## Time varying image analysis

Zoran Duric

- Motion detection
- Computing image motion
- Motion estimation
- Egomotion and structure from motion
- Motion classification

Time-varying image analysis- 1

The problems Visual surveillance - stationary camera watches a workspace - find moving objects and alert an operator - moving camera navigates a workspace - find moving objects and alert an operator Image coding - use image motion to perform more efficient coding of images Navigation - camera moves through the world - estimate its trajectory » use this to remove unwanted jitter from image sequence - image stabilization and mosaicking » use this to control the movement of a robot through the world Time-varying image analysis- 2 Zoran Duric

# Surveillance example: Adding an object to the scene



Time-varying image analysis- 3

Zoran Duric

# Image Sequence Smoothing





Time-varying image analysis- 4

## Motion detection

#### ■ Frame differencing

- subtract, on a pixel by pixel basis, consecutive frames in a motion sequence
- high differences indicate change between the frames due to either motion or changes in illumination

#### Problems

- noise in images can give high differences where there is no motion
  - » compare neighborhoods rather than points
- as objects move, their homogeneous interiors don't result in changing image intensities over short time periods
  - » motion detected only at boundaries
  - » requires subsequent grouping of moving pixels into objects

Time-varying image analysis- 5

Zoran Duric



Image Differencing

Time-varying image analysis- 6





Time-varying image analysis- 7

Zoran Duric



Motion detection algorithms such as these only work if the camera is stationary and objects are moving against a fixed background

Time-varying image analysis- 8

## **Background Subtraction: Results**

Confidence corresponds to gray-level value. High confidence – bright pixels, low confidence – dark pixels.



Time-varying image analysis- 9

Zoran Duric

## Background modeling: color-based

■ At each pixel model colors (*r*,*g*,*b*) or gray-level values *g*. The following equations are used to recursively estimate the mean and the variance at each pixel:

 $\mu_{t+1} = \alpha \mu_t + (1 - \alpha) z_{t+1}$ 

$$\sigma_{t+1}^2 = \alpha (\sigma_t^2 + (\mu_{t+1} - \mu_t)^2) + (1 - \alpha) (z_{t+1} - \mu_{t+1})^2$$

where  $z_{t+1}$  is the current measurement. The mean  $\mu$  and the variance  $\sigma$  can both be time varying. The constant  $\alpha$  is set empirically to control the rate of adaptation (0< $\alpha$ <1).

■ A pixel is marked as foreground if given red value *r* (or for any other measurement, say *g* or *b*) we have

$$|r - \mu_t| > 3 \max(\sigma_r, \sigma_{rcam})$$

Time-varying image analysis- 10

#### Background model

- σ<sub>rcam</sub> is the variance of the camera noise, can be estimated from image differences of any two frames.
- If we compute differences for all channels, we can set a pixel as foreground if any of the differences is above the preset threshold.
- Noise can be cleaned using connected component analysis and ignoring small components.
- Similarly we can model the chromaticity values rc, gc and use them for background subtraction:

 $r_c = r/(r+g+b), g_c = g/(r+g+b)$ 

Time-varying image analysis- 11



## Foreground model

- Use either color histograms (4-bit per color), texture features, edge histograms to model the foreground
- Matching the foreground objects between frames: **tracking**
- Can compare foreground regions directly: shift and subtract. SSD or correlation: *M*, *N* are two foreground regions.

$$SSD = \sum_{i=1}^{n} \sum_{j=1}^{n} [M(i,j) - N(i,j)]^{2}$$
$$C = \frac{\sum_{i=1}^{n} \sum_{j=1}^{n} M(i,j) N(i,j)}{[\sum_{i=1}^{n} \sum_{j=1}^{n} M(i,j)^{2} \sum_{i=1}^{n} \sum_{j=1}^{n} N(i,j)^{2}]^{1/2}}$$

Time-varying image analysis- 13

Zoran Duric



click to start movie

Time-varying image analysis- 14

# Some Intermediate Maps Used in the Method



Time-varying image analysis- 15

Zoran Duric



Time-varying image analysis- 16

#### Results for the sequence



Time-varying image analysis- 17

Zoran Duric

# Using histograms for background modeling

• Use histograms of small regions to model the background:

- Color histograms computed for small regions of the "background" image and the current (new) image (reduced color/ 12 bit bit representation)
- Color edge histograms computed for small regions of the "background" image and the current image (36 bin quantization)

#### **Color Histograms**

Reduced color representation =

C = (R/16) \* 256 + (G/16) \* 16 + (B/16)

(This results in a  $24 \rightarrow 12$  bit color depth reduction)

This results in a 4096 bin histogram

- lowest 4 bits are less useful
- requires less storage
- faster implementation easier to compare histograms

