

# Pattern Recognition Systems – Lab 2

## RANSAC – fitting a line to a set of points

### 1. Objectives

The purpose of this laboratory session is to use the RANSAC method for line fitting.

### 2. Theoretical Background

Random Sample Consensus (RANSAC) is a paradigm for fitting a model to experimental data, introduced by Martin A. Fischler and Robert C. Bolles in 1981.

As stated by Fischler and Bolles [1] "The RANSAC procedure is opposite to that of conventional smoothing techniques: Rather than using as much of the data as possible to obtain an initial solution and then attempting to eliminate the invalid data points, RANSAC uses as small an initial data set as feasible and enlarges this set with consistent data when possible".

The RANSAC algorithm is summarized below [2]:

**Objective:** Robust fit of a model to a data set  $S$  which contains outliers.

**Algorithm:**

1. Randomly select a sample of containing a number of  $s$  data points from  $S$  and instantiate the model from this subset.
2. Determine the set of data points  $S_i$  which is within a distance threshold  $t$  of the model. The set  $S_i$  is the consensus set of the sample and defines the inliers of  $S$ .
3. If the size of  $S_i$  (the number of inliers) is greater than some threshold  $T$ , re-estimate the model using all the points in  $S_i$  and terminate.
4. If the size of  $S_i$  is less than  $T$ , select a new subset and repeat the above.
5. After  $N$  trials the largest consensus set  $S_i$  is selected, and the model is re-estimated using all the points in the subset  $S_i$ .

## 2.1. RANSAC for fitting a line to a set of points

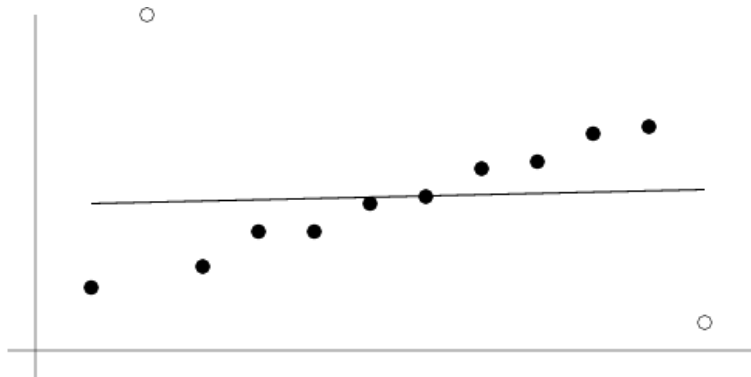


Figure 1-a

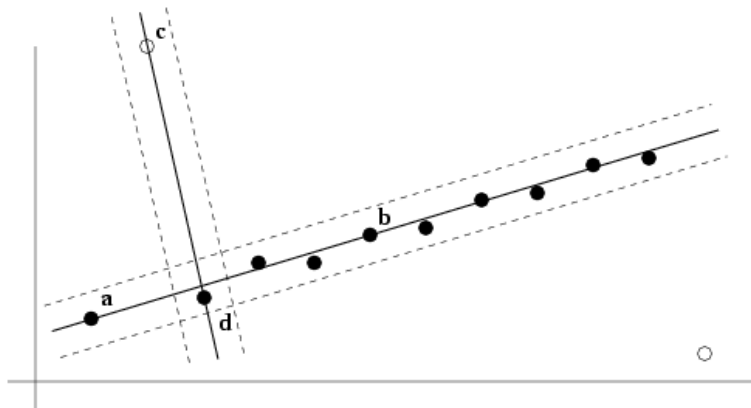


Figure 1-b

The problem, illustrated in Figure 1-a is the following: given a set of 2D data points, find the line which minimizes the sum of squared perpendicular distances (orthogonal regression), subject to the condition that none of the valid points deviates from this line by more than  $t$  units. This is actually two problems: a line fit to the data; and a classification of the data into inliers (valid points) and outliers. The threshold  $t$  is set according to the measurement noise (for example  $t = 3\sigma$ ), and is discussed below.

The idea is the following: two of the points are selected randomly; these points define a line. The *support* for this line is measured by the number of points that lie within a distance threshold. This random selection is repeated a number of times and the line with most support is deemed the robust fit. The points within the threshold distance are the inliers (and constitute the eponymous *consensus* set). The intuition is that if one of the points is an outlier then the line will not gain much support.

Furthermore, scoring a line by its support has the additional advantage of favoring better fits. For example, the line (a, b) in figure 1-b has a support of 10, whereas the line (a, d), where the sample points are neighbors, has a support of only 4. Consequently, even though both samples contain no outliers, the line (a, b) will be selected.

