



TANSZÉKVEZETŐ

DIPLOMATERVEZÉSI FELADAT

Danyi Dávid

Villamosmérnök hallgató részére

Marker alapú helymeghatározás képfeldolgozással

Az információfeldolgozási kapacitás nagymértékű növekedésével párhuzamosan egyre szélesebb körben váltak alkalmazhatóvá (valós időben) képfeldolgozási, gépi látás alapú eljárások. Az ideálistól eltérően a valós képek feldolgozását számos hatás nehezíti, úgy mint zaj, optikai torzítások, megvilágítás- és színbeli különbségek, stb.

A feladat témája aktuális, az egyre gyakrabban alkalmazott mobil robotikai, kiterjesztett és virtuális valóság alkalmazásoknál minden esetben felmerül a „hol vagyok?” kérdés, azaz a nézponti paraméterek becslése, lokalizáció. A SLAM (Simultaneous Localization and Mapping) algoritmus számos szenzorforrás által biztosított információval alkalmazható, így kamerával is. Az általános megoldás tájékozódási pontok egymáshoz képesti elhelyezkedését becsli, a mérési bizonytalanságot folyamatosan csökkentve, valamint méri a tájékozódási pontokhoz képest a megfigyelő pozícióját.

Jelen munka is a fenti témához, kapcsolódik, cél megvizsgálni egy olyan mesterséges, passzív marker megvalósíthatóságát, ami bizonyos paraméterekben (azonosíthatósági tartomány, részleges láthatóság) jobbat kíván biztosítani, mint az elterjedt (ARtag, glyph, stb.) megoldások. A marker egy oldalukon nyitott különböző méretű és alakú négyszögeket tartalmaz. A diplomaterv célja megvizsgálni a javasolt marker jellemzőit és használhatóságát.

A hallgató feladatának a következőkre kell kiterjednie:

- Mutasson be pozícióbecslő algoritmusokat, röviden ismertesse ezek működését!
- Hasonlítsa össze különböző nézőpontbecslési algoritmust síkban elhelyezkedő pontpárok alapján!
- Készítsen egy marker felismerő megoldást!
- Végezzen méréseket (ideális és valós képeken) a marker által meghatározott pozíció pontosságára vonatkozóan!
- Hasonlítsa össze és értékelje az eredményeket!

Tanszéki konzulens: Kovács Viktor, tanársegéd

Budapest, 2017. február 18.

Dr. Charaf Hassan
egyetemi tanár
tanszékvezető





Budapest University of Technology and Economics
Faculty of Electrical Engineering and Informatics
Department of Automation and Applied Informatics

Marker Based Localisation and Pose Estimation Using Image Processing

THESIS

Author

Dávid Danyi

Consultant

Viktor Kovács

May 29, 2017

Contents

List of abbreviations	3
1 Markers	4
1.1 Quad	4
1.1.1 Quad representation	6
1.2 Marker	7
1.2.1 Marker generation	8
1.2.2 Discrete RQIM	9
2 Marker recognition	11
2.1 Preprocessing	12
2.1.1 Segmentation	13
2.1.2 Conditioning	14
2.2 Line fitting	14
2.2.1 Hough transformation	14
2.2.2 Gradient detector	17
2.2.3 Corner detector	17
2.2.4 Skeletoning detector	17
Bibliography	18
Appendix	11

HALLGATÓI NYILATKOZAT

Alulírott *Dávid Danyi*, szigorló hallgató kijelentem, hogy ezt a szakdolgozatot meg nem engedett segítség nélkül, saját magam készítettem, csak a megadott forrásokat (szakirodalom, eszközök stb.) használtam fel. Minden olyan részt, melyet szó szerint, vagy azonos értelemben, de átfogalmazva más forrásból átvettem, egyértelműen, a forrás megadásával megjelöltem.

Hozzájárulok, hogy a jelen munkám alapadatait (szerző(k), cím, angol és magyar nyelvű tartalmi kivonat, készítés éve, konzulens(ek) neve) a BME VIK nyilvánosan hozzáférhető elektronikus formában, a munka teljes szövegét pedig az egyetem belső hálózataán keresztül (vagy autentikált felhasználók számára) közzétegye. Kijelentem, hogy a benyújtott munka és annak elektronikus verziója megegyezik. Dékáni engedéllyel titkosított diplomatervek esetén a dolgozat szövege csak 3 év eltelte után válik hozzáférhetővé.

