

## L3 notes

### Some terms

$d(h * f)$   $h$  convolve with  $f$

Know what the math equations are doing and apply them

Pad the sides with 0 and shift the kernel across the numbers

$\text{floor}(\text{kernel size} / 2)$

$[1, 1, 1]$  symmetric

Linear (shift invariant) = output Is scaled accordingly

Correlation filter to shift pixels in image

Normalised = summed up add to 1

Effects are negligible for large images

Gaussian kernel depends more on the sigma

Salt and pepper = isolated pixels that are problematic

Smooth away the impulse and spread it all around

Impulse is salt without the pepper

Higher signal amplify both the signals and noise

Histogram stretching adjust contrast

Bitmap save the pixels independently

### Laplacian

the process to get the residual is [upsample (fill with 0s)  $\rightarrow$  blur  $\rightarrow$  subtract from original]

(the reason for the blur is to "spread out" the values over the filled-in black pixels)

Q: Does spatial quantization mean the quantisation of values of x and y?  
i.e. x and y can only take discrete integer values?

A: spatial quantization basically happens because we cant represent continuous x and y coordinates in a discrete system

Q: It seems that you always regard intensity as integers, is it a routine for this course?

A: Yes and no

Q: does adjusting abstract and brightness simply refer to the linear mapping or also include the non-linear gamma mapping

A: Non-linear mapping also just the contrast

Photography operation we dont use stretching ( we use images that stretch over the entire range)

Q: In practice, is convolution used more than correlation?

A: Convolution

Q: Is reconstructing the original image using laplacian pyramid 100% loss-less?

A: Up to the difference of quantisation

Q: Can we compensate motion blur by finding the motion blur kernel and use the inverse kernel to demotion blur?

A: Yes, finding the kernel is always hard

Q: Is convolution invertible? In a sense that let's say we have a convolved image, is there always a filter such that we can recover the original image using convolution? e.g, undoing blur convolution using a specific sharpening convolution

A: If we throw away half of the pixels

Q: Does laplacian pyramids take 2\* space than gaussian pyramids simply because laplacian pyramids have 2 sets of the pyramids (gaussian + residual)?

A: Same amount of space

## L4 notes

### Some terms

mean filter does not remove all noise and blurs the image only isolated values 0-255

median filter remove all noise and blurs the image slightly arrange in increasing starting from 0

difference filter sum up all the values in row same value for the entire row

each pixel has a tangent plane

kernel size controls the strength of the filter

border problem output is reduced in size for the pixels along the edge

the image patch and the template always contain positive numbers,  $\cos \Theta \in [0, 1]$ , i.e., the output of normalized cross-correlation is normalized to the interval  $[0, 1]$ , where 0 means no similarity and 1 means a complete similarity.

region of interest can benefit template matching

image has discrete representation we need approximation, positive gradient value when the image change from dark to bright and negative when reversed image

Image sharpening  $g(x,y) = f(x,y) - f(x,y) \circ h(x,y)$

Magnitude =  $\sqrt{gx^2 + gy^2}$

Approximated magnitude =  $|gx| + |gy|$

non maxim suppression if edge has a magnitude too small connected to another pixel above a threshold, dont prune

extract edges to detect start and end edges , useful for 3d to know the shape and geometry

why? resilient and lighting and color useful for recognition

## Template matching

This object is now the template (kernel) and by correlating an image with this template, the output image indicates where the object is. Each pixel in the output image now holds a value, which states the similarity between the template and an image patch (with the same size as the template) centered at this particular pixel position. The brighter a value, the higher the similarity.

Purely correlation since strongest response is when original image exactly matches output

what happens when u zoom the baby's face 2x (refer to image)?

Template is not symmetrical

## Neighborhood processing

Neighbor pixels play a role when determining the output value of a pixel

## Convolution vs Correlation

Convolution (rotated 180 degrees)

$$h(i, j) \cdot f(x - i, y - j)$$

Correlation

$$h(i, j) \cdot f(x + i, y + j)$$

To check if same result, the kernel before and after 180 degrees are same

Intuition of Sobel filter

$$g_x(x, y) = f(x + 1, y) - f(x - 1, y)$$

correlate  $[1, 0, 1]$  with the image

$$[-1, 0, 1]$$

$$[-2, 0, 2] = 2[-1, 0, 1] \text{ (put more weights at center)}$$

$$[-1, 0, -1]$$

Combine the both result using the horizontal and vertical Sobel to get

the final edge image

3x3 kernel is used as we need to include the neighbours as single row or single column kernel is sensitive to noise

