Computer Vision

Computer Science & Information Technology Yung-Yu Chuang 2009/04/10

with slides by Li Zhang, Szymon Rusinkiewicz, Srinivasa Narasimhan

Can computer match human perception?



- Yes and no (but mostly no!)
 - computers can be better at "easy" things
 - humans are much better at "hard" things



• The goal of computer vision is to write computer programs that can interpret images and understand the scene. The holy grail is to mimic human vision system.

OUTPUT

		descriptions	images
5	descriptions		Computer Graphics
INPL	images	Computer Vision	Image Processing

Computer vision vs Human Vision







What we see

What a computer sees





Exposure

- Two main parameters:
 - Aperture (in f stop)





- Shutter speed (in fraction of a second)





Blade (closing) Blade (open) Focal plane (closed)

Focal plane (open)

Effects of shutter speeds



 Slower shutter speed => more light, but more motion blur Slow shutter, speed Slow shutter speed



• Faster shutter speed freezes motion Walking people Running people

Car

1/500





1/250





Fast train

From Photography, London et al.

- Exposure
- Two main parameters:
 - Aperture (in f stop)
 - Shutter speed (in fraction of a second)
- Reciprocity

The same exposure is obtained with an exposure twice as long and an aperture area half as big

- Hence square root of two progression of f stops vs. power of two progression of shutter speed
- Reciprocity can fail for very long exposures



From Photography, London et al.

Depth of field



Changing the aperture size affects depth of field. A smaller aperture increases the range in which the object is approximately in focus



1/125

Depth of field



Changing the aperture size affects depth of field. A smaller aperture increases the range in which the object is approximately in focus





• A digital camera replaces film with a sensor array

scene

lens &

motor

sensor

array

• Each cell in the array is a light-sensitive diode that converts photons to electrons

CCD v.s. CMOS



- CCD is less susceptible to noise (special process, higher fill factor)
- CMOS is more flexible, less expensive (standard process), less power consumption



SLR view finder (3) Prism Your eye Mirror $\overline{(7)}$ (flipped for exposure) Film/sensor 5 Mirror (when viewing) Light from sce 5 (4) (1) 3 lens

SLR (Single-Lens Reflex)



• Not the case for compact cameras

Color

So far, we've only talked about monochrome sensors. Color imaging has been implemented in a number of ways:

- Field sequential
- Multi-chip
- Color filter array
- X3 sensor

Field sequential Field sequential Field sequential Prokudin-Gorskii (early 1900's) Lantern projector

http://www.loc.gov/exhibits/empire/

Digital camera review website

- <u>A cool video of digital camera illustration</u>
- http://www.dpreview.com/

Now, we have images

- We can think of an image as a function, $f: \mathbb{R}^2 \rightarrow \mathbb{R}$:
 - f(x, y) gives the intensity at position (x, y)

	A CONTRACTOR OF						
62	79	23	119	120	105	4	0
10	10	9	62	12	78	34	0
10	58	197	46	46	0	0	48
176	135	5	188	191	68	0	49
2	1	1	29	26	37	0	77
0	89	144	147	187	102	62	208
255	252	0	166	123	62	0	31
166	63	127	17	1	0	99	30
		~					-

• What about color images?

High-level vision

Recognition

"It's a chair!"

Detection

- Edges
- Lines
- Corners

Image filtering

Convolution with a mask

Low-level vision

62	79	23	119	120	105	4	0
10	10	9	62	12	78	34	0
10	58	197	46	46	0	0	48
176	135	5	188	191	68	0	49
2	1	1	29	26	37	0	77
0	89	144	147	187	102	62	208
255	252	0	166	123	62	0	31
166	63	127	17	1	0	99	30

0	0	1	0	0
0	1	1	1	0
1	1	1	1	1
0	1	1	1	0
0	0	1	0	0

Image filtering (motion blur)

• Convolution with a mask

Gaussian filters

One-dimensional Gaussian

$$G_1(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{x^2}{2\sigma^2}}$$

• Two-dimensional Gaussian

$$G_2(x, y) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{x^2 + y^2}{2\sigma^2}}$$

Image filtering (sharpening)

• Convolution with a mask

Demo with PaintShop Pro

Gaussian filters

Gaussian filters

 $Out(x, y) = \sum_{i} \sum_{j} f(i, j) \cdot In(x - i, y - j)$

- If ln is $n \times n$, f is $m \times m$, takes time $O(m^2 n^2)$
- OK for small filter kernels, bad for large ones

