

Perception

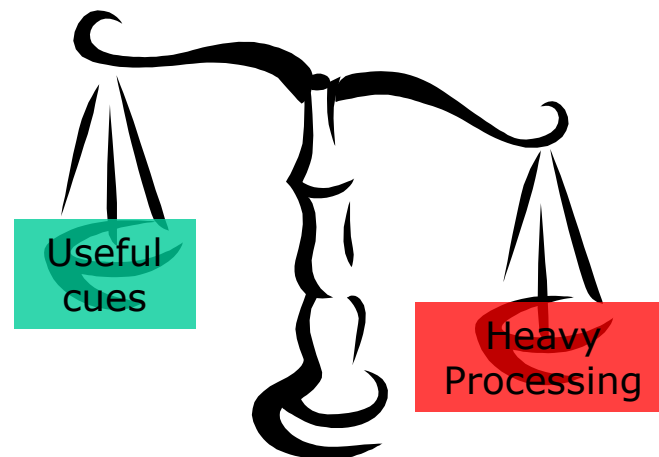
Sensors

Vision

Uncertainties, Line extraction from laser scans

Image Intensities & Data reduction

- Monochrome image \Rightarrow matrix of intensity values
- Typical sizes:
 - 320 x 240 (QVGA)
 - 640 x 480 (VGA)
 - 1280 x 720 (HD)
- Intensities sampled to 256 grey levels \Rightarrow 8 bits
- Images capture a lot of information



- \Rightarrow Reduce the amount of input data:
preserving useful info & discarding redundant info

What is **USEFUL**, What is **REDUNDANT** ?

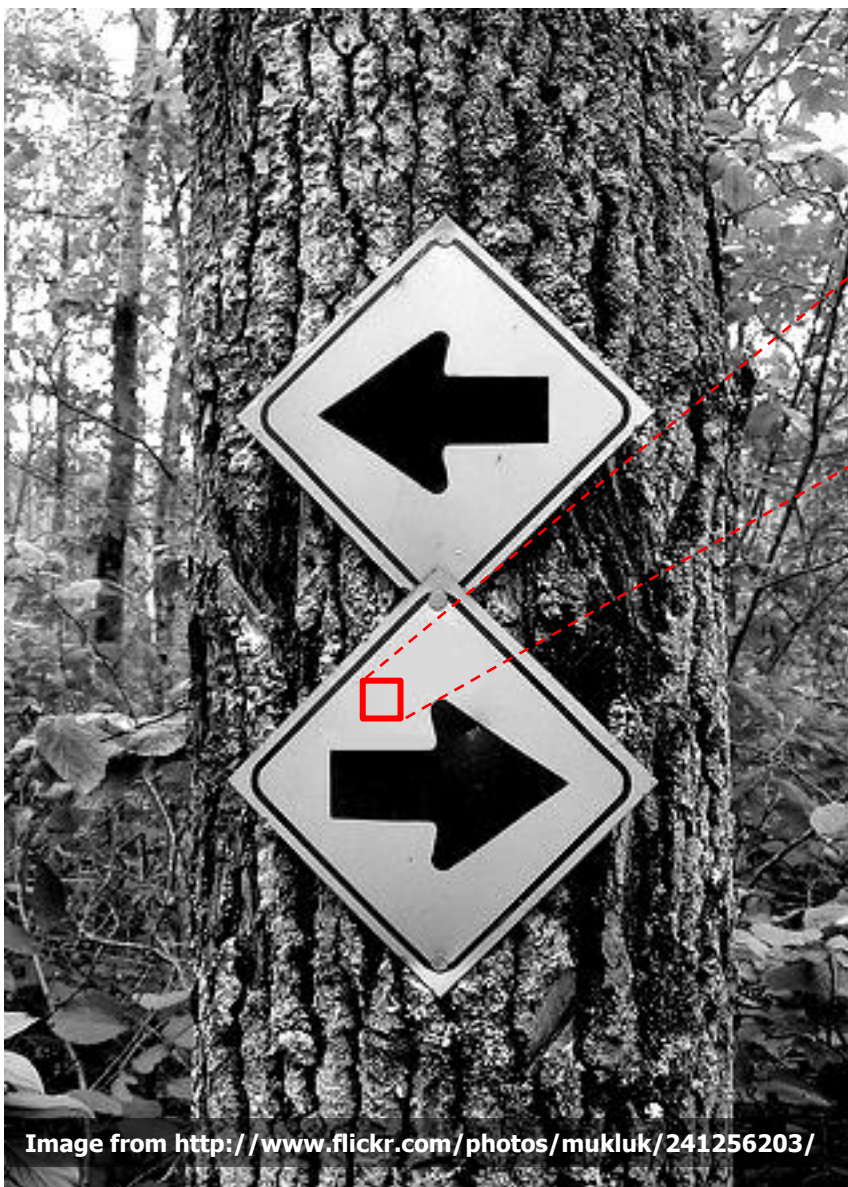
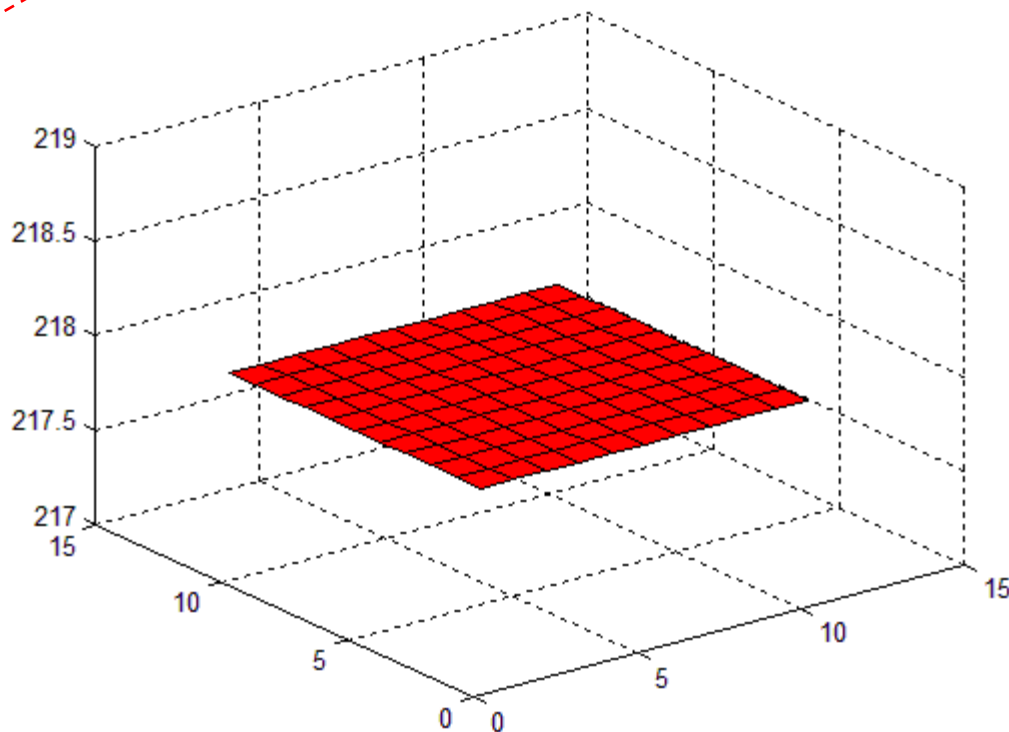


Image from <http://www.flickr.com/photos/mukluk/241256203/>

218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218
 218 218 218 218 218 218 218 218 218 218 218 218



