

TOWARDS A CRITERION FOR EVALUATING THE QUALITY OF 3D RECONSTRUCTIONS

Amit K. Roy Chowdhury, Rama Chellappa

Center for Automation Research,
Dept. of Electrical and Computer Engineering,
University of Maryland, College Park, MD 20742.
{amitrc,rama}@cfar.umd.edu

ABSTRACT

Even though numerous algorithms exist for estimating the structure of a scene from its video, the solutions obtained are often of unacceptable quality. To overcome some of the deficiencies, many application systems rely on processing more information than necessary with the hope that the redundancy will help improve the quality. This raises the question about how the accuracy of the solution is related to the amount of information processed by the algorithm. Can we define the accuracy of the solution precisely enough that we automatically recognize situations where the quality of the data is so bad that even a large number of additional observations will not yield the desired solution? This paper proposes an information theoretic criterion for evaluating the quality of a 3D reconstruction in terms of the statistics of the observed parameters (i.e. the image correspondences). The accuracy of the reconstruction is judged by considering the change in mutual information (or equivalently the conditional differential entropy) between a scene and its reconstructions and its effectiveness is shown through simulations.

1. INTRODUCTION

Obtaining accurate 3D models from video using the structure from motion (SfM) approach [1], [2], is extremely important because of its diverse applications, ranging from multimedia to medical diagnosis. Yet the quality of many of the automatic 3D reconstructions leave much to be desired. This has led many researchers to analyze the sensitivity, robustness and statistical error characterization of the existing algorithms, trying to understand algorithm behavior and the characteristics of the natural phenomenon that is being modeled [3], [4], [5], [6], [7], [8], [9]. To overcome these errors, the tendency has been to add redundancy in the information processed. This raises the question as to how the redundant information affects the quality of the final solution. In this paper, we consider the situation where multiple reconstructions of the same scene are available (called intermediate or individual reconstructions, in this paper), that are combined together to obtain the final estimate (Figure (1)). We compute the incremental mutual information between the unknown 3D structure and increasing numbers of intermediate reconstructions.

Before proceeding to give a detailed description of the idea, we would like to draw the attention of the reader briefly to the area of

Partially supported by NSF grant #0086075; "ITR: Personalized Spatial Audio via Scientific Computing and Computer Vision"

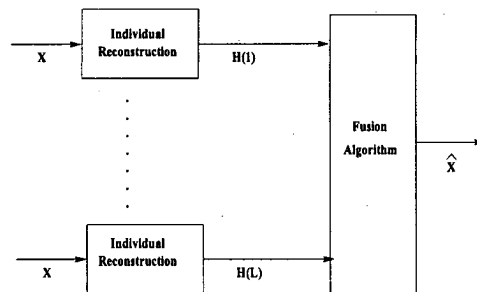


Fig. 1. Block diagram representation of the reconstruction framework. \mathbf{X} is the inverse depth that we want to estimate, $(\mathbf{H}(1), \dots, \mathbf{H}(L))$ are the intermediate reconstructions (e.g. from each individual camera), and $\hat{\mathbf{X}}$ is the final fused estimate.

model selection in statistics (AIC, BIC, MDL etc. [10]). The idea of fitting models to geometric data was formalized by Kanatani using a Geometric Information Criterion (GIC) [11]. However, a large number of SfM algorithms are not model based; they reconstruct individual point features of the scene. Our work tries to define the *quality of reconstruction* from point features in information theoretic terms.

2. AN INFORMATION THEORETIC CRITERION FOR 3D RECONSTRUCTION

2.1. Problem Formulation

We assume that all the depth values are aligned to a common frame of reference. Feature points will be represented by subscripts, separate reconstructions will be within parenthesis. The vector of estimates of the inverse depth¹ $[H_i(1), \dots, H_i(N)]'$ will be denoted by $\mathbf{H}_i^{(N)}$. The boldface notation $\mathbf{H}(i)$ will represent all the features in the i^{th} reconstruction. The final estimate $\hat{\mathbf{X}}$ of $\mathbf{X} = [X_1, \dots, X_M]'$ is obtained by fusing the individual reconstructions $(\mathbf{H}(1), \dots, \mathbf{H}(L))$. To keep the notation simple, the subscript for the feature point will not be mentioned, unless required. We assume that \mathbf{X} is Gaussian with zero mean and variance σ_x^2 .

¹The inverse depth is used throughout this paper since it is the quantity that is estimated from the SfM equations for reconstruction from a video and its statistics can be obtained in an analytic form more easily than for the depth.