Example: smoothing

Original

Smoothed with Gaussian kernel

Canny edge detector

- Smooth
- Find derivative
- Thresholding
- Thinning

Canny edge detector

• First, smooth with a Gaussian of some width σ

Original Image

blurred Image

Canny edge detector

- Next, find "derivative"
- What is derivative in 2D? Gradient:

Horizontal gradient

Vertical gradient

Original Image

Smoothed Gradient Magnitude

Canny edge detector

thresholding

Original Image

Threshold Gradient Magnitude

Canny edge detector

• Thinning

Original Image

Edges

Canny edge detector

- Nonmaximum suppression
 - Eliminate all but local maxima in *magnitude* of gradient
 - At each pixel look along *direction* of gradient: if either neighbor is bigger, set to zero
 - In practice, quantize direction to horizontal, vertical, and two diagonals
 - Result: "thinned edge image"

Detecting lines

- What is the difference between line detection and edge detection?
 - Edges = local
 - Lines = nonlocal
- Line detection usually performed on the output of an edge detector

Canny demo

Hough transform

- General idea: transform from image coordinates to parameter space of feature
 - Need parameterized model of features
 - For each pixel, determine all parameter values that might have given rise to that pixel; vote
 - At end, look for peaks in parameter space

- Generic line: *y* = *ax*+*b*
- Parameters: a and b

Hough transform for lines

- 1. Initialize table of *buckets*, indexed by *a* and *b*, to zero
- 2. For each detected edge pixel (x, y):
 - a. Determine all (a, b) such that y = ax+b
 - b. Increment bucket (a, b)
- 3. Buckets with many votes indicate probable lines

Hough transform for lines

Issues

- Non-uniform sampling of directions
- Can't represent vertical lines
- Angle / distance parameterization
 - Line represented as (r, θ) where $x \cos \theta + y \sin \theta = r$

Hough transform demo Hough transform demo2

Moravec corner detector (1980)

- We should easily recognize the point by looking through a small window
- Shifting a window in *any direction* should give *a large change* in intensity

 Also known as featurs, interesting points, salient points or keypoints. Points that you can easily point out their correspondences in multiple images using only local information.

Moravec corner detector

Moravec corner detector

flat

Moravec corner detector

Moravec corner detector

flat edge Moravec corner detector Change of intensity for the shift [u, v]: $E(u,v) = \sum w(x, y) [I(x+u, y+v) - I(x, y)]^{2}$ window shifted intensity function intensity Window function w(x, y) =1 in window, 0 outside

Four shifts: (u,v) = (1,0), (1,1), (0,1), (-1, 1)Look for local maxima in $min\{E\}$

Problems of Moravec detector

- Noisy response due to a binary window function
- Only a set of shifts at every 45 degree is considered
- Only minimum of E is taken into account
- ⇒Harris corner detector (1988) solves these problems.

Noisy response due to a binary window function ➤ Use a Gaussian function

$$w(x,y) = \exp\left(-\frac{(x^2+y^2)}{2\sigma^2}\right)$$

Window function w(x,y) =

Gaussian

Harris corner detector

Only a set of shifts at every 45 degree is considered ➤ Consider all small shifts by Taylor's expansion

Harris corner detector

Only a set of shifts at every 45 degree is considered

Consider all small shifts by Taylor's expansion

$$E(u, v) = \sum_{x, y} w(x, y) [I(x+u, y+v) - I(x, y)]^2$$

=
$$\sum_{x, y} w(x, y) [I_x u + I_y v + O(u^2, v^2)]^2$$

$$E(u,v) = Au^{2} + 2Cuv + Bv^{2}$$

$$A = \sum_{x,y} w(x,y)I_{x}^{2}(x,y)$$

$$B = \sum_{x,y} w(x,y)I_{y}^{2}(x,y)$$

$$C = \sum_{x,y} w(x,y)I_{x}(x,y)I_{y}(x,y)$$

Harris corner detector

Equivalently, for small shifts [u, v] we have a *bilinear* approximation:

$$E(u,v) \cong \begin{bmatrix} u & v \end{bmatrix} \mathbf{M} \begin{bmatrix} u \\ v \end{bmatrix}$$

, where ${\bf M}$ is a 2×2 matrix computed from image derivatives:

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

Harris corner detector

Only minimum of *E* is taken into account >A new corner measurement by investigating the shape of the error function

 $\mathbf{u}^{T}\mathbf{M}\mathbf{u}$ represents a quadratic function; Thus, we can analyze *E*'s shape by looking at the property of \mathbf{M}

Harris corner detector

High-level idea: what shape of the error function will we prefer for features?

