

Structure from Motion

Paul Heckbert, Nov. 1999

15-869, Image-Based Modeling and Rendering

Approach

Problem:

Reconstruct scene geometry and camera motion from two or more images

Typically, assume

- static scene (model camera motion only)
- diffuse, opaque surfaces (simplifies feature tracking)
- orthographic projection, for starters, then generalize to projective

Steps:

- shoot f frames of video (or sequence of stills)
- track n feature points from frame to frame
- build large matrix of image feature coordinates
- factor this matrix to extract motion and shape
- build 3-D model by connecting the features with triangles

<http://www.ius.cs.cmu.edu/IUS/mbvc0/www/modeling.html>

Image Features

Good features are spots, corners, and other points where the image, regarded as a terrain, has high curvature.

Good features are image windows that can be tracked well [Tomasi-Kanade 92]

Pick good features using the gradient covariance matrix [Lucas-Kanade 81]:

$$\sum_{\text{window}} \nabla I \cdot \nabla I^T = \begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_y I_x & \sum I_y^2 \end{bmatrix}$$

How many large eigenvalues does this matrix have?

2 - good feature (spot)

1 - poor (edge - aperture problem)

0 - poor (flat)

Orthographic Projection

$$\mathbf{u} = \mathbf{\Pi} \mathbf{X} + \mathbf{t}$$

$\begin{matrix} 2 \times 1 & 2 \times 3 & 3 \times 1 & 2 \times 1 \\ \uparrow & \uparrow & \uparrow & \uparrow \\ \text{image point} & \text{projection matrix} & \text{scene point} & \text{image offset} \end{matrix}$

Trick

Choose scene origin to be centroid of 3D points

Choose image origins to be centroid of 2D points

Allows us to drop the camera translation:

$$\mathbf{u} = \mathbf{\Pi} \mathbf{X}$$

$\begin{matrix} 2 \times 1 & 2 \times 3 & 3 \times 1 \end{matrix}$

Tomasi-Kanade Orthographic Method

projection of n features in one image:

$$\begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \end{bmatrix} = \prod_{2 \times n} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}_{3 \times n}$$

projection of n features in f images

$$\begin{bmatrix} \mathbf{u}_1^1 & \mathbf{u}_2^1 & \cdots & \mathbf{u}_n^1 \\ \mathbf{u}_1^2 & \mathbf{u}_2^2 & \cdots & \mathbf{u}_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{u}_1^f & \mathbf{u}_2^f & \cdots & \mathbf{u}_n^f \end{bmatrix} = \begin{bmatrix} \Pi^1 \\ \Pi^2 \\ \vdots \\ \Pi^f \end{bmatrix}_{2f \times 3} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}_{3 \times n}$$

$2f \times n$ $2f \times 3$

W measurement

M motion

S shape

Tomasi-Kanade Orthographic Method

$$\text{known} \rightarrow \boxed{\begin{matrix} \text{W} \\ 2f \times n \end{matrix}} = \boxed{\begin{matrix} \text{M} & \text{S} \\ 2f \times 3 & 3 \times n \end{matrix}} \rightarrow \text{solve for}$$

SVD Factorization Technique

W is at most rank 3 (assuming no noise)

We can use *singular value decomposition* to factor W :

$$\boxed{\begin{matrix} \text{W} \\ 2f \times n \end{matrix}} = \boxed{\begin{matrix} \text{M}' & \text{S}' \\ 2f \times 3 & 3 \times n \end{matrix}}$$

S' differs from S by a linear transformation A :

$$\text{W} = \text{M}'\text{S}' = (\text{MA}^{-1})(\text{AS})$$

Solve for A by enforcing constraints on M

Problems with Basic Orthographic Method, 1

- Measurement matrix \mathbf{W} might have many voids, since occlusion and noise cause features to appear and disappear.

Solution: take linear combinations of rows & columns to “hallucinate” entries.

- Features can be mis-tracked:

Solution: maintain tree of multiple hypotheses

- SVD is slow: $O(fn \cdot \min(f, n))$.

Solution 1: solve bilinear equation by alternating between:

- freeze S and solve $W=MS$ for M
- freeze M and solve $W=MS$ for S

This is faster than SVD.

Solution 2: don't compute full SVD, but partial one

Problems with Basic Orthographic Method, 2

- Feature set too sparse, doesn't yield a good surface model.
Solution: Don't track as a preprocess. Instead solve for correspondence as you solve for motion and shape.
- Orthographic projection is poor approximation for nearby objects.
Solution: model perspective, but this is more complex.
Projective factorization: [Poelman 95] refined orthographic solution using nonlinear optimization (Levenberg-Marquardt)

Triggs' Projective Method

[Triggs CVPR '96] generalized orthographic method by using homogeneous coordinates.

Image coordinates \mathbf{u} have a third, unknown DOF, projective depth. This yields a larger system of equations, but very similar approach.

Steps:

- track features
- find fundamental matrices \mathbf{F} between successive frames
- use \mathbf{F} 's to solve for projective depth
- solve projective factorization equation $\mathbf{W}_{3f \times n} = \mathbf{M}_{3f \times 3} \mathbf{S}_{3 \times n}$