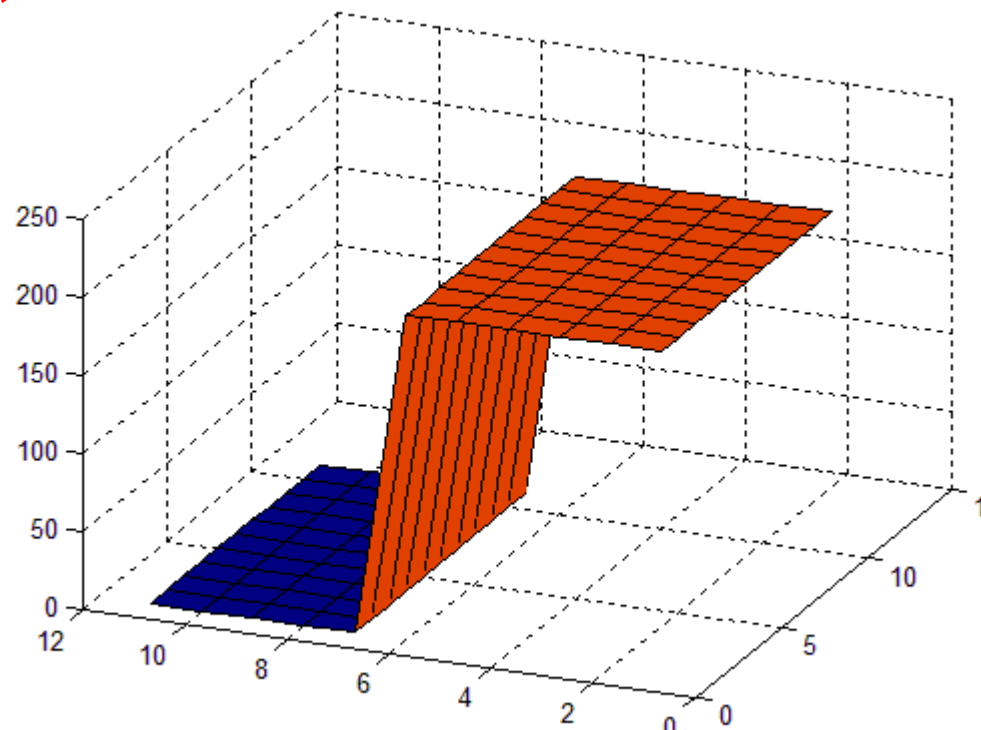
What is **USEFUL**, What is **REDUNDANT** ?



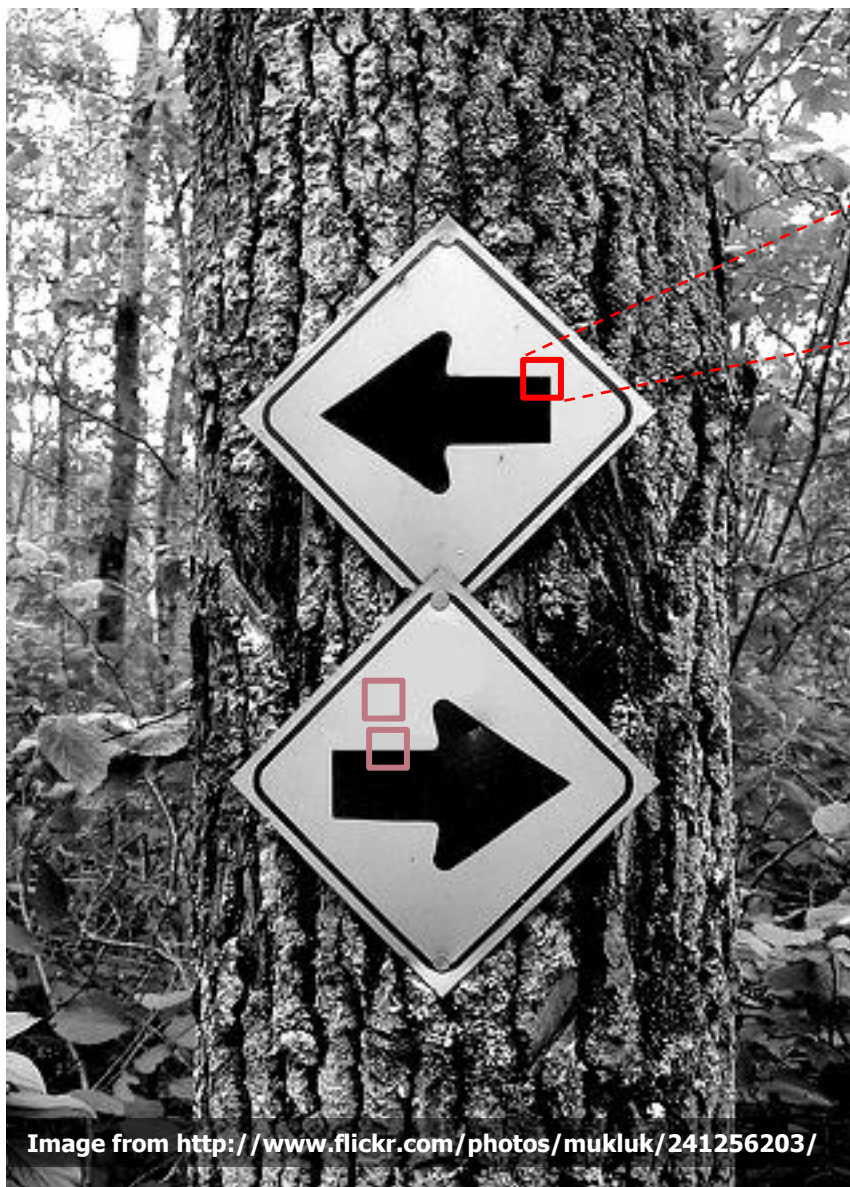
Image from <http://www.flickr.com/photos/mukluk/241256203/>



208	208	208	208	208	208	208	208	208	208	208	208
208	208	208	208	208	208	208	208	208	208	208	208
208	208	208	208	208	208	208	208	208	208	208	208
208	208	208	208	208	208	208	208	208	208	208	208
208	208	208	208	208	208	208	208	208	208	208	208
208	207	208	208	208	208	208	208	208	208	208	208
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0



What is **USEFUL**, What is **REDUNDANT** ?



229	229	229	229	229	229	229	229	229	229	229	229
229	229	229	229	229	229	229	229	229	229	229	229
229	229	229	229	229	229	229	229	229	229	229	229
229	229	229	229	229	229	229	229	229	229	229	229
229	229	229	229	229	230	229	229	229	229	229	229
5	17	31	7	1	0	229	229	229	229	229	229
0	0	1	0	0	0	229	229	229	229	229	229
0	0	0	0	0	0	229	229	229	229	229	229
0	0	0	0	1	4	229	229	229	229	229	229
0	0	0	0	0	11	229	229	229	229	229	229
0	0	0	0	0	5	229	229	229	229	229	229

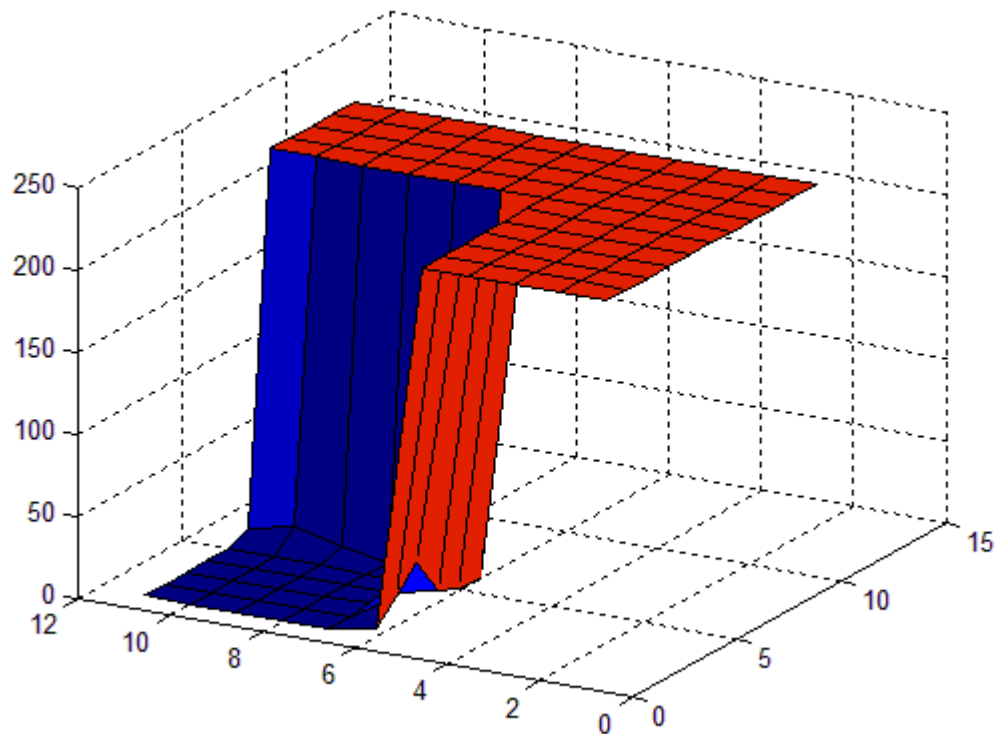


Image from <http://www.flickr.com/photos/mukluk/241256203/>

CS4733 Class Notes, Computer Vision

Sources for online computer vision tutorials and demos - <http://www.dai.ed.ac.uk/HIPR2> and Computer Vision resources online - <http://www.dai.ed.ac.uk/CVonline>