## Derivatives

does not matter the order of the derivative and gaussian blur  
gaussian filter to remove noise before laplacian applied

vertical shows pixel of image row  
horizontal shows the gradient value

$f(x)$  grey level value

$f'(x)$  gradient value

$[1, -2, 1]$

$f''(x)$  gradient of the gradient

$g_{xx}(x, y) \approx f(x-1, y) - 2 \cdot f(x, y) + f(x+1, y)$

$g_{yy}(x, y) \approx f(x, y-1) - 2 \cdot f(x, y) + f(x, y+1)$

## First order derivative

Sobel is a combination of derivative kernel

Sobel results in wide edges as it is first order

First order derivative (thicker edges)

Roberts, Prewitt, Sobel filter

## Second order derivative

Second order derivative can know the exact edge and detect 1 px thin edge

Smoothing is needed

DoG (difference in gaussian)

laplacian is the second order derivative (most sensitive to noise)

-1 white 0 grey 1 black

To approximate the 2nd order derivatives

$g_{xx}(x, y) \approx f(x-1, y) - 2f(x, y) + f(x+1, y)$  represents  $[1; -2; 1]$

$g_{yy}(x, y) \approx f(x, y-1) - 2f(x, y) + f(x, y+1)$  represents  $[1 \ -2 \ 1]$

## Hysteresis

pixels below some value 0

hysteresis join the strong and weak pixels

thresholding depends on the image

cv2.canny is interpolated version

random forest is a learned model from human markups

gradient shift dark to light and light to dark

is the center a local maximum or not if yes keep it as part of the edge  
else discard it

Q: would there be multiple local maxes after performing the non-maximum suppression?

A: no multiple local maxes, strong response right next to each other

Q: is horizontal gradient detection always from left to right? similar for vertical case.

A: left to right

Q: Could you explain the intuition behind why cross correlation isn't commutative/associative, even tho convolution is just flipping the kernel?

A:  $[1, 2, 3]$  cross correlation  $[3, 2, 1]$  is not the same as  $[3, 2, 1]$  cross correlation  $[1, 2, 3]$

$[1, 2, 3]$  cross convolution  $[0, 1, 0]$  is  $[3, 2, 1]$  so it is not identity

Q: From slide 9, how is the 1D derivative filter constructed from the finite differences discrete version equation with  $h = 2$ ?

A: 2 elements we are using

Q: When convolution/cross-correlation is mentioned, is full padding assumed by default?

A: is not full padding, lecture examples are only full padding

Q: is padding to just inserting zero rows and zero columns or that plus blurring/interpolation? It seems to have been used both ways having to do nn interpolation (bilinear, linearly solve between 2 neighbours)0 interspersed in rows and columns

blurring to reduce noise (to avoid enhancing the noise)

Q: I think it was mentioned last lecture that convolutions are generally invertible. How do we reconcile this with the notion that blurring is lossy?

A: convolution is multiplicity in the inverse domain (do division can have problems)

blurring in itself is not lossy only downsampling

for every single elements, composed on several pixels, mixture of weights, same kernel applied to next kernel blurring + downsampling makes it lossy

factor all of this in the system of equations should be able to recover the image

Q: why not take the residual of the original versus upsampled of previous image

A: remove some values these are the values we want to preserve in the residual else will keep some of the detailing

## L6 notes

### Some terms

bandwidth is too narrow not very representative  
we don't know what is the correct/ incorrect value for bandwidth  
sum of all the gaussian forms the new curve  
dirac delta - bandwidth is very small  
number of textons can be smaller/larger than the filter banks  
euclidean distance between histogram ends up compensating for each other  
What are textures used for?  
cue to tell us on the underlying 3d structure  
scale of the envelope, single or many wavelengths

running across all the images where each image is a pixel  
run the cluster where every single pixel and image is a datapoint replace all  
pixels in the email with the texton id, compute the histogram based on the  
samples, each of the texton image is represented with a histogram  
texton id is per pixel  
x and y then the texton would not be purely texture

texture is independent feature to the segment it belongs to, we don't have  
to do segmentation to find the segment in the first place  
the responses is for the whole filter bank  
Do k-means twice  
first k means to form texton dictionary  
second k means NN search (can also apply mean-shift)

histogram is 100 dimension vector also  
Feature Representation  
Filter bank Response  
Texton histogram

Note: Texton histogram as a feature for segmentation will not give a clear  
segmentation, problematic edges, results in a pixelated image for the texture  
(not localised)  
Good localization, pixel from the 1 pixel or all over is the same  
CNN is aggregation of many filters, take a feature response and apply an-  
other kernel to it



blurring each channel independently (2d blurring), now we are blurring the filter response

we have consider the location of the pixel for segmentation

region is bounded by two separate histograms, then u would have separate results

Q: Are we able to somehow use a filter to extract coloraturas as a filter inside the filter bank? So that we can run clustering algo on both texture and colour at the same time?

