

RWTH AACHEN
UNIVERSITY

Computer Vision - Lecture 21

Structure-from-Motion

01.02.2017

Computer Vision WS 15/16

Bastian Leibe
RWTH Aachen
<http://www.vision.rwth-aachen.de>
leibe@vision.rwth-aachen.de

Many slides adapted from Svetlana Lazebnik, Martial Hebert, Steve Seitz

RWTH AACHEN
UNIVERSITY

Announcements

- Exam
 - 1st Date: Friday, 24.02., 09:00 - 12:30h
 - 2nd Date: Thursday, 30.03., 09:30 - 12:30h
 - Closed-book exam, the core exam time will be 2h.
 - We will send around an announcement with the exact starting times and places by email.
- Test exam
 - We will give out a test exam via L2P
 - Purpose: Prepare you for the types of questions you can expect.
- Exchange students
 - If you need a special exam slot due to travel, contact me!

B. Leibe 2

RWTH AACHEN
UNIVERSITY

Announcements (2)

- Last lecture next Monday: Repetition
 - Summary of all topics in the lecture
 - “Big picture” and current research directions
 - Opportunity to ask questions
- Please use this opportunity and prepare questions!

Computer Vision WS 15/16

B. Leibe 3

RWTH AACHEN
UNIVERSITY

Course Outline

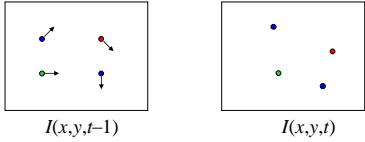
- Image Processing Basics
- Segmentation & Grouping
- Object Recognition
- Local Features & Matching
- Object Categorization
- 3D Reconstruction
 - Epipolar Geometry and Stereo Basics
 - Camera calibration & Uncalibrated Reconstruction
 - Active Stereo
- Motion
 - Motion and Optical Flow
- 3D Reconstruction (Reprise)
 - Structure-from-Motion

Computer Vision WS 15/16

4

RWTH AACHEN
UNIVERSITY

Recap: Estimating Optical Flow



- Given two subsequent frames, estimate the apparent motion field $u(x,y)$ and $v(x,y)$ between them.
- Key assumptions
 - Brightness constancy: projection of the same point looks the same in every frame.
 - Small motion: points do not move very far.
 - Spatial coherence: points move like their neighbors.

Computer Vision WS 15/16

B. Leibe 5

Slide credit: Svetlana Lazebnik

RWTH AACHEN
UNIVERSITY

Recap: Lucas-Kanade Optical Flow

- Use all pixels in a $K \times K$ window to get more equations.
- Least squares problem:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix} \quad A \quad d = b$$

25×2
 2×1
 25×1
- Minimum least squares solution given by solution of

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$
 $A^T b$

Recall the Harris detector!

Computer Vision WS 15/16

B. Leibe 6

Slide adapted from Svetlana Lazebnik

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Recap: Iterative Refinement

- Estimate velocity at each pixel using one iteration of LK estimation.
- Warp one image toward the other using the estimated flow field.
- Refine estimate by repeating the process.
- Iterative procedure
 - Results in subpixel accurate localization.
 - Converges for small displacements.

Slide adapted from Steve Seitz. B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Recap: Coarse-to-fine Estimation

Slide credit: Steve Seitz. B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Recap: Coarse-to-fine Estimation

Slide credit: Steve Seitz. B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Structure from Motion (SfM)
 - Motivation
 - Ambiguity
- Affine SfM
 - Affine cameras
 - Affine factorization
 - Euclidean upgrade
 - Dealing with missing data
- Projective SfM
 - Two-camera case
 - Projective factorization
 - Bundle adjustment
 - Practical considerations
- Applications

Computer Vision WS 15/16 B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Structure from Motion

- Given: m images of n fixed 3D points

$$x_{ij} = P_i X_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$
- Problem: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}

Slide credit: Svetlana Lazebnik. B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

What Can We Use This For?

- E.g. movie special effects

Video

Computer Vision WS 15/16 B. Leibe Video Credit: Stefan Hafensser

RWTH AACHEN UNIVERSITY

Structure from Motion 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 in the image remain exactly the same:

$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left(\frac{1}{k}\mathbf{P}\right)(k\mathbf{X})$$

⇒ It is impossible to recover the absolute scale of the scene!

