

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Reconstruction Project (COL 780)

Ashutosh Agarwal, Shivansh Chandra Tripathi

Abstract

Structure from motion (SfM) is a technique to estimate the 3d structure of objects using their corresponding 2d images taken from multiple views. In the most general case, we've only the images as input, and we are required to predict Camera's intrinsic parameters, and then construct the 3d view. In this project, we'll see a 3d reconstruction model for India Gate. Starting from the base model, we'll see the step by step improvements that we've built on it.

1. Introduction

Structure from motion (SfM) is the process of reconstructing 3D structure from its projections into a series of images taken from different viewpoints. It is a process of estimating camera pose and retrieving a sparse reconstruction simultaneously. Let us now understand the various steps involved to solve the problem :

- Detect Features:** The first stage is correspondence search which finds scene overlap in the input images. It identifies projections of the same points in overlapping images. This step detects features in the images, using any of the feature descriptors (HAHOG, SIFT, SURF, ORB, etc.)
- Match Features:** This step matches feature points between images. It first determines the list of image pairs to run, and then run the matching process for each pair to find corresponding feature points. Since there are a lot of possible image pairs, the process can be very slow. It can be sped up by restricting the list of pairs to match.
- Geometric Verification:** The third stage verifies the potentially overlapping image pairs C. Since matching is based solely on appearance, it is not guaranteed that corresponding features actually map to the same scene point. Therefore, SfM verifies the matches by trying to estimate a transformation that maps feature points between images using projective geometry. Depending on the spatial configuration of an image pair,

different mappings describe their geometric relation. A homography H describes the transformation of a purely rotating or a moving camera capturing a planar scene. If a valid transformation maps a sufficient number of features between the images, they are considered geometrically verified. Since the correspondences from matching are often outlier-contaminated, robust estimation techniques, such as RANSAC, are required. The output of this stage is a set of geometrically verified image pairs C, their associated inlier correspondences M.

4. **Reconstruct:** This step runs the incremental reconstruction process. The goal of the reconstruction process is to find the 3D position of tracks (the structure) together with the position of the cameras (the motion). **Incremental reconstruction algorithm** For this project, we will be using OpenSfM python library to build the reconstruction pipeline. OpenSfM implements an incremental structure from motion algorithm. This is reconstruction algorithm that starts building a reconstruction of a single image pair and then iteratively add the other images to the reconstruction one at a time. The algorithm has three main steps:

- **Finding good initial pairs**

To compute the initial reconstruction using two images, there needs to be enough parallax between them. That is, the camera should have been displaced between the two shots, and the displacement needs to be large enough compared to the distance to the scene.

To compute whether there is enough parallax, we start by trying to fit a rotation only camera model to the two images. We only consider image pairs that have a significant portion of the correspondences that can not be explained by the rotation model. We compute the number of outliers of the model and accept it only if the portion of outliers is larger than 30

The accepted image pairs are sorted by the number of outliers of the rotation only model.

- **Bootstrapping the reconstruction**

108 To bootstrap the reconstruction, we use the first
 109 image pair. If initialization fails we try with the
 110 next on the list. If the initialization works, we
 111 pass it to the next step to grow it with more im-
 112 ages.
 113

114 The reconstruction from two views can be done
 115 by two algorithms depending on the geometry of
 116 the scene. If the scene is flat, a plane-based ini-
 117 tialization is used, if it is not flat, then the five-
 118 point algorithm is used. Since we do not know a
 119 priori if the scene is flat, both initializations are
 120 computed and the one that produces more points
 121 is retained.
 122

123 If the pair gives enough inliers we initialize a re-
 124 construction with the corresponding poses, trian-
 125 gulate the matches and bundle adjust it.
 126

- **Growing the reconstruction**

127 Given the initial reconstruction with two im-
 128 ages, more images are added one by one starting
 129 with the one that sees more of the reconstructed
 130 points.
 131

132 To add an image it needs first needs to be aligned
 133 to the reconstruction. This is done by finding the
 134 camera position that makes the reconstructed 3D
 135 points project to the corresponding position in the
 136 new image. The process is called resectioning.
 137

138 If resectioning works, the image is added to the
 139 reconstruction. After adding it, all features of
 140 the new image that are also seen in other recon-
 141 structed images are triangulated. If needed, the
 142 reconstruction is then bundle adjusted and even-
 143 tually all features are re-triangulated.
 144

