

## CZ4003 Computer Vision

Camera system Imaging device: Screen, optical system (len), Aperture, Optical axis  
 $\frac{1}{z} + \frac{1}{z'} = \frac{1}{f}$

Depth of field: range of depths that scene can be at to be acceptable

Smaller aperture: larger depth of field, more light

Scene radiance at  $p$ :  $L(p) = \rho I^\top n$

Image irradiance at  $p$ :  $E = L \frac{\pi \cos(\alpha)}{4(z')^2} \cos^4 \alpha$



Exposure = irradiance  $\times$  time

Spatial Enhancement: Point processing, Spatial Filtering

PP:  $s = Tcr$

Image Negative:  $s = L - r$ , Contrast Stretching:  $s = \begin{cases} 0 & r \leq r_{\min} \\ \frac{255(r - r_{\min})}{r_{\max} - r_{\min}} & (r_{\min} < r < r_{\max}) \\ 255 & \text{else} \end{cases}$

Power-Law:  $s = cr^\gamma$

Histogram: each bin:  $\frac{MN}{L-1}$ , num of pixels  $k \leq k$ ,  $\sum_{j=0}^k n_j$

bins can be filled  $(\sum_{j=0}^k n_j) \div \frac{MN}{L-1} = \frac{L-1}{MN} (\sum_{j=0}^k n_j) = (L-1) (\sum_{j=0}^k p_j)$

SF: weighted sum:  $g(x,y) = \sum_u \sum_v f(x+u, y+v) W(u,v)$  (correlation)  
 $= \sum_s \sum_t f(x-s, y-t) h(s,t)$   $h(s,t) = W(-s, -t)$  convolution

Average filter:  $\frac{1}{9} \times \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$

Gaussian filter:  $\frac{1}{2\pi\sigma^2} \cdot e^{-\frac{x^2+y^2}{2\sigma^2}}$

Median filter

Frequency Enhancement: 
$$F(u, v) = \sum_{x=0}^{N-1} \sum_{y=0}^{M-1} f(x, y) e^{-j2\pi(\frac{ux}{N} + \frac{vy}{M})}$$
  

$$f(x, y) = \frac{1}{NM} \sum_{u=0}^{N-1} \sum_{v=0}^{M-1} F(u, v) e^{j2\pi(\frac{ux}{N} + \frac{vy}{M})}$$

$F$ : How much,  $|F|$ : How strong,  $\angle F$ : How is shifted.

small and ~~high~~ large  $(u, v)$ : low frequency, median  $(u, v)$ : high  $f$

After fftshift: small/large: high  $f$ ; median: low  $f$


Lowpass filter: remove white noise

Gaussian lowpass filter:  $e^{-\frac{u^2+v^2}{2\sigma^2}}$   $\sigma$ : standard deviation, control cutoff  $f$

Highpass filter: image sharpening, edge filtering

Gaussian highpass:  $1 - e^{-\frac{u^2+v^2}{2\sigma^2}}$

Notch rejects a certain  $f$ , Bandreject rejects band



$$x = f \frac{X}{2} \quad y = f \frac{Y}{2}$$

$$\begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} kx \\ ky \\ k \end{bmatrix}$$

real 3D  
image 2D

Rigid Transform: 
$$P = R P_w + T$$

$\uparrow$  rotate matrix     $\uparrow$  world vector     $\uparrow$  linear transform vector

Neater rigid transform:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & T_x \\ r_{21} & r_{22} & r_{23} & T_y \\ r_{31} & r_{32} & r_{33} & T_z \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

Camera to image: 
$$\begin{cases} X_{im} = -X/s_x + O_x \\ Y_{im} = -Y/s_y + O_y \end{cases}$$

$$\begin{bmatrix} X_{im} \\ Y_{im} \\ 1 \end{bmatrix} = \begin{bmatrix} -1/s_x & 0 & O_x \\ 0 & -1/s_y & O_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

Put all together:

$$\begin{bmatrix} kX_{im} \\ kY_{im} \\ k \end{bmatrix} = \begin{bmatrix} -1/s_x & 0 & O_x \\ 0 & -1/s_y & O_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & T_x \\ r_{21} & r_{22} & r_{23} & T_y \\ r_{31} & r_{32} & r_{33} & T_z \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

image coordinates

camera to image frame

perspective



plane to plane: 
$$\begin{bmatrix} kx_{im} \\ ky_{im} \\ k \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} \\ m_{21} & m_{22} & m_{23} \\ m_{31} & m_{32} & m_{33} \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} kx_{im} \\ ky_{im} \\ k \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} \quad \text{set } m_{34} = 1$$

There's a scale  $k$ , so we can set one of  $m_{ij} = 1$ .

