

Thesis or Article

JeffGWood@mavs.uta.edu

July 11, 2016

Contents

1	Introduction	2
2	Background	4
2.1	Change of Reference	4
2.2	Points and Lines in the Image Plane	4
2.3	Epipolar Geometry	5
2.4	Fundamental Matrix	6
2.5	Intrinsic Calibration Matrix	6
2.6	Essential Matrix	6
3	Point Interpolation	7
4	Image segmentation	9
5	Process	10

Chapter 1

Introduction

Through the development of applications such as augmented and virtual reality, object / scene reconstruction and visual effects, the process of generating images from an arbitrary vantage point can be found in a variety of applications. In this Thesis (or Article) I will discuss various methods for Image Creation from an arbitrary vantage point, which can be accomplished by two main methodologies of Geometric Construction and Image Synthesis. While both methods use stereo correspondance of multiple images, they differ in the way information is stored and used.

Geometric Construction (GC) contains information about the real-world spatial properties (Coordinates in space, Color), thus viewing results are non-constrained in vantage point. Image Synthesis (IS) relies on image properties (pixel displacement) and is thus viewing results are imited in the possible vantage points.

Symbols and Notation

Symbol	Description
\mathbf{v}	Vectors in <i>lowercase</i> bold
v_a	a -component of vector \mathbf{v}
\mathbf{M}	Matrices in <i>uppercase</i> bold
$M_{r,c}$	Entry in row r and column c of matrix \mathbf{M}
\mathbf{m}_c	Vector occurring in column c of matrix \mathbf{M}
\mathbf{x}	Generic 3-dimensional spatial coordinate
$\tilde{\mathbf{x}}$	Generic 3-dimensional spatial coordinate (expressed <i>homogeneously</i>)
\mathbf{y}	Generic 2-dimensionals image coordinate
$\tilde{\mathbf{y}}$	Generic 2-dimensional image coordinate (expressed <i>homogeneously</i>)
\mathbf{u}	Pixelized 2-dimensional image coordinate
$\tilde{\mathbf{u}}$	Pixelized 2-dimensional image coordinate (expressed <i>homogeneously</i>)
${}^A\mathbf{x}$	Generic 3-dimensional spatial coordinate in reference frame A
${}^A\tilde{\mathbf{x}}$	Generic 3-dimensional spatial coordinate (expressed <i>homogeneously</i>) in reference frame A
${}_B^C\tilde{\mathbf{M}}$	Change from of reference frame B to reference frame C
s	Normalizing factor applied to <i>homogeneous</i> vector $\tilde{\mathbf{v}}$ such that original $\mathbf{v} = s \cdot \tilde{\mathbf{v}}$ is recovered
${}^D\mathbb{S}$	Spatial reference frame D
$[\mathbf{x}]_{\times}$	Skew-symmetric matrix version of vector \mathbf{x} used as <i>left</i> -operand in the <i>cross</i> -product such that $[\mathbf{x}]_{\times} \cdot \mathbf{y} = \mathbf{x} \times \mathbf{y}$
l	Epipolar line
\mathbb{P}	Ray (or <i>pencil</i>) of all possible vectors \mathbf{x} where $\mathbf{x} = s \cdot \tilde{\mathbf{x}}$ for some value of s

Chapter 2

Background

Ordinarily, real-world data contains 3-dimensions. Because standard images only include 2-dimensional data, information regarding depth is lost (i.e. it is often difficult to judge distance from a single image without visual cues). *Stereovision* attempts to resolved this by finding the same point in both *stereoscopic* images (known as a *corresponding point*), and recovering the depth information. An elementary example of this occurs in stereoscopic images with relatively low distance between cameras (i.e they are right next to each other). Objects that are *farther* away from the observer occur closer together in the stereo images, whereas objects *closer* to the camera appear appear farther apart in the stereo-images.

2.1 Change of Reference

Each view from a pair of stereo-images encompasses its own *frame of reference* (i.e. the directions of *forward* or *backward* are unique to image and may differ considerably depending on camera displacement). This requires expressing points from different frames of reference (traditionally referred to *left* and *right*) in a single reference frame. As such it is necessary to be able to express coordinates in a given reference frame in any other reference frame.

Coordinates given in ${}^A\mathbf{x}$ can be expressed in ${}^B\mathbf{x}$ by the geometric transformation:

$${}^B\mathbf{x} = {}^B\mathbf{R} \cdot {}^A\mathbf{x} + {}^B\mathbf{t}$$

or

$$\begin{aligned} {}^B\tilde{\mathbf{x}} &= \left[\begin{array}{c|c} {}^B\mathbf{R} & {}^B\mathbf{t} \\ \hline 0 & 1 \end{array} \right] \cdot {}^A\tilde{\mathbf{x}} \\ &= {}^B\mathbf{M} \cdot {}^A\tilde{\mathbf{x}} \end{aligned}$$

where ${}^B\mathbf{M}$ is also the geometric transformation necessary to transform ${}^B\mathbb{S}$ into ${}^A\mathbb{S}$.