1 Vision Sensing

- The fundamental relationship in imaging a surface is: $I(X, Y) = F(R, G, L)$, where I =intensity of the image at pixel (X, Y) , R =Reflectance of the surface, G =Geometry of the surface, and L =Lighting
- Given the image intensities I , we would like to recover the surfaces we have imaged (i.e. depth and orientation at each point on the surface). There are 2 main problems in inverting this equation:
- Mapping is projection from 3-D to 2-D, which means the inverse is multi-valued (each visible point projects to a unique image point, but each image point “back projects” to a line in space.
- The effects of R , G , L on intensity of the image are coupled. They can not be easily separated out.
- To make vision systems work, we need to add constraints. Without constraints, vision problem is too hard and too ill-posed to solve

2 Machine Vision

- Why is it machine vision so hard when we can “see” with so little conscious effort?
 - Matrix is not a retina; variable resolution in retina
 - Biological systems use active vision. High level of coordination between eye movements and processing
 - Biological vision is robust to lighting changes, surface reflectance changes, color changes, resolution changes
- Robot Vision systems are characterized by:
 - Images tend to be binary, not gray scale
 - Resolution reduced to enable real-time processing
 - Lighting is controlled
 - Objects usually in known position and orientation
 - 2-D methods prevail; 3-D methods typically require more computation

- The process of acquiring an image, processing it and understanding its content (i.e. perception) can be thought of as a “Signals to Symbols” paradigm.
- Low-level: image acquisition, noise reduction, enhancement, edge detection.
- Middle-level: Segmentation and region labeling. Surface recovery - depth and orientation (2 1/2-D sketch). Analysis of texture, motion, color, shading etc.
- High-level: Labeling of images with 3-D components, object recognition, functional analysis

LEVELS OF MACHINE VISION		
LOW	MIDDLE	HIGH
Digitization	Shape From Methods	Scene Understanding
Compression	-texture	3-D Preception
Enhancement	-motion	Object Recognition
Morphology	-shading	Model Building
Features	-stereo	
edges, corners	Segmentation	
	2 $\frac{1}{2}$ -D Sketch	

- The hardest problem in using machine vision is getting the low and high levels integrated.

3 Low Level Vision

We first acquire images digitally. An Image is a continuous signal that is sampled at discrete spacings called pixels. Each pixel is typically quantized to 8 bits of resolution monochrome (256 gray levels) or 24 bits for color (8 bits each for the 3 color channels Red, Blue and Green).

Low-level vision is a series of weak methods to understand simple scenes. Many common low-level processes use the following idea:

- For each image, construct a new filtered image.
- The filtered image will consist of a weighted sum of the pixels surrounding each pixel in the image. Every pixel gets combined *locally* with the same set of weights.

3.1 Filtering

- Images are subject to noise. Common filters include median filter to reduce spike noise, averaging and Gaussian smoothing filters to remove high frequency components. Filtering can be done in the spatial domain with convolutions or in the frequency domain using Fourier techniques.

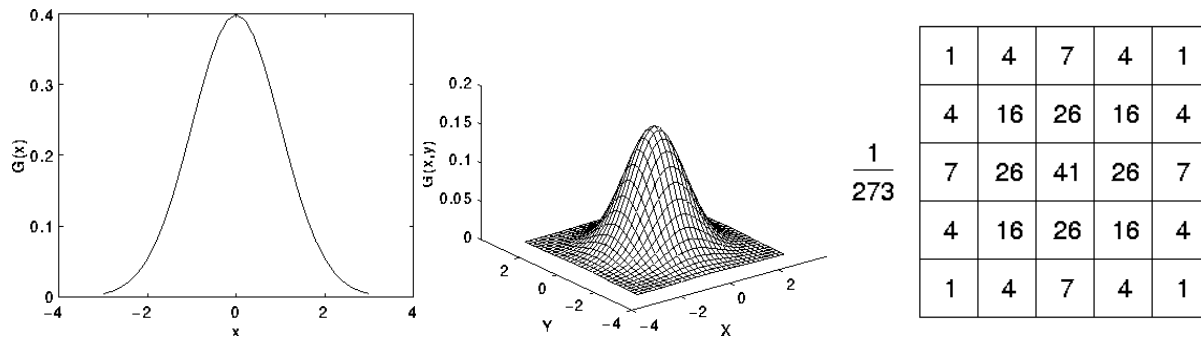


Figure 1: Gaussian filter. Left: 1-D Gaussian with mean=0 and $\sigma = 1$. Middle: 2-D Gaussian with mean=0 and $\sigma = 1$. Right: 5×5 convolution mask for Gaussian smoothing with mean=0 and $\sigma = 1$

- *Mean Averaging Filter*: This filter just averages the pixel values in a neighborhood around a pixel. Neighborhood sizes are variable, depending upon the spatial extent of the filter needed. Common sizes are 3×3 , 5×5 , 7×7 etc. A 3×3 mean filter uses the following set of local weights:

$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$

- *Gaussian Smoothing Filter*: Another smoothing filter is the Gaussian filter, which uses a neighborhood that approximates the fall-off of a Gaussian centered on the pixel of interest. This filter has larger weights for the central pixels and nearest neighbors rather than the mean filter which treats all pixels in the neighborhood with equal weights. See figure 1 above.

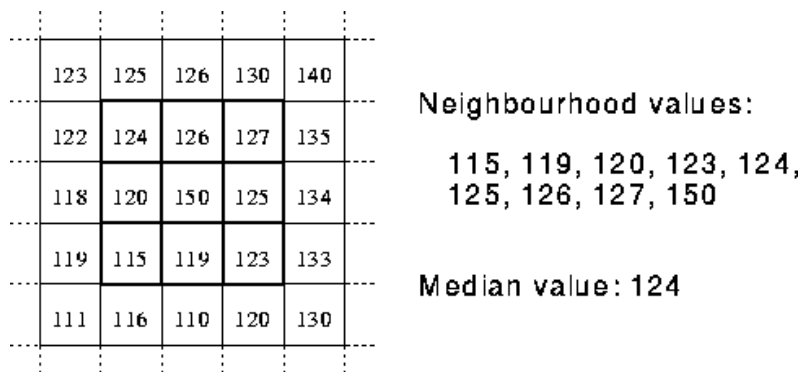


Figure 2: Median filter. Noisy pixel in center (150) is removed by median of its neighborhood.

- *Median Filter*: This filter is used to remove outlier noise values in a region. It is based upon *order statistics*, and is a non-linear filter. In this filter, pixels in a neighborhood are sorted by value, and the *median* value of the pixel's in the neighborhood is taken to be the filter's response. If the pixel being processed is an outlier, it will be replaced by the median value. This filter is useful for "shot" or "salt-and-pepper" noise. See figure 2.

3.2 Enhancement

Often, most of the image values will be centered within a limited range of the full 256 gray levels of an image. *Contrast stretching* performs a linear remapping from the gray level range (I_{low}, I_{high}) to $(0, 255)$, effectively “stretching” the contrast in the image. See figure 3. Before the stretching can be performed it is necessary to specify the upper and lower pixel value limits over which the image is to be normalized. Often these limits will just be the minimum and maximum pixel values in the image. For example for 8-bit graylevel images the lower and upper limits might be 0 and 255. Call the lower and the upper limits a and b respectively.

The simplest sort of normalization then scans the image to find the lowest and highest pixel values currently present in the image. Call these c and d . Then each pixel P is scaled using the following function: $P_{out} = (P_{in} - c) \left(\frac{b-a}{d-c} \right) + a$



Figure 3: Contrast stretching. Original image and histogram and stretched image and histogram.

