Computer Vision and Machine Learning Winter School ENS Lyon 2010

Camera geometry and image alignment

Josef Sivic

http://www.di.ens.fr/~josef INRIA, WILLOW, ENS/INRIA/CNRS UMR 8548 Laboratoire d'Informatique, Ecole Normale Supérieure, Paris

With slides from: O. Chum, K. Grauman, S. Lazebnik, B. Leibe, D. Lowe, J. Philbin, J. Ponce, D. Nister, C. Schmid, N. Snavely, A. Zisserman

Outline

Part I - Camera geometry – image formation

- Perspective projection
- Affine projection
- Projection of planes

Part II - Image matching and recognition with local features

- Correspondence
- Semi-local and global geometric relations
- Robust estimation RANSAC and Hough Transform

Motivation: Stitching panoramas







Extract features



Extract features

Compute *putative matches*



Extract features

Compute *putative matches*

Loop:

• *Hypothesize* transformation *T* (small group of putative matches that are related by *T*)



Extract features

Compute *putative matches*

Loop:

- *Hypothesize* transformation *T* (small group of putative matches that are related by *T*)
- *Verify* transformation (search for other matches consistent with *T*)



Extract features

Compute *putative matches*

Loop:

- *Hypothesize* transformation *T* (small group of putative matches that are related by *T*)
- *Verify* transformation (search for other matches consistent with *T*)

2D transformation models



Why these transformations ???

Camera geometry





Images are two-dimensional patterns of brightness values.



Pinhole perspective projection: Brunelleschi, XVth Century. Camera obscura: XVIth Century.





Pompei painting, 2000 years ago.



Van Eyk, XIVth Century

Brunelleschi, 1415





Massaccio's Trinity, 1425

Pinhole Perspective Equation



Affine projection models: Weak perspective projection



When the scene relief is small compared its distance from the Camera, *m* can be taken constant: weak perspective projection.

Affine projection models: Orthographic projection



$$\begin{cases} x' = x \\ y' = y \end{cases}$$

When the camera is at a (roughly constant) distance from the scene, take *m*=1.



Strong perspective: Angles are not preserved The projections of parallel lines intersect at one point



From Zisserman & Hartley

Strong perspective: Angles are not preserved The projections of parallel lines intersect at one point

Weak perspective: Angles are better preserved The projections of parallel lines are (almost) parallel



Beyond pinhole camera model: Geometric Distortion





Rectification

Radial Distortion Model



Perspective Projection	$x' = f \frac{x}{z}$ $y' = f \frac{y}{z}$	<i>x</i>,<i>y</i>: World coordinates<i>x</i>',<i>y</i>': Image coordinates<i>f</i>: pinhole-to-retina distance
Weak-Perspective Projection (Affine)	$\begin{array}{l} x' \approx -mx \\ y' \approx -my \end{array} m = -\frac{f}{\overline{z}} \end{array}$	<i>x</i>,<i>y</i>: World coordinates<i>x</i>',<i>y</i>': Image coordinates<i>m</i>: magnification
Orthographic Projection (Affine)	$\begin{array}{llllllllllllllllllllllllllllllllllll$	<i>x</i> , <i>y</i> : World coordinates <i>x</i> ', <i>y</i> ': Image coordinates
Common distortion model	$x'' = \frac{1}{\lambda} x'$ $y'' = \frac{1}{\lambda} y'$ $\lambda = 1 + k_1 r^2 + k_2 r^4 + \cdots$	x',y': Ideal image coordinates x",y": Actual image coordinates

Cameras and their parameters



Images from M. Pollefeys

The Intrinsic Parameters of a Camera



Coordinates

$$\begin{cases} u = kf\frac{x}{z} \\ v = lf\frac{y}{z} \end{cases} \rightarrow \begin{cases} u = \alpha\frac{x}{z} + u_0 \\ v = \beta\frac{y}{z} + v_0 \end{cases} \rightarrow \begin{cases} u = \alpha\frac{x}{z} - \alpha \cot\theta\frac{y}{z} + u_0 \\ v = \frac{\beta}{\sin\theta}\frac{y}{z} + v_0 \end{cases}$$

