we knew the correct depths $z_{j,n}$ that we used to form C . And as above, from the form of (19) it is clear that C is, at most, rank 4 (ignoring noise).

$$C = MP. \quad (19)$$

The perspective projection of the scene point \vec{P}_n onto the j^{th} image plane (1) can be written as

$$z_{j,n}\vec{p}_{j,n} = M_j \vec{P}_n,$$

where M_j is now 3×4 , and \vec{P}_n is 4×1 .

By stacking up the camera matrices, M_j , to form the $3J \times 4$ matrix M , and letting $P = (\vec{P}_1, \dots, \vec{P}_N)$ be the $4 \times N$ shape matrix, we obtain a factorization of the data matrix:

$$C = MP.$$

SfM in Photo Tourism: correspondence-finding

92

SfM in Photo Tourism: correspondence-finding

91

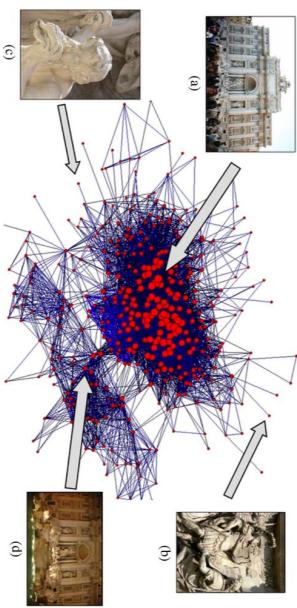


Fig. 6. Image connectivity graph for the Trevi Fountain. This graph contains a node (red dot) for each image in a set of photos of the Trevi Fountain, and an edge between each pair of photos with matching points. The size of a node is proportional to its degree. There are two dominant clusters corresponding to daytime photos (e.g., image (a)) and nighttime photos (image (d)). Similar views of the facade are clustered together in the center of the graph; white nodes in the periphery of the graph; e.g., (b) and (c), are more unusual (often closeup) views. An interactive version of this graph can be found at <http://phototour.cs.washington.edu/imagegraphs/Trevi/>.

SfM in Photo Tourism: 2-view shape & motion

94

3. Processing $\mathcal{I}_1, \mathcal{I}_2$:

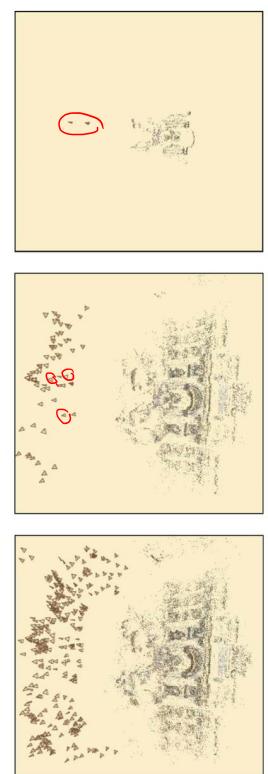
- compute projection matrices W_1, W_2
- compute the Euclidean projection matrices and the reconstructed point coordinates

SfM in Photo Tourism: incremental SfM

96

4. Current state:

- M points with already-estimated 3D coords
- Euclidean projection matrices for images $\mathcal{I}_1, \dots, \mathcal{I}_k$



95

Adding a new image \mathcal{I}_{k+1} :

- ① initialize the projection matrix from the projection of previously-reconstructed points that also project to \mathcal{I}_{k+1} (using DLT algorithm)

Fig. 7. Incremental SfM. Our incremental SfM approach reconstructs the scene a few cameras at a time. This sequence of images shows the Treff data set at three different stages of incremental reconstruction. Left: the initial two-frame reconstruction. Middle: an intermediate stage, after 60 images have been added. Right: the final reconstruction with 360 photos.

