

Estimating the viewpoint position from a three-dimensional image

J.M. Sanchiz †, R.B. Fisher §*

†Departament d'Informàtica, Universitat Jaume I, E-12071 Castelló, Spain

§Institute of Perception, Action and Behavior

University of Edinburgh, 5 Forrest Hill, Edinburgh EH1 2QL, UK

Abstract. We present a method for estimating the viewpoint from where a 3D image has been taken using a central-projection range sensor. We assume we have the 3D coordinates of the points, organized with a known topology, but considerable noise is present in the data. At points in the scene where there are surface discontinuities we estimate step rays through a linear interpolation. The viewpoint is found as the point of minimum distance to the set of step rays. To cope with noise, we define an unbiased distance measure. The minimization of the sum of distances provides the viewpoint. We present results of several experiments carried out with 3D images of an old church.

1 Introduction

With the introduction of long range, wide angle, laser-based range sensors, three-dimensional images are becoming more and more available to the scientific community and to the general public. On the Internet one can already find collections of three-dimensional models made of polygons, or raw three-dimensional images consisting of a matrix of 3D points.

3D images consist of a *point cloud* or can be arranged as a matrix, but usually there is no information available about the sensor: field of view, angular resolution, viewpoint, etc.

Furthermore, even if one has the range sensor, it is possible that it does not provide all this information. The coordinate system origin of the 3D points is not necessarily the sensor origin. Or even if the extrinsic parameters of the sensor are known (the exact position and orientation of the tripod where the sensor is mounted), the transformation from the extrinsic coordinate system to the true sensor origin (viewpoint) is often unknown.

This latter information, the viewpoint, is probably the most important one since, as we assume the 3D coordinates of the point cloud and neighborhood relationships between points in the 3D image are known, knowing the viewpoint allows us to infer other information, such as the aperture of the field of view, sensor orientation or angular resolution. It also allows deduction of occlusion relationships, rejection of outliers, etc.

3D images, like intensity images, are noisy [1]. For example the 3D coordinates of the points may be computed by reading in information from the line scanning device (which deflects the laser beam horizontally by rotating a mirror), or from the tilt head (which modifies the azimuth angle of the rotating mirror). If the readings from the motor encoders are mistaken, it may result in a big angular drift of the points being scanned. The depth may be quite correct, but the point location in space is not.

Time-of-flight or phase-based range sensors also include depth error along the line of sight of each scanned point, usually categorized as Gaussian noise of zero mean and a certain standard deviation. The deviation, of the order of cm for some sensors [2], can be comparable to the scene structure at some areas.

Triangulation sensors have both range and direction errors from sensor noise and imprecision in locating corresponding features.

If we know the viewpoint, we can predict the pan (tilt) angle of each column (row) and discard or correct outliers, that is, put drifted points back in place.

In this work we assume the 3D image has been taken with a central-projection range sensor with geometry as shown in Figure 1. As can be seen, in this geometry all rays start from the center of the mirror. We also assume that enough structure (surface discontinuities) exists in the scene to allow us to estimate view directions. So, the method would not work if the scene is just a plane (a wall, for example).

Made explicit, our assumptions are:

* Email: sanchiz@uji.es, rbf@dai.ed.ac.uk

- Enough depth discontinuities exist, and are in random positions.
- The depth discontinuity lines of sight intersect at the viewpoint.

The question we address is: **is it possible to deduce the viewpoint given a 3D image ?**, and we show that it is possible to do this accurately. We know of no previous work addressing this problem, hence there are only a few references in this paper.

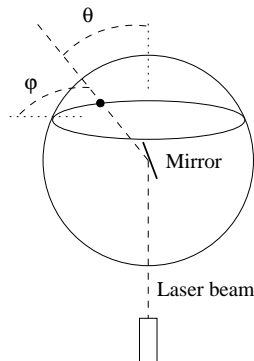


Fig. 1.: Pan-tilt range sensor geometry

2 Viewpoint detection

If we knew the line-of-sight ray (*view ray*) from the center of scanning to some scanned points, then the viewpoint could be easily computed from a number of rays. In theory two rays are enough to compute it, since all rays intersect at the viewpoint, but if we acquire more than two rays, the viewpoint can be computed by minimization as the point of minimum distance to all these rays.

The problem is then to estimate the view rays at a number of points in the 3D image, by some geometric-based technique. Once we acquire a set of estimated view rays, we have the problem of fitting a model (the position of the viewpoint) to a set of noisy data (the estimated view rays) in the presence of outliers. Several fitting methodologies exist to address this problem [3] [4] [5] [6]. We have applied (and discussed) one of them, the *Random Sample Consensus (RanSaC)* [3]. Here, the model fit is the estimated viewpoint and the set of consensual data are the rays that agree with the model.