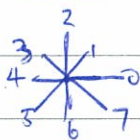
(Actually we can set  $k=1$  ... make it more intuitive.

$n$  points  $\rightarrow 2n$  equations. Need 6 points to solve.

To calibrate cameras: Use calibration chart or cube.

Edgel sequence

Edge processing Path encoding: 1. store the pixels  $(x, y)$  in array or linked list  
2. Chain codes: 0-7 8 directions in total



Edgel sequence: Cost more space

Chain codes: Cost more time to get coordinates

Prewitt filter:

$$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

sobel filter:

$$\begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

Absolute gradient magnitude:  $\nabla f = \sqrt{G_x^2 + G_y^2} \approx |G_x| + |G_y|$

$$h(r) = -e^{-\frac{r^2}{2\sigma^2}} \quad (\text{inversed}) \quad \nabla^2 h(r) = -\left[\frac{r^2 - \sigma^2}{\sigma^4}\right] e^{-\frac{r^2}{2\sigma^2}}$$

Laplacian of Gaussian (LoG) filter: find the zero-crossing in output

Less sensitive to noise due to  $G$  window

Canny Edge Detector:

1. LoG filtering

2. Non-maximal suppression: suppress the point if not local maxima in gradient direction

3. Hysteresis Thresholding:  $t_H$ : set to 1 if higher than  $t_H$

$t_L$ : set to 1 if between  $t_L, t_H$  and have neighbouring pixel

~~set to 0~~ perpendicular to gradient

set to 1

## Polygonal Line Fitting:

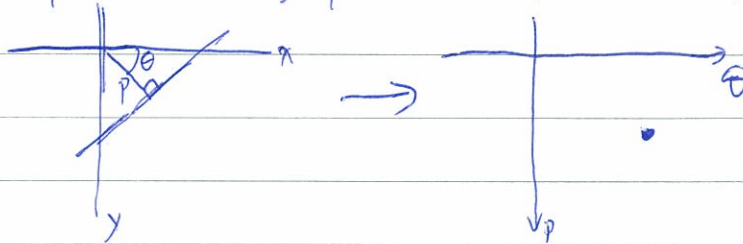
Recursive subdivision algorithm:

1. Connect straight line section connecting end points
2. Find edge with largest distance to any lines
3. Stop when distance is less than threshold
4. Draw new lines

~~Box~~

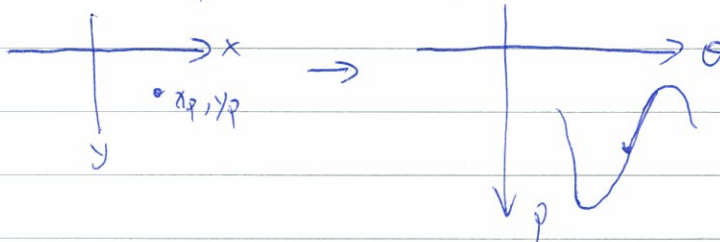
## Hough Transform

~~from~~ Map a line in  $x, y$  plane to a point in  $\theta, \rho$  plane:



Line equation:  $\rho = x_p \cos \theta + y_p \sin \theta = \sqrt{x_p^2 + y_p^2} \sin(\theta + \tan^{-1} \frac{x_p}{y_p})$

Map all lines intersect  $x_p, y_p$  to curve in  $\theta, \rho$  plane:



Hough transform: Transform all edge points into sin curve in  $\rho, \theta$  plane. Find the points in  $\rho, \theta$  plane where most of curves intersect. It corresponds to the line in  $x, y$  plane



Image Region

Processing:

Region encoding

- Labeled mask overlay
- Quadtree representation
- Boundary encoding

Labeled mask overlay: Map the pixel in original image to a colorful mask

Quadtree: Recursively divides a square region into 4 quadrants. Stop ~~when~~ when all pixels in ~~each~~ subregion has common label

Boundary encoding: Define the boundary of each region

Mean - average pixel gray-level:  $\mu = \frac{1}{N} \sum f(x,y) = \sum_{r=0}^{255} r \cdot p(r)$

