

# Image Processing and Computer Vision - Notes

Dom Hutchinson

December 11, 2019

## Contents

<b>1</b>	<b>Image Acquisition</b>	<b>2</b>
<b>2</b>	<b>Image Representation</b>	<b>2</b>
<b>3</b>	<b>Frequency Domains &amp; Image Transforms</b>	<b>5</b>
<b>4</b>	<b>Edges &amp; Shapes</b>	<b>7</b>
4.1	Edge Detection . . . . .	7
4.2	Shapes . . . . .	9
<b>5</b>	<b>Image Segmentation</b>	<b>11</b>
<b>6</b>	<b>Classical Object Detection</b>	<b>13</b>
<b>7</b>	<b>Motion</b>	<b>16</b>
7.1	Motion Estimation . . . . .	18
<b>8</b>	<b>Stereo Vision</b>	<b>20</b>
8.0.1	Correspondence Matching . . . . .	23
<b>0</b>	<b>Reference</b>	<b>26</b>
0.1	Definitions . . . . .	26

# 1 Image Acquisition

## Proposition 1.1 - Common Challenges with Image Acquisition

Below are some common challenges that are faced/produced by image acquisition

Viewpoint Variation	Several images may be taken of the same object but will vary the angle
Illumination	Images may be taken in low/high light
Occlusion	Object may be partly obscured
Scale	Objects may look vary different when placed next to other objects due to their relative scale
Deformation	Objects may have slight variations on the perfect form
Background Clutter	Lots happening behind an object may work to obscure it
Object Intra-Class Variation	Some objects in the same class can vary a lot in shape ( <i>e.g.</i> chairs)
Local Ambiguity	Certain regions of an image can be missinturpred without the rest of the scene being accounted for
World Behind the Image	Depth may need to be accounted for to make sense of an image.

## Definition 1.1 - Dirac Delta-Function, $\delta$

The *Dirac Delta-Function* is used to map continuous distributions to discrete distributions by sampling at particular intervals. Intuitively

$$\delta(t) = \begin{cases} 1 & , t = 0 \\ 0 & , t \neq 0 \end{cases} \implies \delta(t - \alpha) = \begin{cases} 1 & , t = \alpha \\ 0 & , t \neq \alpha \end{cases}$$

## Definition 1.2 - Sifting Property

We can apply the *Dirac Delta-Function* to a function to sample a particular value

$$\int_{-\infty}^{\infty} f(t)\delta(t)dt = f(0) \implies \int_{-\infty}^{\infty} f(t)\delta(t - \alpha)dt = f(\alpha)$$

This can be applied to 2D objects (such as images) as

$$\int_{-\infty}^{\infty} f(a,b)\delta(a - x, b - y)dadb = f(x, y)$$

## Definition 1.3 - Point Spread-Function

A *Point Spread-Function* is applied after sampling an image. It takes the value of a pixel & transforms pixels around it using this value in some way.

*e.g.*  (Should be a white dot on black background but ink).

# 2 Image Representation

## Definition 2.1 - Colour Space

*Colour Space* are different techniques for representing colours. These are generally made up of 3D vectors.

Colour Space	Vector Description
RGB	(Red $\in [0, 255]$ , Green $\in [0, 255]$ , Blue $\in [0, 255]$ )
HSI	(Hue $\in [0, 360]$ , Saturation $\in [0, 1]$ , Intensity $\in [0, 1]$ ) Hue gives the colour in degrees
YUV	(Brightness $\in [0, 255]$ , Blue Projection $\in [0, 255]$ , Red Projection $\in [0, 255]$ )
La*b*	(Luminance $\in [0, 100]$ , Red/Green $\in \{-a, +b\}$ , Blue/Yellow $\in \{+b, -b\}$ )

**Remark 2.1 - Representing Video**

To represent video each fixel is given a third parameter, *time* so we now have

$$f(x, y, t) \mapsto (R, G, B)$$

or any other *Colour Space*.

**Definition 2.2 - Quantisation**

*Quantisation* is representing a continuous single channel function with discrete single channel function that groups the continuous values into a set number of levels.

**Example 2.1 - Quantisation**

16 levels



6 levels



2 levels

**Definition 2.3 - Aliasing**

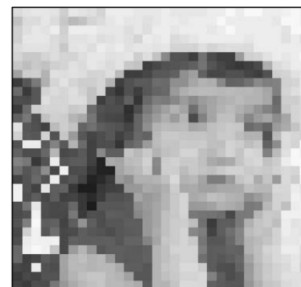
*Aliasing* is the result of sparse sampling since single pixels represent to large an area to get any detail out of it.

**Example 2.2 - Aliasing**

256 x 256



64x64



32x32

**Definition 2.4 - Anti-Aliasing**

*Anti-Aliasing* is the process for avoiding *Aliasing*. This can be achieved by using a sampling rate which is a critical limit defined by the *Shannon-Nyquist Theorem*.

**Theorem 2.1 - Shannon-Nyquist Theorem**

An analogue signal with maximum frequency  $x$ Hz may be completely reconstructed if regular samples are taken with frequency  $2x$ Hz.

**Definition 2.5 - Convolution**

*Convolution* is an operation which takes two functions & produces a third which describes how the shape of one of the two functions is changed by the other.

For functions  $f$  &  $g$

$$(f * g)(x) := \int_{-\infty}^{\infty} f(x - t)h(x)\partial t$$

*N.B.*  $*$  is the symbol for convolution.

**Remark 2.2 - Convolution in Image Representation**

Suppose you have a system, represented by kernel  $g(x)$ , & an input signal, represented by  $f(x)$ . Then  $f * g(x)$  describes the effect of the system on the input signal. The resulting image is called the *Response of  $f$  to the kernel  $h$* .

**Proposition 2.1 - 2D Discrete Convolution**

Since images are represented by discrete 2D functions  $f : \mathbb{N} \times \mathbb{N} \rightarrow (\mathbb{N} \times \mathbb{N} \times \mathbb{N})$  it is pertinent to understand *2D Discrete Convolution*.

$$h(x, y) = \sum_{i \in I} \sum_{j \in J} f(x - i, y - j)g(i, j)$$

Often the kernel,  $g(x, y)$ , has negative indices so the pixel being acted upto is equivalent to the middle pixel in the matrix representation of  $g(x, y)$ .

*N.B.* A convolution whose kernel is symmetric on 180 degree rotation is called a *Correlation*.

**Example 2.3 - 2D Discrete Convolution**

Below is a representation of a grayscale image,  $f(x, y)$ , on the left & a kernel  $g(x, y)$  on the

		$y-1$	$y$	$y+1$		
$x-1$		43	12	61		
$x$		44	45	60		
$x+1$		43	50	61		

	$-1$	$0$	$1$
$-1$	-1	0	1
$0$	-2	0	2
$1$	-1	0	1

right.

$$(f * g)(x, y) = f(x+1, y+1)g(-1, -1) + f(x+1, y)g(0-1, 0) + \dots + f(x-1, y-1)g(1, 1) = -68.$$

**Example 2.4 - Kernels**

Kernels can be defined with specific outcomes in mind.

Operation	Matrix
Identity	$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$
Edge Detection	$\begin{pmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{pmatrix}$ $\begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 1 \end{pmatrix}$ $\begin{pmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{pmatrix}$
Sharpen	$\begin{pmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{pmatrix}$
Box Blur	$\frac{1}{9} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$
Gaussian Blur $3 \times 3$	$\frac{1}{16} \begin{pmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{pmatrix}$

Gaussian Blur $5 \times 5$	$\frac{1}{256} \begin{pmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{pmatrix}$
Unsharp Masking $5 \times 5$	$\frac{-1}{256} \begin{pmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & -476 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{pmatrix}$

### 3 Frequency Domains & Image Transforms

#### Definition 3.1 - Image Transform

An *Image Transform* is deriving a new representation of the input data by encoding the image using another parameter space (e.g. Fourier, DCT, Wavelet, etc.).

#### Remark 3.1 - Purpose of Image Transforms

*Image Transforms* can be used in

- i) Image Filtering;
- ii) Image Compression;
- iii) Feature Extraction;
- iv) etc.