It is easy to estimate view rays at points where the scanned surface has a depth discontinuity, we call these rays *step rays*. Figure 2 illustrates how a depth discontinuity can be detected. The distance d between two adjacent scanned points is $d = \frac{a}{\cos \gamma}$, where a can be approximated, for small angular resolutions, by the arc length at a point between the two scanned points (Figure 2 (top)). $a = \alpha r$, where α is the angle between the lines of sight of the two scanned points, and r is the average distance to the viewpoint. Then $d = \frac{\alpha r}{\cos \gamma}$. A common characteristic of range sensors is that they cannot observe surfaces viewed with an angle between the line of sight and the surface normal bigger than a limiting value (γ_{max}). So, if the two scanned points belong to the same surface, which is locally smooth in the area between them, the maximum distance between these points will be $d_{max} = \frac{\alpha_{max} r_{max}}{\cos \gamma_{max}}$. Figure 2 (bottom) shows two scanned points belonging to different surfaces; the distance between them is now expected to be much bigger.

In order to set a threshold h to select step rays we have to compute d_{max} . This can be done estimating the maximum angular resolution of the sensor, α_{max} , the maximum distance to the surfaces in the scene, r_{max} , and the maximum view angle of the sensor, γ_{max} . For this task, knowledge of the sensor and of the scene will help.

For example, assuming a maximum distance r_{max} of 20 meters, a maximum viewing angle γ_{max} of 45 degrees, and an angular resolution α_{max} of 20 steps per degree, we have

$$d_{max} = \left(\frac{\pi}{180 \cdot 20} \text{ rad}\right) (20 \text{ metres}) \frac{1}{\cos(45 \text{ deg})} = 0.024 \text{ metres} = 2.4 \text{ cm}.$$

So we can use a distance threshold of $h = 24 \text{ cm}$ (one order of magnitude above d_{max} , for example).

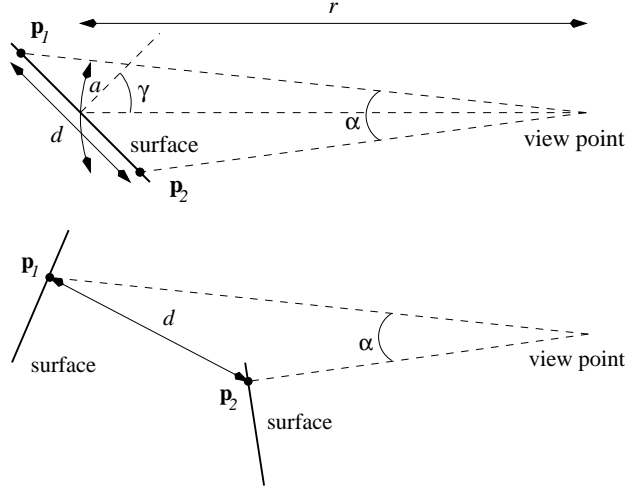


Fig. 2.: Two consecutive points \mathbf{p}_1 , \mathbf{p}_2 , scanned on the same surface (top), and on different surfaces (bottom)

Assume ² that a 3D image is a matrix $[\mathbf{p}_{row,col}]$ ($row \in [0..M-1]$, $col \in [0..N-1]$) where $\mathbf{p}_{i,j} = (x_{i,j}, y_{i,j}, z_{i,j})'$. The method to detect step rays is:

- Traverse the 3D image row by row and select a *horizontal* step ray if the distance from $\mathbf{p}_{i,j}$ to $\mathbf{p}_{i,j+1}$ is bigger than a threshold h .
- Traverse the 3D image column by column and select a *vertical* step ray if the distance from $\mathbf{p}_{i,j}$ to $\mathbf{p}_{i+1,j}$ is bigger than a threshold h .

The direction of a step ray is computed by estimating the point where the surface would have been scanned if a discontinuity had not occurred. Figure 3 illustrates how this point is computed by a linear interpolation. As shown, a horizontal step ray $\mathbf{h}_{i,j+1}$ at point $\mathbf{p}_{i,j+1}$ is computed as:

$$\mathbf{h}_{i,j+1} = \mathbf{p}_{i,j+1} + \lambda(\hat{\mathbf{p}}_{i,j+1} - \mathbf{p}_{i,j+1}) \quad (1)$$

where λ is the parameter that expands the ray, and $\hat{\mathbf{p}}_{i,j+1}$ is the interpolated point, computed as $\hat{\mathbf{p}}_{i,j+1} = \mathbf{p}_{i,j} + (\mathbf{p}_{i,j} - \mathbf{p}_{i,j-1})$.

Similarly, a vertical step ray $\mathbf{v}_{i+1,j}$ at point $\mathbf{p}_{i+1,j}$ is computed as:

$$\mathbf{v}_{i+1,j} = \mathbf{p}_{i+1,j} + \lambda(\hat{\mathbf{p}}_{i+1,j} - \mathbf{p}_{i+1,j}) \quad (2)$$

where now $\hat{\mathbf{p}}_{i+1,j} = \mathbf{p}_{i,j} + (\mathbf{p}_{i,j} - \mathbf{p}_{i-1,j})$.

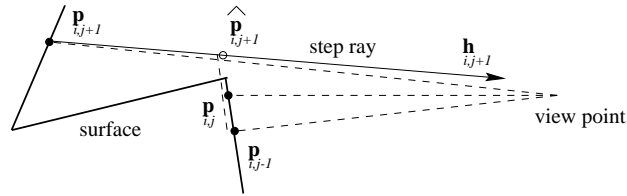


Fig. 3.: A step ray computed by linear interpolation

From the set of step rays (horizontal and vertical), the viewpoint is computed as the point of minimum distance to a subset S of rays, such that the size of S is bigger than a threshold n , and the distance from each ray in S to the viewpoint is smaller than a threshold t . These are the consensual rays.