The individual estimates are modeled as

$$H(i) = X + V(i) \quad (1)$$

where $X \sim \mathcal{N}(0, \sigma_x^2 = P_X)$ and $\{V(i), i = 1, \dots, N\}$ is a sequence of independent random variables distributed as $\mathcal{N}(0, \sigma_{V(i)}^2)$. Let $P_V = \text{diag}[P_V(i)]_{i=1, \dots, N} = \text{diag}[\sigma_{V(1)}^2, \dots, \sigma_{V(N)}^2]^2$.

2.2. Main Result

We will now present an information theoretic measure for evaluating the quality of a 3D reconstruction algorithm by analyzing the contribution of each of the individual reconstructions. Our entire analysis is for a particular point and thus the subscript will be dropped, unless required for clarity. From (1), $E[H(i)] = 0$ and

$$\begin{aligned} E[H(i)H(j)] &= E[(X + V(i))(X + V(j))] \\ &= P_X + P_V(i)\delta_{ij}, \end{aligned} \quad (2)$$

where δ_{ij} is a Kronecker delta function. Thus the covariance of $\mathbf{H}^{(N)}$ is $P_{\mathbf{H}^{(N)}} = P_V^{(N)} + \mathbf{1}_N P_X \mathbf{1}_N^T$, where $\mathbf{1}_N$ is a vector of N ones. Then the mutual information between X and $H(i)$,

$$\begin{aligned} I(X; H(i)) &= h(H(i)) - h(H(i)|X) \\ &= \frac{1}{2} \log \left(1 + \frac{P_X}{P_V(i)} \right). \end{aligned} \quad (3)$$

Next, consider the mutual information between the unknown X and the vector of observations $\mathbf{H}^{(N)}$. We will denote by $|K|$ the determinant of a matrix K .

$$\begin{aligned} I(X; \mathbf{H}^{(N)}) &= h(\mathbf{H}^{(N)}) - h(\mathbf{H}^{(N)}|X) \\ &\stackrel{(a)}{=} h(\mathbf{H}^{(N)}) - \sum_{i=1}^N \frac{1}{2} \log(2\pi e P_V(i)) \\ &\stackrel{(b)}{=} \frac{1}{2} \log \left(\frac{|P_V + \mathbf{1}_N P_X \mathbf{1}_N^T|}{|P_V|} \right). \end{aligned} \quad (4)$$

(a) is a result of applying the chain rule of entropy and substituting the expression for the differential entropy of a Gaussian random variable [12]; (b) is due to the fact that $|P_V| = \prod_{i=1}^N P_V(i) = \prod_{i=1}^N \sigma_{V(i)}^2$. Using the method of induction and the properties of determinants, it can be shown that $|P_V + \mathbf{1}_N P_X \mathbf{1}_N^T| = \prod_{i=1}^N \sigma_{V(i)}^2 + \sigma_x^2 \sum_{i=1}^N \prod_{j=1, j \neq i}^N \sigma_{V(j)}^2$. Then from (4), the expression for the mutual information becomes

$$I(X; \mathbf{H}^{(N)}) = \frac{1}{2} \log \left(1 + \sum_{i=1}^N \frac{\sigma_x^2}{\sigma_{V(i)}^2} \right). \quad (5)$$

Let us compute the difference in the mutual information for the two sets of observations, $\mathbf{H}^{(N)}$ and $\mathbf{H}^{(N-1)}$. We shall call this the

²Where necessary to distinguish a particular feature point, we will use the notation $\sigma_{x_j}^2$ and $P_{V_j}(i)$ or $\sigma_{V_j(i)}^2$ for the j^{th} point.

incremental mutual information, ΔI . Thus,

$$\begin{aligned} \Delta I &= I(X; \mathbf{H}^{(N)}) - I(X; \mathbf{H}^{(N-1)}) \\ &= \frac{1}{2} \log \left(\frac{|P_{V^{(N)}} + \mathbf{1}_N P_X \mathbf{1}_N^T|}{|P_{V^{(N-1)}} + \mathbf{1}_{N-1} P_X \mathbf{1}_{N-1}^T|} \cdot \frac{|P_{V^{(N-1)}}|}{|P_{V^{(N)}}|} \right) \\ &= \frac{1}{2} \log \left(\frac{\prod_{i=1}^N \sigma_{V(i)}^2 + \sigma_x^2 \sum_{i=1}^N \prod_{j=1, j \neq i}^N \sigma_{V(j)}^2}{\prod_{i=1}^{N-1} \sigma_{V(i)}^2 + \sigma_x^2 \sum_{i=1}^{N-1} \prod_{j=1, j \neq i}^{N-1} \sigma_{V(j)}^2} \right) \\ &= \frac{1}{2} \log \left(1 + \frac{1/\sigma_{V(N)}^2}{\frac{1}{\sigma_x^2} + \sum_{i=1}^{N-1} \frac{1}{\sigma_{V(i)}^2}} \right) \\ &= \frac{1}{2} \log \left(1 + \frac{1/P_V(N)}{\frac{1}{\sigma_x^2} + \sum_{i=1}^{N-1} \frac{1}{P_V(i)}} \right). \end{aligned} \quad (6)$$