#### Definition 3.2 - Properties of a Signal

A *Signal* is a sinusoidal function over continuous time. They have the following properties

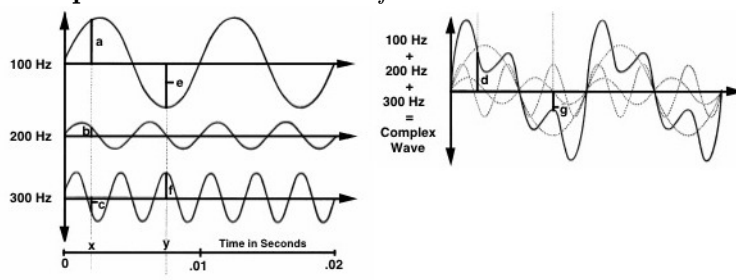
- i) Frequency - Number of cycles per second, Hz;
- ii) Period - Length of a cycle, s;
- iii) Amplitude - Peak intensity of the signal;
- iv) Phase - The shift of the trig wave from its default position,  $\pi$ .

#### Theorem 3.1 - Fourier's Theorem

All periodic functions over continuous time can be expressed as a sum of sin & cos terms, each with their own amplitude & shift.

$$f(t) = a_0 \sin(t + \theta_0) + a_1 \cos(t + \theta_1) + \dots$$

#### Example 3.1 - Fourier Transform



**Proposition 3.1 - Frequency in Images**

*Frequency in Images* is measured as the rate of change in intensity along a given line on the image.

**Remark 3.2 - Fourier Transform on Frequency in Images**

If we read the intensity values along a single row or column we can produce a sinusoidal wave which generalises the distribution & then perform a *Fourier Transform*.

**Definition 3.3 - 2D Discrete-Space Fourier Transform**

Images can be considered as 2-Dimensional discrete space. Let  $f(x, y)$  be the intensity of the pixel at position  $(x, y)$ . 2D Discrete-Space Fourier Transforms have two variables:  $u \in [-\pi, \pi)$  for the vertical frequency; and,  $v \in [-\pi, \pi)$  for the horizontal frequency.

$$\begin{aligned} \underbrace{F(u, v)}_{\text{Fourier Space}} &= \sum_{y=0}^{m-1} \sum_{x=0}^{n-1} f(x, y) e^{i(ux+vy)} \\ &= \sum_{y=0}^{m-1} \sum_{x=0}^{n-1} f(x, y) [\cos(ux + vy) + i \sin(ux + vy)] \end{aligned}$$

*N.B.*  $F(u, v)$  is a complex number.

**Proposition 3.2 - Interpretations of 2D Fourier Transform**

Since  $F(u, v)$  is a complex number we cannot plot it exactly. Thus we consider

i) Magnitude,  $|F(u, v)| := \sqrt{F_r(u, v)^2 + F_i(u, v)^2}$ , and

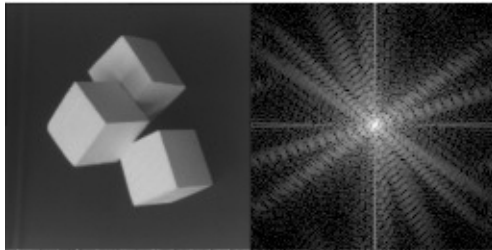
ii) Phase Angles,  $\theta(u, v) := \tan^{-1} \left( \frac{F_i(u, v)}{F_r(u, v)} \right)$

**Remark 3.3 - Expressing  $F(u, v)$  in Polar Coordinates**

$$F(u, v) = |F(u, v)| e^{i\theta(u, v)}$$

**Remark 3.4 - Plotting Magnitude,  $|F(u, v)|$** 

On the left we have a gray scale image & on the right we have the magnitude of a fourier transform on this image. On the right hand image the  $y$ -axis is  $u \in [-\pi, \pi)$  and the  $x$ -axis is  $v \in [-\pi, \pi)$ . We see lots of straight lines since a linear transformation on  $F(u, v) = F(au, av) \forall a \in \mathbb{R}$ . Each line can be interpreted as the frequency of intensity for lines in the left hand image which are parallel to it.

**Theorem 3.2 - Convolution Theorem**

Let  $f$  be an image,  $g$  be a kernel,  $F$  be the result of a fourier transform on  $f$  and  $G$  be a kernel. Then

$$h = f * g \iff H = FG$$

**Proof 3.1 - Convolution Theorem**

$$\begin{aligned}
h(x) &= f(x) * g(x) \\
&= \sum_y f(x-y)g(y) \\
H(u) &= \sum_x \left( \sum_y f(x-y)g(y) \right) e^{iux} \\
&= \sum_x g(y) \sum_x f(x-y) e^{iux} \\
&= \sum_y g(y) (F(u) e^{iuy}) \\
&= \sum_y g(y) e^{iuy} F(u) \\
&= G(u) \cdot F(u) \\
&= F(u) \cdot G(u)
\end{aligned}$$

**Definition 3.4 - Butterworth's Low Pass Filter**

*Butterworth's Low Pass Filter* is a *Signal Processing Filter* designed to have a frequency response which is as flat as possible. It appears to soften an image

$$H(u, v) = \frac{1}{1 + \left( \frac{r(u, v)}{r_0} \right)^{2n}} \text{ of order } n$$

**Definition 3.5 - Butterworth's High Pass Filter**

*Butterworth's High Pass Filter* is a *Signal Processing Filter* designed to have a frequency response which is as flat as possible. It appears to sharpen an image

$$H(u, v) = \frac{1}{1 + \left( \frac{r_0}{r(u, v)} \right)^{2n}} \text{ of order } n$$

## 4 Edges & Shapes

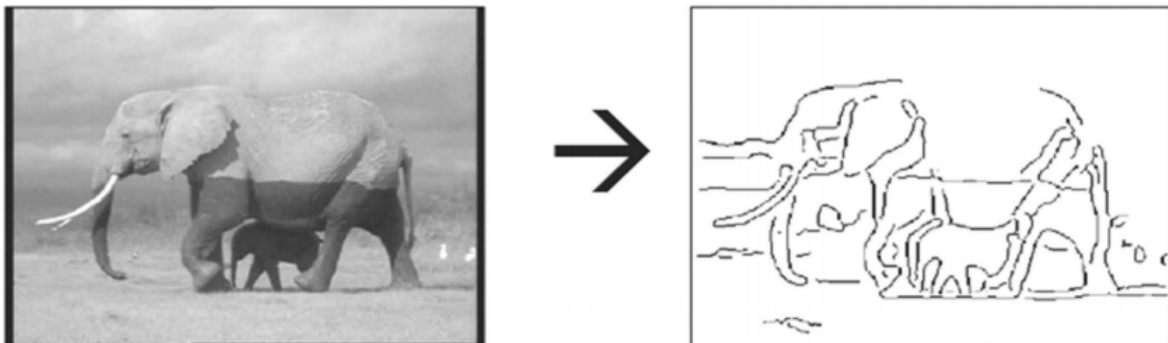
### 4.1 Edge Detection

**Remark 4.1 - Edges are Useful**

Edges are one of the best features to identify an object with. This *Edge-Detection* is a useful thing to be able to do.

**Definition 4.1 - Edge**

An *Edge* is a sharp change in image brightness. Edges are produced by object boundaries, patterns & shadows, this can lead to us finding *Nuisance Edges* which do not help achieve our goal.

**Example 4.1 - Edge Detection****Remark 4.2 - Uses of Edges**

*Edges* are used for

- i) Segmentation - Finding object boundaries;
- ii) Recognition - Extracting patterns;
- iii) Motion Analysis - Finding reliable tracking regions.

**Remark 4.3 - Edge Detection Strategy**

In order to find edges we want to determine the *rate of change* in a pixel's neighbourhood

**Definition 4.2 - Image Gradient**

*Image Gradient* is a vector which gives the direction of greatest change in intensity at a specific pixel. For a pixel at  $(x, y)$  with  $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  mapping to the intensity the *Image Gradient* is

$$\nabla f(x, y) := \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix}$$

N.B. Generally denoted by  $\Psi$ .