We can discuss three main issues that result from the presented algorithm:

1. **What is the distance threshold?** We would like to choose the distance threshold,  $t$ , such that with a probability  $a$  the point is an inlier. This calculation requires the probability distribution for the distance of an inlier from the model. In practice the distance threshold is usually chosen empirically. However, if it is assumed that the measurement error is Gaussian with zero mean and standard deviation  $\sigma$ , then a value for  $t$  may be computed.
2. **How many trials?** It is often computationally infeasible and unnecessary to try every possible sample. Instead the number of samples  $N$  is chosen sufficiently high to ensure with a probability,  $p$ , that at least one of the random samples of  $s$  points is free from outliers. Usually  $p$  is chosen at 0.99. Suppose  $q$  is the probability that any selected data point is an inlier, and thus  $q^s$  is the probability that all  $s$  points are inliers. The complementary event is that there is at least an outlier among the  $s$  points with probability  $1 - q^s$ . Then the probability for  $N$  selections to each have at least 1 outlier is  $(1 - q^s)^N$  which must be equal to  $1 - p$ . We can express  $N$  as  $= \log(1 - p) / \log(1 - q^s)$ .
3. **How large is an acceptable consensus set?** A rule of thumb is to terminate if the size of the consensus set is similar to the number of inliers believed to be in the data set, given the assumed proportion of outliers, i.e. for  $n$  data points  $T = qn$ . For the line-fitting example of Figure 1, a conservative estimate of  $q = 0.8$ , so that  $T = 0.8 \cdot 12 = 9.6$ .

### 3. Mathematical background

The equation of a line through two distinct points  $(x_1, y_1)$  and  $(x_2, y_2)$  is given by:

$$(y_1 - y_2)X + (x_2 - x_1)Y + x_1y_2 - x_2y_1 = 0$$

The distance from a point  $(x_0, y_0)$  to a line given by  $aX + bY + c = 0$  is:

$$d = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

### 4. RANSAC for lines

In this laboratory session our model will be a line. Hence, we consider the example of fitting a line to a set of points. The minimum number of points that define a line is two.

The steps for fitting a line to a set of points are the following:

1. Initialize the parameters of the algorithm.  
For the image in our example (*line.bmp*) we can estimate:  
 $t = 10$ ;  
 $p = 0.99$ ;  
 $q = 0.8$ ;  
 $s = 2$ ;

- Having these parameters we can compute  $N$  and  $T$ .
2. Repeat for  $N$  iterations:
    - a. Randomly select two points from the dataset.
    - b. Compute the parameters  $(a, b, c)$  that define the line passing through those two points.
    - c. Count the number of inliers for the current line.
    - d. If the number of inliers is greater or equal to  $T$  terminate step 2.
    - e. Keep the line which has the maximum number of inliers.
  3. Draw the line having the maximum number of inliers.

## 5. Practical background

Opening an image as grayscale:

```
Mat img = imread("filename", CV_LOAD_IMAGE_GRAYSCALE);
```

Creating a grayscale image:

```
Mat dst(height, width, CV_8UC1); //8bit unsigned 1 channel
```

Accessing the pixel at position row  $i$  and column  $j$ :

```
uchar pixel = img.at<uchar>(i,j); //unsigned char type
```

Modifying the pixel at position row  $i$  and column  $j$ :

```
img.at<uchar>(i,j) = 255; //white
```

Draw a line between two points:

```
line(img, Point(x1, y1), Point(x2, y2), Scalar(B,G,R));
```

Viewing the image:

```
imshow("title", img);  
waitKey();
```

## 6. Practical work

1. Open the image and store points which are represented as black pixels.
2. Calculate the parameters  $N$  and  $T$  from the suggested values for  $t, p, q$  and  $s$ .
3. Apply the RANSAC method:
  - a. choose two different points;
  - b. determine the line equation passing through the selected points;
  - c. Find the distances of each point to the line;
  - d. Count the number of inliers;
  - e. Write the correct termination conditions (based on the size of the consensus set and the maximum number of iterations).
4. Draw the optimal line found by the method.

## **7. References**

- [1] Robert C. Bolles, Martin A. Fischler: *A RANSAC-Based Approach to Model Fitting and Its Application to Finding Cylinders in Range Data*, 1981
- [2] Richard Hartley, Andrew Zisserman: *Multiple View Geometry in Computer Vision*, 2003