Histogram equalization is used to change the response over the entire range of gray values. Often, it is used to create a *uniform* histogram that has all gray values used at the same frequency. This may or may not be useful: large homogeneous regions can get remapped into many gray levels, introducing texture (see figure 4). If an image has R rows and C columns, and there are N gray levels $z_1, z_2, z_3, \dots, z_n$ total (e.g. 256) then uniform histogram equalization requires each gray value to occur $q = \frac{R \times C}{N}$ times. Using the original histogram, we define $H_{in}[i]$ as the number of pixels in the original image having gray level z_i . The first gray level threshold t_1 is found by advancing i in the input image histogram until q pixels are used. All input image pixels with gray level $< t_1$ will be mapped to gray level z_1 in the output image:

$$\sum_{i=1}^{t_1-1} H_{in}[i] \leq q < \sum_{i=1}^{t_1} H_{in}[i]$$

This means that t_1 is the smallest gray level such that the original histogram contains no more than q pixels with lower gray values. The k th threshold t_k is defined by continuing the iteration:

$$\sum_{i=1}^{t_k-1} H_{in}[i] \leq k \cdot q < \sum_{i=1}^{t_k} H_{in}[i]$$

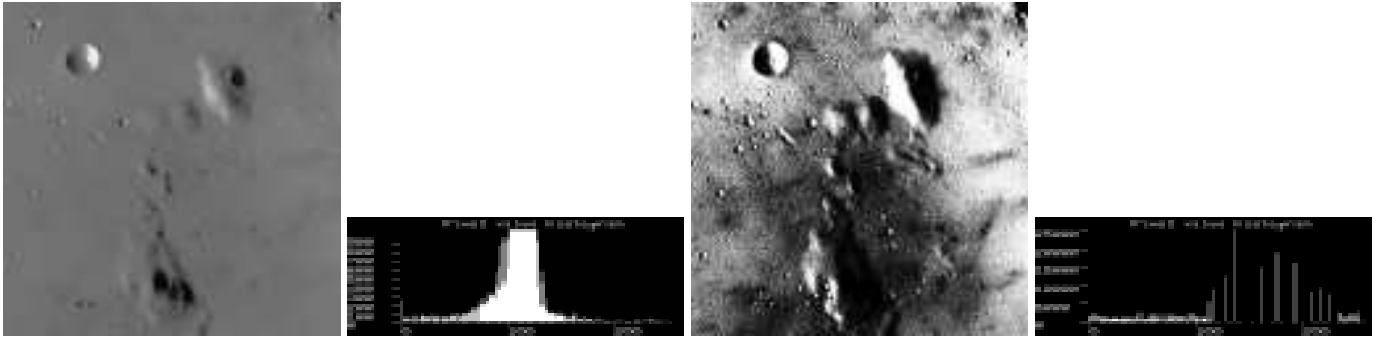


Figure 4: Histogram Equalization. Original image and histogram and equalized image and histogram. See <http://www.dai.ed.ac.uk/HIPR2/histeq.htm>.

3.3 Edge Detection

Find the gradients at each pixel in the image using a gradient operator. Common edge detection masks look for a derivative of the image intensity values in a certain direction. Derivatives are found by differencing the intensity values. The simplest edge detector masks are:

$$VerticalOrientedEdge : \begin{bmatrix} -1 & 1 \end{bmatrix} \quad HorizontalOrientedEdge : \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Each edge detector essentially generates a gradient in the X and Y directions, G_x, G_y . We can calculate the gradient magnitude of the filter's response as:

$$\|G\| = \sqrt{G_x^2 + G_y^2} \quad or \quad ||G|| = |G_x| + |G_y|$$

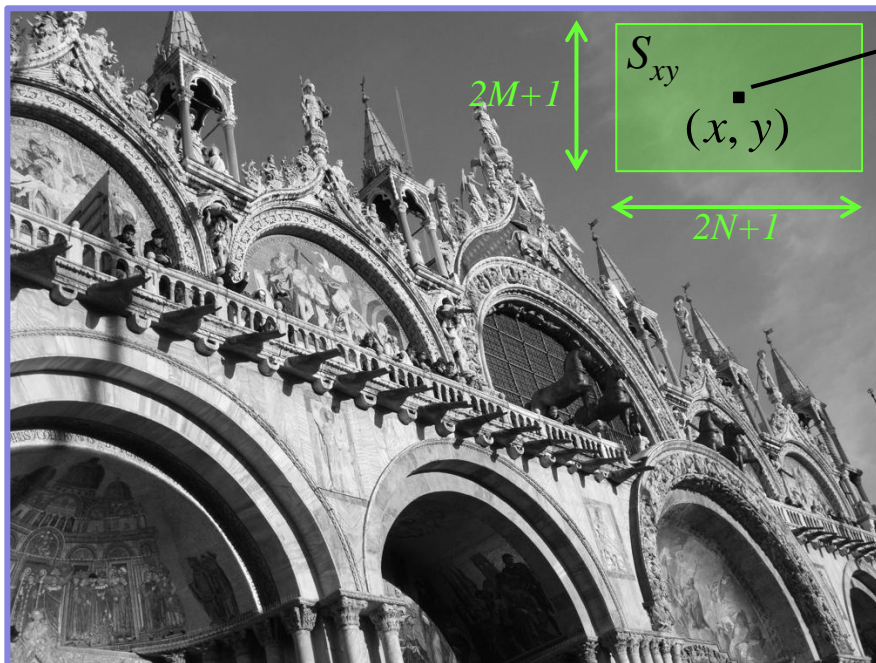
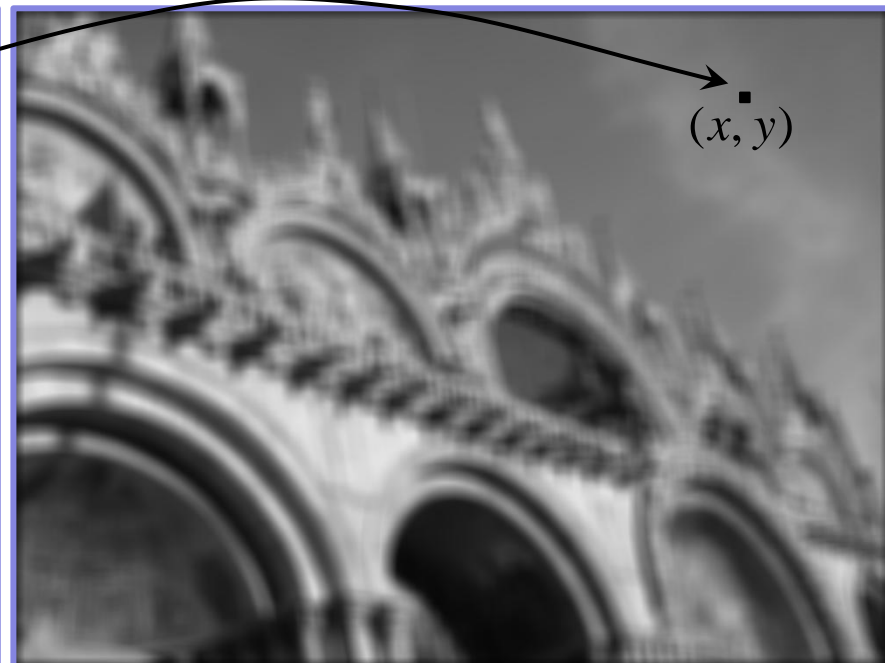
and the edge's orientation (direction) will be $\theta = atan2(G_y, G_x)$.

More sophisticated masks include the Sobel Operators:

$$Vertical : \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \quad Horizontal : \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

Spatial filters

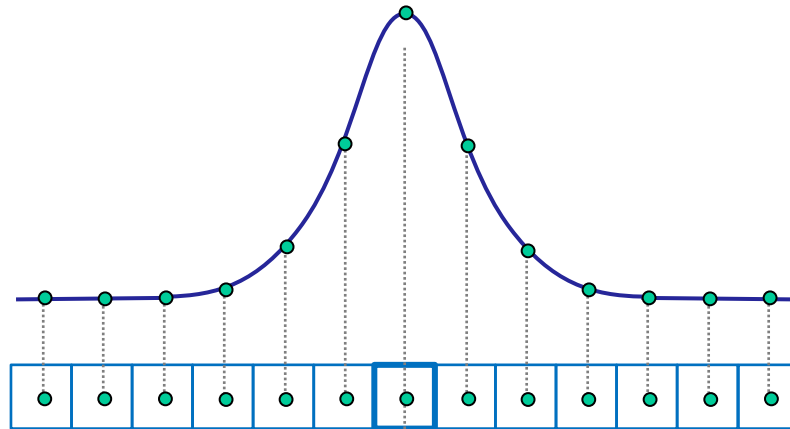
- S_{xy} : neighborhood of pixels around the point (x,y) in an image I
- Spatial filtering operates on S_{xy} to generate a new value for the corresponding pixel at output image J

Image I Filtered Image $J = F(I)$

- For example, an averaging filter is:
$$J(x, y) = \frac{\sum_{(r,c) \in S_{xy}} I(r, c)}{(2M + 1)(2N + 1)}$$

Constructing Filter from a Continuous Fn

- Common practice for image smoothing: use a Gaussian $G(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$



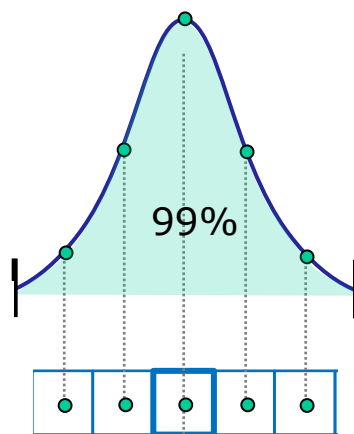
$$\mu = 0$$

σ : controls the amount of smoothing

- Near-by pixels have a bigger influence on the averaged value rather than more distant ones