Budapest, May 29, 2017

Dávid Danyi
hallgató

Abbreviations

This is a complete list of the abbreviations used in this paper.

DOF Degrees of freedom

RQIM Random Quad Image Marker

SHT Simple Hough Transformation

RHT Randomised Hough Transformation

PPHT Progressive Probabilistic Hough Transform

Chapter 1

Markers

One of the goals of this project was to design a fiducial marker with advantageous properties for use in pose estimation. In a typical scenario the marker may be seen from largely varying viewpoints, therefore it has to have some level of scale invariability. If the observer is far from the marker, the smaller details may be lost due to the limited resolution of the camera. If the same observer moves closer to the marker, it may fill the whole field of view and some features may even slip off the image. This leads to another feature the marker needs to have: redundancy. If the observer gets too close to the marker or some obstacle partially blocks the view, the localisation still needs to provide usable results.

The intended use of the markers is spatial localisation and pose estimation. In other words: approximating the observers 3D coordinates (x, y, z) and orientation (ϕ, θ, ψ) with respect to the marker. It is supposed that the observer uses a single camera system for navigation (e.g. smartphone or robotic application with limited resources). This means the marker needs at least 6 degree of freedom.

To sum up the above discussed specifications, a suitable marker would have to:

- have at least 6 DOF
- be (to some degree) scale invariant
- have redundancy

In the following sections will be a recommendation for a marker conforming for the listed specifications. It is based on 3 connected line segments forming a quad with one missing side. The whole marker is built from quads with different side lengths and angles.

1.1 Quad

A marker is put together from quads. Figure 1.1. shows two examples. One side of the quads is left out: they are put together from three joint line segments. The middle segment, with

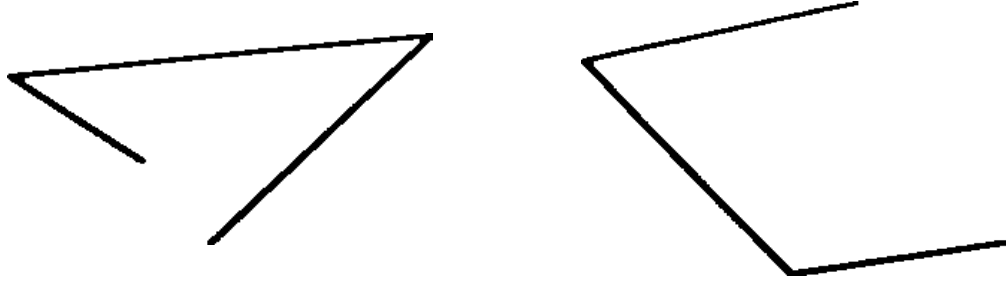


Figure 1.1: *Example for different quads*

two adjoining lines, will be referred to as the 'base' of the quad. The outer segments are going to be called 'arms'.

A quad has 6 degrees of freedom. There are 3 independent distance parameters: the length of the base and the two arm segments. There are also 3 unrelated angle parameters: the angles between each arm and the base, and the orientation of the quad.