**Definition 4.3 - Angle of Gradient**

$$\Psi := \tan^{-1} \left( \frac{\partial f / \partial y}{\partial f / \partial x} \right)$$

**Definition 4.4 - Edge Direction**

*Edge Direction* for a specific pixel is the direction of an edge, it is perpendicular/orthogonal to the *Image Gradient*.

N.B. Generally denoted by  $\Phi$ .

$$\Phi := \Psi - \frac{\pi}{2} = \tan^{-1} \left( \frac{\partial f / \partial y}{\partial f / \partial x} \right) - \frac{\pi}{2}$$

**Definition 4.5 - Image Gradient Magnitude**

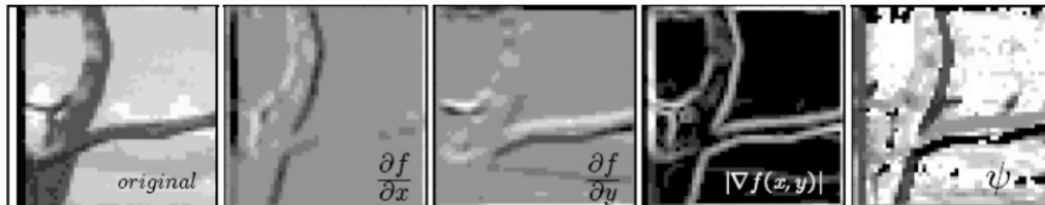
*Magnitude* measures the magnitude of the growth.

$$|\nabla f(x, y)| := \sqrt{\left( \frac{\partial f}{\partial x} \right)^2 + \left( \frac{\partial f}{\partial y} \right)^2}$$

**Proposition 4.1 - Estimating  $\nabla f$**

We can define  $\frac{\partial f}{\partial x}$  and  $\frac{\partial f}{\partial y}$  to be matrices, rather than typical derivatives, in order to produce estimates for the *Image Gradient* which can be used to analyse an image by convolution. We can then combine the results in order to analyse *Image Gradient Angle* and *Magnitude*. Consider

$$\frac{\partial f}{\partial x} = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix} \text{ and } \frac{\partial f}{\partial y} = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$



**Definition 4.6 - Prewitt Operator**

*Prewitt Operators* are a set of  $3 \times 3$  kernels used for estimating the *Image Gradient*. There are 8 *Prewitt Operators*, each in a different direction



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

**Remark 4.4** - *Alternative estimates of  $\nabla f$*

Instead of weighting all pixels in a given direction equally, we may want to weight central ones more heavily. This is since the central pixels are geometrically closer to the pixel which is being acted upon & such should better describe it.

$$\frac{\partial f}{\partial x} = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix} \quad \frac{\partial f}{\partial y} = \begin{pmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{pmatrix}$$

*N.B.* This is called the *Sobel Filter*.

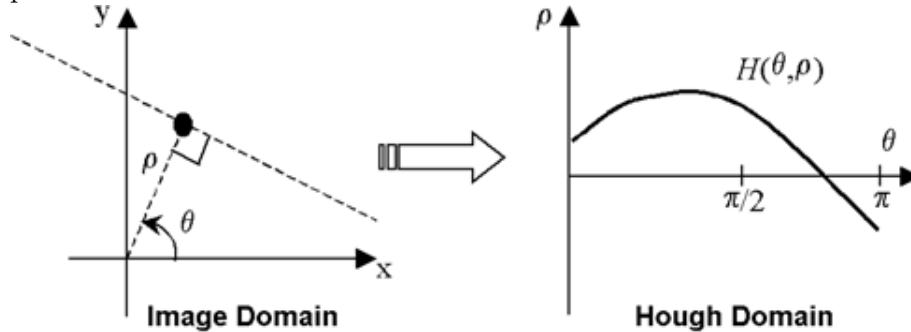
## 4.2 Shapes

**Definition 4.7** - *Hough Transform*

A *Hough Transform* is used to find lines which best explain the set of edge points given to it.

*Hough Transforms* are centred around the premise that lines can be described by  $\rho = x \cos \theta + y \sin \theta$  where either  $(x, y)$  is fixed, or  $(\rho, \theta)$  is fixed.

Suppose we are given an edge point  $(x_0, y_0)$ , we are now dealing with the case where  $(\rho, \theta)$  are variable. Plot all the combinations of  $(\rho, \theta)$  than produce lines which pass through  $(x_0, y_0)$  produces a sinusoidal wave.



**Proposition 4.2** - *Line Detection Algorithm*

Below is an algorithm which uses the *Hough Transform* on a set of edge points in order to detect lines in an image

- i) Make available an  $n = 2$  dimensional array,  $H(\rho, \theta)$ , to be used for the parameter space;
- ii) Find the *Gradient Image*,  $G(x, y)$ ;
- iii) For any pixel  $(x_0, y_0)$  where  $|G(x_0, y_0)| > T_s$  (some threshold value) increment all values of  $H(\rho, \theta)$  where  $(\rho, \theta)$  satisfy  $\rho = x_0 \cos \theta_0 + y_0 \sin \theta_0$ .

$$\forall \rho, \theta \text{ where } \rho = x_0 \cos \theta_0 + y_0 \sin \theta_0 \text{ do } H(\rho, \theta) + 1$$

- iv) Any  $(\rho, \theta)$  where  $H(\rho, \theta) \geq T_h$  (another threshold value) represent a straight line which has been detected in the image.

**Proposition 4.3** - *Circle Detection Algorithm*

- i) Make available an  $n = 3$  dimensional array,  $H(x, y, r)$ , to be used for the parameter space;

- ii) Find the *Gradient Image*,  $G(x, y)$ ;
- iii) For any pixel  $(x_0, y_0)$  where  $|G(x_0, y_0)| > T_s$  increment all  $H(x, y, r)$  which satisfy

$$\forall r \text{ where } x_0 = x + r \cos \Phi \text{ and } y_0 = y + r \sin \Phi$$

- iv) Any  $(x, y, r)$  where  $H(x, y, r) > T_h$  represents a circle with radius  $r$  and centre  $(x_0, y_0)$

**Definition 4.8 - General Hough Transform**

A *General Hough Transform* is used to describe general shapes.

- i) Find the *Gradient Image*,  $G(x, y)$ .
- ii) Define a set of points  $\phi_i := \frac{p_i}{i}$  for  $i \in [1, k]$  and create a table using  $\phi_1, \dots, \phi_k$  as indexes.
- iii) Define a reference point  $(x_c, y_c)$  to act as a *centre of mass*.
- iv) For any given edge point  $(x_0, y_0)$  find

$$r = \sqrt{(x_0 - x_c)^2 + (y_0 - y_c)^2}, \beta = \tan^{-1} \left( \frac{y_0 - y_c}{x_0 - x_c} \right) \text{ and } \Phi_{(x_0, y_0)} \text{ in } G(x, y)$$

- v) Round  $\Phi_{(x_0, y_0)}$  to the nearest  $\phi_i$  and insert  $(r, \beta)$  into the table at  $\phi_i$

**Proposition 4.4 - General Shape Detection Algorithm**

- i) Prepare 2 dimensional array  $H(x_c, y_c)$  for the parameter space (Set of all possible centre of masses).
- ii)  $\forall (x_0, y_0)$  where  $|G(x_0, y_0)| > T_s$  find the table entry,  $\phi_i$ , closest to  $\Phi_{(x_0, y_0)}$ .
- iii)  $\forall (r_j, \beta_j)$  in the table entry find

$$x_c = x + r \cos \beta \text{ and } y_c = y + r \sin \beta$$

Increment  $H(x_c, y_c)$ .

- iv) All  $(x_c, y_c)$  where  $H(x_c, y_c) > T_h$  represents the locations of centre of masses for occurrences of the shape, in the image.

**Remark 4.5 - Issues with Proposition 4.4**

The algorithm defined in **Proposition 4.4** only detects shapes of the same *scale* and *orientation* as the original image. To counter this we introduce two new variables to the parameter space:  $S$ , a scaling factor; and,  $\theta$ , an orientation factor.