A: you can mix the colors if you want, create a gabor filter for each channel

Q: can I clarify whether textons are generated based on clustering done over all filter results within a filterbank? Or is it done using only some sort of 'most differentiating' filter result?

A: filter bank only to apply to all images

apply 2d convolution

Q: how does using the same filter bank ensure the bins are the same

A: having very different feature response

1 image per type then maybe

Q: Why 3 clusters?

A: there is nothing to distinguish the skin from the background

## Quiz 1

Red and green considered equally, but blue is not considered at all  
Maximum shades of grey will not be the shade

## L7 notes

### Some terms

SSD

keypoints are also corners or interest points

low error will look dark, high error will look bright

shape is going to affect the ellipse, size does not matters

$H = \begin{bmatrix} 0 & 0 \\ 0 & C \end{bmatrix}$  gradient is x direction is 0

$H = \begin{bmatrix} A & 0 \\ 0 & 0 \end{bmatrix}$  gradient is y direction is 0

greedy approach non-maximum suppression has keypoints approximately

in the same location as it looks for areas that has high contrast  
automatic scale selection, harris is equivariant to changes of scale

## Weights of derivative

gaussian before we do the summation for the weights, more weight at the center, less weights at the edge

Instead of applying the effects, we are looking at two viewpoints under two different lighting conditions

highest response has a signal has characteristic scale that has the same as gaussian, width of signal corresponds in signal characteristic

Harris operator is more efficient compared to eigen value decomposition  
If the region is unusual then it is difficult for matching, approximated distinctiveness with SSD error

Linearize with 2nd moment matrix

Quantify distinctiveness with eigenvalues of H

R is the cornerness is defined as

$$R = \det(H) - k \text{trace}^2(H)$$

Values can be obtained from the matrix itself

LoG finds blobs (keypoint detector that tries to find roughly circular regions) maximum scale is already built-in (maximum size of the kernel wrt to image)

Auto correlation: Create a template of the patch, shift it around where I use the local patch

Q: Is template the benchmark for comparison? A: strongest peak where it matches to itself Refer to template matching

No matter where u shift the template, will get a strong response at the sky

Third dimension is the intensity

We need 3d How do we analyse the points without looking at 3d

$\lambda_{min}$  and  $\lambda_{max}$  are perpendicular to each other

SSD error  $E(u, v)$

Error will be 0 or greater

H is always positive definite

Semi definite where there is 0 case, flat region only in synthetic images

Cross section is approximately circular ( a lot of change in all directions)

Direction of the fastest change perpendicular to the slowest change will be approximately equal

Which direction I change in it is going to be different

Contrast causes the threshold

$\lambda_{min}$  would be small as  $(\lambda_{min})^{-\frac{1}{2}}$  where there is no change on the straight line

Q: Clustering the responses of the textons

A: The response are the data points, each feature response is 1 dimension, each data point cluster in 38 dimensions

Dimensionality: Refer to L6 supplements

One blue dot is 1 pixel, x location is 1 feature value and y location is 1 feature value

Eigenvalue decomposition of pixel then u would get every pixel also axis length is very long

Q: Texture and template matching are they the same?

A: Correlate the long, what is the pattern vs the individual image. Detecting many texture in 1 image, so the feature response will be different in different part of images. Only a particular filter is tuned to that signature response. Normally we don't have to hand-create the filters. Deep learning to learn whole bunch of features, or gabor filter, can tune the filters. Over many dimensions can get a unique response to that example.

Q: Filter bank to find the size

A: Large enough filter bank or some way to count

Binary classification if the sheet is defect or not

Q: Why is the skin and background cannot be segmented?

A: Texture of the background is very large, approximate scale of the marble is on par with the arm. Doing agnostic segmentation, we don't model the content in the image, only based on filter response. There is no reason why

it should end up in certain segmentations.