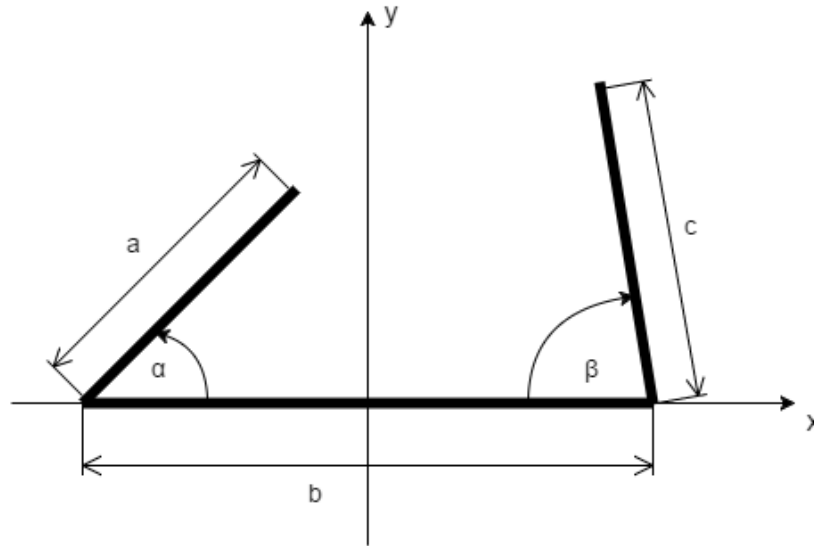


Figure 1.2: *Quad parameters*

Figure 1.2. shows the free parameters of a quad (the orientation is not shown on the image). The following notation is used:

- a : The length of one arm
- b : The length of the base
- c : The length of the other arm
- α : The angle between one arm and the base
- β : The angle between the other arm and the base
- γ : The angle with which the whole quad is rotated

For the sake of simplicity, figure 1.2. does not show the rotation with γ . The quad would be rotated around the origin of it's coordinate system.

The values of the length parameters are given in pixels, although they can be expressed in any unit of distance. The angles can be given in degrees or radians (in the implementation degrees are used for easier human readability).

$$a \in (0, a_{max}] \quad (1.1)$$

$$b \in (0, b_{max}] \quad (1.2)$$

$$c \in (0, c_{max}] \quad (1.3)$$

$$\alpha \in (0, 180^\circ) \quad (1.4)$$

$$\beta \in (0, 180^\circ) \quad (1.5)$$

$$\gamma \in [0, 360^\circ) \quad (1.6)$$

Equations (1.1) through (1.6) specify the range of each parameter. The maximum of the distance parameters are set by the space left on the image for the given marker, there is no theoretical limit for them. There is also no constraint for the resolution of the parameters. From the applications point of view, there are quads with continuous¹ and discrete parameter spaces.

1.1.1 Quad representation

There are several ways to represent quads, each with different advantageous properties. For this work multiple considerations were made in that regard. The most straightforward is to simply store the above mentioned parameters. This is simple and easy for human reading, which is great help in the development process.

A step forward from this is to norm the a, b and c parameter of the quad with the base segment's length. Then the following parameters are used:

s : marker size, the same as the base length

m_a : 'a' multiplier. $m_a = a/b$

m_c : 'c' multiplier. $m_c = c/b$

The α, β, γ angle parameters are not changed. This gives a scale or size parameter for the quad, which is useful for marker generation. These two representations are good for development and marker generation, but not so much for calculations.

¹That is, only limited by the computational precision

A third option is to store the endpoints of the line segments. It requires the storage of 4 points: two endpoints (E_1, E_2) and two inner points (I_1, I_2). Equations (1.7) through (1.10) define the points' coordinates before rotating with γ , using figure 1.2.'s notation.

$$I'_1 = (-\frac{b}{2}, 0) \quad (1.7)$$

$$I'_2 = (\frac{b}{2}, 0) \quad (1.8)$$

$$E'_1 = (-\frac{b}{2} + b * \cos(\alpha), a * \sin(\alpha)) \quad (1.9)$$

$$E'_2 = (\frac{b}{2} - c * \cos(\beta), c * \sin(\beta)) \quad (1.10)$$

$$Rot(\gamma) = \begin{pmatrix} \cos(\gamma) & -\sin(\gamma) \\ \sin(\gamma) & \cos(\gamma) \end{pmatrix} \quad (1.11)$$

The point E_1, E_2, I_1, I_2 can be obtained from E'_1, E'_2, I'_1, I'_2 by a multiplication with the rotational matrix $Rot(\gamma)$.

That method is redundant for storage: it uses 8 parameters instead of 6. However this poses no practical problem in the scope of the project. The one outstanding benefit of this method is it's efficiency in calculations. Because it is based on points in a Euclidean space, linear algebraic methods (matrix multiplications) can be used for calculating the projective transformations.

In this project the second and the third options are used. The first, naive method is omitted because it has no considerable advantage over the other two. The second method, using the size, multiplier and angle parameters is used in marker generation. The third is used during the calculations and the recognition process.

