## Structure from Motion

**Computer Vision**
CSE P 576, Spring 2011
Richard Szeliski
Microsoft Research

## Today's lecture

Geometric camera calibration
- camera matrix (Direct Linear Transform)
- non-linear least squares
- separating intrinsics and extrinsics
- focal length and optic center

CSE 576, Spring 2008          Structure from Motion          2

## Today's lecture

Structure from Motion
- triangulation and pose
- two-frame methods
- factorization
- bundle adjustment
- robust statistics

Photo Tourism

CSE 576, Spring 2008          Structure from Motion          3
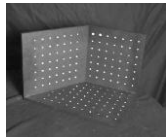
## Camera Calibration

## Camera calibration

Determine camera parameters from *known* 3D points or calibration object(s)

1. *internal* or *intrinsic* parameters such as focal length, optical center, aspect ratio:
   *what kind of camera?*
2. *external* or *extrinsic* (pose) parameters:
   *where is the camera?*

How can we do this?

## Camera calibration – approaches

Possible approaches:

1. linear regression (least squares)
2. non-linear optimization
3. vanishing points
4. multiple planar patterns
5. panoramas (rotational motion)

## Image formation equations

$(X_c, Y_c, Z_c)$

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = \begin{bmatrix} \mathbf{R} \end{bmatrix}_{3\times3} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} + \mathbf{t}$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \sim \begin{bmatrix} U \\ V \\ W \end{bmatrix} = \begin{bmatrix} f & 0 & u_c \\ 0 & f & v_c \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

## Calibration matrix

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \sim \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = \mathbf{K}\ \mathbf{X}_c$$

Is this form of K good enough?
- non-square pixels (digital video)
- skew
- radial distortion

$$\mathbf{K} = \begin{bmatrix} fa & s & u_c \\ 0 & f & v_c \\ 0 & 0 & 1 \end{bmatrix}$$

## Camera matrix

Fold *intrinsic* calibration matrix **K** and *extrinsic* pose parameters (**R**,**t**) together into a *camera matrix*

$$\mathbf{M} = \mathbf{K}\,[\mathbf{R} \mid \mathbf{t}\,]$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \sim \begin{bmatrix} m_{00} & m_{01} & m_{02} & m_{03} \\ m_{10} & m_{11} & m_{12} & m_{13} \\ m_{20} & m_{21} & m_{22} & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

(put 1 in lower r.h. corner for 11 d.o.f.)

CSE 576, Spring 2008          Structure from Motion          10

## Camera matrix calibration

Directly estimate 11 unknowns in the **M** matrix using known 3D points ($X_i, Y_i, Z_i$) and measured feature positions ($u_i, v_i$)

$$u_i = \frac{m_{00}X_i + m_{01}Y_i + m_{02}Z_i + m_{03}}{m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1}$$

$$v_i = \frac{m_{10}X_i + m_{11}Y_i + m_{12}Z_i + m_{13}}{m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1}$$

CSE 576, Spring 2008          Structure from Motion          11

## Camera matrix calibration

Linear regression:

- Bring denominator over, solve set of (over-determined) linear equations.  How?

$$u_i(m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1) =$$
$$m_{00}X_i + m_{01}Y_i + m_{02}Z_i + m_{03}$$
$$v_i(m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1) =$$
$$m_{10}X_i + m_{11}Y_i + m_{12}Z_i + m_{13}$$

- Least squares (pseudo-inverse)
- Is this good enough?

CSE 576, Spring 2008          Structure from Motion          12

## Levenberg-Marquardt

Iterative non-linear least squares [Press'92]

- Linearize measurement equations

$$\hat{u}_i = f(\mathbf{m}, \mathbf{x}_i) + \frac{\partial f}{\partial \mathbf{m}}\Delta\mathbf{m}$$
$$\hat{v}_i = g(\mathbf{m}, \mathbf{x}_i) + \frac{\partial g}{\partial \mathbf{m}}\Delta\mathbf{m}$$

- Substitute into log-likelihood equation:  quadratic cost function in **Δm**

$$\sum_i \sigma_i^{-2}(\hat{u}_i - u_i + \frac{\partial f}{\partial \mathbf{m}}\Delta\mathbf{m})^2 + \cdots$$

CSE 576, Spring 2008          Structure from Motion          15

## Levenberg-Marquardt

Iterative non-linear least squares [Press'92]