² It is not necessary to have a regular array to apply this method. All that is required is knowing the topological or neighborhood relations between points. A regular matrix is the most common topology.

For selecting the subset S we use *RanSaC* [3], an algorithm for general model fitting and consensus selection in noisy data.

Among the most established fitting and data association methodologies, *RanSaC* is a robust and simple one. Its main criticism may be that it is a *classify-then-fit* approach, while it has been reported that a *classify-while-fit* approach [6] is more well suited for multiple-model problems. In this case, we only have one model, and the data has to be classified as belonging or not to the model. The model is found by minimization using all the data that has been previously classified by a minimum amount of randomly selected sample data.

Despite the years passed since *RanSaC* was reported, its approach remains valid for many applications. Its simplicity is also something to take into account. Other alternative methods include the ones based on *Robust Statistics* [7] [4] [5]. Here all the data is used to find the model by minimization using a measure of distance that tends to minimize the influence of outliers (*i.e.* data that are at a relatively far distance from the model). As stated in [6], all these methods trade robustness for accuracy, and fail when the number of outliers increases to more than 50%, which could really be the case in our problem.

With *RanSaC* the user can fix the minimum number of correct data to produce a valid model, as well as the accuracy. We believe this is an advantage in the present problem.

A short review of *RanSaC* is as follows: let A be a set of data, and m be the minimum number of data needed to compute a model.

```

Repeat for a maximum number of times
  Select  $m$  data at random from  $A$ 
  Compute an exact model from the  $m$  data
  Put the elements in  $A$  with distance to the model smaller than a threshold  $t$  into subset  $S$ 
  If  $\text{cardinal}(S)$  is bigger than a threshold  $n$ , exit loop
End of Repeat
Compute the model from all the elements in  $S$  by minimization

```

In our case the minimum number of data to compute the model (viewpoint) is two rays, $m = 2$. From two rays the model is also computed by minimization, as it is from any number of rays. Once a viewpoint is estimated from two sample rays, to check if other rays agree with the model, we compute the distance from each step ray to the viewpoint. The threshold for this distance, t , can be set to a value of a few (3 for example) orders of magnitude smaller than the scene dimensions. If the scene extends to several meters, t can be set to a few mm.

The threshold n , that indicates how many data are considered sufficient to give a final solution to the model, depends on the number of outliers present in the image. It is convenient to express it as a fraction of the data size (step rays). In our experiments we have obtained satisfactory results with a value of 60% of the total number of step rays detected. This figure has been set after a few initial tests. It should be bigger for images with few outliers, and smaller for more noisy images.

3 Error analysis

Although the viewpoint is found by minimization of a distance measure from the set of step rays, we may be introducing some errors in the way these step rays are estimated, which may bias the further estimation of the viewpoint.

As can be seen in Figure 3, if the surface is not perpendicular to the view direction, the linear interpolation of $\hat{\mathbf{p}}_{i,j+1}$ from points $\mathbf{p}_{i,j-1}$ and $\mathbf{p}_{i,j}$ will fall in the next ray, $\mathbf{h}_{i,j+1}$, only if the viewpoint is at infinity. Otherwise $\hat{\mathbf{p}}_{i,j+1}$ will fall before or after the real ray (from $\mathbf{p}_{i,j+1}$ to the viewpoint) depending on the surface orientation, resulting in step rays estimated with too much or too little inclination. This error is unavoidable, since we do not know the position of the viewpoint a priori, but we can consider that the errors tend to cancel if the scene contains a big number of surfaces with random orientations.

Let \mathbf{r}_i ($i \in [1..N]$) be a set of rays, $\mathbf{r}_i = \mathbf{c}_i + \lambda_i \mathbf{n}_i$ ($i \in [1..N]$), where \mathbf{c}_i is the starting point, \mathbf{n}_i is the direction vector, and λ_i is the parameter that expands the ray.

In order to find the point of minimum distance to the set, $\mathbf{v} = (x, y, z)'$, the standard procedure is to minimize the distance from a point \mathbf{v} to a line with respect to x , y and z .

The distance from a point to a line is expressed as:

$$d_i = \frac{\|\mathbf{n}_i \times (\mathbf{v} - \mathbf{c}_i)\|}{\|\mathbf{n}_i\|} \quad (3)$$

The addition of the squared distances to all rays is:

$$D^2 = \sum_{i=1}^N d_i^2 \quad (4)$$

The point of minimum distance is that for which

$$\partial D^2 / \partial x = \partial D^2 / \partial y = \partial D^2 / \partial z = 0 \quad (5)$$

If we assume that the 3D image points used to estimate the step rays include Gaussian noise of zero mean and variance σ_0^2 on each component, then \mathbf{c}_i is a random Gaussian vector of a certain mean and covariance $\sigma_0^2 \mathbf{I}$ (\mathbf{I} being a 3×3 identity matrix). \mathbf{n}_i is a random Gaussian vector of a certain mean and covariance $\sigma_1^2 \mathbf{I}$ ($\sigma_1^2 = 6\sigma_0^2$), since, from (1, 2), $\mathbf{n}_i = -\mathbf{p}_{i,j+1} + 2\mathbf{p}_{i,j} - \mathbf{p}_{i,j-1}$ for horizontal step rays, or $\mathbf{n}_i = -\mathbf{p}_{i+1,j} + 2\mathbf{p}_{i,j} - \mathbf{p}_{i-1,j}$ for vertical step rays.