145 **5. Undistort:** This step creates undistorted version of
 146 the reconstruction, tracks and images. The undistorted
 147 version can later be used for computing depth maps.
 148

149 **6. Compute Depthmaps:** This step computes a dense
 150 point cloud of the scene by computing and merging
 151 depthmaps. It requires an undistorted reconstructions.
 152

2. Dataset Collection

153 The monument we've chosen is India Gate. There
 154 should not be too small disparity between different im-
 155 ages else the depthmap computation will be less pre-
 156 cise. Keeping these points in mind, our initial dataset
 157 was comprised of high resolution images collected
 158 from google and flickr. Later we made an improved
 159 dataset of high resolution images by collecting large
 160 number of frames from a drone video of India Gate.
 161 We also collected data using screenshots of 360 degree
 162 photos on youtube. At last, to get the top views we
 163

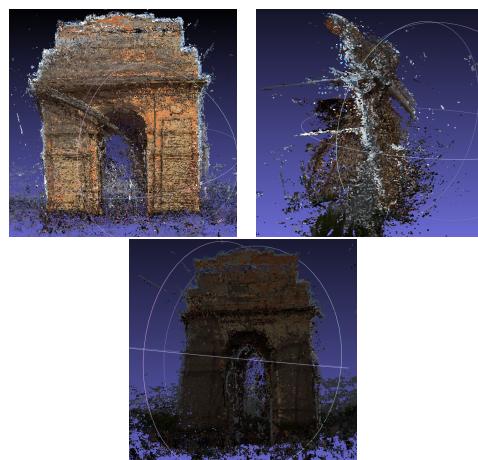
164 took screenshots from the satellite view using google
 165 earth.
 166

3. Reconstruction Models

3.1. Model 1

167 Figure 1 shows the result of opensfm pipeline on a set of
 168 image dataset of set ~ 100 from google and flickr. As we
 169 can clearly see from the figure, only the front facing view
 170 for India Gate is being reconstructed from the image.
 171

172 In these 100 images, we used all 4 view images of the
 173 monument, but still we were able to reconstruct only the
 174 front face because in most of the images includes the
 175 monument details and not its surroundings. Since India gate
 176 is symmetrical, we were getting only the front view of
 177 the monument because the front face and back face corre-
 178 sponded to the same depth point and our model couldn't
 179 distinguish between them.
 180



181 Figure 1. Results of Model 1
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195

3.2. Model 2

196 Lets discuss solution to the above problem. To break
 197 the symmetrical mappings for front-back and both sides, we
 198 can include the surrounding details for the image which will
 199 break the symmetry between the front and back views per se
 200 as we get more context for image points to differentiate be-
 201 tween both view points. We increased the data set to ~ 250
 202 to using screenshots of the 360 degree photos on YouTube
 203 to include surrounding regions of the monument as well.
 204 Figure 2 shows the reconstruction results for this model.
 205

206 What problems do we see in this model?
 207

- 208 1. Illumination variation in the reconstruction
 209
- 210 2. Lack of edges on the face of monument
 211
- 212 3. Sparse Reconstruction
 213
- 214
- 215



Figure 2. Results of Model 2

3.3. Model 3

The problems in the model above simply indicates problem with our data set. To overcome the illumination, we performed data set cleaning. This is one of the standard practice for any learning task, since a noisy input can affect the entire output. We cleaned those images of the monument, which were taking in the night conditions (since they had lighting in the monument) or had quite different illumination in the day time in comparison to the other images. Next, to resolve the lack of details and sparse reconstruction we increased the data set to ~ 500 , by extracting frames from the video and the result using the default parameters is shown in figure 3.



Figure 3. Results of Model 3

4. Scope of Improvements

Let us list the following problems in the above Model 3.

1. Sparse reconstruction of the top of model.

2. The green points along the edges of front of India Gate. 270
3. Noisy points on the top. 271
4. Lack of details like written India Gate on the top of 272 monument and the structure on the side view. 273

The above problems are shown in Figure 4.