Computer Vision WS 15/16 13

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Structure from Motion 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 in the image remain exactly the same.
- More generally: if we transform the scene using a transformation Q and apply the inverse transformation to the camera matrices, then the images do not change

$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}^{-1})\mathbf{Q}\mathbf{X}$$

Computer Vision WS 15/16 14

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Reconstruction Ambiguity: Similarity

Similarity

$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}_S^{-1})\mathbf{Q}_S\mathbf{X}$$

Computer Vision WS 15/16 15

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

RWTH AACHEN UNIVERSITY

Reconstruction Ambiguity: Affine

Affine

$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}_A^{-1})\mathbf{Q}_A\mathbf{X}$$

Computer Vision WS 15/16 16

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

RWTH AACHEN UNIVERSITY

Reconstruction Ambiguity: Projective

Projective

$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}_P^{-1})\mathbf{Q}_P\mathbf{X}$$

Computer Vision WS 15/16 17

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

RWTH AACHEN UNIVERSITY

Projective Ambiguity

Computer Vision WS 15/16 18

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

RWTH AACHEN UNIVERSITY

From Projective to Affine

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

19

RWTH AACHEN UNIVERSITY

From Affine to Similarity

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

20

RWTH AACHEN UNIVERSITY

Hierarchy of 3D Transformations

Projective 15dof	$\begin{bmatrix} A & t \\ v^T & v \end{bmatrix}$		Preserves intersection and tangency
Affine 12dof	$\begin{bmatrix} A & t \\ 0^T & 1 \end{bmatrix}$		Preserves parallelism, volume ratios
Similarity 7dof	$\begin{bmatrix} sR & t \\ 0^T & 1 \end{bmatrix}$		Preserves angles, ratios of length
Euclidean 6dof	$\begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix}$		Preserves angles, lengths

- With no constraints on the camera calibration matrix or on the scene, we get a *projective reconstruction*.
- Need additional information to *upgrade* the reconstruction to affine, similarity, or Euclidean.

Computer Vision WS 15/16 Slide credit: Svetlana Lazebnik B. Leibe

21

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Structure from Motion (SfM)
 - Motivation
 - Ambiguity
- **Affine SfM**
 - Affine cameras
 - Affine factorization
 - Euclidean upgrade
 - Dealing with missing data
- Projective SfM
 - Two-camera case
 - Projective factorization
 - Bundle adjustment
 - Practical considerations
- Applications

Computer Vision WS 15/16 B. Leibe

22

RWTH AACHEN UNIVERSITY

Structure from Motion

- Let's start with *affine cameras* (the math is easier)

Computer Vision WS 15/16 Slide credit: Svetlana Lazebnik B. Leibe Images from Hartley & Zisserman

23

RWTH AACHEN UNIVERSITY

Orthographic Projection

- Special case of perspective projection
 - Distance from center of projection to image plane is infinite

- Projection matrix:

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \Rightarrow (x, y)$$

Computer Vision WS 15/16 Slide credit: Steve Seitz B. Leibe

24

RWTH AACHEN UNIVERSITY

Affine Cameras

Orthographic Projection

Parallel Projection

Computer Vision WS 15/16 25

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Affine Cameras

- A general affine camera combines the effects of an affine transformation of the 3D space, orthographic projection, and an affine transformation of the image:

$$P = [3 \times 3 \text{ affine}] \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} [4 \times 4 \text{ affine}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} A & b \\ 0 & 1 \end{bmatrix}$$

- Affine projection is a linear mapping + translation in inhomogeneous coordinates

$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = A\mathbf{X} + \mathbf{b}$$

Computer Vision WS 15/16 26

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Affine Structure from Motion

- Given: m images of n fixed 3D points:
 - $\mathbf{x}_{ij} = A_i \mathbf{X}_j + \mathbf{b}_i, \quad i = 1, \dots, m, j = 1, \dots, n$
- Problem: use the mn correspondences \mathbf{x}_{ij} to estimate m projection matrices A_i and translation vectors \mathbf{b}_i , and n points \mathbf{X}_j
- The reconstruction is defined up to an arbitrary affine transformation Q (12 degrees of freedom):

$$\begin{bmatrix} A & b \\ 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} A & b \\ 0 & 1 \end{bmatrix} Q^{-1}, \quad \begin{pmatrix} X \\ 1 \end{pmatrix} \rightarrow Q \begin{pmatrix} X \\ 1 \end{pmatrix}$$
- We have $2mn$ knowns and $8m + 3n$ unknowns (minus 12 dof for affine ambiguity).
 - Thus, we must have $2mn \geq 8m + 3n - 12$.
 - For two views, we need four point correspondences.