- Solve for minimum $\frac{\partial C}{\partial \mathbf{m}} = 0$

$$\mathbf{A}\triangle\mathbf{m} = \mathbf{b}$$

Hessian $\quad \mathbf{A} = \left[\sum_i \sigma_i^{-2} \frac{\partial f}{\partial \mathbf{m}} \left(\frac{\partial f}{\partial \mathbf{m}}\right)^T + \cdots\right]$

error: $\quad \mathbf{b} = \left[\sum_i \sigma_i^{-2} \frac{\partial f}{\partial \mathbf{m}} (u_i - \hat{u}_i) + \cdots\right]$

CSE 576, Spring 2008 　　　Structure from Motion 　　　16

## Camera matrix calibration

Advantages:
- very simple to formulate and solve
- can recover **K** [**R** | **t**] from **M** using QR decomposition [Golub & VanLoan 96]

Disadvantages:
- doesn't compute internal parameters
- more unknowns than true degrees of freedom
- need a separate camera matrix for each new view

CSE 576, Spring 2008 　　　Structure from Motion 　　　18

## Separate intrinsics / extrinsics

New feature measurement equations
$$\hat{u}_{ij} = f(\mathbf{K}, \mathbf{R}_j, \mathbf{t}_j, \mathbf{x}_i)$$
$$\hat{v}_{ij} = g(\mathbf{K}, \mathbf{R}_j, \mathbf{t}_j, \mathbf{x}_i)$$

Use non-linear minimization

Standard technique in photogrammetry, computer vision, computer graphics

- [Tsai 87] – also estimates $\kappa_1$ (freeware @ CMU)
  http://www.cs.cmu.edu/afs/cs/project/cil/ftp/html/v-source.html
- [Bogart 91] – *View Correlation*

CSE 576, Spring 2008 　　　Structure from Motion 　　　19

## Intrinsic/extrinsic calibration

Advantages:
- can solve for more than one camera pose at a time
- potentially fewer degrees of freedom

Disadvantages:
- more complex update rules
- need a good initialization (recover **K** [**R** | **t**] from **M**)
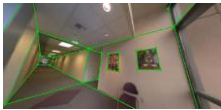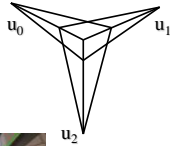
CSE 576, Spring 2008 　　　Structure from Motion 　　　20

## Vanishing Points

Determine focal length *f* and optical center ($u_c$, $v_c$) from image of cube's (or building's) *vanishing points*
[Caprile '90][Antone & Teller '00]

$u_0$   $u_1$

$u_2$

## Vanishing point calibration

Advantages:
- only need to see vanishing points (e.g., architecture, table, …)

Disadvantages:
- not that accurate
- need rectihedral object(s) in scene

## Multi-plane calibration

Use several images of planar target held at *unknown* orientations [Zhang 99]
- Compute plane homographies

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix} \sim \mathrm{K} \begin{bmatrix} r_1 & r_2 & t \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} \sim \mathrm{HX}$$

- Solve for $K^{-T}K^{-1}$ from $H_k$'s
  - 1plane if only *f* unknown
  - 2 planes if (*f*, $u_c$, $v_c$) unknown
  - 3+ planes for full *K*
- Code available from Zhang and OpenCV

## Rotational motion

Use pure rotation (large scene) to estimate *f*
1. estimate *f* from pairwise homographies
2. re-estimate *f* from 360º "gap"
3. optimize over all {*K*, *R*$_j$} parameters
   [Stein 95; Hartley '97; Shum & Szeliski '00; Kang & Weiss '99]

f=510     f=468

Most accurate way to get *f*, short of surveying distant points

## Pose estimation and triangulation

## Pose estimation

Once the internal camera parameters are
known, can compute camera pose

$$\hat{u}_{ij} = f(K, \boxed{R_j, t_j}, x_i)$$
$$\hat{v}_{ij} = g(K, \boxed{R_j, t_j}, x_i)$$

[Tsai87] [Bogart91]

Application: superimpose 3D graphics onto
video

How do we initialize (**R**,**t**)?