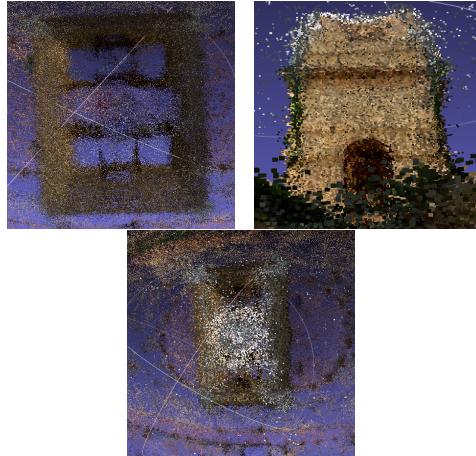


Figure 4. Sparse reconstruction, Green points on edges, Noisy top

4.1. Sparse reconstruction of the top of model

We thought that this problem exists due to lack of images from the top view from all of our sources, videos, internet scraping and 360 view photos. To overcome this caveat, we clicked some screenshots from google Earth for top view. Since, google earth gives satellite views, we were able to obtain ~ 10 images for the top view. We included those images the data set and performed reconstruction pipeline once again. Figure 5 shows slight improvement on the top view, but we were still not able to achieve much dense reconstruction of top view due to lack of images.

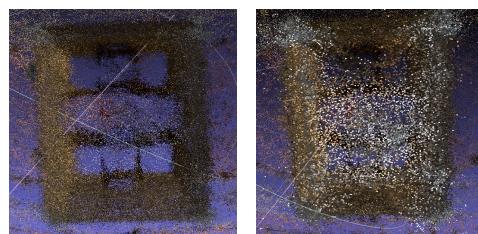


Figure 5. Sparse reconstruction, Dense reconstruction

4.2. Green points along the edges

We reasoned that this may be caused due to problems in feature matching and feature descriptors for images in the opensfm pipeline. So, this section we will experiment with different feature descriptors and matches.

324 **4.2.1 Feature Descriptors**
 325
 326 We experiment with 4 descriptors in the openSfM pipeline,
 327 SIFT, SURF, ORB and HAHOG.
 328
 329
 330
 331
 332
 333
 334
 335

336 Figure 6. With ORB as feature detector
 337

338 It can be from 6 seen that ORB mapped side and front
 339 to depths which couldn't preserve the monument geometry,
 340 next we use SIFT and the results are shown in figure 7(Note
 341 that we didn't mention SURF because it's outputs were simi-
 342 lar to SIFT, but we choose SIFT because it performed better
 343 in further improvements with its parameter tweaked). On
 344 the other hand HAHOG gives a sparse construction.
 345
 346 **4.2.2 Parameters for SIFT**
 347

348 Next, we tried to tune the parameters of the SIFT descrip-
 349 tors. We performed 3 experiments here. The SIFT detector
 350 is controlled mainly by two parameters: the peak thresh-
 351 old and the (non) edge threshold. The peak threshold filters
 352 peaks of the DoG scale space that are too small (in absolute
 353 value). Increasing it means more features, so may be better
 354 matching. The edge threshold eliminates peaks of the DoG
 355 scale space whose curvature is too small (such peaks yield
 356 badly localized frames). Decreasing it means more features,
 357 so may be better matching.
 358 The default parameters for SIFT in opensfm are:
 359 sift peak threshold: 0.01
 360 sift edge threshold: 10

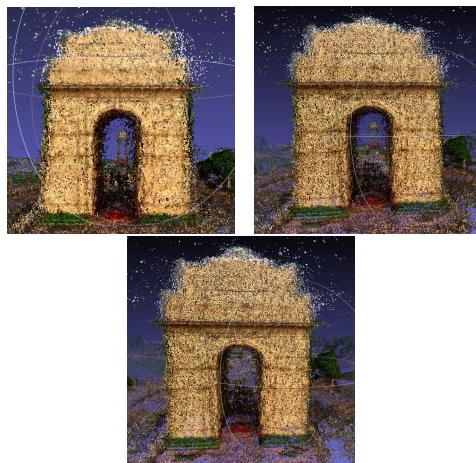
Peak threshold	edge threshold
0.01	10
0.001	10
0.01	100
0.001	100

367 It can be seen in Figure 7 that first one is better of the
 368 three. So, we fix the parameters of sift as

369 sift peak threshold: 0.001
 370 sift edge threshold: 10

373 **4.2.3 Feature Matchers**

374 For matching, we used FLANN matcher which is the de-
 375 fault parameter in opensfm. We tried using BoW matcher
 376 too, but since it only works with HAHOG descriptors and



392 Figure 6. With ORB as feature detector
 393

394 Figure 7. SIFT with (peak,edge) threshold as (0.001,10),
 395 (0.01,100) and (0.001,100) respectively.
 396

397 HAHOG descriptors gave bad results, we used matcher as
 398 FLANN. FLANN stands for Fast Library for Approximate
 399 Nearest Neighbors. It contains a collection of algorithms
 400 optimized for fast nearest neighbor search in large data sets
 401 and for high dimensional features. It works faster than BF-
 402 Matcher for large data sets. We will see the second example
 403 with FLANN based matcher. Also, the second reason for
 404 using this the time constraints, since we found out FLANN
 405 based matching was faster than BoW and of course brute
 406 force, we decided to go ahead with it.