1.2 Marker

Quads are 6 DOF shapes: in theory it would be enough to use only one of them for localisation and pose estimation. However that method would have very low error tolerance and questionable accuracy even in a best case scenario. To comply with the specifications written in the beginning of this chapter, the markers are put together from multiple quads. By placing quads with different orientations and sizes the error tolerance and accuracy can be greatly improved.

An intrinsic positive quality of using multiple quads with varying sizes is the scale invariance. As mentioned, even a single quad is sufficient for the task at hand. If the smaller quads become unrecognisable because of the low resolution or too large distance, a successful measurement is still possible. The same is true on the other end of the spectrum: if the observer is too close to the marker and the larger ones leave the field of view, the position and orientation can be calculated from the smaller quads.

Figure 1.3. shows an example for a marker. It is generated with the simple algorithm described in the next section, and is not optimal in many ways. Nonetheless it is functional, even if only a fraction of the quads are registered for the measurement.

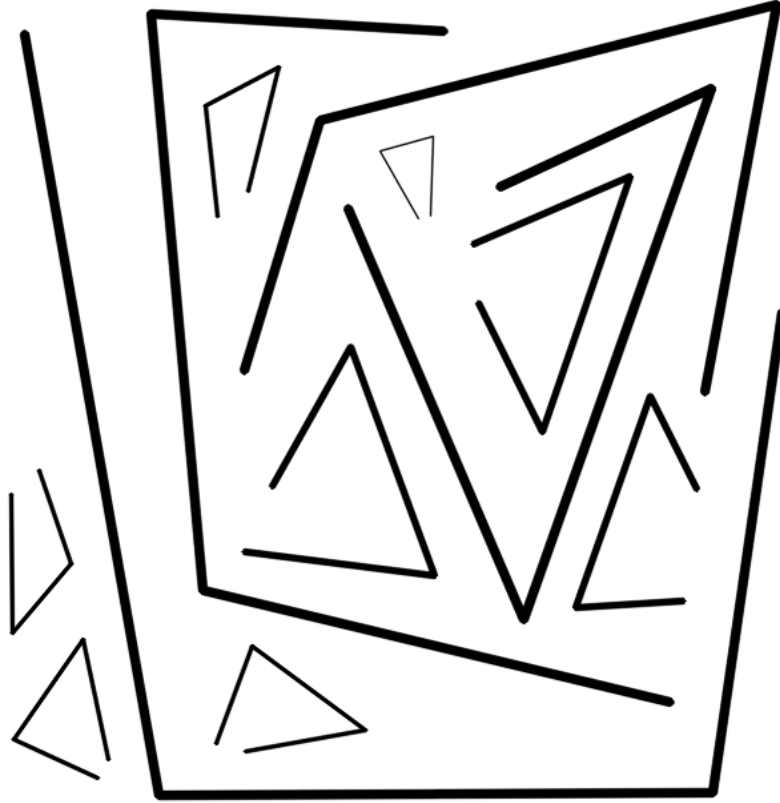


Figure 1.3: *An example for a marker*

The markers are going to be referred to as RQIM, which means Random Quad Image Marker. As the name suggests, the quads are randomly generated and placed on the markers.

Unless otherwise specified, RQIMs use quads with continuous parameter spaces. In this chapter there will also be a small introduction to discrete parameter markers and their potential applications.

1.2.1 Marker generation

In the current state of the project, markers are randomly generated using a simple algorithm. The generator routine receives the number of quads to be used in the current RQIM. The core concept is to create the desired number of random quads and place them on the image.

Let the number of quads to generate be n . First, the quad sizes are picked. There is an upper and a lower limit for them, given in percent of the image size. The generated sizes

follow an exponential distribution:

$$s = e^{-x*f} \quad (1.12)$$

Where s is the quad size, x is random number between 0 and 1 with uniform distribution, and f is a scale factor. Then the n sizes are ordered in descending order.

After the scale factors are picked, the whole quads are generated by the following method. A random quad is created with the first (the largest) scale and placed on the image. Then another quad is created with the next largest size. After every new quad a check is performed whether or not it can be placed on the marker. If it cannot, then a new quad is generated with the same scale factor until it can be placed or the algorithm reached the limit of retries.