Computer Vision WS 15/16 27

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Affine Structure from Motion

- Centering: subtract the centroid of the image points

$$\hat{\mathbf{x}}_{ij} = \mathbf{x}_{ij} - \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik} = A_i \mathbf{X}_j + \mathbf{b}_i - \frac{1}{n} \sum_{k=1}^n (A_i \mathbf{X}_k + \mathbf{b}_i)$$

$$= A_i \left(\mathbf{X}_j - \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \right) = A_i \hat{\mathbf{X}}_j$$
- For simplicity, assume that the origin of the world coordinate system is at the centroid of the 3D points.
- After centering, each normalized point $\hat{\mathbf{x}}_{ij}$ is related to the 3D point \mathbf{X}_j by

$$\hat{\mathbf{x}}_{ij} = A_i \mathbf{X}_j$$

Computer Vision WS 15/16 28

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Affine Structure from Motion

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

$$D = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \dots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \dots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \dots & \hat{\mathbf{x}}_{mn} \end{bmatrix}$$

↓ Cameras (2m)
→ Points (n)

C. Tomasi and T. Kanade, *Shape and motion from image streams under orthography: A factorization method*, IJCV, 9(2):137-154, November 1992.

Computer Vision WS 15/16 29

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Affine Structure from Motion

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

$$D = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \dots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \dots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \dots & \hat{\mathbf{x}}_{mn} \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \dots & \mathbf{X}_n \end{bmatrix}$$

↓ Cameras (2m x 3)
→ Points (3 x n)

- The measurement matrix $D = MS$ must have rank 3!

C. Tomasi and T. Kanade, *Shape and motion from image streams under orthography: A factorization method*, IJCV, 9(2):137-154, November 1992.

Computer Vision WS 15/16 30

Slide credit: Svetlana Lazebnik B. Leibe

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Factorizing the Measurement Matrix

Measurements = Motion \times Shape

$$D = MS$$

Slide credit: Martial Hebert B. Leibe 31

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Factorizing the Measurement Matrix

- Singular value decomposition of D:

Slide credit: Martial Hebert 32

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Factorizing the Measurement Matrix

- Singular value decomposition of D:

To reduce to rank 3, we just need to set all the singular values to 0 except for the first 3

Slide credit: Martial Hebert 33

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Factorizing the Measurement Matrix

- Obtaining a factorization from SVD:

Slide credit: Martial Hebert 34

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Factorizing the Measurement Matrix

- Obtaining a factorization from SVD:

Possible decomposition:
 $M = U_3 W_3^{1/2}$ $S = W_3^{1/2} V_3^T$

This decomposition minimizes $|D-MS|^2$

Slide credit: Martial Hebert B. Leibe 35

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Affine Ambiguity

- The decomposition is not unique. We get the same D by using any 3×3 matrix C and applying the transformations $M \rightarrow MC$, $S \rightarrow C^{-1}S$.
- That is because we have only an affine transformation and we have not enforced any Euclidean constraints (like forcing the image axes to be perpendicular, for example). We need a *Euclidean upgrade*.

Slide credit: Martial Hebert B. Leibe 36

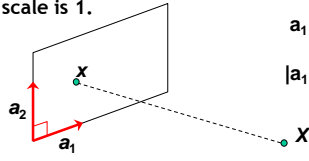
RWTH AACHEN UNIVERSITY

Estimating the Euclidean Upgrade

- Orthographic assumption: image axes are perpendicular and scale is 1.

$$\mathbf{a}_1 \cdot \mathbf{a}_2 = 0$$

$$|\mathbf{a}_1|^2 = |\mathbf{a}_2|^2 = 1$$



- This can be converted into a system of $3m$ equations:

$$\begin{cases} \hat{a}_{i1} \cdot \hat{a}_{i2} = 0 \\ |\hat{a}_{i1}| = 1 \\ |\hat{a}_{i2}| = 1 \end{cases} \Leftrightarrow \begin{cases} \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i2} = 0 \\ \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i1} = 1, \\ \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i2} = 1 \end{cases} \quad i = 1, \dots, m$$