Harris corner detector

Intensity change in shifting window: eigenvalue analysis

$$E(u,v) \cong \begin{bmatrix} u, v \end{bmatrix} \mathbf{M} \begin{bmatrix} u \\ v \end{bmatrix}$$
 λ_1, λ_2 - eigenvalues of \mathbf{M}

Ellipse E(u, v) = const

direction of the fastest change direction of the slowest change

Summary of Harris detector

4. Define the matrix at each pixel

$M(x, y) = \begin{bmatrix} S_{x^{2}}(x, y) & S_{xy}(x, y) \\ S_{xy}(x, y) & S_{y^{2}}(x, y) \end{bmatrix}$

- 5. Compute the response of the detector at each pixel $R = \det M - k(\operatorname{trace} M)^2$
- 6. Threshold on value of R; compute nonmax suppression.

Threshold on R

Local maximum of R

Mid-level vision

Harris corner detector

Segmentation and clustering

- Defining regions
 - Should they be compact? Smooth boundary?
- Defining similarity
 - Color, texture, motion, ...
- Defining similarity of regions
 - Minimum distance, mean, maximum

Clustering based on color

- Let's make a few concrete choices:
 - Arbitrary regions
 - Similarity based on color only
 - Similarity of regions = distance between mean colors

k-means Clustering

- 1. Pick number of clusters k
- 2. Randomly scatter k "cluster centers" in color space
- 3. Repeat:
 - a. Assign each data point to its closest cluster center
 - b. Move each cluster center to the mean of the points assigned to it

k-means Clustering

k-means Clustering

k-means Clustering

k-means Clustering

k-means Clustering

Results of Clustering

Original Image

k-means, k=5

k-means, k=11

k-means Clustering

Results of Clustering

Sample clusters with *k*-means clustering based on color

Other Distance Measures

- Suppose we want to have compact regions
- New feature space: 5D
 (2 spatial coordinates, 3 color components)
- Points close in this space are close both in color and in actual proximity

Interactive segmentation

video1

video2

Matting

High-level vision

Recognition

Recognition problems

- What is it?
 - Object detection
- Who is it?
 - Recognizing identity
- What are they doing?
 Activities
- All of these are classification problems
 - Choose one class from a list of possible candidates

Face detection

• How to tell if a face is present?

One simple method: skin detection

- Skin pixels have a distinctive range of colors
 - Corresponds to region(s) in RGB color space

• for visualization, only R and G components are shown above

Skin classifier

- A pixel X = (R,G,B) is skin if it is in the skin region
- But how to find this region?

Skin detection

- Learn the skin region from examples
 - Manually label pixels in one or more "training images" as skin or not skin
 - Plot the training data in RGB space
 - skin pixels shown in orange, non-skin pixels shown in blue
 - some skin pixels may be outside the region, non-skin pixels inside. Why?

Skin classifier

• Given X = (R,G,B): how to determine if it is skin or not?

Skin classification techniques

Skin classifier

- Given X = (R,G,B): how to determine if it is skin or not?
- · Nearest neighbor
 - find labeled pixel closest to X
 - choose the label for that pixel
- Data modeling
 - fit a model (curve, surface, or volume) to each class
- · Probabilistic data modeling
 - fit a probability model to each class

Probability

- Basic probability
 - X is a random variable
 - P(X) is the probability that X achieves a certain value

Conditional probability: P(X | Y)
probability of X given that we already know Y

Probabilistic skin classification

- Now we can model uncertainty
 - Each pixel has a probability of being skin or not skin
 - $P(\sim \text{skin}|R) = 1 P(\text{skin}|R)$

Skin classifier

- Given X = (R,G,B): how to determine if it is skin or not?
- Choose interpretation of highest probability
 - set X to be a skin pixel if and only if $R_1 < X \leq R_2$

Where do we get P(skin|R) and $P(\sim skin|R)$?

Learning conditional PDF's

- We can calculate P(R | skin) from a set of training images
 - It is simply a histogram over the pixels in the training images
 - each bin R_i contains the proportion of skin pixels with color R_i

This doesn't work as well in higher-dimensional spaces. Why not?