Withough calculating any new quantities, rearranging allows us to express coordinates in ${}^B\mathbf{x}$ in the ${}^A\mathbf{x}$ reference frame as:

$${}^B\mathbf{R}^\top \cdot ({}^B\mathbf{x} - {}^B\mathbf{t}) = {}^A\mathbf{x}$$

and similarly transforms ${}^A\mathbb{S}$ into ${}^B\mathbb{S}$.

2.2 Points and Lines in the Image Plane

Points in *world-space* of \mathbb{R}^3 are converted to points in the *image-plane* of \mathbb{R}^2 by *homogenization*. This occurs when a *world-coordinate* of $\mathbf{x} = [x_1, x_2, x_3]^\top$ is mapped to a *homogeneous image coordinate* of

$\tilde{\mathbf{y}} = [y_1, y_2, 1]^\top = [x_1/x_3, x_2/x_3, x_3/x_3]^\top$ or a *non-homogeneous image coordinate* of $\mathbf{y} = [y_1, y_2]^\top = [x_1/x_3, x_2/x_3]^\top$. Points of the form $\tilde{\mathbf{y}} = [y_1, y_2, 0]^\top$ are special case of homogeneous point referred to as a *point at infinity*.

Lines in \mathbb{R}^2 can be represented in different contexts. The *vector offset* method calculates a line $\mathbf{s}(t)$ between points \mathbf{y}_1 and \mathbf{y}_2 as

$$\begin{aligned}\mathbf{s}(t) &= (1 - t) \cdot \mathbf{y}_1 + t \cdot \mathbf{y}_2 \\ &= \mathbf{y}_1 + t \cdot (\mathbf{y}_2 - \mathbf{y}_1)\end{aligned}$$

in which the line is parallel to the vector $\mathbf{y}_2 - \mathbf{y}_1$ and offset from the origin by the vector \mathbf{y}_1 . Lines are also represented by their coefficients as $\mathbf{l} = [a, b, c]^\top$ where

$$\begin{aligned}\mathbf{l}^\top \cdot \tilde{\mathbf{y}} &= \begin{bmatrix} a & b & c \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ 1 \end{bmatrix} \\ &= a \cdot y_1 + b \cdot y_2 + c \cdot 1 \\ &= 0\end{aligned}$$

This definition lets us say $\tilde{\mathbf{y}}$ is located on line \mathbf{l} *if and only if* $\mathbf{l}^\top \cdot \tilde{\mathbf{y}} = 0$. The line \mathbf{l} joining two *homogeneous image coordinates* $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$ is then calculated as the cross product of $\mathbf{l} = \tilde{\mathbf{y}}_1 \times \tilde{\mathbf{y}}_2$.

2.3 Epipolar Geometry

Each point of interest (also referred to as a *feature*) in a single image occurs in a 2-dimensional space at location $\tilde{\mathbf{y}} = [x, y, 1]^\top$. The same point in space when viewed from an image at a similar (though different) angle is referred to as a *corresponding point* with location of $\tilde{\mathbf{y}}' = [x', y', 1]^\top$. This set of infinitely many points form a 1-dimensional subspace (also known as a *pencil*) of the 3-dimensional world space.

The pencil, when viewed from an image at a different angled-position, appears as a line $\mathbf{l}' = [A', B', C']^\top$, known as the *epipolar line*. The fact that the corresponding point (in the *angled image*) of $\tilde{\mathbf{y}}' = [x', y', 1]^\top$ occurs on this epipolar line is referred to as the *epipolar constraint*. It is formalized, using the previously given *line-point equality* of $\mathbf{l}'^\top \cdot \tilde{\mathbf{y}}' = 0$ for the *angled image*. Similarly, the corresponding point of $\tilde{\mathbf{y}}' = [x', y', 1]^\top$ produces an epipolar line in the *original image* of $\mathbf{l} = [A, B, C]^\top$. The original point of $\tilde{\mathbf{y}} = [x, y, 1]^\top$ must lie located on this epipolar line as required by the epipolar constraint, resulting in the *line-point equality* of $\mathbf{l}^\top \cdot \tilde{\mathbf{y}} = 0$ for the *original image*.