The Intrinsic Parameters of a Camera



Calibration Matrix

$$oldsymbol{p} = \mathcal{K}\hat{oldsymbol{p}}, ext{ where } oldsymbol{p} = egin{pmatrix} u \ v \ 1 \end{pmatrix} ext{ and } \mathcal{K} \stackrel{ ext{def}}{=} egin{pmatrix} lpha & -lpha \cot heta & u_0 \ 0 & rac{eta}{\sin heta} & v_0 \ 0 & rac{\sin heta}{\sin heta} & v_0 \ 0 & 0 & 1 \end{pmatrix}$$

The Perspective $p = \frac{1}{z}MP$, where $M \stackrel{\text{def}}{=} (\mathcal{K} \ \mathbf{0})$ Projection Equation

Notation



Euclidean Geometry

Recall: Coordinate Changes and Rigid Transformations



$$\begin{bmatrix} B \\ 1 \end{bmatrix} = \begin{bmatrix} B \\ A \\ 0 \end{bmatrix} \begin{bmatrix} B \\ 0 \end{bmatrix} \begin{bmatrix} B \\ 0 \end{bmatrix} \begin{bmatrix} A \\ 0$$

The Extrinsic Parameters of a Camera

• When the camera frame (C) is different from the world frame (W), $\binom{^{C}P}{1} = \binom{^{C}C}{\mathbf{0}^{T}} \binom{^{C}O_{W}}{1} \binom{^{W}P}{1}.$

• Thus,

$$\boldsymbol{p} = \frac{1}{z} \mathcal{M} \boldsymbol{P}, \quad \text{where} \quad \begin{cases} \mathcal{M} = \mathcal{K} (\mathcal{R} \quad \boldsymbol{t}), \\ \mathcal{R} = {}_{W}^{C} \mathcal{R}, \\ \boldsymbol{t} = {}_{W}^{C} \mathcal{R}, \\ \boldsymbol{t} = {}_{W}^{C} \mathcal{O}_{W}, \\ \boldsymbol{t} = {}_{Z}^{C} \mathcal{O}_{W}, \\ \boldsymbol{t} = {}_{Z}^{M} P_{1} \sum_{\boldsymbol{m}_{12}} {}_{\boldsymbol{m}_{13}} {}_{\boldsymbol{m}_{24}} {}_{\boldsymbol{m}_{34}} {}_{\boldsymbol{m}_{34}} \\ \begin{pmatrix} {}_{W}^{w} \boldsymbol{x} \\ {}_{W}^{w} \boldsymbol{y} \\ {}_{W}^{w} \boldsymbol{z} \\ 1 \end{pmatrix} \quad \begin{pmatrix} {}_{W} \mathcal{R} \\ \boldsymbol{w} \boldsymbol{y} \\ \boldsymbol{w} \boldsymbol{z} \\ 1 \end{pmatrix} \quad \boldsymbol{W} \boldsymbol{P} = \begin{pmatrix} {}_{W}^{W} \boldsymbol{x} \\ {}_{W}^{w} \boldsymbol{y} \\ {}_{W}^{w} \boldsymbol{z} \end{pmatrix}$$

• Note: z is *not* independent of \mathcal{M} and \mathbf{P} :

$$\mathcal{M} = egin{pmatrix} oldsymbol{m}_1^T \ oldsymbol{m}_2^T \ oldsymbol{m}_3^T \end{pmatrix} \Longrightarrow z = oldsymbol{m}_3 \cdot oldsymbol{P}, \quad ext{or} \quad \left\{egin{array}{c} u = rac{oldsymbol{m}_1 \cdot oldsymbol{P}}{oldsymbol{m}_3 \cdot oldsymbol{P}}, \ v = rac{oldsymbol{m}_2 \cdot oldsymbol{P}}{oldsymbol{m}_3 \cdot oldsymbol{P}}. \end{array}
ight.$$

Explicit Form of the Projection Matrix

$$\mathcal{M} = \begin{pmatrix} \alpha \boldsymbol{r}_1^T - \alpha \cot \theta \boldsymbol{r}_2^T + u_0 \boldsymbol{r}_3^T & \alpha t_x - \alpha \cot \theta t_y + u_0 t_z \\ \frac{\beta}{\sin \theta} \boldsymbol{r}_2^T + v_0 \boldsymbol{r}_3^T & \frac{\beta}{\sin \theta} t_y + v_0 t_z \\ \boldsymbol{r}_3^T & \boldsymbol{t}_z \end{pmatrix}$$

Note: If $\mathcal{M} = (\mathcal{A} \ \mathbf{b})$ then $|\mathbf{a}_3| = 1$.