SfM in Photo Tourism: initial image pair selection

93

2. Initial pair selection

- find pair $\mathcal{I}_1, \mathcal{I}_2$ with
 - ① large # of successful matches
 - ② matches not consistent with a plane homography

Topic 04:

98

SfM in Photo Tourism: Incremental SfM

97

- Adding a new image I_{k+1}

- initialize the projection matrix from the projection of previously-reconstructed points that also project to I_{k+1} (using DLT algorithm)
- Add new matches to the reconstruction
- Run a non-linear optimizer to optimize projection matrices & point positions

- SfM in Photo Tourism
 - projective matrices & Euclidean upgrade from F
 - incremental SfM

projection matrices from the F-matrix

100

projection matrices from the F-matrix

99

- compute projection matrices W_1, W_2 from fundamental matrix F

- ② It is possible to show that W_2 can be chosen as (see Hartley & Zisserman, 2000, p. 237)

$$* \quad W_2 = \begin{bmatrix} [e_2]_\times F \\ e_2 \end{bmatrix}$$

left epipole, i.e. $e_2^\top F = 0$

- ① Choose W_1 so that camera 1 defines the projective frame of reference

$$* \quad W_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Euclidean upgrade: constraints from rotations

Topic 04:

Structure from Motion (SfM)

- big picture
- Euclidean & affine structure
- the factorization algorithm
- projective structure
- projective factorization
- SfM in Photo Tourism
 - projective matrices & Euclidean upgrade from F
 - incremental SfM

adding cameras: exterior orientation & intrinsics

108

4. Adding a new image \mathcal{I}_{k+1} :

- ① Initialize the projection matrix from
the projection of previously-reconstructed
points that also project to \mathcal{I}_{k+1}
(using DLT algorithm)

$$W_i \begin{bmatrix} Q \\ q \end{bmatrix} = \begin{bmatrix} f_i & 0 & 0 \\ 0 & f_i & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} R_i & t_i \end{bmatrix}$$

projective projection

Euclidean upgrade

unknown rotation

unknown focal length

107

SfM in Photo Tourism: 2-view shape & motion

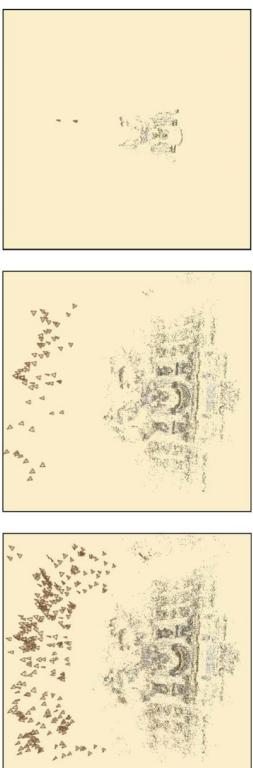


Fig. 7. Incremental SfM. Our incremental SfM approach reconstructs the scene a few cameras at a time. This sequence of images shows the Treff data set at three different stages of incremental reconstruction. Left: the initial two-frame reconstruction. Middle: an intermediate stage, after 60 images have been added. Right: the final reconstruction with 360 photos.

- ② Add new matches to the reconstruction
- ③ Run a non-linear optimizer to optimize projection matrices & point positions

⇒ We can compute Q^∞ from 2 views
⇒ Once Q^∞ is known, we can compute the upgrade matrix using diagonalization of Q^∞

(see notes p.23 and Forsyth & Ponce p. 293)

Once Q is known, we can choose an arbitrary q and thus fully determine R_i and t_i .

unknown origin of world coordinate system (free parameter)

intrinsic calibration matrix unknown translation

3x4 W_i $\begin{bmatrix} Q \\ q \end{bmatrix}$ $\begin{bmatrix} f_i & 0 & 0 \\ 0 & f_i & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} R_i & t_i \end{bmatrix}$

DLT algorithm for projection matrix estimation

110

Given $M \geq 6$ correspondences between world points and image points, compute projection matrix Π

Basic relation:

$$x_i \approx \Pi X_i$$

known homogeneous \rightarrow \uparrow known homogeneous 3D world
2D coords equality up + to scale

