
Slum Segmentation and Change Detection : A Deep Learning Approach

Shishira R Maiya*

Robert Bosch Center for Cyber Physical Systems
Indian Institute of Science
Bangalore, Karnataka 560054
shishirar@iisc.ac.in

Sudharshan Chandra Babu*

Department of Computer Science
Indian Institute of Technology Bombay
Mumbai, Maharashtra 400076
cbsudu@gmail.com

Abstract

More than one billion people live in slums around the world. In some developing countries, slum residents make up for more than half of the population and lack reliable sanitation services, clean water, electricity, other basic services. Thus, slum rehabilitation and improvement is an important global challenge, and a significant amount of effort and resources have been put into this endeavor. These initiatives rely heavily on slum mapping and monitoring, and it is essential to have robust and efficient methods for mapping and monitoring existing slum settlements. In this work, we introduce an approach to segment and map individual slums from satellite imagery, leveraging regional convolutional neural networks for instance segmentation using transfer learning. In addition, we also introduce a method to perform change detection and monitor slum change over time. We show that our approach effectively learns slum shape and appearance, and demonstrates strong quantitative results, resulting in a maximum AP of 80.0.

1 Introduction

Currently, about one-quarter of the world’s urban population live in slums [1]. These slum residents lack basic resources such as clean water, proper sanitation, electricity, and other necessary basic services. Various initiatives have been undertaken by international organizations and world governments in the past few decades towards slum improvement and rehabilitation. These initiatives rely heavily on the information provided by slum mapping and monitoring, such as scale, boundaries and slum growth, crucial for slum policy planning and development. Thus it is essential to have automated, robust, and efficient methods for slum mapping and monitoring. Slums differ greatly in terms of shape and appearance. Current approaches to slum segmentation and change detection [2–4] are limited and do not adapt very well to the variance in shape and texture. In addition, these approaches classify all detected instances of slums as one entity, given a satellite image, and do not recognize individual slums in a satellite image, which is essential for developing rehabilitation strategies for individual slums.

We concentrate our work on the slums in Mumbai–Dharavi, The Mankhurd-Govandi belt, Kurla-Ghatkopar belt, Dindoshi and The Bhandup-Mulund slums. The number of slum-dwellers in Mumbai is estimated to be around 9 million, up from 6 million in 2001 that is, 62% of of Mumbai live in informal slums [5].

In our work we introduce the following contributions:

1. We propose an instance segmentation based approach to the problem of slum mapping, that recognizes each slum in a given image, leveraging transfer learning, without the need for a

*Joint First authors.

large dataset. Our approach is based on the Mask R-CNN [6] framework, and automatically recognizes and segments individual slums from satellite imagery. In order to study this, we curate a custom dataset, consisting of satellite images of slums along with their polygon masks.

2. We introduce a method to monitor slum size increase or decrease over time and perform change detection.

We show that our method effectively identifies individual slums and demonstrates good qualitative and quantitative results.

2 Related Work

Early work on slum segmentation and mapping, include those based on object-based image analysis (OBIA) and texture-based methods [7–9]. In texture-based methods, the co-occurrence matrix (GLCM) is commonly used [2–4]. They’re limited by the fact that the extraction of a specific feature depends on the technique used. In addition, they have parameters that need to be optimized through trial and error [10–12]. In addition, these approaches classify all the slums as one entity in a single satellite image, and do not identify individual instances of slums.

Deep convolutional networks for segmentation have demonstrated strong results over other approaches [13], and can work well in the problem of slum mapping. Change detection is the process of identifying differences by observing images at different times [14–16]. Change detection on satellite imagery can provide important insights on urban development, and growth of informal and formal settlements [15].

Our approach identifies individual instances of slum in a given image using the Mask R-CNN model and we show that it has sufficient capacity to learn the visual and spatial features about slum settlements in satellite imagery. The model captures the inherent visual distribution sufficiently, and overcomes the above limitations of other approaches. Our approach to change detection is straightforward and builds upon other popular approaches [15, 16].