With this simple logic n quads are placed on the RQIM and the creation process is finished. Below is the pseudo-code of the algorithm.

```
n_max = number of quads to create
f = scale factor for exponential distribution
lowlim = lower size limit
uplim = upper size limit
n = 0
while n < n_max
    size = exp(-rand() * f)
    if size > lowlim and size < uplim
        store size
        n = n+1
    endif
end
sort(sizes, descending)
n = 0
while n < n_max
    while quad placed or max tries
        quad = create_random_quad(sizes(n))
        if quad can be placed
            place quad
            n = n+1
        end
    end
end
return marker
```

This method is not optimal and is based on trial and error, but it gives usable markers for the development process.

1.2.2 Discrete RQIM

There are experiments in progress with discrete parameter space quads. It may be advantageous to quantize the parameter space in order to decrease the error probability in the pose estimation process.

Quads with finite possible states can be stored using much less resources than their continuous counterpart. As an example let us take a look at the following quantisation.

Angles: $15^\circ, 30^\circ, 45^\circ, 60^\circ, 75^\circ, 90^\circ, 105^\circ, 120^\circ$

Multipliers: 0.40, 0.60, 0.80, 1.0, 1.25, 1.50, 1.75, 2.0

Orientations: $0^\circ, 22.5^\circ, 45^\circ, 67.5^\circ \dots 270^\circ, 292.5^\circ, 315^\circ, 337.5^\circ$

Sizes: 1, 0.8, 0.6, 0.5, 0.4, 0.3, 0.25, 0.2, 0.1, 0.08, 0.06, 0.05, 0.04, 0.03, 0.025, 0.02

In this example, there are 8 possible values for the angle parameters, also 8 for the multipliers, 16 for the orientation and also 16 for the sizes. If the possibilities are stored in a lookup table, it is enough for the quad to store the index at which the value is accessible. A quad is defined by two angle parameters, two multipliers, an orientation and a size. The angles (in this case) require at least 3 bits each, the multipliers also. The orientation and the size need 4 bits each. That gives a sum of 20 bits per quad, which is significantly less than the space required to store 6 floating point numbers per quad.

This 20 bit word is also usable as an ID for the quad. It may be possible to code information in these ID-s, so the marker could provide additional information. This information could be related to and used by the localisation process, or be totally unrelated, general data. These possibilities have not yet been extensively researched.

The discrete RQIMs are usually less dense than the continuous ones, due to the limited angle possibilities. This means fewer quads per marker, which leads to decreasing redundancy. An optimum must be found between the number of quads per RQIM and distance between quads in the parameter space.

Chapter 2

Marker recognition

In this chapter there will be a summary of the image processing algorithms tried and used for the recognition of the fiducial markers. From a computer vision point of view the task is to detect joint line segments. This is a relatively easy task in image processing, there are many well tried algorithms for it.

In this chapter will be a short summary of the algorithms used for testing and performance comparison. The preprocessing steps used for preparing the images for the line fitters (segmentation, thresholding, filtering etc.) will also be discussed.

The general flow of processing is the same for every line fitting solution. The input is the raw image taken¹ by the observer. The first problem is finding the RQIM on the picture. When the marker area is located, it is necessary to discard the only partially visible and/or unrecognisable quads. At this point there is an image or set of images containing potentially good quads. The process here diverges depending on which line fitting algorithm is used. They all need differently conditioned input images for optimal performance. The line fitter routines not necessarily have the same output format², so conversion may be needed. This is the end of the marker recognition phase. This step of the process takes the raw input image and initiates quad structures based on the observed picture.

Four separate line fitters are profiled in this experiment.

- Hough-transformation
- Corner detection
- Skeletoning detector
- Gradient detector (enhanced Hough transformation)

¹In the development phase there were rendered pictures used for better repeatability

²Some return line segments defined by their endpoint, others use the polar representation of a line etc.

The first one uses the classic Hough-transformation for line detection. In the OpenCV framework there is another, probabilistic implementation of the transformation. It will also be tested.

The second detector is based on corner recognition. There are more variants of this method to try out, too. The corner metrics of a feature can be calculated differently with (Harris metric, eigenvalues, etc.) varying results. It is also needed for the solution to be scale invariant, which also can be achieved in a number of ways.