CSE 576, Spring 2008          Structure from Motion          32

## Pose estimation

Previous initialization techniques:
- vanishing points [Caprile 90]
- planar pattern [Zhang 99]

Other possibilities
- *Through-the-Lens Camera Control* [Gleicher92]:
  differential update
- 3+ point "linear methods":
  [DeMenthon 95][Quan 99][Ameller 00]

CSE 576, Spring 2008          Structure from Motion          33

## Triangulation

Problem:  Given some points in
*correspondence* across two or more images
(taken from calibrated cameras), $\{(u_j, v_j)\}$,
compute the 3D location **X**

CSE 576, Spring 2008          Structure from Motion          35
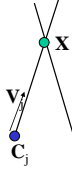
## Triangulation

**Method I**: intersect viewing rays in 3D, minimize:

$$\arg \min_{\mathbf{X}} \sum_j \|\mathbf{C}_j + s\mathbf{V}_j - \mathbf{X}\|$$

- **X** is the unknown 3D point
- **C**$_j$ is the optical center of camera $j$
- **V**$_j$ is the *viewing ray* for pixel ($u_j$,$v_j$)
- $s_j$ is unknown distance along **V**$_j$

Advantage: geometrically intuitive

CSE 576, Spring 2008          Structure from Motion          36

## Triangulation

**Method II**: solve linear equations in **X**
- advantage: very simple

$$u_i = \frac{m_{00}X_i + m_{01}Y_i + m_{02}Z_i + m_{03}}{m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1}$$

$$v_i = \frac{m_{10}X_i + m_{11}Y_i + m_{12}Z_i + m_{13}}{m_{20}X_i + m_{21}Y_i + m_{22}Z_i + 1}$$

**Method III**: non-linear minimization

- advantage: most accurate (image plane error)

CSE 576, Spring 2008          Structure from Motion          37

## Structure from Motion

## Today's lecture

Structure from Motion
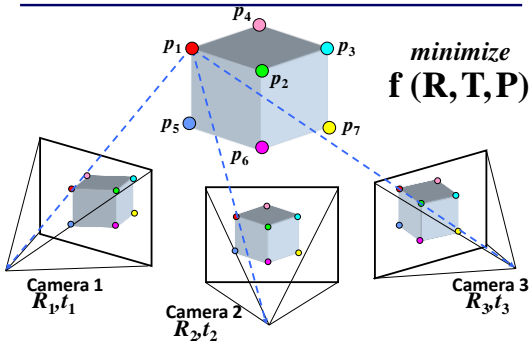- two-frame methods
- factorization
- bundle adjustment
- robust statistics

CSE 576, Spring 2008          Structure from Motion          39

## Structure from motion



$p_4$  
$p_1$  $p_3$  
$p_2$  
$p_5$  $p_7$  
$p_6$

*minimize*

$$\mathbf{f}\,(\mathbf{R}, \mathbf{T}, \mathbf{P})$$

**Camera 1**  
$R_1, t_1$

**Camera 2**  
$R_2, t_2$

**Camera 3**  
$R_3, t_3$

## Structure from motion

Given many points in *correspondence* across several images, $\{(u_{ij}, v_{ij})\}$, simultaneously compute the 3D location $\mathbf{x}_i$ and camera (or *motion*) parameters (**K**, $\mathbf{R}_j$, $\mathbf{t}_j$)

$$\hat{u}_{ij} = f(\mathrm{K}, \mathrm{R}_j, \mathrm{t}_j, \mathrm{x}_i)$$
$$\hat{v}_{ij} = g(\mathrm{K}, \mathrm{R}_j, \mathrm{t}_j, \mathrm{x}_i)$$

Two main variants: calibrated, and uncalibrated (sometimes associated with Euclidean and projective reconstructions)

CSE 576, Spring 2008        Structure from Motion        41

## Structure from motion

$$\hat{u}_{ij} = f(\mathrm{K}, \mathrm{R}_j, \mathrm{t}_j, \mathrm{x}_i)$$
$$\hat{v}_{ij} = g(\mathrm{K}, \mathrm{R}_j, \mathrm{t}_j, \mathrm{x}_i)$$

How many points do we need to match?