- i) Prepare 4 dimensional array  $H(x_c, y_c, S, \theta)$  for the parameter space (Set of all possible centre of masses, scales & orientations).
- ii)  $\forall (x_0, y_0)$  where  $|G(x_0, y_0)| > T_s$  find the table entry,  $\phi_i$ , closest to  $\Phi_{(x_0, y_0)}$ .
- iii)  $\forall (r_j, \beta_j)$  in the table entry find

$$x_c = x + rS \cos(\beta + \theta) \text{ and } y_c = y + rS \sin(\beta + \theta)$$

Increment  $H(x_c, y_c, S, \theta)$ .

- iv) All  $(x_c, y_c)$  where  $H(x_c, y_c, S, \theta) > T_h$  represents the locations of centre of masses, scaling factors & orientations for occurrences of the shape in the image.

## 5 Image Segmentation

### Definition 5.1 - Image Segmentation

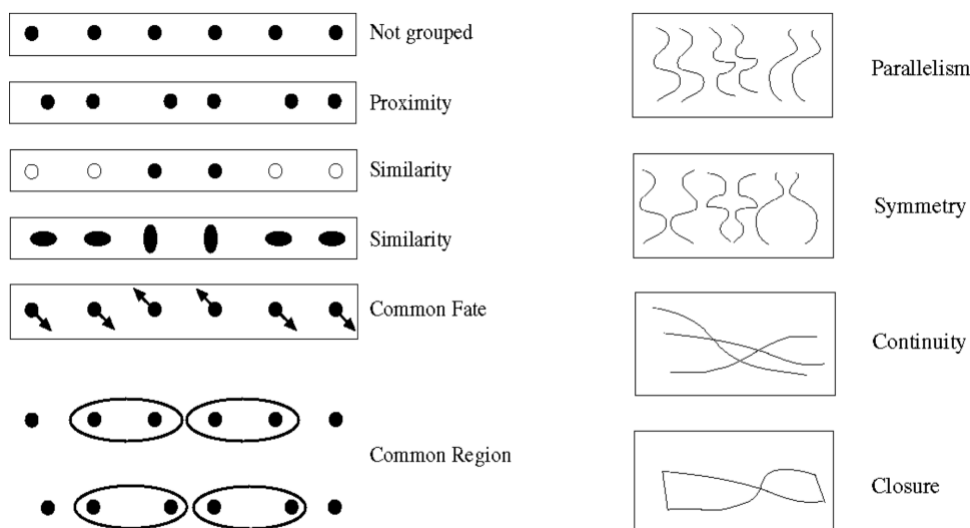
*Image Segmentation* is the process of partitioning the pixels in an image into homogenous regions. This may wrt colour, texture, object.

### Remark 5.1 - Usefulness of Image Segmentation

*Image Segmentation* is useful since it simplifies an image from millions of pixels into a few regions making computation much easier.

### Definition 5.2 - Gestalt Rules for Grouping & Segmentation

It is useful to consider the different ways in which we can group objects.



### Remark 5.2 - Perfect Image Segmentation is hard

It is hard to achieve perfect *Image Segmentation* since pixels may straddle the boundary of objects & thus not truly belong to any object. Noise, non-uniform illumination, occlusions etc. will also cause problems.

### Definition 5.3 - Over-Segmentation

*Over-Segmentation* is when an image is segmented into too many regions & loses the efficiencies gained by the simplification of *Image Segmentation*.

### Example 5.1 - Over-Segmentation



### Definition 5.4 - Under-Segmentation

*Under-Segmentation* is when an image is segmented into too few regions & too much information is lost.

### Example 5.2 - Under-Segmentation



**Proposition 5.1 - Segmentation Technique Types**

There are several groups of techniques that can be used for *Image Segmentation*

- i) Thresholding. Categorise pixels wrt intensity.
- ii) Edge-Based. Region boundaries are produced using an edge map.
- iii) Region-Based. Regions are grown from seed pixels or using split-merge techniques.
- iv) Clustering & Statistical. Global partitioning, often based around histograms. (*e.g.* K-means).
- v) Topographic (OUT OF SCOPE). Stepwise simplifications that take spatially wider image configurations into account.

**Proposition 5.2 - Thresholding Image Segmentation**

*Thresholding* is a good technique when trying to distinguish dark objects(background) from bright objects (foreground). One possible algorithm is

- i) Choose a threshold value  $T$ .
- ii) For each pixel:
  - (a) If the brightness is less than  $T$  map to 0 (black).
  - (b) else, map to 255 (white).

**Remark 5.3 - Choosing Threshold Value**

Choosing a value for *Thresholding* is important since if it is too high then the background pixels are classified as the foreground. And visa-versa for too low.

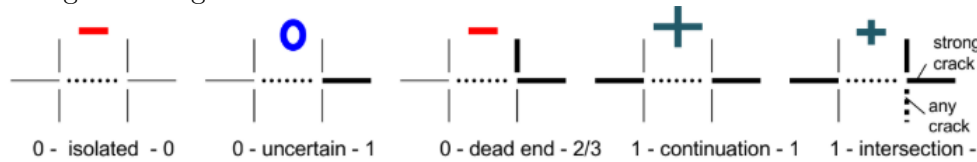
**Proposition 5.3 - Choosing Threshold Value**

There are several techniques used for choosing *Threshold Value*. By consider the histogram of an image's intensity we can produce regions for classification.

- i) Make an initial estimate for threshold  $T$ .
- ii) Segment the image using  $T$ : Let  $G_1$  be the pixels with intensity  $> T$ ; and  $G_2$  be the pixels with intensity  $< T$ .
- iii) Compute the average intensity values  $m_1$  &  $m_2$  for  $G_1$  &  $G_2$  respectively.
- iv) Compute a new threshold value  $T = \frac{1}{2}(m_1 + m_2)$ .
- v) Repeat **ii)-iv)** until convergence.

**Proposition 5.4 - Edge-Based Segmentation**

When performing *Edge-Based Image Segmentation* we want to detect strong edges that are relevant to the object we wish to detect. There are several techniques for assessing the strength of an edge, but generally we wish to remove edges that are isolated or irrelevant to an object. One valid technique is edge relaxation where we change the strength of an edge depending upon its 6 neighbour edges.

**Proposition 5.5 - Region Growing-Seed Pixel**

A function is required to analyse the homogeneity of two pixels  $(s, p) \rightarrow \{0, 1\}$ .

- i) Choose an initial seed pixel;
- ii) Consider the neighbouring pixels:
  - (a) If it is close to seed pixel, add to region & to queue to be analysed.
  - (b) Else, do nothing with it.
- iii) Repeat ii) until there are no pixels left in the queue to be analysed.

*N.B.* You may wish to add a final step to remove small regions.

**Proposition 5.6 - Split & Merge Segmentation**

Let  $H$  be a function to analyse the homogeneity of a region  $(R) \rightarrow \{0, 1\}$

- i) Let  $R_0$  represent the entire image.
- ii) If  $H(R_i) = 0$  (*i.e.* is inhomogenous) then split it into four regions.
- iii) Repeat ii) until all regions are homogenous.
- iv) Merge all subregions that satisfy  $H(R_i \cup R_j) = 1$  (*i.e.* are homogeneous)

**Proposition 5.7 - Clustering Segmentation**

If we map the RGB values of an image into the 3D real space we can use clustering algorithms to produce  $k$  clusters. We then map back to the pixel space using each cluster as a segment.

*N.B.* Look at k-means clustering.

## 6 Classical Object Detection

**Remark 6.1 - Classical Object Detection**

In this course we only cover *Classical Object Detection*. This technique came to prominence with the 2001 paper *First Real-Time Detection Method: Viola & Jones' - 2001* and was the industry standard for 10+ years. Now there are methods which use deep learning techniques.

**Definition 6.1 - Object Detection**

*Object Detection* bridges the semantic gap between pixel values & meaningful objects by grouping pixels and classifying. This is a difficult problem.