Equation (6) gives us a measure of the extra information that would be obtained by including an additional observation into the fusion process. Also, since

$$I(X; \mathbf{H}^{(N)}) - I(X; \mathbf{H}^{(N-1)}) = h(X|\mathbf{H}^{(N-1)}) - h(X|\mathbf{H}^{(N)}), \quad (7)$$

the quantity defined as the incremental mutual information can also be referred to as the incremental conditional entropy. Thus we are measuring the reduction in the uncertainty of the solution as we consider an extra observation. The difference in the differential entropy determines the decrease in the coding length of the scene structure as the number of observations increases [12].

The above calculation requires computing the variances of the intermediate reconstructions. Any method to compute them is perfectly suitable. In an earlier work [13], we have shown how to do this for the case of 3D reconstruction using optical flow. It should be remembered that all the geometric quantities have to be with respect to a particular frame of reference; hence it may be necessary to transform the variances appropriately.

An Estimation Theoretic Interpretation: We will now present an alternative interpretation of the result in (6) from an estimation theoretic perspective. The mean squared distortion is defined as

$$D(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{M} \sum_{j=1}^M E[(X_j - \hat{X}_j)^2]. \quad (8)$$

Let $p(X_j, H_j(1), \dots, H_j(N))$ denote the joint density function of the parameter and observations. The mean square error estimator \hat{X}_j of X_j , obtained from $\mathbf{H}^{(N)}$, is $\hat{X}_j(N) = E[X_j | H_j^{(N)}]$. From the Cramer-Rao lower bound (CRLB) we can write the following set of inequalities.

$$\begin{aligned} D &\geq \frac{1}{M} \sum_{j=1}^M \frac{1}{E \left[-\frac{\partial^2}{\partial X^2} \log p(X_j, H_j(1), \dots, H_j(N)) \right]} \\ &= \frac{1}{M} \sum_{j=1}^M \frac{1}{\frac{1}{\sigma_{x_j}^2} + \sum_{i=1}^N E \left[-\frac{\partial^2}{\partial X^2} \log p(H_j(i)|X) \right]} \\ &\geq \frac{1}{\frac{1}{M} \sum_{j=1}^M \left(\frac{1}{\sigma_{x_j}^2} + \sum_{i=1}^N \frac{1}{P_{V_j}(i)} \right)} \\ &\triangleq \frac{1}{\frac{1}{M} \sum_{j=1}^M \frac{1}{D_j(N)}}. \end{aligned} \quad (9)$$

The last step is a result of the application of Jensen's inequality [14] and the fact that $E\left[-\frac{\partial^2}{\partial X^2} \log p(H_j(i)|X)\right] = \frac{1}{P_{V_j(i)}}$. Recalling that (6) is for a particular feature point where the subscript has been suppressed for clarity of notation, let us denote $\Delta I_j \triangleq I(X_j; \mathbf{H}_j^{(N)}) - I(X_j; \mathbf{H}_j^{(N-1)})$. Then from (9) and the last expression of (6), we get

$$\Delta I_j = \frac{1}{2} \log \left(\frac{D_j(N-1)}{D_j(N)} \right). \quad (10)$$

Alternatively, the innovations at the N^{th} stage, $\gamma_N = X_N - \hat{X}_N$. Then following the standard derivation for the Kalman filter [14], it can be shown that variance of the innovations

$$P_{\gamma_N} = \sigma_{V(N)}^2 \left(1 + \frac{1/\sigma_{V(N)}^2}{\frac{1}{\sigma_z^2} + \sum_{i=1}^{N-1} \frac{1}{\sigma_{V(i)}^2}} \right), \quad (11)$$

which shows that, for each feature point, the incremental mutual information is related to P_{γ_N} as

$$\Delta I = \frac{1}{2} \log \left(\frac{P_{\gamma_N}}{\sigma_{V(N)}^2} \right). \quad (12)$$

These relationships provide an alternative estimation theoretic interpretation to our result. Taken together (6), (10) and (12) demonstrate the use of statistical evaluation techniques to the SfM problem, when it is suitably formulated.