Replacing \mathcal{M} by $\lambda \mathcal{M}$ in

$$\left\{ egin{array}{l} u = \displaystyle rac{oldsymbol{m}_1 \cdot oldsymbol{P}}{oldsymbol{m}_3 \cdot oldsymbol{P}} \ v = \displaystyle rac{oldsymbol{m}_2 \cdot oldsymbol{P}}{oldsymbol{m}_3 \cdot oldsymbol{P}} \end{array}
ight.$$

does not change u and v.

M is only defined up to scale in this setting!!

Weak perspective (affine) camera $z_r = m_3^T P = \text{const.}$

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \frac{1}{z_r} \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix} P = \begin{pmatrix} m_1^T P / m_3^T P \\ m_2^T P / m_3^T P \\ m_3^T P / m_3^T P \end{pmatrix}$$

$$\begin{pmatrix} u \\ v \end{pmatrix} = \frac{1}{z_r} \begin{bmatrix} m_1^T \\ m_2^T \\ m_2^T \end{bmatrix} P$$

$$\begin{pmatrix} u \\ v \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ A_{2\times 3} & b_{2\times 1} \end{bmatrix} \begin{pmatrix} w \\ w \\ w \\ v \\ w \\ z \\ 1 \end{pmatrix} = A^W P + b$$

Re-cap: imaging and camera geometry (with a slight change of notation)

- perspective projection
- camera centre, image point and scene point are collinear
- an image point back projects to a ray in 3-space



• depth of the scene point is unknown

The camera model for perspective projection is a linear map between homogeneous point coordinates

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{bmatrix} P (3 \times 4) \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

Image Point

Scene Point

$$\lambda \mathbf{x} = \mathbf{P} \mathbf{X}$$

e.g. if P = [I|0] then

$$x = \frac{X}{Z}$$
 $y = \frac{Y}{Z}$

• P has 11 degrees of freedom (essential parameters).

How a "scene plane" projects into an image?


Plane projective transformations



Choose the world coordinate system such that the plane of the points has zero z coordinate. Then the 3×4 matrix P reduces to

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{pmatrix} \mathsf{x} \\ \mathsf{y} \\ \mathsf{0} \\ \mathsf{1} \end{pmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{14} \\ p_{21} & p_{22} & p_{24} \\ p_{31} & p_{32} & p_{34} \end{bmatrix} \begin{pmatrix} \mathsf{x} \\ \mathsf{y} \\ \mathsf{1} \end{pmatrix}$$

which is a 3×3 matrix representing a general plane to plane projective transformation.

Projective transformations continued

$$\begin{pmatrix} x'_1 \\ x'_2 \\ x'_3 \end{pmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

or $\mathbf{x}' = \mathbf{H}\mathbf{x}$, where \mathbf{H} is a 3 × 3 non-singular homogeneous matrix.

- This is the most general transformation between the world and image plane under imaging by a perspective camera.
- It is often only the 3×3 form of the matrix that is important in establishing properties of this transformation.
- A projective transformation is also called a ``homography" and a ``collineation".
- H has 8 degrees of freedom.

Planes under affine projection

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \end{bmatrix} \begin{pmatrix} x \\ y \\ 0 \\ 1 \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = A_{2 \times 2} P + b_{2 \times 1}$$

Points on a world plane map with a 2D affine geometric transformation (6 parameters)

Summary

• Affine projections induce affine transformations from planes onto their images.

• Perspective projections induce projective transformations from planes onto their images.





2D transformation models



When is homography a valid transformation model?