Clustering is done over all of the pixels over all of the images. Creating the local histogram is trial and error

## **L8 notes**

### **Some terms**

Descriptors are parameter independent

SIFT - scaling iterative feature response

Pixel wise difference (gradients) - more sensitive to relative location but sensitive to deformations

Color histogram - 3d histogram

What size of the cells do we want, fine grids, overlap etc Bin sizes - more finely tuned, a small shift will affect the histogram

Compare histogram - chi square

Orientation normalization

### **Averaging**

Bimodal distribution average end up with dominant orientation with no gradient

MOPS or GIST is no longer used as a descriptor

### **Mode - SIFT**

SIFT - understand the implication of doing this

Can be combined with other keypoint detector

as the threshold is larger, it removes more keypoints, depends on illumination difference

poorly localised if it sits on an edge and not a corner

corners have high curvature, edges have low curvature

needs to be clear where the keypoint is matched to when moving in a straight

line

ratio of eigenvalues to be below the edge thresh high threshold more tolerant of edges

smaller bar to be corner, longer bar is edge

white on black and black on white difference

black bar cannot be keypoints as we also have the surrounding neighbors

We don't assign gradients in original orientation - dominant orientation

Taking the original histogram with some wrapping operation, form of normalisation wrt dominant orientation

Normalize to unit length

If it exceeds the threshold set the threshold value, rectify and assign it. Then renormalize.

Clamping is to encourage invariance to small changes to illumination, over count in certain orientations.

can use in the dense form where it can use in any rotation

High robust

## Feature matching

### Image transformation

Filtering

Only changes on the actual image intensities

$G(x) = h \{ F(x) \}$  Warping

$G(x) = F \{ h(x) \}$  changes the location of the pixels

## L9 notes

### 2d projective transformation

affine: parallel lines are still preserved projective: also known as a homography, try to solve for the parameters for a projective transformation

Not an arbitrary matrix with 9 degrees of freedom (3x3)

Solve for H using direct linear transform

Apply homography to align two images together

## SVD

$$A = U \Sigma V^T$$

A -  $8 \times 9$

U -  $8 \times 9$

$\sigma$  -  $9 \times 9$

V -  $9 \times 9$

V is a singular matrix

singular values found on diagonal values of sigma

columns of V are single vectors

N number of correspondences

Normalization such that the centroid is (0, 0) of the similarity transformation is a whitening step

RANSAC deal with outliers, sample from the set of noisy observations

Even though we can sample more, we usually s to be equal to the number of the outliers

## L10 notes

Hessian is the 2nd order derivative

## Optical flow

The bigger the arrow, the stronger the flow

Small motion assumption - Linear wrt to the location

Constant flow

All flow in the same region will have a flow

$Ax = b$  A is not a square matrix

Have different motions to avoid aperture problem

motion is not small, we end up with wrong matches and estimate shift are actually not smaller than the actuals shift

→ reduce the resolution

smoothing the flow vectors are

warp and upsample, compute the flow until we reach the highest resolution

Line and smooth location, we cannot take the inverse, can only get sparse motion flow with Lucas kanade  
Until we get the closest match  
Mode represent location in the image space, feature representation most closest to the template

Issues with computing flows at object edges

How to associate from 1 frame to another frame

With optimal flow, we can see if people are moving in or out  
sum of  $x$  is the projective range of the template, only for the template not warping the whole image to find the template

When  $p = 0$ , there is no translation, so there is only the change of intensity over time

Do warp on observation, make the comparison between the warped observation and increment on top of it  
Template should stay the same from start to finish, depending on the direction of the warp

$I_t$  if  $p = 0$  as we are updating the whole image after every step, warp image on template

$$T(x) - I(W(x; p))$$

Warping from template to observation image or vice versa  
Number of pixels in the blob would give an indication of how many people are together

Suitable for large images

### **Horn Schunck**

Assumption: smooth flow field  
Is  $\lambda$  is small that the flow goes to 0, then it gives a degenerate flow  
Then we have the smoothest possible flow, i.e. a constant flow field.  
We can also pick  $\lambda$  to be extremely large, in which case there will be no smoothing at all.  
As a min problem,  $u$  and  $v$  to be small

Each pixel flow can be different  
Difference in flow between adjacent pixels are small  
uniformly flow field will have the lowest flow field compared to the random flow field  
We also consider the right and top neighbour so we won't have to double count everything  
 $u_{i,j}$  appears 4 times as it is being referenced by the neighbours  
Only suitable when in the video frames movements are small  
Does not mean that there is motion in the underlying object  
For every model, go back and question our assumptions

Motions are large how to solve this?  
Coarse to fine estimation using LK  
Or non linear assumptions like Horn Schuncks

## **L11 notes**

### **Tracking**

Tracking the speed of the action  
Estimate volume of heart chamber  
We can't rely only on optimal flow for tracking