When viewed in their respective images, each point ($\tilde{\mathbf{y}}$ and $\tilde{\mathbf{y}}'$) has a pencil that coincides with that point. Since the pencils act as *directional*-vectors in 3-dimensional space, there is a unique 2-dimensional plane which contain both of these vectors, known as the *epipolar plane*. It is the intersection of the epipolar plane with the *original image*-plane and the *angled image*-plane that results in the epipolar lines of \mathbf{l} and \mathbf{l}' , respectively. In fact, the *epipolar plane* (in each image's *coordinate systems*)² has the same vector form as its epipolar line. Specifically, in the *original image* reference frame $\mathbf{l} = \mathbf{P} = [A, B, C]^\top$, and in the *angled image* reference frame $\mathbf{l}' = \mathbf{P}' = [A', B', C']^\top$. This results from the fact that any *world-point* \mathbf{x} lying on the *epipolar plane* \mathbf{P} will result in a *homogeneous image-point* $\tilde{\mathbf{y}}$ that also lies on the plane \mathbf{P} . Specifically, when $\mathbf{x} = s \cdot \tilde{\mathbf{y}}$ for some non-zero value of s , then $\mathbf{P}^\top \cdot \mathbf{x} = 0$ implies $\mathbf{P}^\top \cdot \mathbf{x} = \mathbf{P}^\top \cdot (s \cdot \tilde{\mathbf{y}}) = 0$. Since $s \neq 0$, its true that $\mathbf{P}^\top \cdot \tilde{\mathbf{y}} = 0$.

In the majority of images, the sets of epipolar lines will converge at a point known as an *epipole*, denoted as \mathbf{e} in the *original image* \mathbf{e}' in the *angled image*.

¹A *change of reference* is implied between coordinates $\tilde{\mathbf{y}} = [x, y, 1]^\top$ and $\tilde{\mathbf{y}}' = [x', y', 1]^\top$. The majority of corresponding points do not occur at the same *image coordinates* between images (i.e $\tilde{\mathbf{y}} \neq \tilde{\mathbf{y}}'$). The only way a single *world coordinate* can yield different *image coordinates*, is if a *change of reference* occurs in *world space* each time the *image coordinates* are obtained by dividing by z_{world} .

²There is a single *epipolar plane* for each pair of corresponding points $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{y}}'$. However, the single plane can be parameterized infinitely many ways, depending on the *frame of reference*.

2.4 Fundamental Matrix

In stereo vision, points ($\tilde{\mathbf{x}}$) in one image I are related to the epipolar line (l') that contain the corresponding point ($\tilde{\mathbf{x}}'$) by the *Fundamental Matrix* (\mathbf{F}).

$$l' = \mathbf{F} \cdot \tilde{\mathbf{x}}$$

2.5 Intrinsic Calibration Matrix

A point \mathbf{x} in the *camera-coordinate system* of \mathbb{R}^3 is projected to the point $\tilde{\mathbf{y}}$ in \mathbb{R}^2 by means of the *pinhole camera model*. The set of all $\tilde{\mathbf{y}}$ are the result of *rays* passing through the *image plane* located at $z = f$, and converging at the *optical center* as shown in the figure below: The location of $\tilde{\mathbf{y}}$ is determined by utilizing

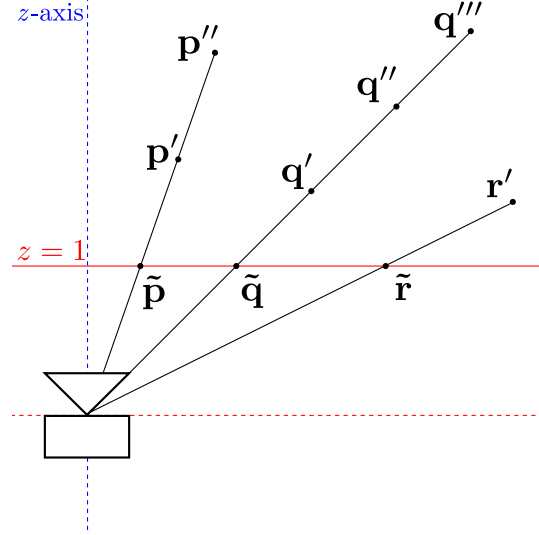


Figure 2.1: Pinhole Camera model

the *similarity of triangles* between \mathbf{x} and $\tilde{\mathbf{y}}$. Specifically, we see that $y_1/f = x_1/x_3$ and $y_2/f = x_2/x_3$ lets us express the *image coordinate* $\tilde{\mathbf{y}}$ as $y_1 = f \cdot x_1/x_3$ and $y_2 = f \cdot x_2/x_3$. The point in the *image plane* of $\tilde{\mathbf{y}}$ is derived from the point \mathbf{x} in *camera space* by means of the *Camera Projection Matrix* \mathbf{P} such that

$$\begin{aligned} \mathbf{P} \cdot \tilde{\mathbf{x}} &= \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ 1 \end{bmatrix} = \begin{bmatrix} f \cdot x_1 \\ f \cdot x_2 \\ x_3 \end{bmatrix} \\ &= x_3 \cdot \begin{bmatrix} f \cdot x_1/x_3 \\ f \cdot x_2/x_3 \\ 1 \end{bmatrix} = x_3 \cdot \tilde{\mathbf{y}} \end{aligned}$$