Case I: Plane projective transformations



Choose the world coordinate system such that the plane of the points has zero z coordinate. Then the 3×4 matrix P reduces to

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{pmatrix} \mathsf{x} \\ \mathsf{y} \\ \mathsf{0} \\ \mathsf{1} \end{pmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{14} \\ p_{21} & p_{22} & p_{24} \\ p_{31} & p_{32} & p_{34} \end{bmatrix} \begin{pmatrix} \mathsf{x} \\ \mathsf{y} \\ \mathsf{1} \end{pmatrix}$$

which is a 3×3 matrix representing a general plane to plane projective transformation.

Case II: Cameras rotating about their centre



- The two image planes are related by a homography H
- H depends only on the relation between the image planes and camera centre, C, not on the 3D structure

Case II: Cameras rotating about their centre



Slide credit: A. Zisserman

Outline – the rest of the lecture

Part 2. Image matching and recognition with local features Correspondence Semi-local and global geometric relations Robust estimation – RANSAC and Hough Transform

Image matching and recognition with local features

The goal: establish correspondence between two or more

images



Image points x and x' are in correspondence if they are projections of the same 3D scene point X.

Example I: <u>Wide baseline matching</u>

Establish correspondence between two (or more) images.

Useful in visual geometry: Camera calibration, 3D reconstruction, Structure and motion estimation, ...

Scale/affine – invariant regions: SIFT, Harris-Laplace, etc.



Example II: Object recognition

Establish correspondence between the target image and (multiple) images in the model database.



[D. Lowe, 1999]

Example III: Visual search

Given a query image, find images depicting the same place / object in a large unordered image collection.







Find these landmarks

... in these images and 1M more

Establish correspondence between the query image and all images from the database depicting the same object / scene.



Database image(s)

Why is it difficult?

Want to establish correspondence despite possibly large changes in scale, viewpoint, lighting and partial occlusion



Scale



Viewpoint





... and the image collection can be very large (e.g. 1M images)

Approach

Pre-processing (last lecture):

- Detect local features.
- Extract descriptor for each feature.

Matching:

- 1. Establish tentative (putative) correspondences based on local appearance of individual features (their descriptors).
- 2. Verify matches based on semi-local / global geometric relations.

Example I: Two images -"Where is the Graffiti?"





Step 1. Establish tentative correspondence

Establish tentative correspondences between object model image and target image by nearest neighbour matching on SIFT vectors



Need to solve some variant of the "nearest neighbor problem" for all feature vectors, $\mathbf{x}_i \in \mathcal{R}^{128}$, in the query image:

$$\forall j \ NN(j) = \arg\min_i ||\mathbf{x}_i - \mathbf{x}_j||,$$

where, $\mathbf{x}_i \in \mathcal{R}^{128}$, are features in the target image.

Can take a long time if many target images are considered.

Problem with matching on local descriptors alone



- too much individual invariance
- each region can affine deform independently (by different amounts)
- Locally appearance can be ambiguous

Solution: use semi-local and global spatial relations to verify matches.

Example I: Two images -"Where is the Graffiti?"

Initial matches

Nearest-neighbor search based on appearance descriptors alone.



After spatial verification



Step 2: Spatial verification (now)

a. Semi-local constraints

Constraints on spatially close-by matches

b. Global geometric relations

Require a consistent global relationship between all matches

Semi-local constraints: Example I. – neighbourhood consensus



Fig. 4. Semi-local constraints : neighbours of the point have to match and angles have to correspond. Note that not all neighbours have to be matched correctly.

[Schmid&Mohr, PAMI 1997]

Semi-local constraints: Example I. – neighbourhood consensus

[Schaffalitzky & Zisserman, CIVR 2004]



After neighbourhood consensus

Semi-local constraints: Example II.







Figure 5: Surface contiguity filter. a) the pattern of intersection between neighboring correct region matches is preserved by transformations between the model and the test images, because the surface is contiguous and smooth. b) the filter evaluates this property by testing the conservation of the area ratios.