for the transformation matrix $C \Rightarrow$ goal: estimate C

37

Computer Vision WS 15/16
Slide adapted from S. Lazebnik, M. Hebert, B. Leibe

RWTH AACHEN UNIVERSITY

Estimating the Euclidean Upgrade

- System of $3m$ equations:

$$\begin{cases} \hat{a}_{i1} \cdot \hat{a}_{i2} = 0 \\ |\hat{a}_{i1}| = 1 \\ |\hat{a}_{i2}| = 1 \end{cases} \Leftrightarrow \begin{cases} \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i2} = 0 \\ \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i1} = 1, \\ \mathbf{a}_i^T \mathbf{C} \mathbf{C}^T \mathbf{a}_{i2} = 1 \end{cases} \quad i = 1, \dots, m$$
- Let $L = \mathbf{C} \mathbf{C}^T$ $A_i = \begin{bmatrix} \mathbf{a}_{i1}^T \\ \mathbf{a}_{i2}^T \end{bmatrix}$, $i = 1, \dots, m$
- Then this translates to $3m$ equations in L

$$A_i L A_i^T = I, \quad i = 1, \dots, m$$
 - Solve for L
 - Recover C from L by Cholesky decomposition: $L = \mathbf{C} \mathbf{C}^T$
 - Update M and S : $M = \mathbf{M} \mathbf{C}$, $S = \mathbf{C}^{-1} \mathbf{S}$

38

Computer Vision WS 15/16
Slide adapted from S. Lazebnik, M. Hebert, B. Leibe

RWTH AACHEN UNIVERSITY

Algorithm Summary

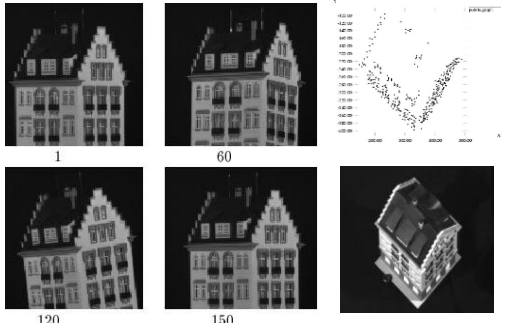
- Given: m images and n features x_{ij}
- For each image i , center the feature coordinates.
- Construct a $2m \times n$ measurement matrix D :
 - Column j contains the projection of point j in all views
 - Row i contains one coordinate of the projections of all the n points in image i
- Factorize D :
 - Compute SVD: $D = U W V^T$
 - Create U_3 by taking the first 3 columns of U
 - Create V_3 by taking the first 3 columns of V
 - Create W_3 by taking the upper left 3×3 block of W
- Create the motion and shape matrices:
 - $M = U_3 W_3^{1/2}$ and $S = W_3^{1/2} V_3^T$ (or $M = U_3$ and $S = W_3 V_3^T$)
- Eliminate affine ambiguity

39

Computer Vision WS 15/16
Slide credit: Martial Hebert

RWTH AACHEN UNIVERSITY

Reconstruction Results



1 60 120 150

C. Tomasi and T. Kanade. *Shape and motion from image streams under orthography: A factorization method.* *IJCV*, 9(2):137-154, November 1992.

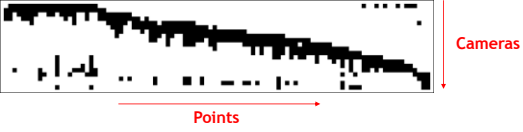
40

Computer Vision WS 15/16
Slide credit: Svetlana Lazebnik, B. Leibe, Image Source: Tomasi & Kanade

RWTH AACHEN UNIVERSITY

Dealing with Missing Data

- So far, we have assumed that all points are visible in all views
- In reality, the measurement matrix typically looks something like this:



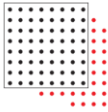
41

Computer Vision WS 15/16
Slide credit: Svetlana Lazebnik, B. Leibe

RWTH AACHEN UNIVERSITY

Dealing with Missing Data

- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
 - Finding dense maximal sub-blocks of the matrix is NP-complete (equivalent to finding maximal cliques in a graph)
- Incremental bilinear refinement



- Perform factorization on a dense sub-block

F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce. *Segmenting, Modeling, and Matching Video Clips Containing Multiple Moving Objects.* *PAMI* 2007.