The third alternative uses skeletoning for line detection. It is a RANSAC-like method, which makes it very robust, though quite computationally expensive. The core concept is to thin the observed marker (band) down to single pixel lines and try to find two point with the most inliers.

The gradient detector is a bit Hough-like in it's nature. It uses the image gradient to calculate the angle in the Hough-space, and the pixels only vote in their distance parameter.

A typical marker shot with partial visibility is shown figure 2.1..

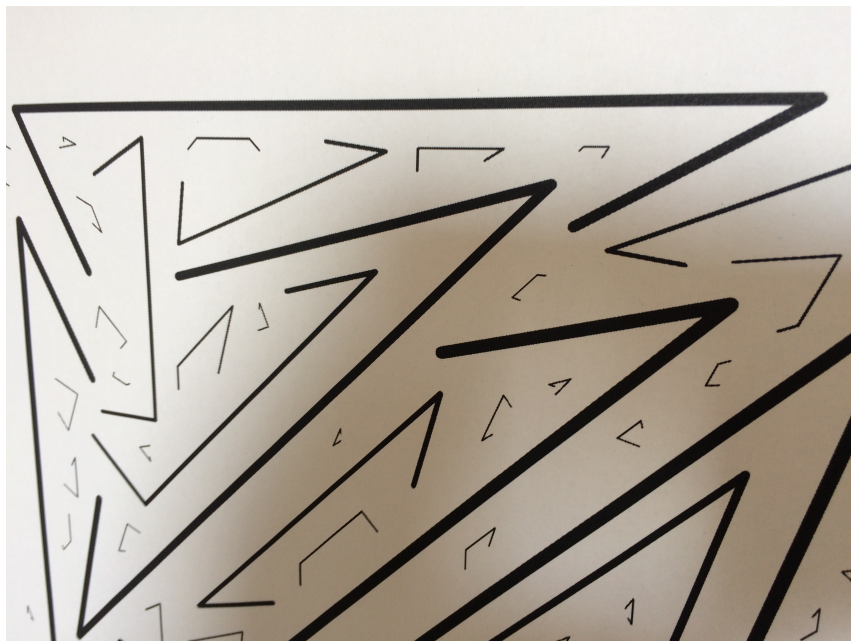


Figure 2.1: *Partially visible marker (taken with commercial smartphone)*

2.1 Preprocessing

The preprocessing is done in two stages. The first is the segmentation, when quad-like blobs are found. The second step is the preparation of the aforementioned blobs for the line fitting algorithms' needs.

2.1.1 Segmentation

The segmentation process is carried out on images roughly like the one shown in figure 2.1.. First the photos are converted to binary format by applying a threshold. The image is inverted in the process, because it makes more sense for the objects to be marked with non-zero elements than vice-versa. The threshold's value is determined using Otsu's method, which maximises the inter-class variance of the clusters³. The implementation is provided by the OpenCV framework.

Afterwards, the binary image is conditioned with a *close* morphology operator. The closing removes the gaps from the large connected areas (possible quads) and removes the *salt and pepper*-like noise. In the current implementation the kernel size of the morphology operator is constant, however it could be beneficial to calculate it from the global or local image parameters⁴.

The segmentation is based on finding continuous contours on the binary image. The OpenCV framework provides great functionality for this. The implementation is based on calculating the 8-neighbour chain code for the binary blobs on the image. The functions returns a list of list of points for the borders os each distinct contour.

The next step is the filtering of the found blobs. First the surely partial quads are discarded. This is done by calculating the bounding box of the contours, and if one of it's sides are touching the image border, the blob is marked as partial. With this approach it is possible that some fully visible quads that only touch the image border with one of their corner are lost. This problem can be easily fixed by checking the neighbourhood of the contact point, but this is not yet implemented.



Figure 2.2: *Quad candidates after segmentation*

The next blob-filtering step is to filter out the false-positive contours. These false hits can be caused by the light conditions or the scene around the marker. For this purpose a simple metric is used to measure how likely a blob is to be a quad. This metric is the ratio of a blob's area and circumference. By experimentation this ratio for quads is found to be in the range of 10 and 50. The contours with ratios outside these limits are discarded.