**Remark 6.2 - Failings of Shape Detection for Object Detection**

Typically using classical *Shape Detection & Segmentation* does not work for real world applications of *Object Detection* due

- i) High Intra-Class & Low Inter-Class variance;
- ii) Classes are rarely well defined;
- iii) Variations in illumination, scale, pose, deformation, occlusion etc.

**Remark 6.3 - Shift-Scale Invariance**

We require *Object Detection* algorithms to be shift & scale invariant so it can detect objects of any size at any point in the image.

**Definition 6.2 - Sliding Window Detectors**

*Sliding Window Detectors* ensure shift & scale invariance by setting a minimum size to check (say  $20 \times 20$ px) and shifts the region it checks by 1px right-to-left, top-to-bottom until the whole is checked. They then increase the size of the region it checks (say 10%) and repeats the shifting. It does this until the size of the region it checks exceeds the size of the image.

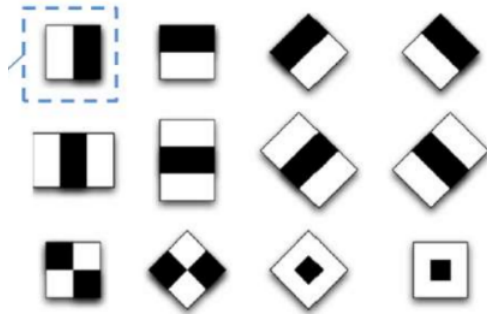
*N.B.* This technique does not work in real-time due to its high computation complexity.

**Definition 6.3 - Haar-Like Features**

*Haar-Like Features* are features of digital images used in *Object Detection*. They define a set of adjacent rectangles and assigns them to either be *light* or *dark* regions. It then returns a value, the average gray value of the light region less the average gray value of the dark region. This value is then compared to a threshold value to determine whether the feature occurs in the image.

**Proposition 6.1 - Haar-Like Feature Types**

There are a finite set of *Haar-Like Features Types*



**Definition 6.4 - Image Integral**

The *Integral* of an image  $I$  is the image  $J$  where  $J(x, y) = \sum_{i=0}^x \sum_{j=0}^y I(i, j)$ . *i.e.* Each pixel value is the sum of all pixels between it and the origin, inclusive.

**Proposition 6.2 - Calculating Image Integral**

Let  $I$  be an image then we can calculate the *Integral* of  $J$  in *Linear Time* using the following formulae

$$\begin{aligned}
 A(x, -1) &= 0 \\
 A(x, y) &= A(x, y - 1) + I(x, y) \\
 J(-1, y) &= 0 \\
 J(x, y) &= J(x - 1, y) + A(x, y)
 \end{aligned}$$

**Theorem 6.1 - Integration Rule of Convolution**

**Proposition 6.3 - Calculating the Average Pixel Value of a Large Block**

Calculating the average pixel value of a large block can be very slow if the region is very large. Here is a technique that can do it by only querying 4 values.

- i) Create a matrix which is the same size as the image.

- ii) Encode the matrix st its values are 1 in the region we want to check & 0 everywhere else.
- iii) Double Differentiate this matrix to produce  $I''$

(a) First, horizontally by applying  $\begin{pmatrix} 0 & -1 & 1 \end{pmatrix}$  to every pixel.

(b) Then, vertically by applying  $\begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}$  to every pixel.

This leaves us with four non-zeros values (two 1s and two  $-1$ s).<sup>7</sup>

- iv) Calculate the image integral  $J$ .

v) Return  $\sum_i \sum_j I''(i, j) \times J(i, j)$

*N.B.* This can be simplified to just be the four non-zero terms in  $I''$ .

#### Proposition 6.4 - Training Data

Two categories of *Training Data* should be provided: *Positive Samples*, images of the object you wish to classify' and, *Negative Samples*, images which do not contain the object. These should be labeled  $(x_i, y_i)$  with  $y_i = 0, 1$  for *Negative* and *Positive Samples* respectively.

#### Definition 6.5 - Adaboost

*Adaboost* (Short for *Adaptive Boosting*) is a machine learning algorithm.

Let  $(x_1, y_1), \dots, (x_n, y_n)$  with  $y_i = 0, 1$  for negative & positive samples respectively

- i) Initialise weights  $w_{1,i} = \begin{cases} \frac{1}{2^m}, & \text{if } y_i = 0 \\ \frac{1}{2^l}, & \text{if } y_i = 1 \end{cases}$  where  $m$  is the number of negative samples &  $l$  is the number of positive samples.

- ii) For  $t = 1, \dots, T$  where  $T$  is the number of features you wish to find

- (a) Normalise the weights st  $w_t$  is a probability distribution

$$w_{t,i} := \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

- (b) For each feature,  $j$ , train a classifier  $h_j$  which is restricted to using a single feature (Haar-Like Feature). Calculate the error wrt  $w_t$

$$e_j = \sum_i w_i |h_j(x_i) - y_i|$$

- (c) Choose the classifier,  $h_t$ , with the lowest error  $e_t$ .

- (d) Update the weights

$$w_{t+1,i} = w_{t,i} \left( \frac{e_t}{1 - e_t} \right)^{1 - e_i}$$

where  $e_i = 0$  if sample  $x_i$  is classified correctly and  $e_i = 1$  otherwise

- iii) The final strong classifier is

$$h(x) = \begin{cases} 1 & \sum_{i=1}^T \log \left( \frac{e_t}{1 - e_t} \right) h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \log \left( \frac{1 - e_t}{e_t} \right) \\ 0 & \text{otherwise} \end{cases}$$

## 7 Motion

### Definition 7.1 - Video Sequences

A video can be considered as a series of matrices, each representing an image at a different point in time. The time difference between consecutive matrices will be constant,  $\Delta t$  seconds.

Define  $I(x, y, t)$  to be the intensity of the pixel at position  $(x, y)$  at time  $t$ .

### Definition 7.2 - Frame Rate

*Frame Rate* is the number of images per second in a video. Typically measured in FPS (*Frames Per Second*), this is equivalent to Hz.

$$\text{Frame Rate} = \frac{1}{\Delta T}$$

*N.B.* Typical frame rates are  $24fps$ ,  $30fps$  &  $60fps$ .

### Definition 7.3 - 2D Tracking

*2D Tracking* is tracking how an object moves in a two dimensional plane. The object may really be moving a three dimensional plane, but we only care about its movement in a specified 2D plane.

This means the object has 3 attributes: x position, y position & orientation. Scale can be considered too.

### Example 7.1 - 2D Tracking

In the below image the path the pen has taken in the plane parallel to the camera has been tracked. The pen is actually moving in 3-dimensions, but we are ignoring depth in this example.



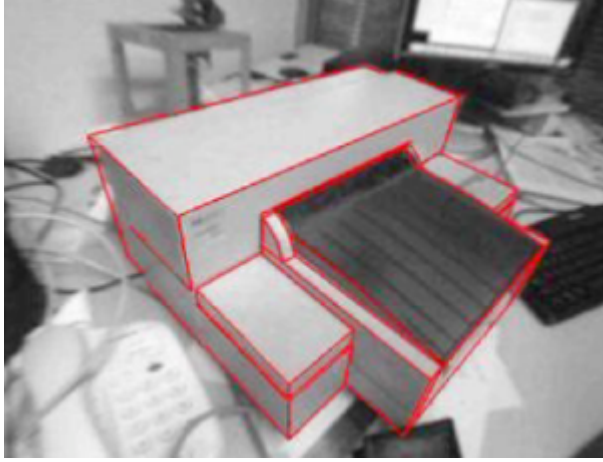
### Definition 7.4 - 3D Tracking

*3D Tracking* is used to track an object's movement in a three dimensional space. This means an object has six factors: x,y,z positions and rotation in the x,y,z directions. Scale can be considered too.

### Example 7.2 - 3D Tracking

In the below image a printer is tracker so that a skeleton of it can be mapped to it.




**Definition 7.5 - Simultaneous Localisation and Mapping**

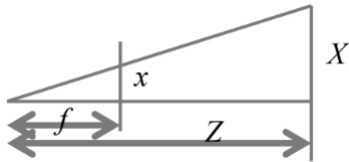
*Simultaneous Localisation and Mapping*, SLAM, is the process of constructing and updating a map of an unknown environment while tracking the camera's position within it.