To see if the expression in (3) is biased, we find its expectation, resulting in:

$$E[d_i^2] = \frac{\|\bar{\mathbf{n}}_i \times (\mathbf{v} - \bar{\mathbf{c}}_i)\|^2 + 2\sigma_0^2 \|\bar{\mathbf{n}}_i\|^2 + 6\sigma_0^2 \sigma_1^2 + 2\sigma_1^2 \|\mathbf{v} - \bar{\mathbf{c}}_i\|^2}{\|\bar{\mathbf{n}}_i\|^2 + 3\sigma_1^2} \quad (6)$$

where the mean values $\bar{\mathbf{c}}_i$, $\bar{\mathbf{n}}_i$ can be estimated by the measured values, \mathbf{c}_i , \mathbf{n}_i (*Maximum Likelihood* or *Least Squares* criterion [8]).

So, instead of using expression (3), we minimize the unbiased distance measure:

$$d_{unbiased, i}^2 = \frac{\|\mathbf{n}_i \times (\mathbf{v} - \mathbf{c}_i)\|^2 - 2\sigma_0^2 \|\mathbf{n}_i\|^2 - 6\sigma_0^2 \sigma_1^2 - 2\sigma_1^2 \|\mathbf{v} - \mathbf{c}_i\|^2}{\|\mathbf{n}_i\|^2 - 3\sigma_1^2} \quad (7)$$

Taking the partial derivatives (5) of (7), we obtain a linear system $\mathbf{A}\mathbf{v} = \mathbf{b}$, where

$$\mathbf{A} = \begin{pmatrix} \sum_{i=1}^N (n_{iy}^2 + n_{iz}^2 - 2\sigma_1^2) - \sum_{i=1}^N (n_{ix}n_{iy}) & -\sum_{i=1}^N (n_{ix}n_{iz}) \\ -\sum_{i=1}^N (n_{ix}n_{iy}) & \sum_{i=1}^N (n_{ix}^2 + n_{iz}^2 - 2\sigma_1^2) - \sum_{i=1}^N (n_{iy}n_{iz}) \\ -\sum_{i=1}^N (n_{ix}n_{iz}) & -\sum_{i=1}^N (n_{iy}n_{iz}) & \sum_{i=1}^N (n_{ix}^2 + n_{iy}^2 - 2\sigma_1^2) \end{pmatrix} \quad (8)$$

and

$$\mathbf{b} = \begin{pmatrix} \sum_{i=1}^N c_{ix}(n_{iy}^2 + n_{iz}^2 - 2\sigma_1^2) - \sum_{i=1}^N (c_{iy}n_{ix}n_{iy}) - \sum_{i=1}^N (c_{iz}n_{ix}n_{iz}) \\ -\sum_{i=1}^N (c_{ix}n_{ix}n_{iy}) + \sum_{i=1}^N c_{iy}(n_{ix}^2 + n_{iz}^2 - 2\sigma_1^2) - \sum_{i=1}^N (c_{iz}n_{iy}n_{iz}) \\ -\sum_{i=1}^N (c_{ix}n_{ix}n_{iz}) - \sum_{i=1}^N (c_{iy}n_{iy}n_{iz}) + \sum_{i=1}^N c_{iz}(n_{ix}^2 + n_{iy}^2 - 2\sigma_1^2) \end{pmatrix} \quad (9)$$

Solving the system we obtain the point $\mathbf{v} = (x, y, z)'$ of minimum distance to the set of rays.

4 Results

We have checked the approach with ten 3D panoramic images taken with a range scanner (LARA 12600, Zoller&Frohlich) [9] inside an old church in Bornholm, Denmark. Each image is of 8000×1400 points, the ground-truth viewpoint is known, $(0,0,0)'$, and the extent of these images is of several meters.

Although we only have ten images, the volume of data they contain is considerable. The number of points used in the experiments has been of 1.12×10^8 and, by subsampling and extracting subimages of smaller angular extent, we have used more than 600 images in the tests.

Figure 4 shows an overhead view of one of the test images, with the consensual step rays after applying *RanSaC*. The image covers quite a big area, and it is difficult to appreciate details in it. As can be seen, the rays point at the viewpoint, in the center of the church. Similar results can be seen in Figure 5, but for a more detailed partial view of the church containing some chairs. In fact the method can be used not only with panoramic images, but with whatever the aperture of the field of view.