³Foreground and background

⁴e.g. image size, area of the connected region, etc.

The segmentation processes output is available as a single image with colour-coded⁵ blobs or as a list of separate images each containing a quad candidate. Figure 2.2. shows an example for the output of the segmentation process.

2.1.2 Conditioning

Before running the line detection algorithms some conditioning steps are done in order to improve their effectiveness. These processes are not uniform - every fitting method needs it's own.

The Hough transformation traditionally works best on thin lines. The easiest way is to generate an edge image with high-pass filtering. The OpenCV framework offers a wide variety of features for this task. The best results are obtained by using the Canny edge detector.

The skeletoning detector does not need a conditioning step, as it performs the band thinning on it's own.

The methods based on image gradients and corner detection both require some level of smoothing on the picture. In case of the corner detection the smoothing is useful for removing the false positive matches caused by the jagged edges, or in case of JPEG images the artefacts caused by the compression. The gradient detector simply gives a more spread-out and easier to analyse result on a smoothly changing gradient than it would on a strict edge. The smoothing is also implemented using the OpenCV framework, which provides easy access to Gaussian filtering. The OpenCV implementation is based on convolution with a configurable Gaussian kernel. The kernel size and the deviation in x and y direction can be set. The Gaussian kernel and the convolution itself is handled in the framework.

2.2 Line fitting

2.2.1 Hough transformation

One of the most commonly used methods for line detection on images is the Hough transform. Over its long history many publications have been made about its applications, performance and improvements.

Originally it was developed by Paul Hough in 1959 and later patented in 1962[3]. It was intended to be used for machine analysis of bubble chamber photographs. In its modern form (with the $\theta - \rho$ parametrisation) was introduced in 1972 by Duda and Hart[2]. The transformation became popular in the image processing community after Ballard's article[1] about generalising the algorithm for detection of arbitrary shapes. There were many optimised and improved variants of the transformation, however the basic concept

⁵Gray level, to be exact

remained the same. In 1990 a publication[4] introduced the Randomized Hough Transform, which was a fundamentally new approach to the algorithm with notable merits. As opposed to the one-to-many mapping of the simple Hough transform, the randomised version uses a convergent many-to-one mapping when creating the parameter space.

After this short historical overview the theory and the application of the transformation will be discussed.

Theory

The transformation is used to find instances of a model on digital images. The models are usually simple geometric shapes like lines, circles or ellipses. The curves are described by their parameters, e.g. slope and intercept for a line, centre point and radius for a circle etc.. Every non-zero pixel⁶ votes for the features it could be part of. The number of votes is stored for every possible parameter combination. Then a threshold is applied to the stored votes, and the remaining parameters are accepted as model instances.

At first Hough described the algorithm to lines, but later the method would be generalised to any analytic⁷ curve or shape. This theoretical overview is based on the example of line detection. The process is the same for every analytic curve, the only difference is the parameter space's number of dimensions. The original patent[3] used the slope-intercept representation of lines.

$$y = m * x + b \quad (2.1)$$

In this case, the *parameter space* is 2 dimensional and it's axes are m and b . Every point in the parameter space represent an image space line. With this representation every non-zero pixel in the image space transforms into a line in the parameter space. For a given (x_0, y_0) pair (2.2) gives the line in the parameter space.

$$b = -x_0 * m + y_0 \quad (2.2)$$

Collinear points in the image show up in the parameter space as intersecting lines. The more lines intersect in a given (m_0, b_0) , the more likely it is the image contains the $y = m_0 * x + b_0$ line. The problem with this parametrisation is that the parameter space is unbounded along both axes. Both intersect and slope can have values in the range of $(-\infty, \infty)$. Duda and Hart[2] proposed an alternative parametrisation, which turned out to be better for application. They used the *normal parametrisation* of a line, shown in (2.3).

$$\rho = x * \cos(\theta) + y * \sin(\theta) \quad (2.3)$$

In (2.3) ρ means the distance of the line from the image plane's origin. θ is angle of the normal vector of the line. If the *normal parametrisation* is used the parameter space