*N.B.* This can be made easier by using a depth sensor.

**Definition 7.6 - Perspective Projection Equations**

Suppose a 3D point  $\mathbf{P} = (X, Y, Z)$ , which sits on the surface of an object, is projected to a point in two dimensions,  $\mathbf{p} = (x, y, f)$  ( $f$  is the distance of the 2D plane from the camera). Then, by similar triangles

$$x = \frac{fX}{Z} \text{ and } y = \frac{fY}{Z}$$


**Definition 7.7 - Rotation Matrices**

A *Rotation Matrix* can be applied to a position, by matrix-multiplication, to perform a rotation on that position. A *Rotation Matrix*,  $R$  can be broken down into *Rotation Matrices* which represent the rotation about each axis. In three-dimensions this means  $R = R_X R_Y R_Z$ .

*N.B.* *Rotation Matrices* are  $n \times n$  in  $n$ -dimensional space.

**Remark 7.1 - 3D Rotation Matrices**

$$R_X(\theta) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta \\ 0 & \sin \theta & \cos \theta \end{pmatrix} \quad R_Y(\theta) = \begin{pmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{pmatrix} \quad R_Z(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

**Proposition 7.1 - Rotation Matrices for small  $\theta$** 

When  $\theta$  is small we can approximate  $\sin \theta \approx \theta$  and  $\cos \theta \approx 1$ . Thus

$$R_X(\theta) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -\theta \\ 0 & \theta & 1 \end{pmatrix} \quad R_Y(\theta) = \begin{pmatrix} 1 & 0 & \theta \\ 0 & 1 & 0 \\ -\theta & 0 & 1 \end{pmatrix} \quad R_Z(\theta) = \begin{pmatrix} 1 & -\theta & 0 \\ \theta & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

**Definition 7.8 - Rigid Object**

A *Rigid Object* is one which does not change in shape, only scale & rotation.

**Proposition 7.2 - 3D Rigid Motion**

Consider a *Rigid Object* & a point on that object  $\mathbf{P}$ .

If the object is moved st point  $\mathbf{P}$  is now at position  $\mathbf{P}'$  then we can describe  $\mathbf{P}'$  in terms of  $\mathbf{P}$  by defining a translation matrix  $T$  and a rotation matrix  $R$ .

$$\mathbf{P}' = R\mathbf{P} + T$$

N.B.  $\mathbf{P}$ ,  $\mathbf{P}'$  &  $T$  are  $3 \times 1$  matrices and  $R$  is a  $3 \times 3$  matrix.

### Proposition 7.3 - 3D Motion Field

Consider a *Rigid Object* & a point on that object  $\mathbf{P}$ .

Let  $V$  be the vector which represents the translation from  $\mathbf{P}$  to  $\mathbf{P}'$ . We have

$$V = \lim_{\Delta t \rightarrow 0} [\mathbf{P}' - \mathbf{P}] = \lim_{\Delta t \rightarrow 0} [(R - I)\mathbf{P} + T]$$

We note that for small  $\theta_X, \theta_Y$  &  $\theta_Z$

$$R = \begin{pmatrix} 1 & -\theta_Z & \theta_Y \\ \theta_Z & 1 & -\theta_X \\ -\theta_Y & \theta_X & 1 \end{pmatrix}$$

Thus

$$\begin{aligned} \frac{dX}{dt} &= V_X = \theta_Y Z - \theta_Z Y + T_X \\ \frac{dY}{dt} &= V_Y = \theta_Z X - \theta_X Z + T_Z \\ \frac{dZ}{dt} &= V_Z = \theta_X Y - \theta_Y X + T_Z \end{aligned}$$

N.B. Here  $(\theta_X, \theta_Y, \theta_Z)$  is called the *Angular Velocity* and  $(T_X, T_Y, T_Z)$  is called the *Rectilinear Velocity* since we are taking the limit as  $\Delta t$  tends to 0.

### Proposition 7.4 - 2D Motion Field Equations

Consider a point  $\mathbf{p} = (x, y, f)$  which represents a 2D projection of a point  $\mathbf{P} = (X, Y, Z)$ .

Then

$$\begin{aligned} v_x &= \frac{d}{dt} x \\ &= \frac{d}{dt} \left( \frac{fX}{Z} \right) \\ &= f \frac{V_X Z - X V_Z}{Z^2} \text{ by quotient rule} \end{aligned}$$

Substituting  $V_X, V_Y, V_Z$  for the definitions in **Proposition 7.3** we get

$$v_x = \frac{fT_X - xT_Z}{Z} + f\theta_Y - \theta_Z y - \frac{\theta_X xy - \theta_Y x^2}{f} \text{ and } v_y = \frac{fT_Y - yT_Z}{Z} - f\theta_X + \theta_Z x + \frac{\theta_Y xy - \theta_X y^2}{f}$$

We can recognise two components of these equations: Translation, terms involving  $T$ ; and Rotation, terms involving  $\theta$ .

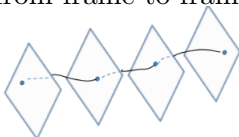
## 7.1 Motion Estimation

### Definition 7.9 - Motion Estimation

*Motion Estimation* is the process of using spatial & temporal variation in pixel values to detet amount of movement in an image region.

### Definition 7.10 - Optical Flow

*Optical Flow* is the movement of pixel values between frames. We assume that pixels follow a trajectory through frames. *i.e.* There is a sequence  $(x(t), y(t))$  that shows how a pixel moves from frame to frame.



*N.B.* Also known as *Apparent Motion*.

**Remark 7.2 - Optical Flow & True Motion**

In some cases the *Optical Flow* in an image is not equivalent to the *True Motion* of the object. Consider a diagonal rod moving across the screen, it is generally very hard to tell if it is moving horizontally or vertically or any other direction.

*N.B.* This is sometimes the result of the *Aperture Problem*.

**Definition 7.11 - Optical Flow Equation**

Let  $I(x, y, t)$  be the grayscale value of pixel  $(x, y)$  at time  $t$ . If we wish to follow a particular pixel between frames we find that  $I(x, y, t)$  is fixed &  $(x, y)$  depend on  $t$ . Further

$$\frac{d}{dt}I(x, y, t) = 0$$

By the chain rule we have

$$\frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = 0 \implies I_x v_x + I_y v_y + I_t = 0$$

*N.B.* Gradient =  $(I_x, I_y, I_t)$  and Motion =  $(v_x, v_y)$ .

**Remark 7.3 - Estimating the Optical Flow Equation**

The *Optical Flow Equation* has two unknowns  $v_x$  &  $v_y$ . This means we cannot solve it for a single pixel, as it would be underconstrained. Hence we have to *estimate* its value by using multiple pixels.

Often a  $2 \times 2$  square of pixels is used to estimate the *Normal Flow* of the central point.

*N.B.* The number of pixels used to *estimate* depends on the complexity of  $v_x$  &  $v_y$ . If  $v_x$  &  $v_y$  are linear in  $x$  &  $y$  with 3 parameters (e.g.  $v_x = ax + by + c$ ) then four pixels is sufficient.

**Proposition 7.5 - Constant Velocity Model**

Suppose we have a region  $R$  of the image (often taken as the whole image) we can assume it has constant velocity  $(v_x, v_y)$  accross it. Then our problem is to find  $v_x, v_y$  which minimise

$$\varepsilon(v_x, v_y) = \sum_R (I_x v_x + I_y v_y + I_t)^2$$

*N.B.* This is a minimisation problem since, ideally,  $I_x v_x + I_y v_y + I_t = 0$ .

**Theorem 7.1 - Lucas-Kanade Algorithm**

The *Lucas-Kanade Algorithm* aims to solve the problem in **Proposition 7.5**.