[Ferrari et al., IJCV 2005]



Model image



Matched image



Matched image

Geometric verification with global constraints

- All matches must be consistent with a global geometric relation / transformation.
- Need to simultaneously (i) estimate the geometric relation / transformation and (ii) the set of consistent matches





Tentative matches

Matches consistent with an affine transformation

Epipolar geometry (not considered here)

In general, two views of a 3D scene are related by the epipolar constraint.



- A point in one view "generates" an epipolar line in the other view
- The corresponding point lies on this line.

Slide credit: A. Zisserman

Epipolar geometry (not considered here)

Epipolar geometry is a consequence of the coplanarity of the camera centres and scene point



The camera centres, corresponding points and scene point lie in a single plane, known as the epipolar plane

Epipolar geometry (not considered here)

Algebraically, the epipolar constraint can be expressed as



where

- x, x' are homogeneous coordinates (3-vectors) of corresponding image points.
- F is a 3x3, rank 2 homogeneous matrix with 7 degrees of freedom, called the **fundamental matrix**.

3D constraint: example (not considered here)

Matches must be consistent with a 3D model



3D constraint: example (not considered here)

With a given 3D model (set of known X's) and a set of measured image points x, the goal is to find find camera matrix P and a set of geometrically consistent correspondences x <> X.



 $\mathbf{x} = \mathbf{P}\mathbf{X}$

- $P: 3 \times 4$ matrix
- x : 4-vector
- x : 3-vector

2D transformation models



Example: estimating 2D affine transformation

- Simple fitting procedure (linear least squares)
- Approximates viewpoint changes for roughly planar objects and roughly orthographic cameras
- Can be used to initialize fitting for more complex models



Example: estimating 2D affine transformation

- Simple fitting procedure (linear least squares)
- Approximates viewpoint changes for roughly planar objects and roughly orthographic cameras
- Can be used to initialize fitting for more complex models



Matches consistent with an affine transformation

Fitting an affine transformation

Assume we know the correspondences, how do we get the transformation?


Fitting an affine transformation



Linear system with six unknowns

Each match gives us two linearly independent equations: need at least three to solve for the transformation parameters Dealing with outliers

The set of putative matches may contain a high percentage (e.g. 90%) of outliers

How do we fit a geometric transformation to a small subset of all possible matches?

Possible strategies:

- RANSAC
- Hough transform

Strategy 1: RANSAC

RANSAC loop (Fischler & Bolles, 1981):

- Randomly select a *seed group* of matches
- Compute transformation from seed group
- Find *inliers* to this transformation
- If the number of inliers is sufficiently large, re-compute least-squares estimate of transformation on all of the inliers
- Keep the transformation with the largest number of inliers

Example: Robust line estimation - RANSAC

Fit a line to 2D data containing outliers



There are two problems

- 1. a line fit which minimizes perpendicular distance
- a classification into inliers (valid points) and outliers
 Solution: use robust statistical estimation algorithm RANSAC
 (RANdom Sample Consensus) [Fishler & Bolles, 1981]

RANSAC robust line estimation

Repeat

- 1. Select random sample of 2 points
- 2. Compute the line through these points
- 3. Measure support (number of points within threshold distance of the line)

Choose the line with the largest number of inliers

• Compute least squares fit of line to inliers (regression)



















Algorithm summary – RANSAC robust estimation of 2D affine transformation

Repeat

- 1. Select 3 point to point correspondences
- 2. Compute H (2x2 matrix) + t (2x1) vector for translation
- 3. Measure support (number of inliers within threshold distance, i.e. $d_{transfer}^2 < t$) $d_{transfer}^2 = d(\mathbf{x}, \mathbf{H}^{-1}\mathbf{x}')^2 + d(\mathbf{x}', \mathbf{H}\mathbf{x})^2$



Choose the (H,t) with the largest number of inliers (Re-estimate (H,t) from all inliers)

How many samples?

Number of samples *N*

- Choose *N* so that, with probability *p*, at least one random sample is free from outliers
- e.g.:
 - > p=0.99
 - > outlier ratio: e

Probability a randomly picked point is an inlier

$$\left(1-\left(1-e\right)^{s}\right)^{v}=1-p$$

Probability of all points in a sample (of size s) are inliers

How many samples?