⁶The transformation works on binary images

⁷The Generalised Hough Transform even extends to arbitrary shapes

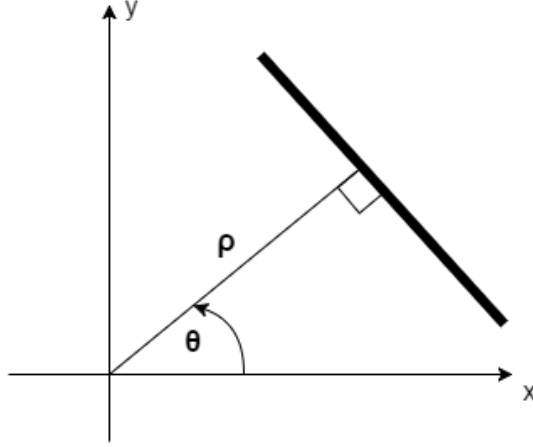


Figure 2.3: *Normal line parameters*

becomes finite in both dimensions. θ is in the range of $(0, 2\pi)$, ρ is bounded by the image size. In this case the image points define sinusoid curves in the parameter plane, and the line detection is done by searching for their intersections.

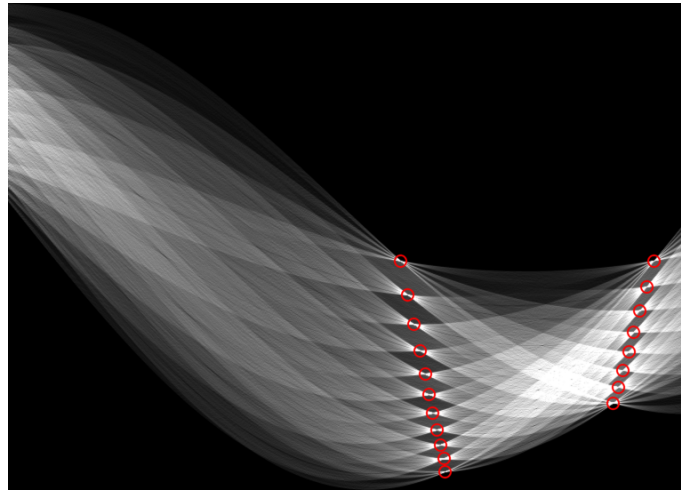


Figure 2.4: *Hough-transform of a chessboard pattern*

As mentioned before, the line detection is based on a voting scheme. The parameter space (in this case a 2 dimensional plane) is divided into *bins*. ρ and θ are quantised in the desired resolution. The discrete (ρ, θ) pairs define the bins. Every bin has accumulator. When a given (ρ_i, θ_i) pair gets a vote it's corresponding accumulator is incremented by 1. The **SHT** (Standard Hough Transform) uses one-to-many divergent mapping. This means that every non-zero pixel votes for every possible parameter pair it could belong to. The above mentioned sinusoid is calculated with the desired resolution for the pixel, and the corresponding accumulators are updated.

When the accumulation phase is completed for the whole image, the local maxima of the accumulators are found. Usually a threshold is applied in order to reduce noise and eliminate too short line segments. The radius of the non-maxima suppression also has

impact on the results of the line fitting, it must be chosen carefully. After this step the parameters for the most likely line candidates are available.

As the **SHT** does not provide the endpoints of the line, they must be found by examining the original binary image. This can be done by simple checking every pixel along the line with the given parameters and deciding whether or not it is part of the feature. If it is desired, lines with gaps can also be accepted with this method. For more accurate fitting, a Least Squares approximation can also be done on the pixels belonging to the line.

Application

2.2.2 Gradient detector

2.2.3 Corner detector

2.2.4 Skeletoning detector

Bibliography

- [1] D.H. Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern Recognition*, 13(2):111 – 122, 1981.
- [2] R. O. Duda and P. E. Hart. Use of the hough transformation to detect lines and curves in pictures,. *Comm. ACM*, pages 11–15, January 1972.
- [3] P.V.C. Hough. Method and means for recognizing complex patterns. *U.S. Patent 3,069,654*, Dec. 1962.
- [4] Lei Xu, Erkki Oja, and Pekka Kultanen. A new curve detection method: Randomized hough transform (rht). *Pattern Recognition Letters*, 11(5):331 – 338, 1990.