Figure 6 (left) shows the estimated view points for the test images. The mean distances from these centers to the ground truth $(0,0,0)'$ is 1.26 mm and the standard deviation is 1.43 mm. Figure 6 (right)

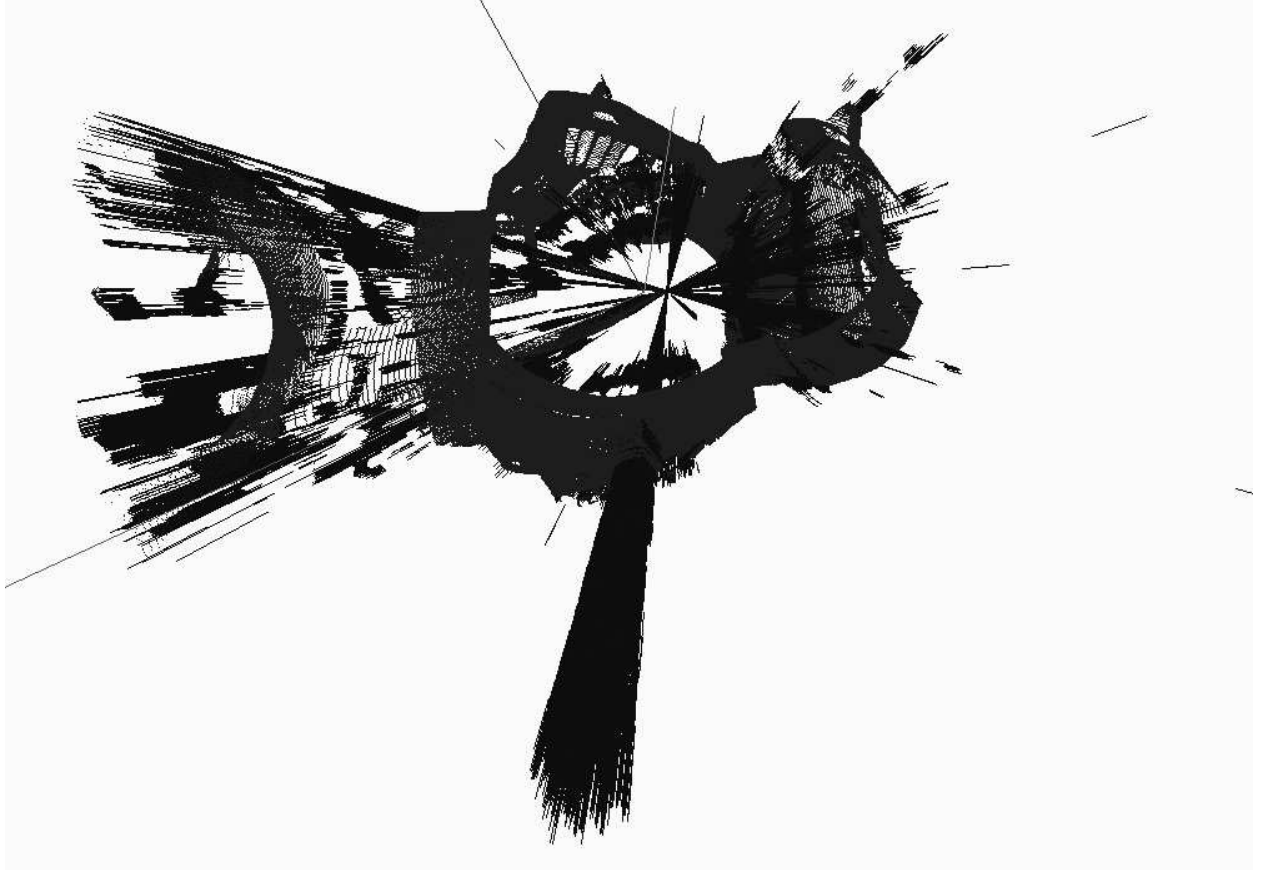


Fig. 4: Full 3D image (8000×1400 points) of a church seen from above, with the consensual step rays after applying *RanSaC*, converging on the viewpoint in the centre of the church.

shows plots with the mean distance from all of the consensual rays to the estimated view point, and the standard deviation of these distances. Both figures are of the order of few mm.

Figure 7 shows the number of iterations of the *RanSaC* algorithm until the number of consensual rays exceeds the threshold (60% of the total number of step rays), and the percentage of consensual data found. The threshold value of 60% was chosen running *RanSaC* first with a threshold of 100% (all data should fit the model), and observing the different percentages of consensual data obtained at each iteration.

To check the dependence of the approach with the amount of data used, we ran two more experiments focused on using only part of the data. In the first experiment, we used different sized horizontal sectors of panoramic images, but with the same image resolution. Figure 8 shows the mean distance from the consensual rays to the estimated centroid, and their standard deviation taking different sector sizes from the panoramic image in Figure 4. It also shows the distance from the estimated viewpoint to the ground-truth viewpoint. The method works well until the sector width becomes quite small (about 20 degrees).