3. ANALYSIS AND EXPERIMENTS

3.1. Analysis:

Present methods to evaluate the quality of a reconstruction involve computing the distortion in (8). For a fusion algorithm, this means that we need to compute (8) at every stage of the fusion and decide when to stop. This is computationally intensive, distortion measures are not always very useful in practical experiments since the choice of an acceptable threshold is often arbitrary and the source of the error (whether in the intermediate reconstructions or in the fusion algorithm) is difficult to identify. In our approach, (6) gives a direct way to measure the contribution of the intermediate solutions and the accuracy of the final solution as the algorithm progresses. The statistics of the error can be computed using the SfM equations and its solutions, as described in [13]. If the solution is far from its desired values, the error would be larger than if the solution is close to its true value. When the error in the intermediate reconstructions is small, D_j is small and hence the difference in the mutual information is small. Ideally, this difference should go to zero as we include more and more observations. If the error is large, D_j would be large and ΔI_j would not decrease appreciably with the number of observations. Another salient feature of our method is that we measure the information content between the true structure and the reconstructions *before* the fusion. This allows us to understand the source of the error better since the effect of intermediate reconstructions and fusion algorithm are separated.

One scenario where this idea can be applied is reconstruction from a video sequence where intermediate reconstructions, $\mathbf{H}(1), \dots, \mathbf{H}(L)$, obtained from a few frames (two or three) are

combined together. Another application would be where partial reconstructions have been obtained from multiple cameras³. These partial models would have common overlapping regions which can be combined together to form the single estimate. In this case, $\mathbf{H}(1), \dots, \mathbf{H}(L)$ would represent these common sub-regions from L separate reconstructions.

The statistical assumptions of independence and Gaussianity are necessary in order to derive closed form expressions for the quantities of interest. The independence of the intermediate estimates $H(1), \dots, H(L)$ may be valid when these are obtained from separate imaging systems and then combined. When the same camera is used, the intermediate reconstructions should be obtained with non-overlapping frames; otherwise the common frames increase the dependencies. Regarding the Gaussianity assumptions, it has been pointed out by Zhang in [7] that the correspondence errors in SfM are usually normally distributed, if we can get rid of the outliers in the matches.

3.2. Experiments:

Experiment 1: A set of 3D points were generated so that we know their true positions. The perspective projections of these points were generated and Gaussian noise with zero mean and known variance was added to these 2D locations. The projections were taken for different positions of the camera, so that in the end a set of tracked features was obtained. From every pair of such tracked features, the positions of the original 3D points were estimated, which results in a set of 3D reconstructions. The first plot of Figure (2) shows the true value of the 3D points and their estimated reconstruction from all the frames over which the features could be tracked.⁴ The second diagram in Figure (2) plots the decrease in the incremental mutual information with the increasing number of intermediate reconstructions.

Experiment 2: As in the previous simulation, a set of features were tracked over a number of frames. However, the level of noise added to the feature positions was higher and it led to a mismatch of some of the features. The 3D positions of the points were estimated using the SfM algorithm and the results were erroneous as is clear from the first plot of Figure (3). The second plot of Figure (3) depicts this case where the incremental mutual information remains large and does not follow any trend.

Experiment 3: We will now present our result on a real video sequence. The video consists of a person moving his head in front of a static camera. The aim was to reconstruct the model of the head of the person from this video. The focal length of the camera was known. Figure (4)(a) represents an image from the video along with some of the feature points which were tracked. Figure (4)(b) represents the change in the incremental mutual information between the unknown 3D structure and the intermediate reconstructions from every pair of frames. Based on this measure, the 3D model was reconstructed using 25 frames and Figure (4)(c) shows one particular view of this model.

³This is the set-up in the "Eye Vision" technology developed by Carnegie Mellon University (CMU) and CBS Television (<http://www.ri.cmu.edu/events/sb35/tksuperbowl.html>).

⁴The first point was used to set the scale of the reconstruction, so that the geometric indeterminacies do not affect the result.

4. CONCLUSION

In this paper, we have introduced a method to evaluate the quality of 3D reconstruction from a video sequence. Existing methods rely on computing the distortion between the projections of the reconstructions and the original images and deciding that the reconstruction is of acceptable quality when the distortion is below a certain empirically chosen threshold. In this paper, we have shown that it is possible to evaluate the quality of the 3D structure estimate as the algorithm proceeds by computing the incremental mutual information, which determines the importance of considering an additional observation. It is related to the decrease in the coding length of the actual structure conditioned on the increasing number of observations. Finally, experimental results have been provided to justify these claims.