2.6 Essential Matrix

When coordinates from a reference frame are expressed as *normalized image coordinates* the range of possible NIC values in the corresponding image are given by the

Chapter 3

Point Interpolation

Pixels from image a and image b can be used to create a new images. This is done by interpolating the pixel positions (\mathbf{p}_{uv}^a and \mathbf{p}_{uv}^b) of corresponding points between frames. Because not all pixels are established as corresponding points, pixel correspondances *between* corresponding points (\mathbf{p}_{uv}) are calculated through bi-linear interpolation of 4 established corresponding points:

$$\mathbf{P}_{uv} = \mathbf{P}_{00} \cdot (1 - u) \cdot (1 - v) + \mathbf{P}_{10} \cdot u \cdot (1 - v) + \mathbf{P}_{01} \cdot (1 - u) \cdot v + \mathbf{P}_{11} \cdot u \cdot v$$

This is done through the following series of linear equations

$$\begin{aligned} x_{uv} &= \begin{bmatrix} u & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} v \\ 1 \end{bmatrix} \\ y_{uv} &= \begin{bmatrix} u & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} y_{00} & y_{01} \\ y_{10} & y_{11} \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} v \\ 1 \end{bmatrix} \end{aligned}$$

or as a single matrix equation of

$$\begin{bmatrix} x_{uv} & 0 \\ 0 & y_{uv} \end{bmatrix} = \begin{bmatrix} \mathbf{u} & \mathbf{0} \\ \mathbf{0} & \mathbf{u} \end{bmatrix}^T \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{M} \end{bmatrix}^T \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{Y} \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{M} \end{bmatrix} \begin{bmatrix} \mathbf{v} & \mathbf{0} \\ \mathbf{0} & \mathbf{v} \end{bmatrix}$$

where

$$\mathbf{u} = \begin{bmatrix} u \\ 1 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} v \\ 1 \end{bmatrix}, \mathbf{X} = \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \end{bmatrix}, \mathbf{Y} = \begin{bmatrix} y_{00} & y_{01} \\ y_{10} & y_{11} \end{bmatrix}, \text{ and } \mathbf{M} = \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix}$$

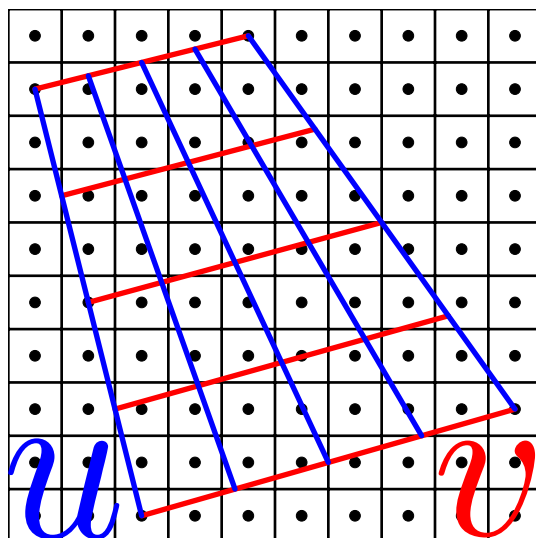


Figure 3.1: Bi-Linear Point Correspondance

Chapter 4

Image segmentation

Central to our need to localize corresponding points in stereo images is the ability to partition images by similar texture or planar attributes compared to the image at hand. Such techniques are referred to as *image segmentation*.

Image segmentation of regions of similar color or textures region is often approached from a graph-theory standpoint, in which individual pixels form the nodes of the graph. Edges are formed by a number of methods, the simplest of which is for each pixel to have 4 equally weighted edges connecting with the 4 immediate adjoining pixels in *North*, *East*, *South* and *West* vicinities (referred to as the **4-neighborhood region**). A common variation of this is to *also* include the next 4 closest adjoining pixels in the *Northeast*, *Southeast*, *Southwest* and *Northwest* vicinities (referred to as the **8-neighborhood region**). More sophisticated methods assign edge weightings proportional to the difference in color values (*scalar gray values* or *euclidean distance of color vectors*) between each pixel-pair.

Binary segmentation (partitioning into two regions) can be accomplished through min-cut / max-flow algorithms

Chapter 5

Process

The system in question contains 3 main components

1. Image Acquisition System

- Webcam / Kinect set-up
- If Webcam should also contain Image-Processing module for:
 - Feature Identification
 - Point-correspondance
 - Sub-Pixel interpolation

2. Point Cloud Processing

- Should take inputs
- Should produce point-clouds as one of the output
- (Possible) Options for Surface Reconstruction include:
 - Calculation of surface Normal through PCA
 - Mesh construction through Delaunay triangulation
 - Parametrization of Bezier surface through linear-least squares.