Constructing Filter from a Continuous Fn

- Common practice for image smoothing: use a Gaussian $G(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$



Normalize filter so that values always add up to 1

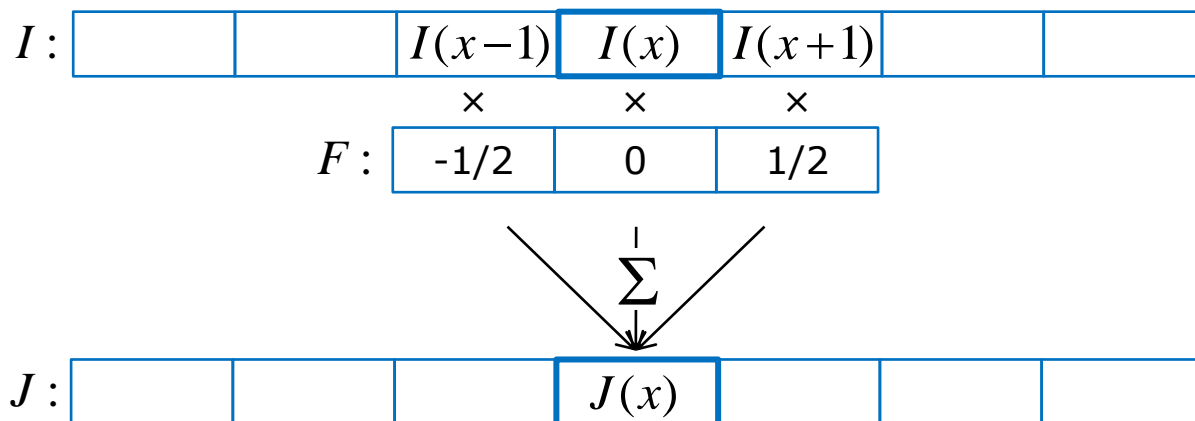
$$\mu = 0$$

σ : controls the amount of smoothing

- Near-by pixels have a bigger influence on the averaged value rather than more distant ones

Taking derivatives with Correlation

- Derivative of an image: quantifies how quickly intensities change (along the direction of the derivative)
- Approximate a derivative operator:



Matching using Correlation

- Find locations in an image that are similar to a **template**
- Filter = template

3	8	3
---	---	---

 \Rightarrow test it against all image locations

I :

3	2	4	1	3	8	4	0	3	8	7	7
---	---	---	---	---	---	---	---	---	---	---	---

- Similarity measure: Sum of Squared Differences (SSD)

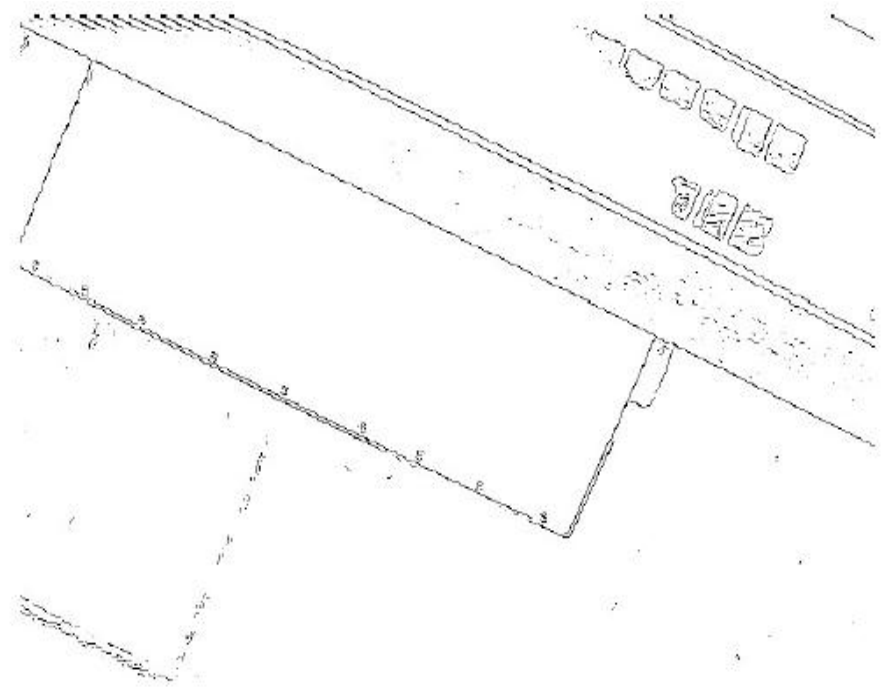
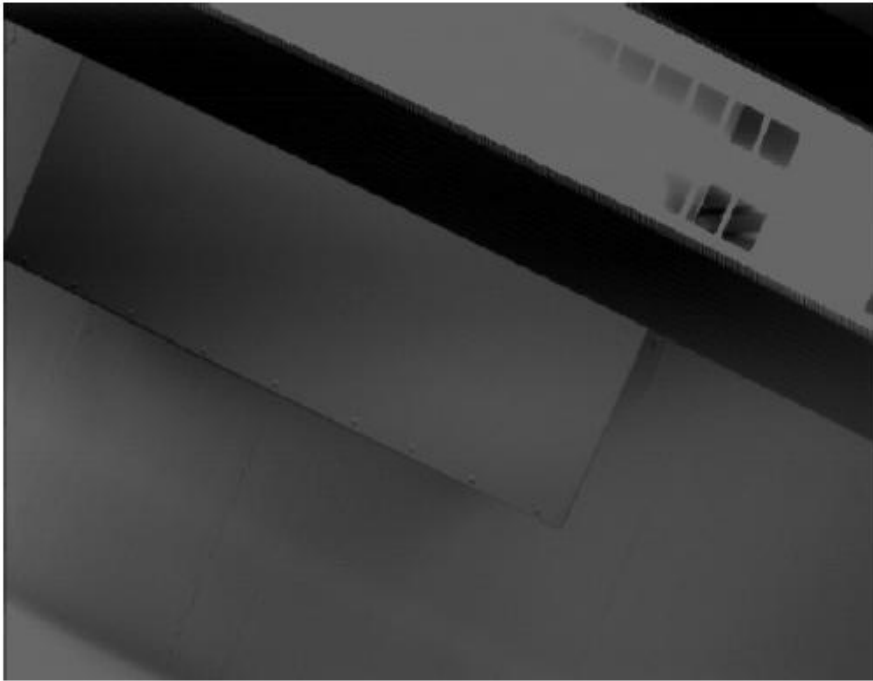
$$\sum_{i=-N}^N (F(i) - I(x+i))^2$$

J :

26	37	21	50	54	1	50	65	59	16	42	17
----	----	----	----	----	---	----	----	----	----	----	----

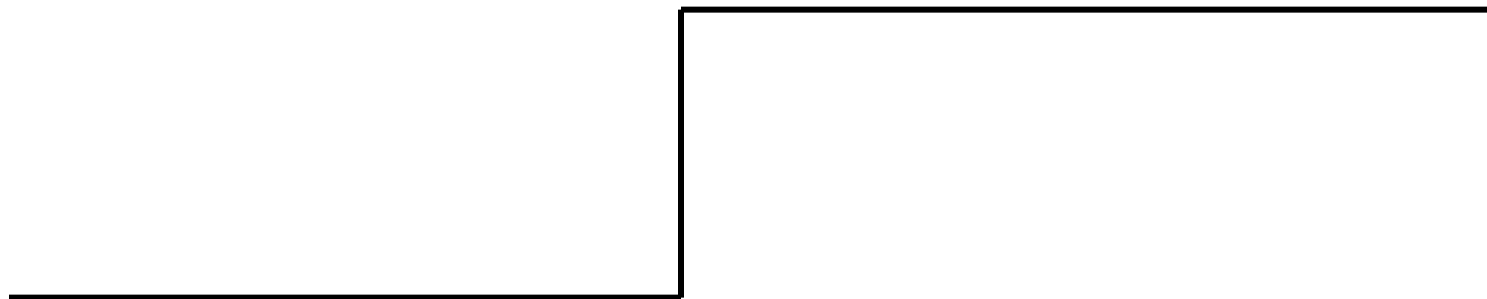
Edge Detection

- Ultimate goal of edge detection: an idealized line drawing.
- Edge contours in the image correspond to important scene contours.