In the second experiment we sub-sampled the test images using one of every number of rows and one of every number of columns. Since the step rays estimation (1, 2) is based on assuming that points $\mathbf{p}_{i,j}$ and $\mathbf{p}_{i,j-1}$, or $\mathbf{p}_{i,j}$ and $\mathbf{p}_{i-1,j}$, belong to the same surface, the method will start to fail when the sub-sampling rate is so big that this assumption does not hold any more. The experiment used the same thresholds as the previous experiments, but we ran *RanSaC* one thousand times to determine, for every sub-sampling rate, the biggest number of consensual rays that we could use. Figure 9 corresponds to the sub-sampling of the panoramic image in Figure 4. It shows how the mean and standard deviation of the distances from the consensual rays to the estimated viewpoint gradually increase with the sub-sampling rate. Figure 9 also shows how the number of consensual rays decreases with the sub-sampling rate since, the more we sub-sample, the less accurate the step ray estimation is. This is probably the most indicative graph because, for problems where we only have one model to fit (the viewpoint), we would like to acquire the highest amount of useful data to compute the model.

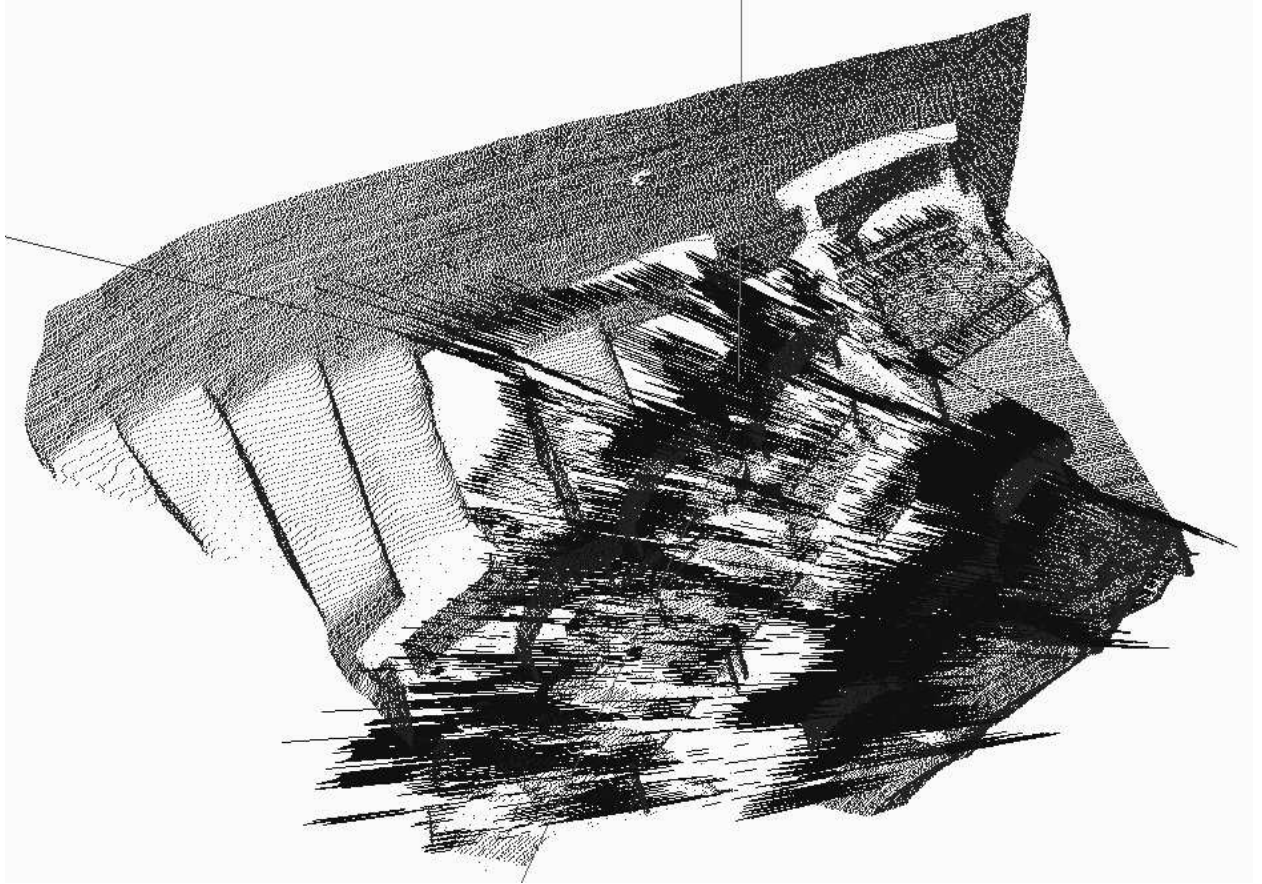


Fig. 5.: Detailed 3D subimage (1000×700 points), and consensual step rays after applying *RanSaC*, the true viewpoint is to the right.