Variance -  $\sigma^2 = \frac{1}{N} \sum (f(x,y) - \mu)^2 = \sum_{r=0}^{255} (r - \mu)^2 \cdot p(r)$

Naïve Gray-level thresholding: use average value or 127

Otsu Method: Minimizing within group variance thresholding

Within-group variance:  $\sigma_w^2 = q_L(t) \sigma_L^2(t) + q_H(t) \sigma_H^2(t)$   
 $q_L(t) = \sum_{r=0}^t p(r)$     $q_H(t) = \sum_{r=t+1}^{255} p(r)$

Texture representation: structural and ~~statistical~~ statistical

Structural: hard to extract texture

Statistical: neighborhood mean/variance

Co-occurrence Matrix  $C_d$ : dimensions = number of gray-levels

$d$ : displacement vector indicating pairwise relationship

$C_d[r,c]$  = number of a pixel with gray-level  $r$  occurs ~~together~~ together with a pixel with gray-level  $c$  at  $(x+dx, y+dy)$

Normalized  $C_d$ :  ~~$N_d$~~   $N_d[r,c] = \frac{C_d[r,c]}{\sum_r \sum_c C_d[r,c]}$

Energy =  $\sum_r \sum_c N_d^2[r,c]$ , Entropy =  $\sum_r \sum_c N_d[r,c] \log_2 N_d[r,c]$

Contrast =  $\sum_r \sum_c (r-c)^2 N_d[r,c]$ , Homogeneity =  $\sum_r \sum_c \frac{N_d[r,c]}{1 + |r-c|}$

Gabor filters bank to get different statistical features and texture vector

Euclidean distance  $D(T_1, T_2) = \|T_1 - T_2\|$

K-means and Quadtree Decomposition for clustering

Covariance of feature vectors: mean  $\mu = \frac{1}{N} \sum T_i$   
 $S = \frac{1}{N} \sum (T_i - \mu)(T_i - \mu)^T$

The largest eigenvalue  $\lambda_{\max}$  of  $S$  is used to measure ~~the~~ distance between vectors.

After splitting, merge!

## Object Recognition :

Levels of recognition: Image categorization, detection/localization, segmentation

Challenges : view point variation, illumination, ~~occl~~ occlusion, scale, deformation, background clutter, intra-class variation

K-means :

select  $k$  initial points

repeat:

Assign all points to closest centroid

Recompute centroid to each cluster

until centroids don't change much

Solution to initial points:

1. multiple runs
2. hierarchical clustering
3. select more than  $k$  centroids, then select  $k$  from them

Evaluation for K-means :

Cohesion, measured by ~~cluster~~ within-cluster sum of square:

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

Separation, measured by between-cluster sum of square:

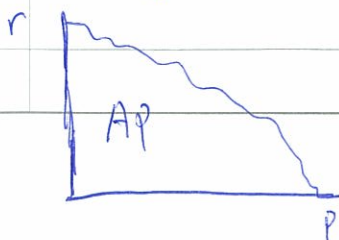
$$BSS = \sum_i |C_i| (m_i - \bar{m})^2$$

Classification :

$$\text{precision} = TP / (TP + FP) \quad \text{recall} = TP / (TP + FN)$$

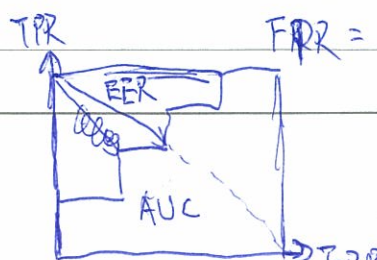
$$F\text{-measure} = 2pr / (p+r)$$

Average Precision



ROC curve:  $TPR = \text{recall} = TP / (TP + FN)$

$$FPR = FP / (FP + TN)$$





(Equal Error Rate)

EER: The point on ROC having same false positive and false negative rate

Area Under Curve (AUC): The higher the better

Bag-of-words model:

Region selection  $\rightarrow$  region description  $\rightarrow$  vector quantization  $\rightarrow$  histogram  $\rightarrow$  classifier

Harris Corner Detector:

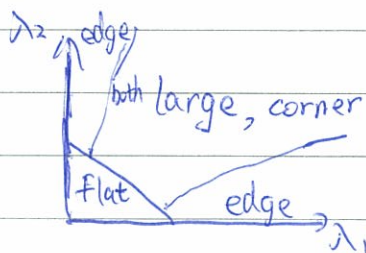
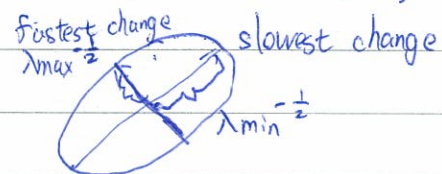
$$E(u,v) = \sum_{x,y} W(x,y) [I(x+u, y+v) - I(x,y)]^2 \quad (W(x,y) = e^{-\frac{r^2}{2\sigma^2}})$$

$$\approx \sum_{x,y} W [u I_x + v I_y]^2$$

$$= [u, v] M [u, v]^T$$

$$M = \sum_{x,y} W(x,y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

$\lambda_1, \lambda_2$  eigenvalues of  $M$ ,



$$R = \lambda_1 \lambda_2 - k (\lambda_1 + \lambda_2)^2 \quad (k = 0.04 - 0.06)$$

large positive = corner      small = flat

large negative = edge

Algorithm: Find points with large  $R$ , take local maxima

Properties: Rotation invariance. Intensity ~~invariance~~ change invariance  
non-invariant to image scale change

SIFT: divide sample image into  $4 \times 4$  subregions.

For each region, generate the histogram of gradient in 8 orientations

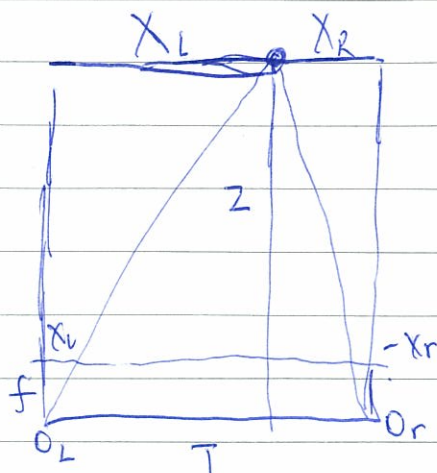
Bag-of-words pros and cons:

flexible to geometry / deformation / view points

compact summary, vector representation, good result

cons: ignore geometry, optimal formation remains unclear

3D stereo vision:



$$x_L = X_L \frac{f}{z}, \quad x_R = X_R \frac{f}{z}, \quad (x_L - x_R) = (X_L - X_R) \frac{f}{z} = \frac{Tf}{z}$$

$$z = \frac{Tf}{d}, \quad \frac{\partial z}{\partial d} = \left| \frac{Tf}{d^2} \right| = \frac{z^2}{Tf} \rightarrow \infty \text{ when } T \rightarrow 0$$

Matching: Appearance based, feature based

~~Appearance~~ Appearance-based: Minimize sum-of-square difference (SSD)

$$\arg \min_{(x,y)} \sum_{u=0}^{M-1} \sum_{v=0}^{M-1} (I(x+u, y+v) - g(u,v))^2$$

$$\arg \max_{(x,y)} 2 \sum_{u=0}^{M-1} \sum_{v=0}^{M-1} I(x+u, y+v) g(u,v) - \sum \sum I(x+u, y+v)^2$$

partition left image into patches, find corresponding patches in right image minimizing SSD

Feature based matching:

Find candidates by special property

Use heuristic to find correct solution

Reconstruction:

$$L \text{ cam: } \begin{bmatrix} kx_L \\ ky_L \\ k \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & a_3 & a_4 \\ a_5 & a_6 & a_7 & a_8 \\ a_9 & a_{10} & a_{11} & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$$R \text{ cam: } \begin{bmatrix} mx_R \\ my_R \\ m \end{bmatrix} = \begin{bmatrix} b_1 & \dots \\ \dots & \dots \\ \dots & \dots & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$



$$\underset{(3 \times 3)}{W} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \underset{(4 \times 1)}{q}$$

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = W^+ q = (W^T W)^{-1} W^T q$$

SVM

Target : find  $y = f(x)$  classify dataset  $\{x_i\}$

Expected risk  $R[f] = \mathbb{E} C(f(x), y) d(x, y)$   $C$  could be squared error

Empirical risk  $R_{\text{emp}}[f] = \frac{1}{N} \sum_{i=1}^N C(f(x_i), y_i)$

~~SVM~~ SVM : define a hyperplane to separate dataset

Margin : minimum distance of a point to the plane, ideally being as large as possible

Larger margin = over learning, cannot tune

minimize  $J(w) = \frac{1}{2} \|w\|^2$  using lagrangian optimization