42

Computer Vision WS 15/16
Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Dealing with Missing Data

- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
 - Finding dense maximal sub-blocks of the matrix is NP-complete (equivalent to finding maximal cliques in a graph)
- Incremental bilinear refinement

(1) Perform factorization on a dense sub-block

(2) Solve for a new 3D point visible by at least two known cameras (linear least squares)

F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce. [Segmenting, Modeling, and Matching Video Clips Containing Multiple Moving Objects](#). PAMI 2007.

43

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Dealing with Missing Data

- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
 - Finding dense maximal sub-blocks of the matrix is NP-complete (equivalent to finding maximal cliques in a graph)
- Incremental bilinear refinement

(1) Perform factorization on a dense sub-block

(2) Solve for a new 3D point visible by at least two known cameras (linear least squares)

(3) Solve for a new camera that sees at least three known 3D points (linear least squares)

F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce. [Segmenting, Modeling, and Matching Video Clips Containing Multiple Moving Objects](#). PAMI 2007.

44

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Comments: Affine SfM

- Affine SfM was historically developed first.
- It is valid under the assumption of *affine cameras*.
 - Which does not hold for real physical cameras...
 - ...but which is still tolerable if the scene points are far away from the camera.
- For good results with real cameras, we typically need projective SfM.
 - Harder problem, more ambiguity
 - Math is a bit more involved... (Here, only basic ideas. If you want to implement it, please look at the H&Z book for details).

45

Computer Vision WS 15/16

B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Structure from Motion (SfM)
 - Motivation
 - Ambiguity
- Affine SfM
 - Affine cameras
 - Affine factorization
 - Euclidean upgrade
 - Dealing with missing data
- Projective SfM
 - Two-camera case
 - Projective factorization
 - Bundle adjustment
 - Practical considerations
- Applications

46

Computer Vision WS 15/16

B. Leibe

RWTH AACHEN UNIVERSITY

Projective Structure from Motion

- Given: m images of n fixed 3D points

$$x_{ij} = P_i X_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$
- Problem: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}

47

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik

B. Leibe

RWTH AACHEN UNIVERSITY

Projective Structure from Motion

- Given: m images of n fixed 3D points
 - $z_{ij} x_{ij} = P_i X_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$
- Problem: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}
- With no calibration info, cameras and points can only be recovered up to a 4×4 projective transformation Q :

$$X \rightarrow QX, \quad P \rightarrow PQ^{-1}$$
- We can solve for structure and motion when

$$2mn \geq 11m + 3n - 15$$
- For two cameras, at least 7 points are needed.

48

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik

B. Leibe

RWTH AACHEN UNIVERSITY

Projective SfM: Two-Camera Case

- Assume fundamental matrix F between the two views
 - First camera matrix: $[I|0]Q^{-1}$
 - Second camera matrix: $[A|b]Q^{-1}$
- Let $\tilde{X} = QX$, then $zAx = [I|0]\tilde{X}$, $z'x' = [A|b]\tilde{X}$
- And

$$z'x' = A[I|0]\tilde{X} + b = zAx + b$$

$$z'x' \times b = zAx \times b$$

$$(z'x' \times b) \cdot x' = (zAx \times b) \cdot x'$$

$$0 = (zAx \times b) \cdot x'$$
- So we have $x'^T [b_\times] Ax = 0$

$$F = [b_\times] A \quad b: \text{epipole } (F^T b = 0), \quad A = -[b_\times] F$$

Computer Vision WS 15/16 49
 Slide adapted from Svetlana Lazebnik B. Leibe FRP sec. 13.3.1

RWTH AACHEN UNIVERSITY

Projective SfM: Two-Camera Case

- This means that if we can compute the fundamental matrix between two cameras, we can directly estimate the two projection matrices from F .
- Once we have the projection matrices, we can compute the 3D position of any point X by triangulation.
- How can we obtain both kinds of information at the same time?

Computer Vision WS 15/16 50
 B. Leibe

RWTH AACHEN UNIVERSITY

Projective Factorization

$$D = \begin{bmatrix} z_{11}x_{11} & z_{12}x_{12} & \dots & z_{1n}x_{1n} \\ z_{21}x_{21} & z_{22}x_{22} & \dots & z_{2n}x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ z_{m1}x_{m1} & z_{m2}x_{m2} & \dots & z_{mn}x_{mn} \end{bmatrix} = \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_m \end{bmatrix} \begin{bmatrix} X_1 & X_2 & \dots & X_n \end{bmatrix}$$

Points ($4 \times n$)

Cameras ($3m \times 4$)

$D = MS$ has rank 4

- If we knew the depths z , we could factorize D to estimate M and S .
- If we knew M and S , we could solve for z .
- Solution: iterative approach (alternate between above two steps).