Acknowledgements

We would like to thank Prof. Adrian Papamarcou for many interesting and helpful suggestions. The first author would also like to thank Gang Qian, Kaushik Chakraborty and Damianos Karakos for many related technical discussions.

5. REFERENCES

- [1] O.D. Faugeras, *Three-Dimensional Computer Vision: A Geometric Viewpoint*, MIT Press, 1993.
- [2] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, Cambridge University Press, 2000.
- [3] J. Oliensis, "A critique of structure from motion algorithms," *NECI Tech. Report*, vol. <http://www.neci.nj.nec.com/homepages/oliensis/>, 2000.
- [4] Y. Ma, J. Kosecka, and S. Sastry, "Linear differential algorithm for motion recovery: A geometric approach," *International Journal of Computer Vision*, vol. 36, pp. 71–89, January 2000.
- [5] J. Weng, N. Ahuja, and T.S. Huang, "Optimal motion and structure estimation," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 15, pp. 864–884, September 1993.
- [6] G.S. Young and R. Chellappa, "Statistical analysis of inherent ambiguities in recovering 3-d motion from a noisy flow field," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 14, pp. 995–1013, October 1992.
- [7] Z.Y. Zhang, "Determining the epipolar geometry and its uncertainty: A review," *International Journal of Computer Vision*, vol. 27, pp. 161–195, March 1998.
- [8] D.D. Morris, K. Kanatani, and T. Kanade, "3d model accuracy and gauge fixing," in *CMU-RI-TR*, 2000.
- [9] Z. Sun, V. Ramesh, and A.M. Tekalp, "Error characterization of the factorization method," *Computer Vision and Image Understanding*, vol. 82, pp. 110–137, May 2001.
- [10] M. Hansen and B. Yu, "Model selection and the principle of minimum description length," *To appear in Journal of the American Statistical Association*.
- [11] K. Kanatani, "Geometric information criterion for model selection," *IJCV*, vol. 26, no. 3, pp. 171–189, February 1998.
- [12] T. Cover and J. Thomas, *Elements of Information Theory*, John Wiley and Sons, 1991.
- [13] A. Roy Chowdhury and R. Chellappa, "Robust estimation of depth and motion using stochastic approximation," in *Proc. IEEE ICIP-01*, 2001.
- [14] H.V. Poor, *An Introduction to Signal Detection and Estimation*, Springer-Verlag, 1988.

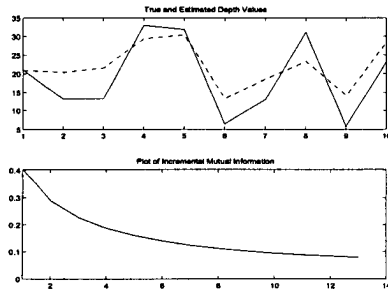


Fig. 2. The upper plot shows the true value of the depth of the 3D points using the solid line and the fused estimate from the intermediate reconstructions from all the frames using the dotted lines. The second diagram plots the decrease in the incremental information with the increasing number of frames.

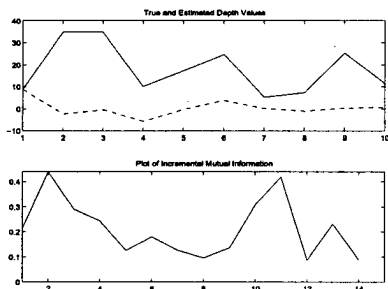


Fig. 3. The upper plot shows the true value of the depth of the 3D points using the solid line and the fused estimate from the intermediate reconstructions from all the frames using the dotted lines. The lower plot is the change in the mutual information with increasing number of frames. This is the case where the estimated reconstruction does not converge to the true value even with increasing observations.

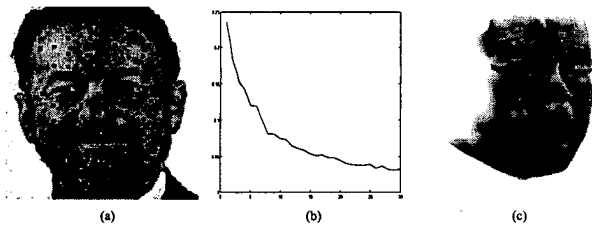


Fig. 4. The above figures represent a 3D reconstruction from video using the method of measuring the incremental mutual information to judge the quality of the result. (a) is one of the images from the video along with the set of tracked features used for the reconstruction. (b) represents the change in the incremental mutual information with the number of images; (c) depicts one view from the reconstructed model.