Learning conditional PDF's

- We can calculate P(R | skin) from a set of training images
 - It is simply a histogram over the pixels in the training images
 - each bin \boldsymbol{R}_i contains the proportion of skin pixels with color \boldsymbol{R}_i

But this isn't quite what we want

- · Why not? How to determine if a pixel is skin?
- We want P(skin | R) not P(R | skin)
- How can we get it?

Bayes rule

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

• In terms of our problem:

normalization term

$P(R) = P(R|skin)P(skin)+P(R| \sim skin)P(\sim skin)$

The prior: P(skin)

Could use domain knowledge

what we want

(posterior)

- P(skin) may be larger if we know the image contains a person
- for a portrait, P(skin) may be higher for pixels in the center
- Could learn the prior from the training set. How?
 - P(skin) may be proportion of skin pixels in training set

- Bayesian estimation
 = minimize probability of misclassification
 - Goal is to choose the label (skin or ~skin) that maximizes the posterior
 - this is called Maximum A Posteriori (MAP) estimation

 $60,000 \times 100 = 6,000,000$ Unique Features

Select 200 by Adaboost

Run-time: 15fps (384x288 pixel image on a 700 Mhz Pentium III)

Application

Smart cameras: auto focus, red eye removal, auto color correction

Face recognition

- Suppose you want to recognize a *particular* face
- How does this face differ from average face
- Consider variation from average face
- Not all variations equally important
 - Variation in a single pixel relatively unimportant
- If image is high-dimensional vector, want to find directions in this space with high variation

Application

Lexus LS600 Driver Monitor System

PCA

• Principal Components Analysis (PCA): approximating a high-dimensional data set with a lower-dimensional subspace

Using PCA for Recognition

• Store each person as coefficients of projection onto first few principal components

image =
$$\sum_{i=0}^{i_{max}} a_i Eigenface_i$$

• Compute projections of target image, compare to database ("nearest neighbor classifier")

 $a_2 \mathbf{v}_2 \ a_3 \mathbf{v}_3 \ a_4 \mathbf{v}_4 \ a_5 \mathbf{v}_5 \ a_6 \mathbf{v}_6 \ a_7 \mathbf{v}_7 \ a_8 \mathbf{v}_8$ $a_1\mathbf{v}_1$

Choosing the dimension K

- NM
- How many eigenfaces to use?
- · Look at the decay of the eigenvalues
 - the eigenvalue tells you the amount of variance "in the direction" of that eigenface
 - ignore eigenfaces with low variance

Kal

High dynamic range imaging/display

Advanced topics

Image warping/morphing

someone not that famous

someone very famous

video

Image warping/morphing

Tracking

Feature tracking

Image stitching

MatchMove

Move matching using scene planes

Matchmove

Move matching using scene planes

Matchmove

Move matching using scene planes

Video matching

Matrix MOCO (Motion control camera)

Photo tourism

Photo Tourism Exploring photo collections in 3D

Video matching

Video matching

Matting and compositing

Titanic

Image manipulation

GraphCut Texture

Matting

Image manipulation

Poisson blending

Image-based modeling

photogrammetric modeling and projective texture-mapping

Image-based modeling

photogrammetric modeling and projective texture-mapping

Image-based modeling

photogrammetric modeling and projective texture-mapping

Image-based modeling

Tour into a picture

Image-based modeling

Tour into a picture

Structured Light and Ranging Scanning

http://graphics.stanford.edu/projects/mich/

3D photography (active)

Cyberware whole body scanner

3D photography (active)

Photometric stereo

Bullet time video

Image-based rendering

Surface lightfield

View interpolation

High-Quality Video View Interpolation

Making face

Gollum

Spacetime face

Video rewrite

Trainable videorealistic speech animation

Inpainting (wire removal)

Inpainting

Texture synthesis/replacement

Texture replacement

Andy Serkis, Gollum, Lord of the Rings

Novel Cameras and Displays

http://www1.cs.columbia.edu/CAVE/projects/cc.htm

Capturing Light Field

Camera Arrays, Graphics Lab, Stanford University

Digital visual effects

Reality?

Bush campaign's TV AD, 2004

Retouching

Texture synthesis and inpainting

This section shows a sampling of the duplication of soldiers.

Production pipeline

Domestic example

Production pipeline

Preproduction

Storyboard

Preproduction

Reference & Research

Preproduction

Artwork

Production

Shooting

Post-production

Visual effects production