Time-varying image analysis- 19



## **Histogram Matching**

Histogram Intersection

$$I(h_{c}, h_{b}) = \frac{\sum_{i} \min\{h_{c}(i), h_{b}(i)\}}{\sum_{i} \max\{h_{c}(i), h_{b}(i)\}}$$

Chi Squared Formula

$$\chi^{2}(h_{c},h_{b}) = \sum_{i} 2 \frac{(h_{c}(i) - h_{b}(i))^{2}}{h_{c}(i) + h_{b}(i)}$$

Time-varying image analysis- 21

Zoran Duric



- Divide each frame into 40x40 pixel blocks
- To make sure that we do not miss objects on grid block boundaries we tile the frame by overlaying two grids, one of which is shifted by 20 pixels in x and y directions



## Criteria for block activation

- On a block by block basis, similarity measures between background and foreground histograms are computed
- For histogram intersection: If the similarity is below a threshold, *T*, then the block contains a foreground object and is activated for display
- For chi squared: If the *X*<sup>2</sup> measure is greater than a threshold, *T*, then the block contains a foreground object and is activated for display

Time-varying image analysis- 23

Zoran Duric



Time-varying image analysis- 24

# Using edge histograms for detection





Zoran Duric

# Moving person in a cluttered scene



Time-varying image analysis- 26



Time-varying image analysis- 27

Zoran Duric



Time-varying image analysis- 28

# Surveillance: dropping an object



Time-varying image analysis- 29

Zoran Duric

## Surveillance: removing an object



Time-varying image analysis- 30

## Surveillance: Interacting people



Time-varying image analysis- 31

Zoran Duric



## A 1-d gradient technique

- Suppose we have a 1-D image that changes over time due to a translation of the image
- Suppose we also assume that the image function is, at least over small neighborhoods, well approximated by a linear function.
  - completely characterized by its value and slope
- Can we estimate the motion of the image by comparing its spatial derivative at a point to its temporal derivative?
  - example: spatial derivative is 10 units/pixel and temporal derivative is 20 units/frame
  - then motion is (20 units/frame) / (10 units/pixel) = 2 pixels/frame



Time-varying image analysis- 33

Zoran Duric



- Assume I(x,y,t) is a continuous and differentiable function of space and time
- Suppose the brightness pattern is locally displaced by a distance dx, dy over time period dt.
  - this means that as the time varying image evolves, the image brightnesses of points don't change (except for digital sampling effects) as they move in the image
  - I(x,y,t) = I(x + dx, y + dy, t + dt)
- We expand I in a Taylor series about (x,y,t) to obtain
  - $I(x + dx, y + dy, t + dt) = I(x,y,t) + dx \partial I/\partial x + dy \partial I/\partial y + dt \partial I/\partial t$  + (higher order terms)
- - valid only if temporal change is due entirely to motion
- Can rewrite this as  $dI/dt = G_x u + G_y v + G_t = 0$ . The G's are derivatives measured from the image sequence, and u and v are the unknown optic flow components in the x and y directions, respectively



• So, the spatial and temporal derivatives at a point in the image only provide a linear constraint on the optic flow

Time-varying image analysis- 35





Time-varying image analysis- 37

Zoran Duric

# Motion Flow Example: Images



## Motion Flow Example: Normal Flow



Time-varying image analysis- 39

Zoran Duric



Time-varying image analysis- 40

### Recovering u and v

- If the constraint lines in a neighborhood are nearly parallel (i.e., the gradient directions are all similar), then the location of the best fitting (u,v) will be very sensitive to errors in estimating gradient directions.
- More generally, one could fit a parametric form to local neighborhoods of constraint lines, finding parameters that bring constraint lines "nearest" to the estimated motion assigned to each pixel.
  - for example, if we assume that the surface we are viewing in any small image neighborhood is well approximated by a plane, then the optical flow will be a quadratic function of image position in that image neighborhood

Time-varying image analysis- 41

Zoran Duric



- therefore, the solution is not unique - how to choose one?

## A regularization approach

- Solution add a priori knowledge that can choose between the solutions
- Formally, suppose we have an ill posed problem of determining z from data y expressed as
  - Az = y, where A is a linear operator (e.g., projection operation in image formation)
- We must choose a quadratic norm || || and a so-called stabilizing functional ||Pz|| and then find the z that minimizes:
  - $||Az-y||^2 + \lambda ||Pz||^2$
  - $-\lambda$  controls the compromise between the degree of regularization and the closeness of the solution to the input data (the first term).
- T. Poggio, V. Torre and C. Koch, Computational vision and regularization theory, *Nature*, **317**, 1984.