After computing Π , we express it in the form $[f \ f^0 \ 0^1] [R \ t]$ for some f and rotation matrix R (see Hartley & Zisserman p.143, p.170)

bundle adjustment: objective function

Objective: minimize total reprojection error

$$\ast \sum_{i=1}^N \sum_{j=1}^M w_{ij} \|x_{ij} - \hat{P}(\Theta_i, X_{ij})\|^2$$

projection operations
 1 if point j appears in image i
 2D image coordinates of point j in image i
 (w_{ij}, f_i, c_i)
 3 unknown rotation angles
 unknown image center
 focal length

Basic relation:

$$x_i \approx \Pi X_i = 0 \quad (\text{linear in the } M \text{ elements of } \Pi)$$

known homogeneous \rightarrow \uparrow known homogeneous 3D world
2D coords equality up + to scale

$\ast \Rightarrow$ 2 linearly-independent eqs per point so ≥ 6 points are needed to get Π

bundle adjustment

- ① Initialize the projection matrix from the projection of previously-reconstructed points that also project to \mathcal{I}_{k+1} (using DLT algorithm)
 - ② Add new matches to the reconstruction
 - ③ Run a non-linear optimizer to optimize projection matrices & point positions
- known
image
center

DLT algorithm for projection matrix estimation

109

Given $M \geq 6$ correspondences between world points and image points, compute projection matrix Π

bundle adjustment: focal length estimation

114

two sources of information:

① image EXIF data focal length entry

② result of DLT algorithm

- prefer ① unless there is a big difference between the two (then use ②)

$$\begin{pmatrix} \vec{\omega}_i & f_i & \vec{c}_i \end{pmatrix}_{\text{unknown}} \quad \begin{pmatrix} \text{camera center} \\ \text{focal length} \\ \text{rotation angles} \end{pmatrix}_{\text{known}}$$

bundle adjustment

Objective: minimize total reprojection error

$$\sum_{i=1}^N \sum_{j=1}^M w_{ij} \| \vec{x}_{ij} - \vec{P}(\Theta_i, \vec{X}_{ij}) \|^2$$

projection operations

\vec{x}_{ij} → 3D coordinates of point j in image i

w_{ij} → 2D image coordinates of point j in image i

$\vec{P}(\Theta_i, \vec{X}_{ij})$ → params image i of camera i

115

bundle adjustment: projection function *

Objective: minimize total reprojection error

$$\begin{pmatrix} \vec{\omega}_i & f_i & \vec{c}_i \end{pmatrix}_{\text{known}} \quad \begin{pmatrix} \text{camera center} \\ \text{focal length} \\ \text{rotation angles} \end{pmatrix}_{\text{unknown}}$$

$$\sum_{i=1}^N \sum_{j=1}^M w_{ij} \| \vec{x}_{ij} - \vec{P}(\Theta_i, \vec{X}_{ij}) \|^2$$

projection operations

\vec{x}_{ij} → 3D coordinates of point j in image i

w_{ij} → 2D image coordinates of point j in image i

$\vec{P}(\Theta_i, \vec{X}_{ij})$ → params image i of camera i

$\vec{\omega}_i, f_i, \vec{c}_i$ → inliers → rotation matrix camera center

115

bundle adjustment: rotation matrix representation

113

we already have an initial estimate \vec{R}_j so we optimize an incremental adjustment

$$R(\theta, \vec{n}) = \left(I + \sin \theta [\vec{n}]_x + (1 - \cos \theta) [\vec{n}]_x^2 \right) R_j$$

with $\vec{n} = \vec{Q}_m \vec{n}_i$ and \vec{n}_i unit vector

*The objective typically optimized using sparse methods.

- Hartley & Zisserman, *Multiple View Geometry in Computer Vision*, 2000
- Triggs et al, "Bundle Adjustment – A modern Synthesis," *Vision Algorithms: Theory and Practice*, LNCS 1883, 2000
- Snavely et al, "Photo Tourism: Exploring Photo Collections in 3D," ACM Siggraph 2006
- Agrawal et al, "Building Rome in a Day," *Proc. ICCV 2009*
- Pollefeys et al, Self-calibration and metric reconstruction in spite of varying and unknown intrinsic camera parameters, *IJCV* 1999