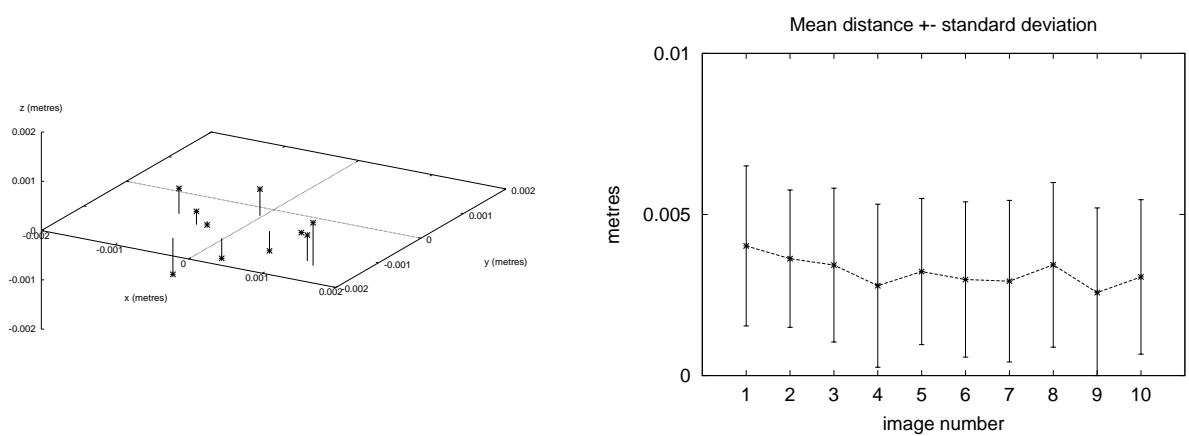


Fig. 6.: *Left*: Estimated centres for 10 test images. *Right*: Statistical analysis of the view point estimation in 10 test images. Dots are the mean distance from the consensual rays to the estimated viewpoint. Bars indicate one standard deviation of the distances, added to and subtracted from the mean.

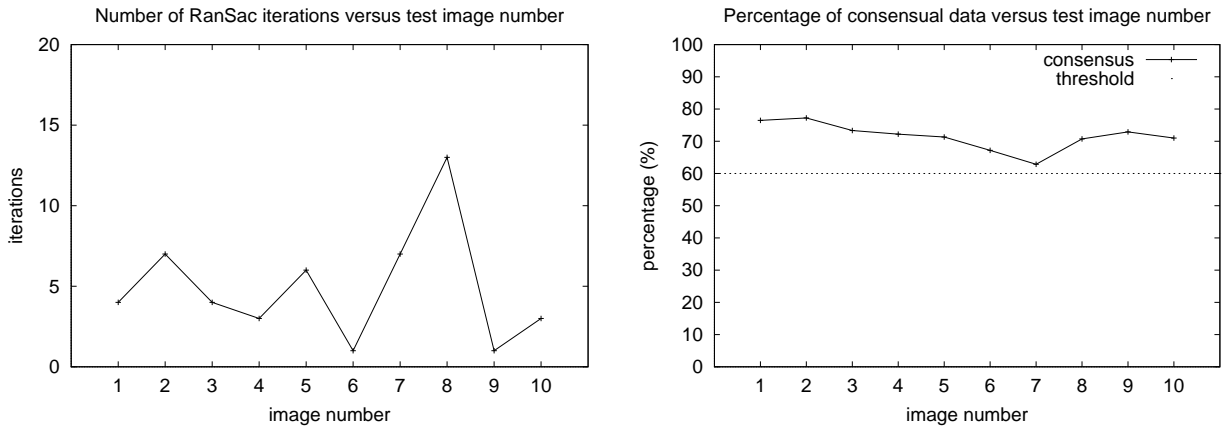


Fig. 7.: *RanSaC* results on the 10 test images. *Left*: number of iterations until the number of consensual rays exceeds the threshold (fixed at 60% of the total number of step rays). *Right*: consensus as a percentage of the total number of step rays. The threshold is pointed out as a dotted line.

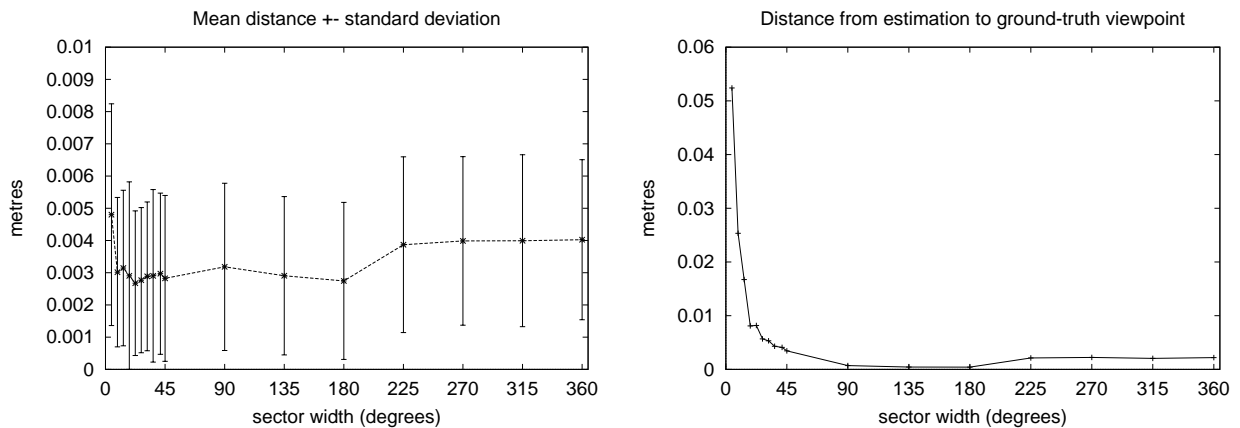


Fig. 8.: Estimation error as a function of the amount of data used to compute the viewpoint, the horizontal axis is the sector angle of the image. *Left*: mean distance and standard deviation (added and subtracted) from the consensual rays to the estimated viewpoint. *Right*: Distance from the estimated viewpoint to the ground truth.