Time-varying image analysis- 43

Zoran Duric



## Token and correlation methods

- Gradient based methods only work when the motion is "small" so that the derivatives can be reliably computed
  - although for "large" motions, one can employ multiresolution methods
- Tracking algorithms can compute motion when the motion is "large"
  - correlation
  - feature tracking
- Correlation
  - choose a kxk window surrounding a pixel, p, in frame i.
  - compare this window against windows in similar positions in frame i+1
  - The window of best match determines the displacement of p from frame i to frame i+1

Time-varying image analysis- 45

Zoran Duric



#### Correlation

- sum of squared gray level differences
- sum of absolute intensity differences
- "robust" versions of these sensitive to outliers

#### Drawbacks of correlation

- matching in the presence of rotation is computationally expensive since all orientations of the window must be matched in frame i+1
- if motion is not constant in the kxk window then the window will be distorted by the motion, so simple correlation methods will fail
  - » this suggests using smaller windows, within which motion will not vary significantly
  - » but smaller windows have less specificity, leading to matches more sensitive to noise



Time-varying image analysis- 47

Zoran Duric

## Motion estimation – token matching

Extract features from each frame (grey level windows, edge detection)

$$S = \begin{pmatrix} \Sigma E_x^2 & \Sigma E_x E_y \\ \Sigma E_x E_y & \Sigma E_y^2 \end{pmatrix} \qquad E_z$$

 $E_x$  and  $E_y$  are x and y components of image gradient

- $-\lambda_1 \ge \lambda_2 \ge 0$  are eigenvalues of M
- If  $\lambda_1 = \lambda_2 = 0$ , mean squared magnitude of the gradient is 0 (flat, unchanging area in the image)
- If  $\lambda_1 > \lambda_2 = 0$ , values do not change in the direction of the corresponding eigenvector (edge)
- If  $\lambda_1 > 0$  and  $\lambda_2 > 0$ , gray values change in multiple directions (corner)
  - »  $\lambda_2 > \tau$ , where  $\tau$  is some threshold

## Motion estimation – token matching

- Match them from frame to frame. Detect tokens in the next frame using lower threshold. Why?
  - Minimize SSD (sum of squared differences) over a neighborhood in the new image. M is a small area around the token (5x5,7x7,11x11)

$$SSD = \sum_{i=1}^{n} \sum_{j=1}^{n} [M(i, j) - N(i, j)]^{2}$$

- Maximize the correlation over a neighborhood in the new image  $\sum_{i=1}^{n} \sum_{j=1}^{n} M(i, j) N(j, j)$ 

$$C = \frac{\sum_{i=1}^{n} \sum_{j=1}^{n} M(i,j) N(i,j)}{\left[\sum_{i=1}^{n} \sum_{j=1}^{n} M(i,j)^{2} \sum_{i=1}^{n} \sum_{j=1}^{n} N(i,j)^{2}\right]^{1/2}}$$

Time-varying image analysis- 49

Zoran Duric

Multiresolution methods

Consider using edges as features for a tracking algorithm for motion estimation. What should the scale of the edge detector be?

- small scale
  - » many edges are detected
  - » easily confused with one another
  - » computationally costly matching problem
- coarse scale
  - » relatively few edges identified
  - » localized only poorly, so motion estimates have high errors
  - » simple matching problem

## Multiresolution methods

- Multiresolution process the image over a range of scales, using the results at coarser scales to guide the analysis at finer scales
  - detect edges at a coarse scale
  - estimate motion by tracking
  - use these estimates as initial conditions for matching edges at next finest scale
- These are also called **focusing** methods or **scale space** methods
  - can also apply to gradient based motion estimators

Time-varying image analysis- 51

Zoran Duric

## 3-D motion and optical flow

- Assume a camera moving in a static environment
- A rigid body motion of the camera can be expressed as a translation and a rotation about an axis through the origin.
- Let
  - **t** be the translational component of the camera motion
  - $\omega$  be the angular velocity
  - $\mathbf{r}$  be the column vector [X Y Z] <sup>T</sup>
- Then the velocity of r with respect to the XYZ coordinate system is
  - $\mathbf{V} = -\mathbf{t} + \mathbf{\omega} \mathbf{x} \mathbf{r}$
- Let the components of
  - **t** = [U V W]<sup>T</sup>
  - $\mathbf{w} = [A B C]^T$