Since it is a minimisation problem we take the derivatives and find

$$\frac{\partial \varepsilon}{\partial v_x} = 2 \sum_R I_x (I_x v_x + I_y v_y + I_t) \quad \text{and} \quad \frac{\partial \varepsilon}{\partial v_y} = 2 \sum_R I_y (I_x v_x + I_y v_y + I_t)$$

Equating these to zero and solving we find

$$\begin{pmatrix} v_x \\ v_y \end{pmatrix} = A^{-1}b \text{ where } A := \sum_R \begin{pmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{pmatrix} \text{ and } b = - \sum_R \begin{pmatrix} I_t I_x \\ I_t I_y \end{pmatrix}$$

**Remark 7.4 - Motion for Object Detection**

*Motion Estimates* can be a good cue for object detection (and segmentation) as it implies depth in the image. However it requires an object to have an approximately constant velocity otherwise it will be partitioned (consider trying to detect the body of a person who is waving their

arm).

**Definition 7.12 - Layered Representation**

If the camera is moving with a constant motion we can infer depth from the image by identifying motion-homogeneous regions. We can use this information to produce a *Layered Representation* of the image, where each layer is a *depth-homogeneous* region.

**Definition 7.13 - Layered Video**

Every frame of an image can be split into three maps

1. Colour Map - The colour of each pixel in the layer.
2. Alpha Map - The pixels which are in the layer.
3. Velocity Map - The *Motion Estimate* for pixels in the layer.

**Proposition 7.6 - Motion Segmentation Process**

1. Perform *Block Motion Fitting*.
  - (a) Split the image into blocks of equal area.
  - (b) Estimate the motion in each block.
2. Detect the dominant motions using clustering algorithms and returning the centroids of each cluster.
3. Perform *Motion Segmentation* by assigning pixels to the dominant motion which is closest to theirs
4. Produce Final Parameteric Motion Fitting by updating our motion estimates using our assigned pixel motions.

**Proposition 7.7 - Region Linking & Alignment**

Once we have used the process in **Proposition 7.6** we can pick a particular region & align it across every frame. By taking the median values of every pixel in this aligned form we can produce the maps mentioned in **Definition 7.13**

## 8 Stereo Vision

**Definition 8.1 - Stereo Vision**

*Stereo Vision* is the process of producing a 3D representation of an enviroment using images taken from different viewpoints. The variation in position of objects between images is used to infer depth.

$$-\frac{1}{\text{depth}} \propto \text{disparity}$$

This shows that if we know disparit & viewpoint positions we can construct a 3D scene. *N.B.* Disparity is the difference in position between two images.

**Proposition 8.1 - Problems of Stereo Vision**

1. **Calibration** - determine the relative position & orientation of the cameras.
2. **Correspondence** - determine matching points in the stereo views.

3. **Reconstruction** - determine the 3D location of matched points in the image, via triangulation.

**Remark 8.1** - *Proposition 8.1*

- The three problems in **Proposition 8.1** are all inter-related as the information from one is used in the other.
- Which problem is hardest depends on the scenario. Generally reconstruction is the easiest as it relies fully on the results of the first two problems, but can be hard when it has a sparse depth matrix.

**Definition 8.2** - *Simple Stereo Case*

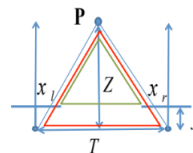
The simplest case where *Stereo Vision* can be implemented is where we have just two cameras which are translations of each other, with no rotation.

Here we can refer to the cameras as left & right.

**Definition 8.3** - *Triangulation - Simple Stereo*

In the *Simple Stereo Case* we can infer depth by considering similar triangles.

Let  $p$  be a similar point;  $x_l$  &  $x_r$  be the  $x$ -coordinate of  $p$  in the left & right image respectively;  $T$  be the translation matrix between the two cameras (left to right); and,  $f$  be the focal length of the cameras (assume equal for left & right). Then

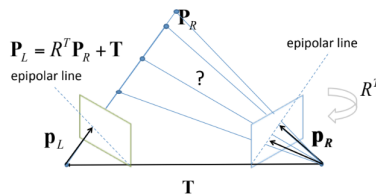


$$\frac{T}{Z} = \frac{T + x_r - x_l}{Z - f} \implies Z = \frac{fT}{x_l - x_r}$$

**Definition 8.4** - *Epipolar Geometry*

*Epipolar Geometry* defines the relationship between two stereo views. Consider a scenario where we have just two cameras which are related by a transformation  $T$  & rotation  $R^T$  (i.e.  $P_L = R^T P_R + T$ ).

Let  $p_l$  be the position of a point in the image of the left camera. We define two *Epipolar Lines* related



1. For  $p_l$  we have an epipolar line which is the line which runs from the left camera's focal point, through  $p_l$  and beyond. This is the set of all possible true positions of  $P$  in the 3D space.
2. For  $p_r$  we have an epipolar line which represents the set of all possible positions of  $P$  on the right camera's image.

Further we define the *Epipolar Plane* to be the plane defined by these two *Epipolar Lines* & *Epipoles* to be the points at which the *Epipolar Lines* intersect  $T$ .

**Proposition 8.2** - *Epipolar Geometry*

Consider the rigid transformation between cameras

$$P_L = R^T P_R + T \iff P_R = R(P_L - T)$$

We can account for perspective projection with

$$P_L = \begin{pmatrix} X_L \\ Y_L \\ Z_L \end{pmatrix} \quad p_L = \begin{pmatrix} x_L \\ y_L \\ f \end{pmatrix} = \frac{f P_L}{Z_L} \quad P_R = \begin{pmatrix} x_R \\ y_R \\ f \end{pmatrix} = \frac{f P_R}{Z_R}$$

TODO - More maths. Not sure the context of it though (ie why we want  $=0$ )

**Proposition 8.3 - Derivation of the Epipolar Line - Right Image**

Consider the *Simple Stereo Case* where the focal points are separated by transformation  $T$  and rotation  $R$ .

Let  $P$  be the point we wish to triangulate.

Let  $P_L$  &  $P_R$  be the vectors to  $P$  from the focal points of the left & right cameras respectively.

We have that  $P_L$  &  $T$  both lie on the epipolar plane, thus  $P_L - T$  is in the epipolar plane and  $T \otimes P_L$  is perpendicular to the plane. Thus

$$(P_L - T)^T(T \otimes P_L) = 0 \text{ since they are perpendicular to each other}$$

We have that  $P_R = R(P_L - T)$ .

Note that  $R^T = R^{-1}$  giving

$$R^T P_R = (P_L - T) \implies P_R^T R = (P_L - T)^T$$

Further we can redefine

$$(T \otimes P_L) = S P_L \text{ where } S := \begin{pmatrix} 0 & -T_Z & T_Y \\ T_Z & 0 & -T_X \\ -T_Y & T_X & 0 \end{pmatrix}$$

Substituting these two results into the first equation we get

$$(P_L - T)(T \otimes P_L) = P_R^T R S P_L = 0$$

For simplicity, we define the *Essential Matrix*  $E := RS$ . Thus

$$P_R^T R S P_L = 0 \implies P_R^T E P_L = 0$$

Let  $p_L$  &  $p_R$  be the vectors from the focal point to the pixel which represents  $P$  in the left & right images respectfully. Then

$$p_L = \frac{f P_L}{Z_L} \text{ and } p_R = \frac{f P_R}{Z_R} \implies \frac{Z_R}{f} p_R^T E \frac{Z_L}{f} p_L = 0 \implies p_R^T E p_L = 0$$

where  $f$  is the focal length and  $Z_R$  &  $Z_L$  are the depth of  $P$  from each image.

Define  $\mathbf{u}_L = E p_L = \begin{pmatrix} u_{L1} \\ u_{L2} \\ u_{L3} \end{pmatrix}$ . Then

$$p_R^T E p_L = p_R^T \mathbf{u}_L = x_R u_{L1} + y_R u_{L2} f u_{L3} = 0$$

This is the equation for the *Epipolar Line* in the right image.

**Definition 8.5 - Fundamental Matrix**

Let  $(\bar{o}_x, \bar{o}_y)$  be the origin (often the centre) of the image &  $(\hat{x}, \hat{y})$  be the row-column coordinates of a pixel of interest.