Computer Vision WS 15/16 51
 Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Sequential Structure from Motion

- Initialize motion from two images using fundamental matrix
- Initialize structure
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image - *calibration*

Computer Vision WS 15/16 52
 Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Sequential Structure from Motion

- Initialize motion from two images using fundamental matrix
- Initialize structure
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image - *calibration*
 - Refine and extend structure: compute new 3D points, re-optimize existing points that are also seen by this camera - *triangulation*

Computer Vision WS 15/16 53
 Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

Sequential Structure from Motion

- Initialize motion from two images using fundamental matrix
- Initialize structure
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image - *calibration*
 - Refine and extend structure: compute new 3D points, re-optimize existing points that are also seen by this camera - *triangulation*
- Refine structure and motion: *bundle adjustment*

Computer Vision WS 15/16 54
 Slide credit: Svetlana Lazebnik B. Leibe

Computer Vision WS 15/16

Bundle Adjustment

- Non-linear method for refining structure and motion
- Minimizing mean-square reprojection error

$$E(\mathbf{P}, \mathbf{X}) = \sum_{i=1}^m \sum_{j=1}^n D(\mathbf{x}_{ij}, \mathbf{P}_i \mathbf{X}_j)^2$$

Slide credit: Svetlana Lazebnik

B. Leibe

55

Computer Vision WS 15/16

Bundle Adjustment

- Seeks the Maximum Likelihood (ML) solution assuming the measurement noise is Gaussian.
- It involves adjusting the bundle of rays between each camera center and the set of 3D points.
- Bundle adjustment should generally be used as the final step of any multi-view reconstruction algorithm.
 - Considerably improves the results.
 - Allows assignment of individual covariances to each measurement.
- However...
 - It needs a good initialization.
 - It can become an extremely large minimization problem.
- Very efficient algorithms available.

B. Leibe

56

Computer Vision WS 15/16

Projective Ambiguity

- If we don't know anything about the camera or the scene, the best we can get with this is a reconstruction up to a projective ambiguity Q .
 - This can already be useful.
 - E.g. we can answer questions like "at what point does a line intersect a plane"?
- If we want to convert this to a "true" reconstruction, we need a **Euclidean upgrade**.
 - Need to put in additional knowledge about the camera (calibration) or about the scene (e.g. from markers).
 - Several methods available (see F&P Chapter 13.5 or H&Z Chapter 19)

Images from Hartley & Zisserman

B. Leibe

57

Computer Vision WS 15/16

Self-Calibration

- Self-calibration (auto-calibration) is the process of determining intrinsic camera parameters directly from uncalibrated images.
- For example, when the images are acquired by a single moving camera, we can use the constraint that the intrinsic parameter matrix remains fixed for all the images.
 - Compute initial projective reconstruction and find 3D projective transformation matrix Q such that all camera matrices are in the form $P_i = K [R_i | t_i]$.
- Can use constraints on the form of the calibration matrix: square pixels, zero skew, fixed focal length, etc.

Slide credit: Svetlana Lazebnik

B. Leibe

58

Computer Vision WS 15/16

Practical Considerations (1)

1. Role of the baseline
 - Small baseline: large depth error
 - Large baseline: difficult search problem
- Solution
 - Track features between frames until baseline is sufficient.

Slide adapted from Steve Seitz

B. Leibe

59

Computer Vision WS 15/16

Practical Considerations (2)

2. There will still be many outliers
 - Incorrect feature matches
 - Moving objects
 ⇒ Apply RANSAC to get robust estimates based on the inlier points.
3. Estimation quality depends on the point configuration
 - Points that are close together in the image produce less stable solutions.
 ⇒ Subdivide image into a grid and try to extract about the same number of features per grid cell.

B. Leibe

60

RWTH AACHEN UNIVERSITY

General Guidelines

- Use calibrated cameras wherever possible.
 - It makes life so much easier, especially for SfM.
- SfM with 2 cameras is *far* more robust than with a single camera.
 - Triangulate feature points in 3D using stereo.
 - Perform 2D-3D matching to recover the motion.
 - More robust to loss of scale (main problem of 1-camera SfM).
- Any constraint on the setup can be useful
 - E.g. square pixels, zero skew, fixed focal length in each camera
 - E.g. fixed baseline in stereo SfM setup
 - E.g. constrained camera motion on a ground plane
 - Making best use of those constraints may require adapting the algorithms (some known results are described in H&Z).