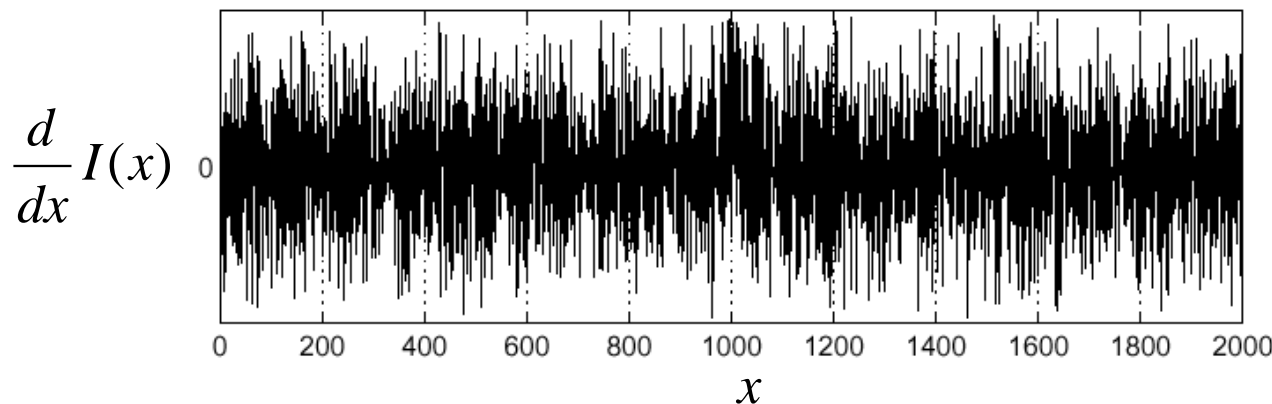
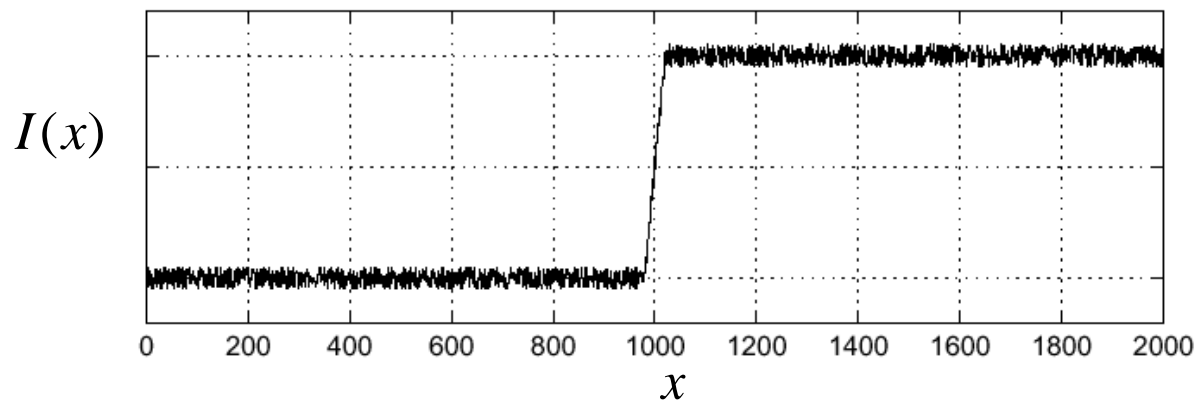
Edge = intensity discontinuity in one direction

- Edges correspond to sharp changes of intensity
- Change is measured by 1st order derivative in 1D
- Big intensity change \Leftrightarrow magnitude of derivative is large
- Or 2nd order derivative is zero.