Define  $x := s_x(\hat{x} - \bar{o}_x)$  and  $y := s_y(\hat{y} - \bar{o}_y)$  for some  $s_x, s_y \in \mathbb{R}$ .

Then

$$p_L = \begin{pmatrix} x_L \\ y_L \\ f \end{pmatrix} = M_L \begin{pmatrix} \hat{x}_L \\ \hat{y}_L \\ f \end{pmatrix} = M_L \hat{p}_L$$

Thus

$$p_R^T E p_L = \hat{p}_R^T M_R^T E M_L \hat{p}_L = 0$$

For simplicity we define the *Fundamental Matrix*  $F = M_R^T E M_L$ .

$$\hat{p}_R^T F \hat{p}_L = 0$$

**Proposition 8.4 - Fundamental Matrix from Correspondences**

Consider a set of *Correspondences*, i.e. a set of pairs of points  $(\hat{p}_R(i), \hat{p}_L(i))$  which represent the same point in the scene, but in the separate images.

We have

$$\hat{p}_R^T(i) F \hat{p}_L(i) = 0 \text{ for } i \in [1, N]$$

This produces a series of simultaneous equations which can be solved for  $F$ .

We can simplify the problem to

$$A\mathbf{f} = 0$$

where  $A$  is an  $N \times 9$  matrix defined by the correspondence vectors and  $\mathbf{f}$  is the components of  $F$ .

N.B.  $A\mathbf{f} = 0$  can be solved for  $\mathbf{f}$  using Singular Value decomposition.

**Proposition 8.5 - 3D Reconstruction**

When considering the two rays (from each image) on which a point  $P$  may lie it is very unlikely that the two lines will intersect thus we consider the point of closest passing on each of these rays.

We predict that  $P$  lies on the midpoint between these two points.

Let  $ap_L$  be the vector the passing point on the left image, then  $bR^T p_R + T$  is a vector the pass point on the ray of the right image wrt the left image.

To find the closest passing points we want to find  $a, b, c > 0$  st

$$ap_L - bR^T p_R - T - \underbrace{c(p_L \otimes R^T p_R)}_{\text{miforces perpendicular}} = 0$$

where  $ap_L$  is the line of possible points for  $P$  given it appears at pixel  $p_L$  in the left image,  $bR^T p_R + T$  is the equivalent for the right image &  $T$  is the translation from left to right image. Given corresponding points we know  $p_L$  &  $p_R$  and given calibrated views we know  $R$  &  $T$  making  $a, b, c$  the only unknowns in this equation.

Once solved we give our estimate for  $P$  as

$$\hat{P} = \frac{1}{2}(ap_L + bR^T p_R + T)$$

### 8.0.1 Correspondence Matching

**Definition 8.6 - Correspondence Matching**

*Correspondence Matching* is the process of finding pairs of features (in different images) which correspond to the same point in the scene.

**Remark 8.2 - Challenges for Correspondence Matching**

1. Noise and lighting effects.
2. Perspective transformations.
3. Occlusions.
4. Repetitive Structures.

**Proposition 8.6 - What to Match?**

When looking for corresponding pairs there are several features we can consider

1. Pixels - Dense but they are not very distinctive.
2. Regions - Dense and are more distinctive than pixels, but must be wary of depth variation in a single region.
3. Feature Points - Sparse but distinctive.
4. Edges - Sparse and not necessarily distinctive as it can be hard to tell what point on an edge you are dealing with.
5. Objects - Sparse and distinctive but hard to identify and match.

*N.B.* Dense means there are many occurrences of this feature in the image (sparse is the opposite).

**Proposition 8.7 - Where to Search?**

Once we have identified a feature in one image we need to consider where the corresponding point may be in the other image. Where we search depends on whether the system is calibrated or not

1. **Uncalibrated** - Corresponding point may lie anywhere, even off the image.
2. **Calibrated** - The fundamental matrix is known thus we can limit our search to the *epipolar line*.

**Proposition 8.8 - How to compare Elements?**

When trying to find a corresponding pair we need a metric to judge similarities between features/regions. Some variables we may consider are

1. Pixel intensities.
2. Local Colour Variation.
3. Local Structure Variation.
4. Position and Angle of Edges.
5. Shape and Pose

*N.B.* The later of these are harder to calculate, but are more robust.

**Proposition 8.9 - Region Matching**

Let  $I_L(\cdot, \cdot)$  &  $I_R(\cdot, \cdot)$  return the intensity at a given coordinate for a pair of stereo images.

Let  $W$  be a region, centred around  $(x, y)$ , which we want to find a corresponding region to in our second image.

$$\hat{\mathbf{d}} := \operatorname{argmin}_{(d_x, d_y)} \sum_{(i, j) \in W} \operatorname{sim}[I_L(x + i, y + j), I_R(x + i + d_x, y + j + d_y)]$$

$\hat{W}$  which is centred around  $(x + d_x, y + d_y)$  is the best corresponding match.

*N.B.* Sometimes we want to find  $\operatorname{argmax}$ , depending upon the similarity function used.

**Proposition 8.10 - Similarity Measures**

Here are a few possible *Similarity Measures*

1. Sum of squared differences

$$\operatorname{Sim}(u, v) = (u - v)^2$$



## 2. Similar Pixel Count

$$\text{Sim}(u, v) = \begin{cases} 1 & , |u - v| < T \\ 0 & \text{otherwise} \end{cases}$$

where  $T$  is a thresholding value.

## 3. Normalised Cross Correlation

$$\text{Sim}(u, v) = \frac{(u - \bar{u})(v - \bar{v})}{N_u N_v} \text{ where } N_x := \sqrt{\sum_{k, l \in W} (x - \bar{x})^2}$$

**Proposition 8.11 - Feature Based Methods**

In *Feature Based Methods for Correspondence Matching* we restrict our search to sparse sets of features.

We want to find salient points in each view and match the points.

We can match points by comparing pixels or image descriptors in local regions about each point. Some examples of *Feature Based Methods* include

1. Salient Points - Harris Corner Detector.
2. Uncalibrated Matching - RANSAC
3. Scale-Invariant Feature Transform

**Proposition 8.12 - Harris Corner Detection**

Let  $W$  be a region we want to compare and  $I_x, I_y$  be the image derivatives in the  $x$  and  $y$  directions respectively.

Calculate the covariance matrix of  $I_x$  &  $I_y$ .

$$A := \sum_{(x,y) \in W} \begin{pmatrix} I_x(x, y)^2 & I_x(x, y)I_y(x, y) \\ I_x(x, y)I_y(x, y) & I_y(x, y)^2 \end{pmatrix}$$

Calculate the eigenvalues,  $\lambda_1, \lambda_2$ , of  $A$ . These represent the spread of gradients in the region. Our similarity metric is

$$\text{Sal}(W) = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2}$$

Large values are better here as they represent a busier region.

**Proposition 8.13 - RANSAC**

Consider a system where we don't know the relationship in position between the two cameras (*i.e.* Uncalibrated).

The following steps describe *RANSAC*

1. Find potential correspondences (*e.g.* by comparing regions around salient points)
2. Select a subset, of at least 7, of these matches at random.
3. Compute the *Fundamental Matrix*,  $F$ , from this subset.
4. Assess how well  $F$  describes the relationship between the unused potential correspondences.
5. Repeat several times then pick the best  $F$ .

*N.B.* *RANSAC* stands for *RAN*dom *SAM*ple *C*onsensus.

**Proposition 8.14 - Scale-Invariant Feature Transform**

TODO

*N.B.* This is known as *SIFT*

## 0 Reference

### 0.1 Definitions

**Definition 0.1 - *Connectivity***

When analysing a pixel there are several ways to consider its adjacent neighbourhood.

- 4-Connectivity. The pixels above, below, left & right.
- 8-Connectivity. Above, below, left, right & single step diagonals.

**Definition 0.2 - *Kernel***

A *Kernel* is a small matrix used in convolution. Typically  $3 \times 3$  or  $5 \times 5$ . *Kernels* can be defined for blurring, sharpening, embossing, edge detection & more *N.B.* This definition only applies to image processing & is different from the definition in linear algebra.