Number of samples N

- Choose *N* so that, with probability *p*, at least one random sample is free from outliers
- e.g.:
 - > *p*=0.99
 - > outlier ratio: e

Probability that all N samples (of
size s) are corrupted (contain an
outlier)



Probability of at least one point in a sample (of size s) is an outlier

$$N = \log(1-p)/\log((1-(1-e)^s))$$

	proportion of outliers <i>e</i>						
S	5%	10%	20%	30%	40%	50%	90%
1	2	2	3	4	5	6	43
2	2	3	5	7	11	17	458
3	3	4	7	11	19	35	4603
4	3	5	9	17	34	72	4.6e4
5	4	6	12	26	57	146	4.6e5
6	4	7	16	37	97	293	4.6e6
7	4	8	20	54	163	588	4.6e7
8	5	9	26	78	272	1177	4.6e8

Source: M. Pollefeys

Example: line fitting

p = 0.99 s = ? e = ?

N = ?



Example: line fitting

p = 0.99s = 2e = 2/10 = 0.2

N = 5



Compare with exhaustively trying all point pairs:

$$\begin{pmatrix} 10 \\ 2 \end{pmatrix} = 10*9 / 2 = 45$$

_								
	proportion of outliers <i>e</i>							
	S	5%	10%	20%	30%	40%	50%	90%
	1	2	2	3	4	5	6	43
	2	2	3	5	7	11	17	458
	3	3	4	7	11	19	35	4603
	4	3	5	9	17	34	72	4.6e4
	5	4	6	12	26	57	146	4.6e5
	6	4	7	16	37	97	293	4.6e6
	7	4	8	20	54	163	588	4.6e7
_	8	5	9	26	78	272	1177	4.6e8

Source: M. Pollefeys

How to reduce the number of samples needed?

- 1. Reduce the proportion of outliers.
- 2. Reduce the sample size
 - use simpler model (e.g. similarity instead of affine tnf.)
 - use local information (e.g. a region to region correspondence is equivalent to (up to) 3 point to point correspondences).



Number of samples N

	proportion of outliers <i>e</i>						
S	5%	10%	20%	30%	40%	50%	90%
1	2	2	3	4	5	6	43
2	2	3	5	7	11	17	458
3	3	4	7	11	19	35	4603
4	3	5	9	17	34	72	4.6e4
5	4	6	12	26	57	146	4.6e5
6	4	7	16	37	97	293	4.6e6
7	4	8	20	54	163	588	4.6e7
8	5	9	26	78	272	1177	4.6e8

RANSAC (references)

- M. Fischler and R. Bolles, "Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography," Comm. ACM, 1981
- R. Hartley and A. Zisserman, Multiple View Geometry in Computer Vision, 2nd ed., 2004.

Extensions:

- B. Tordoff and D. Murray, "Guided Sampling and Consensus for Motion Estimation, ECCV'03
- D. Nister, "Preemptive RANSAC for Live Structure and Motion Estimation, ICCV'03
- Chum, O.; Matas, J. and Obdrzalek, S.: Enhancing RANSAC by Generalized Model Optimization, ACCV'04
- Chum, O.; and Matas, J.: Matching with PROSAC Progressive Sample Consensus , CVPR 2005
- Philbin, J., Chum, O., Isard, M., Sivic, J. and Zisserman, A.: Object retrieval with large vocabularies and fast spatial matching, CVPR'07

Chum, O. and Matas. J.: Optimal Randomized RANSAC, PAMI'08

Strategy 2: Hough Transform

- Origin: Detection of straight lines in cluttered images
- Can be generalized to arbitrary shapes
- Can extract feature groupings from cluttered images in linear time.
- Illustrate on extracting sets of local features consistent with a similarity transformation.

Hough transform for object recognition

Suppose our features are scale- and rotation-covariant

• Then a single feature match provides an alignment hypothesis (translation, scale, orientation)



David G. Lowe. "Distinctive image features from scaleinvariant keypoints", *IJCV* 60 (2), pp. 91-110, 2004.