• 2 frames:  
  $(\boldsymbol{R}, \boldsymbol{t})$: 5 dof + $3n$ point locations $\leq$  
  $4n$ point measurements $\Rightarrow$  
  $n \geq 5$

• $k$ frames:  
  $6(k{-}1){-}1 + 3n \leq 2kn$

• always want to use many more

CSE 576, Spring 2008        Structure from Motion        42

## Two-frame methods

Two main variants:

1. Calibrated: "Essential matrix" $E$  
   use ray directions $(\boldsymbol{x}, \boldsymbol{x}_i{}')$

2. Uncalibrated: "Fundamental matrix" $F$

[Hartley & Zisserman 2000]

CSE 576, Spring 2008        Structure from Motion        43

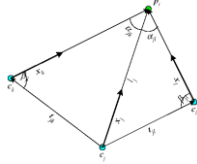## Essential matrix

Co-planarity constraint:

$$x' \approx R\,x + t$$

$$[t]_\times x' \approx [t]_\times R\,x$$

$$x'^T [t]_\times x' \approx x'^{\,T}[t]_\times R\,x$$

$$x'^{\,T} E\,x = 0 \text{ with } E = [t]_\times R$$

- Solve for $E$ using least squares (SVD)
- $t$ is the least singular vector of $E$
- $R$ obtained from the other two sing. vectors

CSE 576, Spring 2008          Structure from Motion          44

## Fundamental matrix

Camera calibrations are unknown

$$x'\,F\,x = 0 \text{ with } F = [e]_\times H = K'[t]_\times R\,K^{-1}$$

- Solve for $F$ using least squares (SVD)
  - re-scale $(x, x_i')$ so that $|x_i| \approx 1/2$ [Hartley]
- $e$ (epipole) is *still* the least singular vector of $F$
- $H$ obtained from the other two s.v.s
- "plane + parallax" (projective) reconstruction
- use self-calibration to determine $K$ [Pollefeys]

CSE 576, Spring 2008          Structure from Motion          45

## Multi-frame Structure from Motion

## Bundle Adjustment

$$\hat{u}_{ij} \;=\; f(K, R_j, t_j, x_i)$$
$$\hat{v}_{ij} \;=\; g(K, R_j, t_j, x_i)$$

What makes this non-linear minimization hard?
- many more parameters: potentially slow
- poorer conditioning (high correlation)
- potentially lots of outliers
- gauge (coordinate) freedom

CSE 576, Spring 2008          Structure from Motion          56

## Levenberg-Marquardt

Iterative non-linear least squares [Press'92]
- Linearize measurement equations

$$\hat{u}_i = f(\mathbf{m}, \mathbf{x}_i) + \frac{\partial f}{\partial \mathbf{m}} \Delta \mathbf{m}$$

$$\hat{v}_i = g(\mathbf{m}, \mathbf{x}_i) + \frac{\partial g}{\partial \mathbf{m}} \Delta \mathbf{m}$$

- Substitute into log-likelihood equation: quadratic cost function in **Δm**

$$\sum_i \sigma_i^{-2} (\hat{u}_i - u_i + \frac{\partial f}{\partial \mathbf{m}} \Delta \mathbf{m})^2 + \cdots$$

CSE 576, Spring 2008          Structure from Motion          57

## Levenberg-Marquardt

Iterative non-linear least squares [Press'92]
- Solve for minimum $\frac{\partial C}{\partial \mathbf{m}} = 0$

$$\mathbf{A} \Delta \mathbf{m} = \mathbf{b}$$

Hessian $\quad \mathbf{A} = \left[ \sum_i \sigma_i^{-2} \frac{\partial f}{\partial \mathbf{m}} \left( \frac{\partial f}{\partial \mathbf{m}} \right)^T + \cdots \right]$

error: $\quad \mathbf{b} = \left[ \sum_i \sigma_i^{-2} \frac{\partial f}{\partial \mathbf{m}} (u_i - \hat{u}_i) + \cdots \right]$

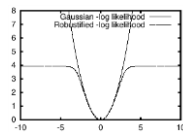CSE 576, Spring 2008          Structure from Motion          58

## Robust error models

Outlier rejection
- use robust penalty applied to each set of joint measurements

$$\sum_i \sigma_i^{-2} \rho \left( \sqrt{(u_i - \hat{u}_i)^2 + (v_i - \hat{v}_i)^2} \right)$$

- for extremely bad data, use random sampling [RANSAC, Fischler & Bolles, CACM'81]

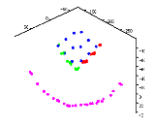CSE 576, Spring 2008          Structure from Motion          63

## Structure from motion: limitations

Very difficult to reliably estimate *metric* structure and motion unless:
- large (*x* or *y*) rotation          *or*
- large field of view and depth variation

Camera calibration important for Euclidean reconstructions

Need good feature tracker

CSE 576, Spring 2008          Structure from Motion          65