## 3-D Motion and Optic Flow

## Rewrite in component form:

X' = -U - BZ + CYY' = V - CY + AZ

$$Y' = -V - CX + AZ$$

 $\mathbf{Z'} = -\mathbf{W} - \mathbf{A}\mathbf{Y} + \mathbf{B}\mathbf{X}$ 

- where the differentiation is with respect to time
- The optic flow at a point (x,y) is (u,v) where

$$u = x', x = fX/Z$$

v = y', y = fY/Z

Differentiating x and y with respect to time, we obtain

- $u = X'/Z XZ'/Z^2 = (-U/Z B + Cy) x(-W/Z Ay + Bx)$
- $v = Y'/Z YZ'/Z^2 = (-V/Z Cx + A) y(-W/Z Ay + Bx)$

Time-varying image analysis- 53

Zoran Duric

#### 3-D Motion and Optic Flow These can be written in the form $u = u_t + u_r$ $\mathbf{v} = \mathbf{v}_{t} + \mathbf{v}_{r}$ $(u_t, v_t)$ denotes the translational component of the optic flow $(u_r, v_r)$ denotes the rotational component of the optic flow $u_t = [-U + xW]/Z$ $v_t = [-V + yW]/Z$ $u_r = Axy - B(x^2 + 1) + Cy$ $v_r = A(y^2 + 1) - Bxy - Cx$ Notice that the rotational part is independent of Z - it just depends on the image location of a point So, all information about the structure of the scene is revealed through the translational component

## Special case of a plane in motion

Suppose we are looking at a plane while the camera moves

 $- Z = Z_0 + pX + qY$ 

- Then for any point on this plane
  - $Z pX qY = Z_0$
  - $1 p(X/Z) p(Y/Z) = Z_0/Z$
  - $1/Z = [1-pX/Z qY/Z]/Z_0 = [1-px qy]/Z_0$
- So, we can rewrite the translational components of motion for a plane as:
  - $u_t = [-U + xW][1 px qy]/Z_0 = [-U/Z_0 + xW/Z_0] [1 px qy]$
  - $v_t = [-V + yW][1- px qy]/Z_0 = [-V/Z_0 + xW/Z_0] [1- px qy]$
- These are quadratic equations in x and y
- So, if we can compute the translational component of the optic flow at "enough" points from a planar surface, then we can recover the translational motion (with unknown scaling) and the orientation of the plane being viewed.

Time-varying image analysis- 55

Zoran Duric



## Pure translation

- Another way to look at it
  - Let  $\Delta t = 1$ , so that the image center at time t moves from (0,0,0) to (U,V,W) at time t+1
  - Think of the two images as a stereo pair
  - The location of the projection of (U,V,W), the lens center at time t+1 (the "right" image), in the image at time t (the left image) is at location (U/W, V/W) = (u,v)
  - All conjugate lines at time t must pass through this point
  - So, given a point (x,y) at time t, the location of its corresponding point at time t+1 in the **original** coordinate system must line on the line connecting (x,y) to (u,v)
- So, if we know the optic flow at two points in the case of pure translation, we can find the focus of expansion
  - in practice want more than two points

Time-varying image analysis- 57



- Can we recover the third component of motion, W?
- No, because the same optic flow field can be generated by two similar surfaces undergoing similar motions (U,V and W always occur in ratio with Z).

## Normal flows and camera motion estimation

- If we can compute optic flow at a point, then the foe is constrained to lie on the extension of the optic flow vector
- But the aperture problem makes it difficult to compute optic flow without making assumptions of smoothness or surface order
- Normal flow (the component of flow in the gradient direction) can be locally computed at a pixel without such assumptions
- Can we recover camera motion from normal flow?



Time-varying image analysis- 59

Zoran Duric

#### Identifying the FOE from normal flow Assume that the foe is within the field of view of the camera For each point, p, in the image For each normal flow vector, **n**, If p lies in the "correct" halfplane of **n**, then score a vote for p The FOE is the centroid of the connected component of highest scoring points (might be a single pixel, but ordinarily will not be). Alternative code - maintain an array of counters in register with the image For each normal flow vector,**n**, Increment the counters corresponding to all pixels in the "correct" halfplane of **n** Search the array of counters for the connected component of highest vote count For an image containing N normal flow vectors and mxm pixels, both algorithms are (m<sup>2</sup>N), but (2) is more efficient

## Identifying the FOE from normal flow

- What if the FOE is outside the field of view of the camera?
- The image plane is a bad place to represent the FOE to begin with
  - FOE indicates the direction of translational motion
  - Pixels in a perspective projection image do not correspond to equal angular samples of directions
    - » in the periphery, a pixel corresponds to a wide range of directions
  - Solution represent the array of accumulators as a sphere, with an equiangular sampling of the surface of the sphere
    - » Each normal vector will then cast votes for all samples in a hemisphere
    - » Simple mathematical relationship between the spherical coordinate system of the array of counters, and the image coordinate system



Time-varying image analysis- 61

Zoran Duric