Hough transform for object recognition

Suppose our features are scale- and rotation-covariant

- Then a single feature match provides an alignment hypothesis (translation, scale, orientation)
- Of course, a hypothesis obtained from a single match is unreliable
- Solution: Coarsely quantize the transformation space. Let each match vote for its hypothesis in the quantized space.



David G. Lowe. "Distinctive image features from scaleinvariant keypoints", *IJCV* 60 (2), pp. 91-110, 2004.

model

Basic algorithm outline

- 1. Initialize accumulator H to all zeros
- 2. For each tentative match compute transformation hypothesis: tx, ty, s, θ H(tx,ty,s,θ) = H(tx,ty,s,θ) + 1 end end



ty

- Find all bins (tx,ty,s,θ) where H(tx,ty,s,θ) has at least three votes
- Correct matches will consistently vote for the same transformation while mismatches will spread votes.
- Cost: Linear scan through the matches (step 2), followed by a linear scan through the accumulator (step 3).

Hough transform details (D. Lowe's system)

- **Training phase:** For each model feature, record 2D location, scale, and orientation of model (relative to normalized feature frame)
- **Test phase:** Let each match between a test and a model feature vote in a 4D Hough space
 - Use broad bin sizes of 30 degrees for orientation, a factor of 2 for scale, and 0.25 times image size for location
 - Vote for two closest bins in each dimension
- Find all bins with at least three votes and perform geometric verification
 - Estimate least squares affine transformation
 - Use stricter thresholds on transformation residual
 - Search for additional features that agree with the alignment

Hough transform in object recognition (references)

- P.V.C. Hough, Machine Analysis of Bubble Chamber Pictures, Proc. Int. Conf. High Energy Accelerators and Instrumentation, 1959
- D. Lowe, "Distinctive image features from scale-invariant keypoints", IJCV 60 (2), 2004.
- H. Jegou, M. Douze, C. Schmid, Hamming embedding and weak geometric consistency for large scale image search, ECCV'2008

Extensions (object category detection):

- B. Leibe, A. Leonardis, and B. Schiele., Combined Object Categorization and Segmentation with an Implicit Shape Model, in ECCV'04 Workshop on Statistical Learning in Computer Vision, Prague, May 2004.
- S. Maji and J. Malik, Object Detection Using a Max-Margin Hough Tranform, CVPR'2009
- A. Lehmann, B. Leibe, L. Van Gool. Fast PRISM: Branch and Bound Hough Transform for Object Class Detection, IJCV (to appear), 2010.
- O. Barinova, V. Lempitsky, P. Kohli, On the Detection of Multiple Object Instances using Hough Transforms, CVPR, 2010

Comparison

Hough Transform

Advantages

- Can handle high percentage of outliers (>95%)
- Extracts groupings from clutter in linear time

Disadvantages

- Quantization issues
- Only practical for small number of dimensions (up to 4)

Improvements available

- Probabilistic Extensions
- Continuous Voting Space
- Can be generalized to arbitrary shapes and objects

RANSAC

Advantages

- General method suited to large range of problems
- Easy to implement
- "Independent" of number of dimensions

Disadvantages

 Basic version only handles moderate number of outliers (<50%)

Many variants available, e.g.

- PROSAC: Progressive RANSAC [Chum05]
- Preemptive RANSAC [Nister05]

Beyond affine transformations

What is the transformation between two views of a planar

surface?



What is the transformation between images from two cameras that share the same center?



Beyond affine transformations

Homography: plane projective transformation (transformation taking a quad to another arbitrary quad)



Case II: Cameras rotating about their centre



- The two image planes are related by a homography H
- H depends only on the relation between the image planes and camera centre, C, not on the 3D structure

Fitting a homography

Recall: homogenenous coordinates $\begin{bmatrix} x \\ y \end{bmatrix}$

$$(x,y) \Rightarrow \begin{bmatrix} x\\ y\\ 1 \end{bmatrix}$$

Converting to homogenenous image coordinates

$$\begin{bmatrix} x \\ y \\ w \end{bmatrix} \Rightarrow (x/w, y/w)$$

Converting from homogenenous image coordinates

Fitting a homography

Recall: homogeneous coordinates $(x,y) \Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \qquad \qquad \begin{bmatrix} x \\ y \\ w \end{bmatrix} \Rightarrow (x/w, y/w)$

$$(x,y) \Rightarrow \begin{bmatrix} x\\ y\\ 1 \end{bmatrix}$$

Converting to homogenenous image coordinates

Equation for homography:

$$\lambda \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Fitting a homography