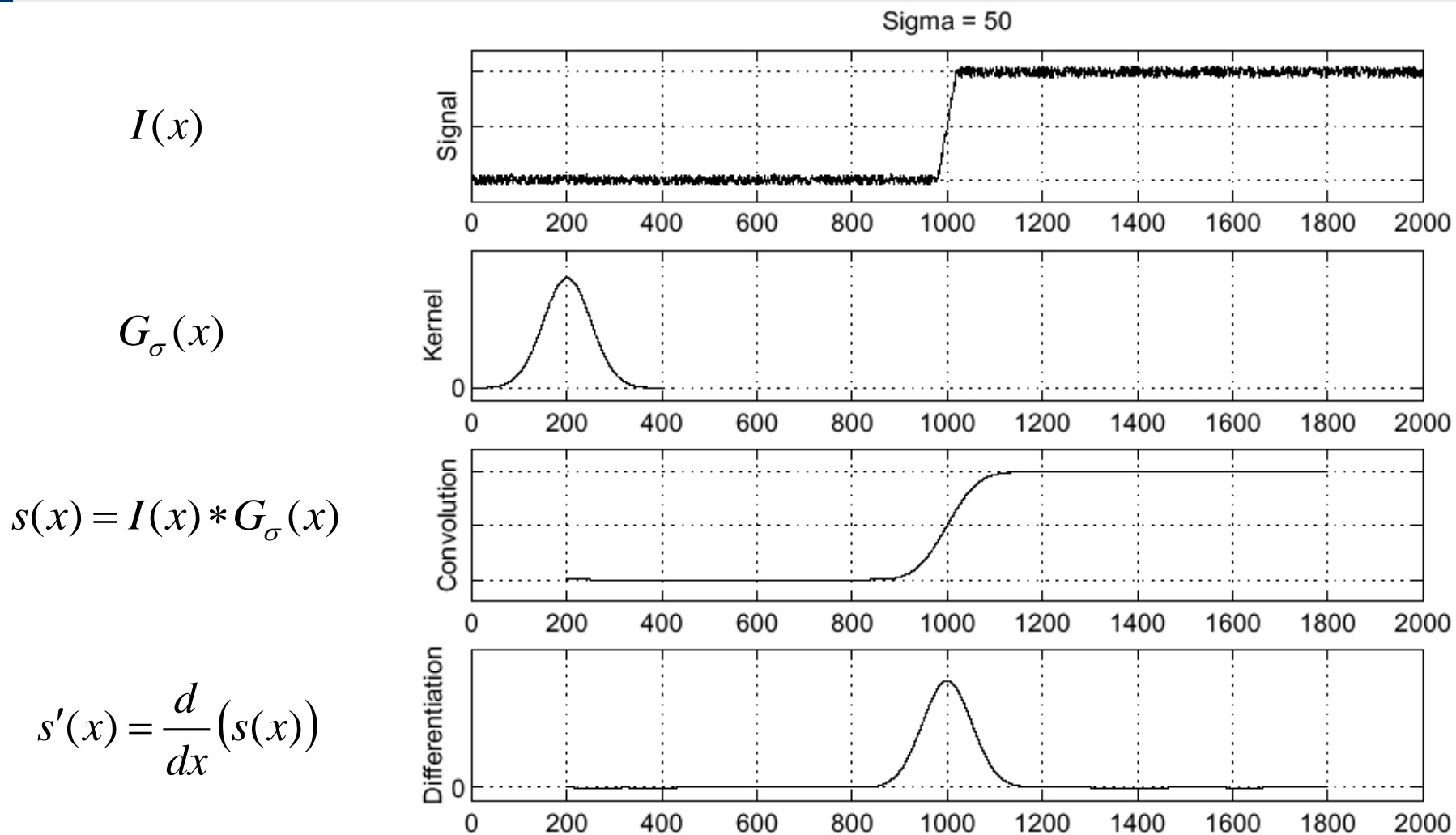
1D Edge detection

- Image intensity shows an obvious change



- Where is the edge? \Rightarrow image noise cannot be ignored

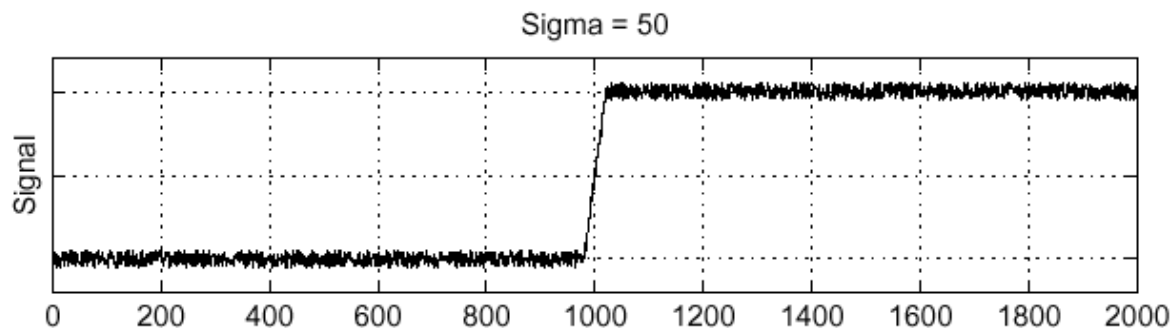
Solution: smooth first



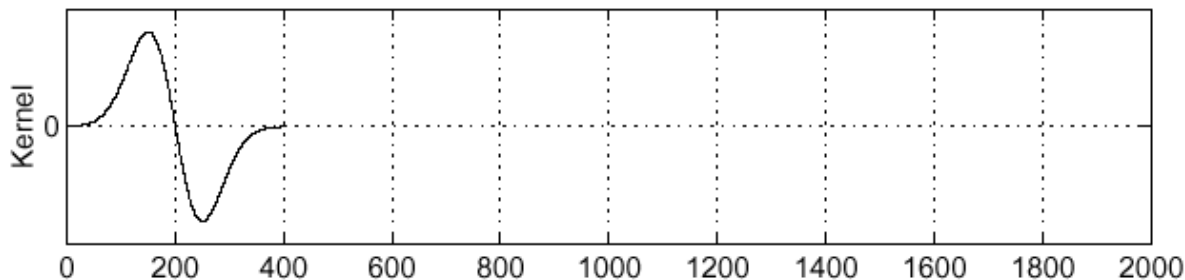
- Where is the edge? At the extrema of $s'(x)$

Derivative theorem of convolution

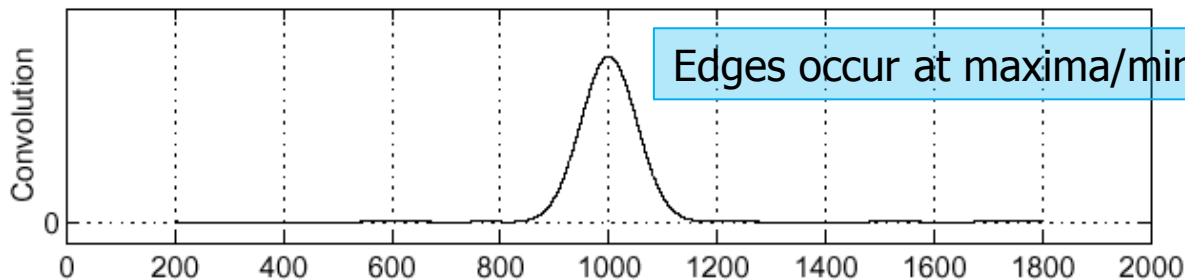
- $$s'(x) = \frac{d}{dx} (G_\sigma(x) * I(x)) = G'_\sigma(x) * I(x)$$
- This saves us one operation.

 $I(x)$


$$G'_\sigma(x) = \frac{d}{dx} G_\sigma(x)$$



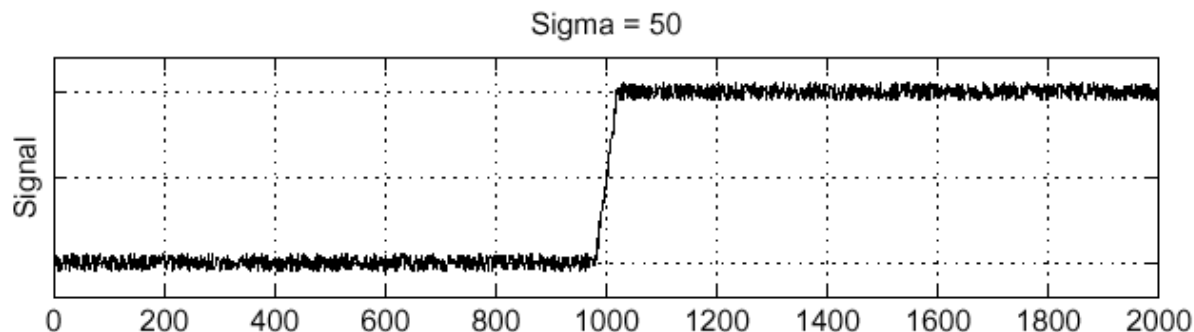
$$s'(x) = G'_\sigma(x) * I(x)$$



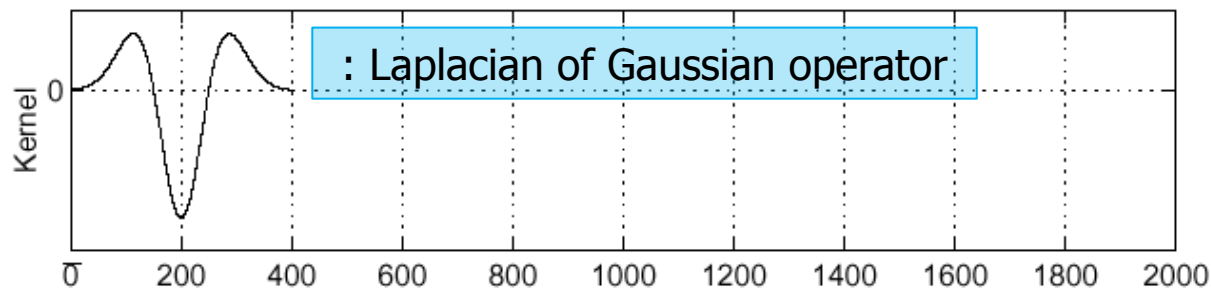
Edges occur at maxima/minima of $s'(x)$

Zero-crossings

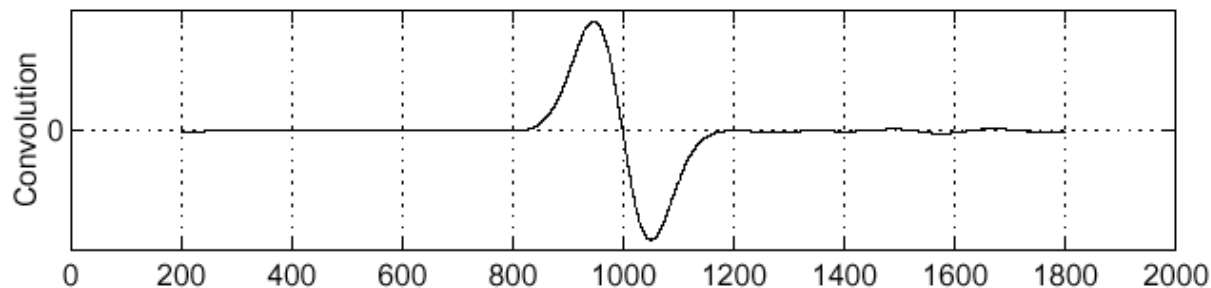
- Locations of Maxima/minima in $s'(x)$ are equivalent to zero-crossings in $s''(x)$

 $I(x)$ 

$$G''_{\sigma}(x) = \frac{d^2}{dx^2} G_{\sigma}(x)$$



$$s''(x) = G''_{\sigma}(x) * I(x)$$



2D Edge detection

- Find gradient of smoothed image in both directions

Usually use a separable filter such that:
 $G_\sigma(x, y) = G_\sigma(x)G_\sigma(y)$

$$\nabla S = \nabla(G_\sigma * I) = \begin{bmatrix} \frac{\partial(G_\sigma * I)}{\partial x} \\ \frac{\partial(G_\sigma * I)}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{\partial G_\sigma}{\partial x} * I \\ \frac{\partial G_\sigma}{\partial y} * I \end{bmatrix} = \begin{bmatrix} G'_\sigma(x)G_\sigma(y) * I \\ G_\sigma(x)G'_\sigma(y) * I \end{bmatrix}$$

- Discard pixels with $|\nabla S|$ below a certain threshold
- Non-maximal suppression:** identify local maxima of $|\nabla S|$ along the directions $\pm|\nabla S|$

2D Edge detection: Example



I : original image (Lena image)

2D Edge detection: Example

$$\nabla S = \nabla(G_\sigma * I)$$



$|\nabla S|$: Edge strength

2D Edge detection: Example



Thresholding $|\nabla S|$

2D Edge detection: Example



Thinning: non-maximal suppression

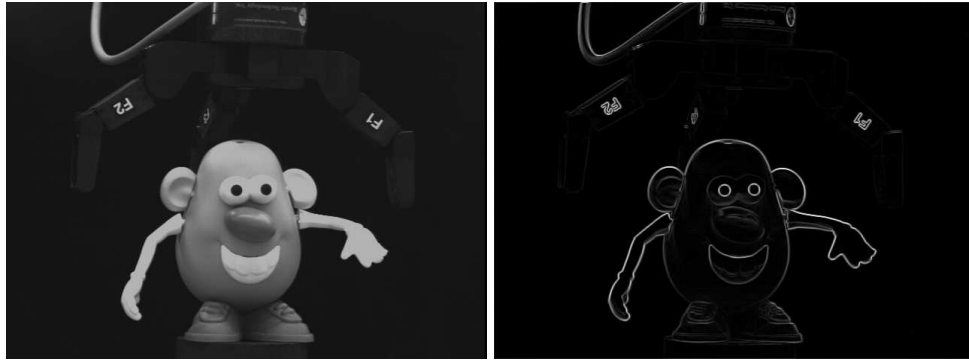


Figure 5: Edge Detection on image. Edge color signifies edge magnitude (brighter == larger magnitude).