61

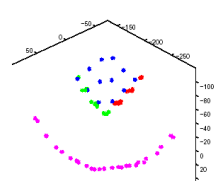
Computer Vision WS 15/16

B. Leibe

RWTH AACHEN UNIVERSITY

Structure-from-Motion: Limitations

- Very difficult to reliably estimate metric SfM unless
 - Large (x or y) motion *or*
 - Large field-of-view and depth variation
- Camera calibration important for Euclidean reconstruction
- Need good feature tracker



62

Computer Vision WS 15/16

B. Leibe

Slide adapted from Steve Seitz

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Structure from Motion (SfM)
 - Motivation
 - Ambiguity
- Affine SfM
 - Affine cameras
 - Affine factorization
 - Euclidean upgrade
 - Dealing with missing data
- Projective SfM
 - Two-camera case
 - Projective factorization
 - Bundle adjustment
 - Practical considerations
- Applications

63

Computer Vision WS 15/16

B. Leibe

RWTH AACHEN UNIVERSITY

Commercial Software Packages

- boujou (<http://www.2d3.com/>)
- PFTrack (<http://www.thepixelfarm.co.uk/>)
- MatchMover (<http://www.realviz.com/>)
- SynthEyes (<http://www.ssontech.com/>)
- Icarus (<http://aig.cs.man.ac.uk/research/reveal/icarus/>)
- Voodoo Camera Tracker (<http://www.digilab.uni-hannover.de/>)

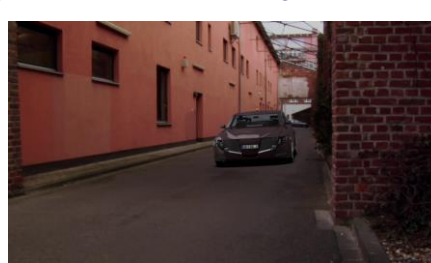
64

Computer Vision WS 15/16

B. Leibe

RWTH AACHEN UNIVERSITY

Applications: Matchmoving



- Putting virtual objects into real-world videos
 - [Original sequence](#) [Tracked features](#)
 - [SfM results](#) [Final video](#)

66

Computer Vision WS 15/16

B. Leibe Videos from Stefan Hafenecker

RWTH AACHEN UNIVERSITY

Applications: Matchmoving



- Putting virtual objects into real-world videos
 - [Original sequence](#) [Tracked features](#)
 - [SfM results](#) [Final video](#)

67

Computer Vision WS 15/16

B. Leibe Videos from Stefan Hafenecker

RWTH AACHEN UNIVERSITY

Applications: Matchmoving

Computer Vision WS 15/16

68

B. Leibe Videos from Stefan Hafeneeger

RWTH AACHEN UNIVERSITY

Applications: Matchmoving

Computer Vision WS 15/16

69

B. Leibe Videos from Stefan Hafeneeger

RWTH AACHEN UNIVERSITY

Applications: Matchmoving

Computer Vision WS 15/16

70

B. Leibe Videos from Stefan Hafeneeger

RWTH AACHEN UNIVERSITY

Applications: Large-Scale SfM from Flickr

Computer Vision WS 15/16

71

B. Leibe


S. Agarwal, N. Snavely, I. Simon, S.M. Seitz, R. Szeliski, [Building Rome in a Day](http://grail.cs.washington.edu/rome/), ICCV'09, 2009. (Video from <http://grail.cs.washington.edu/rome/>)

RWTH AACHEN UNIVERSITY

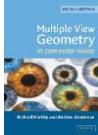
References and Further Reading

- A (relatively short) treatment of affine and projective SfM and the basic ideas and algorithms can be found in Chapters 12 and 13 of

D. Forsyth, J. Ponce,
Computer Vision - A Modern Approach.
Prentice Hall, 2003


- More detailed information (if you really want to implement this) and better explanations can be found in Chapters 10, 18 (factorization) and 19 (self-calibration) of

R. Hartley, A. Zisserman
Multiple View Geometry in Computer Vision
2nd Ed., Cambridge Univ. Press, 2004



Computer Vision WS 15/16

72

B. Leibe