407 **Where we able to solve the green point on the edge prob-**
 408 **lem?** No, but experimenting with this, we we able to get a
 409 denser reconstruction.

410 **4.3. Noisy points on the top**

411 The reason for noisy top will be the projection of noisy
 412 points in the 3d reconstruction. So, to overcome this we
 413 performed an experiment to change the bundle adjustment
 414 loss. A LossFunction is a scalar valued function that is used
 415 to reduce the influence of outliers on the solution of non-
 416 linear least squares problems. In addition to SoftOne loss
 417 we used to more losses to see if we can get rid of noisy
 418 points at the top of the model namely:

419 **Huber Loss** Huber loss is a loss function used in robust
 420 regression, that is less sensitive to outliers in data than the
 421 squared error loss.

$$\rho(s) = \begin{cases} s & s \leq 1 \\ 2\sqrt{s} - 1 & s > 1 \end{cases}$$

422 **Cauchy Loss**

$$\rho(s) = \log(1 + s)$$

423 Figure 8 shows that Huber loss performs better in
 424 comparison to the Cauchy loss and Softone loss that is the de-
 425 fault parameter. Cauchy loss performs similar to the Soft-
 426 one loss. Even though the density reduced but the noise is
 427

432 reduced. As, we can see those dense points were actually
 433 the noisy points.
 434

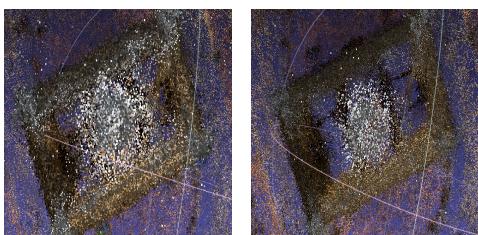


Figure 8. Before, After

445 4.4. Lack of details

446 The reason for noisy top will be the projection of noisy
 447 points in the 3d reconstruction. So, to overcome this we
 448 performed two experiments:

- 449 1. **Changing the triangulation threshold** We thought
 450 that lack of the details was due to sparsity of the
 451 construction. So, to overcome, we increased the trian-
 452 gulation threshold, which indicated how much error can
 453 a point have for it to be present in the construction.
 454 Unfortunately, it backfired which is clearly visible in
 455 Figure 9. The solid on one of the pillars is protruded.
 456



Figure 9. Before, After

- 466 2. **Adding more images corresponding to the detailed
 467 feature** To do this, we took the highest quality India
 468 Gate images, and obtained more images that corre-
 469 sponds to features like India Gate written on top, wheel
 470 and the side structure. These are shown in Figure 10.
 471 After running a reconstruction, we saw small refine-
 472 ments in the small objects. This is shown in Figure
 473 11.

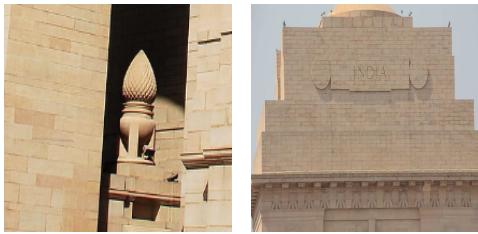


Figure 10. Detailed Featured of Monuments



Figure 11. Before, After

References

- [1] <https://www.opensfm.org>
- [2] <https://medium.com/@kazuyukimiyazawa/revisit-structure-from-motion-revisited-part-1-b1edde150a06>
- [3] <http://www.ceres-solver.org>
- [4] <https://blog.mapillary.com/update/2016/10/31/denser-3d-point-clouds-in-opensfm.html>
- [5] <https://github.com/rajkataria/ReliableResectioning>
- [6] <https://www.google.com>
- [7] <https://www.flickr.com>
- [8] <https://www.youtube.com>

486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539