4 Middle Level Vision

Middle level vision tries to move beyond the pixel level to larger abstractions including shape and geometry.

4.1 Region Labeling: Recursive Region Growing

Recursive region growing is a simple method. Starting from a binary image, it scans the image for any foreground pixels (not black). For each foreground pixel, it labels that pixel with a unique label, “grows” the pixel by coloring any of its non-black 4-neighbors with this unique color label, and pushing these pixels on a queue. The queue is then processed until empty. All 4-connected pixels in the region will be labeled consistently. Recursive method can be slow however, and may need large memory for recursive calls.

Recursive Region Grower

Do the following for every unvisited SEED pixel.....

Input: Binary image – White(255) = foreground, Black(0) = background. Output: labeled regions

Choose a foreground SEED pixel (pixel whose value = White), $I(c,r)$

Enqueue(c,r), and mark (c,r) as Visited

Label = K # random color

While Queue !Empty do

 (c,r) = Dequeue

 Out_image(c,r) = K

 if $I(c-1,r) \neq \text{Visited} \ \&\& \ I(c-1,r) == \text{White}$ # WEST neighbor pixel, hasn't been Visited and is foreground pixel

 enqueue(c-1,r), mark (c-1,r) as Visited

 if $I(c+1,r) \neq \text{Visited} \ \&\& \ I(c+1,r) == \text{White}$ # EAST neighbor pixel, hasn't been Visited and is foreground pixel

 enqueue(c+1,r), mark (c+1,r) as Visited

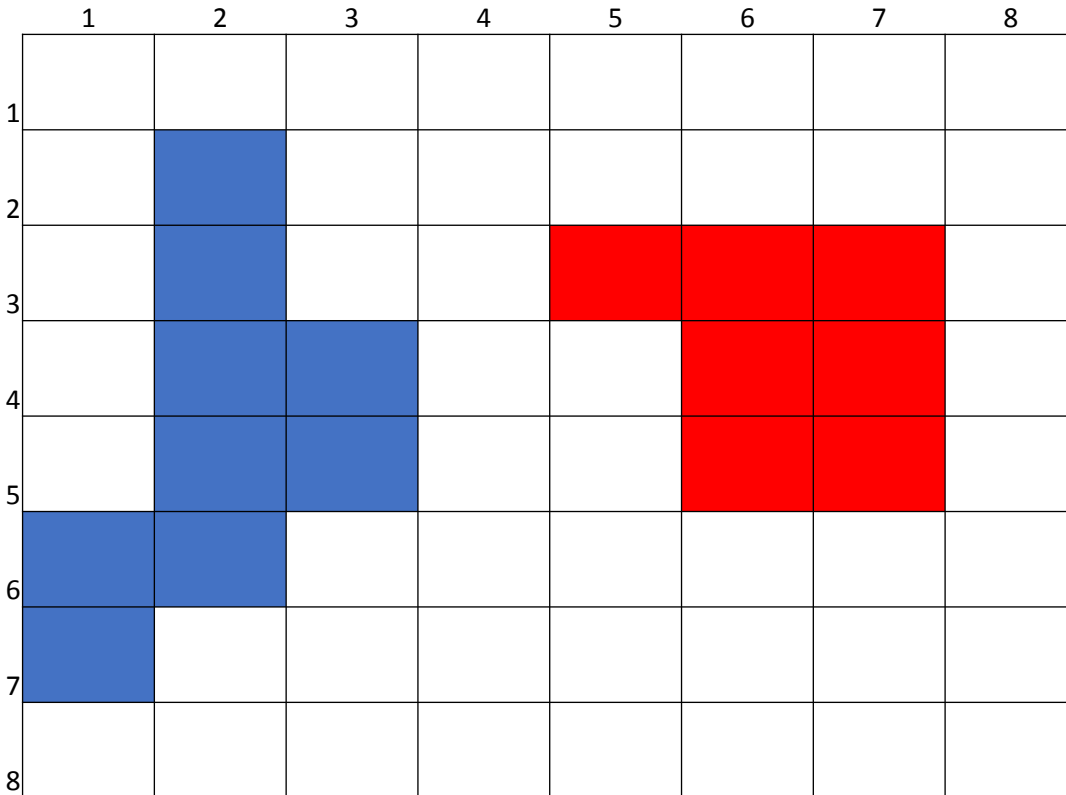
 if $I(c,r-1) \neq \text{Visited} \ \&\& \ I(c,r-1) == \text{White}$ # NORTH neighbor pixel, hasn't been Visited and is foreground pixel

 enqueue(c,r-1), mark (c,r-1) as Visited

 if $I(c,r+1) \neq \text{Visited} \ \&\& \ I(c,r+1) == \text{White}$ # SOUTH neighbor pixel, hasn't been Visited and is foreground pixel

 enqueue(c,r+1), mark (c,r+1) as Visited

Recursive Region Grower - Seed Pixel is (2,2)



Queue	Queue	Queue	Queue
2,2	2,3	2,4	3,4 2,5
Queue	Queue	Queue	Queue
2,5 3,5	3,5 2,6	2,6	1,6
Queue	Queue		
1,7	empty		

4.2 Region Labeling: Blob Coloring

This algorithm uses 2 passes. The first pass labels each pixel and the second pass merges the labels into a consistent labeling.

Let the initial color, $k = init_{color}$, and choose a `color_increment` to change the color each time a new blob is found. Scan the image from left to right and top to bottom. Assign colors to each non-zero pixel in pass 1. In pass2, we merge the regions whose colors are equivalent. To maintain the equivalence table between merged colors, we can use a standard disjoint set Union-Find data structure.

```
If  $I(x_C) = 0$  then continue
else begin

    if  $I(x_U) = 1$  and  $I(x_L) = 0$ 
    then color  $(x_C) :=$  color  $(x_U)$ 
    if  $I(x_L) = 1$  and  $I(x_U) = 0$ 
    then color  $(x_C) :=$  color  $(x_L)$ 
    if  $I(x_L) = 1$  and  $I(x_U) = 1$ 
    then begin /* two colors are equivalent. */
        color  $(x_C) :=$  color  $(x_L)$ 
        color  $(x_L)$  is equivalent to color  $(x_U)$ 
    end

    if  $I(x_L) = 0$  and  $I(x_U) = 0$  /* new color */
    then color  $(x_C) := k$ ;  $k := k +$  color_increment

end
```

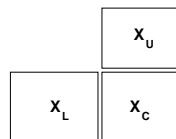


Figure 6: Image topology of x_u , x_c , x_l for region growing



Figure 7: Blob coloring. Left: original binary image. Middle: blob color assignment after first pass. Right: Blob color assignment after merging colors.

Below are 3 ascii images, showing the original test pattern in figure 6, the first pass results, and the final image after region labels are merged. The initial color=80 and the color increment is 50.

```

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

0  0  0  0  0  0  0  0  0
0  80 0  0  130 130 130 130 0
0  80 0  0  0  0  0  130 0
0  80 0  180 0  230 0  130 0
0  80 0  180 0  230 0  130 0
0  80 0  180 180 180 0  130 0
0  80 0  0  0  0  0  130 0
0  80 80 80 80 80 80 80 0

0  0  0  0  0  0  0  0  0
0  130 0  0  130 130 130 130 0
0  130 0  0  0  0  0  130 0
0  130 0  230 0  230 0  130 0
0  130 0  230 0  230 0  130 0
0  130 0  230 230 230 0  130 0
0  130 0  0  0  0  0  130 0
0  130 130 130 130 130 130 130 0

```

5 Simple Shape Matching

- Template Matching: Simple matching of masks (templates) that contain object's image structure
- Object is represented as a region of pixels. Region is compared against all other positions in the image.
- Measure is absolute value of difference between template pixels and image pixels - zero means exact match. Find minimum response for template operator and this is best match
- Problems: Translation, Rotation, Scaling, Lighting changes between image and template
- Translation is handled by applying template everywhere in image
- Rotation handled by using a set of templates oriented every few degrees. Increases cost