5 Conclusions

In this paper we have applied a paradigm for model fitting in the presence of outliers (*RanSaC*) to the problem of estimating the viewpoint from where a 3D image has been taken, assuming a central-projection range sensor.

We have used an approach to step ray detection based on first-order interpolation. Step rays are the data used for model fitting (viewpoint estimation). We have presented clues for setting the parameters involved.

We have used an unbiased distance estimator that takes into account the data noise, modeling it as Gaussian with zero mean and a known (or estimated) standard deviation.

The experiments give an estimate of the viewpoint with an accuracy of three orders of magnitude smaller than the scene's extent.

Two experiments focused on estimating the viewpoint from a reduced volume of data: taking only sectors of a panoramic image, and image sub-sampling. Accuracy is largely independent of sector width provided at least 20 degrees of sampling is obtained. On the other hand, image sub-sampling reduces the confidence of the step ray estimates, resulting in fewer consensual data being found, but only slowly degrades the accuracy of the estimated viewpoint.

From our results we can conclude that the method presented is a valid approach for viewpoint estimation

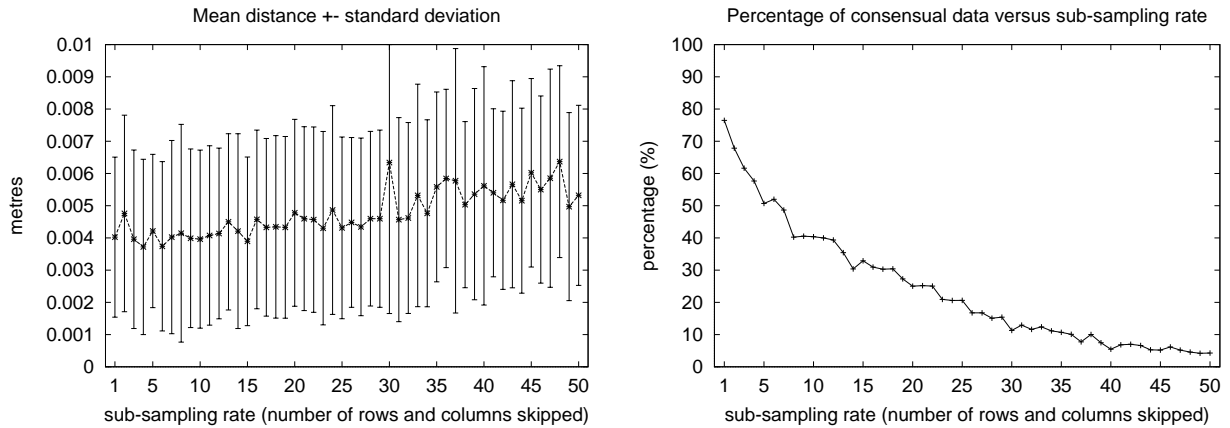


Fig. 9.: Estimation error as a function of the amount of data used to compute the viewpoint, the horizontal axis is the sub-sampling rate (number of rows and columns skipped). *Left*: mean distance and standard deviation (added and subtracted) from the consensual rays to the estimated viewpoint. *Right*: Percentage of consensual data found for every sub-sampling rate. As expected, the method finds less consensual data as the sub-sampling rate increases.

from noisy images.

Acknowledgments

The authors wish to acknowledge the support of the European Union through SMART II (ERB4061PL950841) and CAMERA (ERB FMRX-CT97-0127) research networks, and of the Spanish *Comisión Interministerial de Ciencia y Tecnología* (CICYT), under project TAP1999-0590-C02-01.

References

1. M. Hebert and E. Krotkov. 3d measurements from imaging laser radars: how good are they ? *International Journal of Imaging and Vision Computing*, 10(3):170–178, 1992.
2. G. Noé, M.I. Ribeiro, and J.A. Santos Victor. 3d laser scanner performance evaluation. Technical report, Instituto de Sistemas e Robótica, Lisboa, Portugal, 1999.
3. M.A. Fischler and B.C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
4. P. Rousseeuw and A. Leroy. *Robust Regression and Outlier Detection*. John Wiley & Sons, 1987. Regression and data association.
5. P. Meer, D. Mintz, and A. Rosenfeld. Robust regression methods for computer vision. *International Journal of Computer Vision*, 6(1):59–70, 1991.
6. G. Danuser and M. Stricker. Parametric model fitting: from inlier characterization to outlier detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(2):263–280, 1998.
7. P. Huber. *Robust Statistics*. John Wiley & Sons, 1981. Regression and data association.
8. Y. Bar-Shalom and T.E. Fortmann. *Tracking and Data Association*. Academic Press, 1988. Estimation.
9. G. Dalton. Reverse engineering using laser metrology. *Sensor Review*, 18(2):92–96, 1998.