Equation for homography:

$$\lambda_{i} \begin{bmatrix} x'_{i} \\ y'_{i} \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x_{i} \\ y_{i} \\ 1 \end{bmatrix} \qquad \lambda_{i} \mathbf{x}'_{i} = \mathbf{H} \mathbf{x}_{i} = \begin{bmatrix} \mathbf{h}_{1}^{T} \\ \mathbf{h}_{2}^{T} \\ \mathbf{h}_{3}^{T} \end{bmatrix} \mathbf{x}_{i}$$

9 entries, 8 degrees of freedom
(scale is arbitrary)

$$\mathbf{x}'_{i} \times \mathbf{H} \mathbf{x}_{i} = 0$$
 $\mathbf{x}'_{i} \times \mathbf{H} \mathbf{x}_{i} = \begin{bmatrix} y'_{i} \mathbf{h}_{3}^{T} \mathbf{x}_{i} - \mathbf{h}_{2}^{T} \mathbf{x}_{i} \\ \mathbf{h}_{1}^{T} \mathbf{x}_{i} - x'_{i} \mathbf{h}_{3}^{T} \mathbf{x}_{i} \\ \mathbf{h}_{1}^{T} \mathbf{x}_{i} - y'_{i} \mathbf{h}_{3}^{T} \mathbf{x}_{i} \\ x'_{i} \mathbf{h}_{2}^{T} \mathbf{x}_{i} - y'_{i} \mathbf{h}_{1}^{T} \mathbf{x}_{i} \end{bmatrix}$

$$\begin{bmatrix} \mathbf{0}^T & -\mathbf{x}_i^T & y_i' \mathbf{x}_i^T \\ \mathbf{x}_i^T & \mathbf{0}^T & -x_i' \mathbf{x}_i^T \\ -y_i' \mathbf{x}_i^T & x_i' \mathbf{x}_i^T & \mathbf{0}^T \end{bmatrix} \begin{pmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{pmatrix} = 0 \quad 3 \text{ equations, only 2 linearly}$$
independent

Direct linear transform

$$\begin{bmatrix} \mathbf{0}^{T} & \mathbf{x}_{1}^{T} & -y_{1}' \, \mathbf{x}_{1}^{T} \\ \mathbf{x}_{1}^{T} & \mathbf{0}^{T} & -x_{1}' \, \mathbf{x}_{1}^{T} \\ \cdots & \cdots & \cdots \\ \mathbf{0}^{T} & \mathbf{x}_{n}^{T} & -y_{n}' \, \mathbf{x}_{n}^{T} \\ \mathbf{x}_{n}^{T} & \mathbf{0}^{T} & -x_{n}' \, \mathbf{x}_{n}^{T} \end{bmatrix} \begin{pmatrix} \mathbf{h}_{1} \\ \mathbf{h}_{2} \\ \mathbf{h}_{3} \end{pmatrix} = \mathbf{0} \qquad \mathbf{A} \, \mathbf{h} = \mathbf{0}$$

H has 8 degrees of freedom (9 parameters, but scale is arbitrary)

One match gives us two linearly independent equations

Four matches needed for a minimal solution (null space of 8x9 matrix)

More than four: homogeneous least squares

Application: Panorama stitching



Images courtesy of A. Zisserman.
Recognizing panoramas

Given contents of a camera memory card, automatically figure out which pictures go together and stitch them together into panoramas



M. Brown and D. Lowe, "Recognizing panoramas", ICCV 2003.

1. Estimate homography (RANSAC)



1. Estimate homography (RANSAC)



1. Estimate homography (RANSAC)



2. Find connected sets of images



2. Find connected sets of images











2. Find connected sets of images











3. Stitch and blend the panoramas



Results