First search at coarser scale  
Only need to search on local region based on previous state  
Based on initial location so it has to be good  
Can have other types of transformation, rather than just homography (projective transform)  
no closed form to express  $I$  as it has no parameters for  $p$   
Can only compare as much pixels as in the template  
Do in a matrix instead of a loop form, each element is independent so we can parallelise

### **KLT algorithm**

Refresh population of possible corners as features of tracking

### **Duality of feature tracking vs Optimal flow**



Feature tracking

Follow specific pixels As we solve optical flow problem then we can track features, vice versa

## Mean shift

Depends on the background

Spatial layout - what is in the template layout

Not allowing the template to change in size then it will change the distribution of the color histogram, much less than in original template, has issues with scaling (additional dimension in scale)

How do we bin the colors?

10 bins per color channel, then we have 30, 3d histogram that is compounded for every channel.

10 bins may not be discriminative enough to distinguish two shades of color

Initialize for a specific data point and apply to every point

Every point that end up at the same attraction basin belongs to same centroid

Move the points to the centroid to calculate the mean shift vector for every single point

Location that we end up at are the local maxima called attraction basin

The reason we do for every single point for segmentation is we want to identify the membership for each mode

b is the binning function

Look at m which bin does it go into, if it go to mf bin, then  $b(x_n - m) = 0$  and delta function goes to 1

Contribution to q is 0, does not contribute to descriptor value

We are trying to pick out the points that have the colors

Normalizing with the size of the target

We use Bhattacharyya co-efficient as similarity measure (compare histograms) chosen as our feature, if we were using distance measure choose L2

Bhattacharyya co-efficient, L2 is not very ideal for mean-shift

Find the peak of mode or maximizing over all the possible locations y

Assign to new target and calculate the candidate again

### **Challenges of mean shift**

We assume the bounding box is going to move the same but if the person is moving the hand further from the camera, so it is not a good assumption

Repeating the procedure from frame to frame, we may suffer issue from drift (cannot count each color pixel equally), might need to reinitialise at some point

Target may leave out of the frame, there are a lot of objects in the scene

Stationary vs moving cameras

ID shift, move to another object

Moving background foreground and background are changing

## **L12 notes**

### **Data driven image classification**

Learning classifiers based on data

Train classifiers using KNN or DNN

### **Challenges in CV**

Invariance and generalization

### **Where deep learning comes in**

Every single image in orientation and shape

Extrapolate some form of semantics

### **AI**

Rule based system that uses expert indulge human intelligence

Narrow AI only do a dedicated task

Small sliver of human intelligence

Neurons are non linear processing unit

NN is not close to how the brain even works

Mainly done for statistical pattern recognition

Non linearity (Perceptron) are building blocks, weighting the sum which passes through a non-linear activation function

Features are formed with cascade of functions that transform features

### **Perceptron Mark I**

Original perceptron was created to do image processing

Each photocell is sensitive to light

Wires are creating different combinations of inputs

Fully connected - Every single neuron has a connection to the output layer, weighted summation of every single pixel, each representing a probability of the class

Locally connected instead of fully connected

Flatten the matrix into a vector, each pixel is 1 element in the vector which correspond to the input

Hidden units = input\_width  $\times$  input\_height

Parameters for 1 hidden unit = Hidden unit <sup>2</sup> (+1 is negligible)

### **Stationarity**

Definition of corners does not change

Relative comparison to the local neighbourhood upper left and bottom right

Apply same weights to all neurons colored different

Otherwise we would need different weights

Every neuron is centered at every pixel

Square area is called the receptive field of the perceptron

If it has a reduction in half the dimension, receptive field at  $3 \times 3$  will correspond to  $6 \times 6$  in the original image

Will progressively increase even if kernel stay in same size due to reduce resolution

Perceptrons are using convolution operation - sweep the kernel over image

Kernels whose weights are learned

Filtering extract edges, orientations of edges

E.g. Gabor, Sobel filter

Instead learn weights for the filters

### **Feature maps**

Also called a channel

NN also uses pooling rather than simply convolution  
After pooling the resolution would be very coarse  
Aggregate response into a simple value like downsampling  
ReLU rectified to 0  
Convolution is done by many kernels rather than just 1, using parameter sharing

ImageNet is 3 orders of magnitude larger than previous Nets  
Top 1 = 100 - 1  
Top 5 = As long as u have 5 of the classes in the top of the probabilities, it is also considered correct, top 5 most confident correspond to ground truth

### **Extended Image classification**

Not only find out what it is, also where it is