3 Approach

3.1 Dataset

We curated a dataset containing 3-band (RGB) satellite imagery with 65 cm per pixel resolution collected from Google Earth. Each image has a pixel size of 1280x720. The dataset consists of satellite images in two scales–100 m and 1000 m. These two scales provide different features that are used in small-scale and large-scale slum analysis respectively. The satellite imagery covers most of Mumbai and we include images from 2002 to 2018, to analyze slum change. Variability in resolutions of older images exist, due to the difference in satellites. Each image contains slum(s) along with formal settlements and vegetation. All images in the dataset have a paired list of polygons that describes sum instances. To verify our annotations, we used data provided by the Slum Rehabilitation Authority of Mumbai (SRA). We used 513 images for training, and 97 images for testing, for each scale. An example image at the 100 m scale and it’s ground truth is depicted in Figure 1.



Figure 1: Example ground truth image (left) and segmentation mask (right) at a scale of 100 m. Map data © 2018 Google, DigitalGlobe.

3.2 Slum Segmentation

Our approach to slum segmentation is based on Mask R-CNN, a powerful and flexible instance segmentation model. The Mask R-CNN model consists of two stages—the first stage scans the image and generates region based proposals, and the second stage classifies the proposals and generates bounding boxes and masks. The Mask R-CNN architecture we use is based on a ResNet-101 [17] and a Feature Pyramid Network [18] backbone. We train the Mask R-CNN by optimizing the following multi-task loss function, which combines the classification, localization and segmentation mask loss.

$$\mathcal{L} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{box}} + \mathcal{L}_{\text{mask}}$$

where \mathcal{L}_{cls} and \mathcal{L}_{box} are the classification loss and bounding box regression loss respectively. The mask branch of the network generates a $m \times m$ dimensional mask for each Region of interest (RoI) and for each class, with K classes in total. Thus the resulting output tensor size is $K \cdot m^2$. $\mathcal{L}_{\text{mask}}$ is the average binary cross-entropy loss, which includes the k -th mask in the region is mapped with the ground truth class k .

$$\mathcal{L}_{\text{mask}} = -\frac{1}{m^2} \sum_{1 \leq i, j \leq m} [y_{ij} \log \hat{y}_{ij}^k + (1 - y_{ij}) \log(1 - \hat{y}_{ij}^k)]$$

The input to the model is a satellite image, and the outputs are the bounding boxes, predicted masks and the confidence score. We train two models, one for each scale. We leverage transfer learning and pre-train the network on the COCO dataset [13].

3.3 Slum change detection

The input consists of a pair of satellite images, representing the same location, but at different points of time. For detecting change in the size of slums, we follow a two stage approach—We first pass both the images through the Mask R-CNN and predict masks for each image. We then subtract the binary masks and obtain a percentage increase or decrease.

3.4 Training details

We used an open source implementation of Mask R-CNN [19]. We fine-tuned on the pre-trained Mask R-CNN network for 128 epochs, with a batch size of 2. We used the Adam optimizer [20] with an initial learning rate of 10^{-4} , and we decayed the learning rate at 50 and 120 epochs by a factor of 10. We padded and resized each image to 1024x1024, without changing the aspect ratio. We performed data augmentation by horizontal and vertical flipping, rotation, translation, and variations in hue and saturation. We trained the model for around 4 hours on a Nvidia 1080 Ti GPU. We trained two models for each scale, with minor changes in the training process.

4 Results

4.1 Slum segmentation

We evaluate our approach on two metrics – IoU (Intersection over Union) and AP₅₀ (Average Precision at 50% overlap), commonly used in segmentation. We demonstrate strong quantitative and qualitative results on these metrics on the test set, and recognize individual slum instances, as illustrated in Table 1 and in Figure 2. We find that the model performs better on the images from 100 m scale, due to the larger number of visual features (slum huts, boundaries etc) available in the image. We also discover that the model shows satisfactory results on certain slum regions. These instances are quite heterogeneous and contain multiple buildings and vegetation within the slum.

4.2 Slum change detection

We predict the percentage change of the slums in a given satellite image. Figure 3 shows the change on a test image between 2005 and 2018. (35.25% change)



Figure 2: **Top:** Results on 100 m scale. **Bottom:** Results on 1000 m scale. It can be observed that the 1000 m images have less learnable visual features. Both the rows contain input images from the test set (left), Ground truth masks (center), and predicted masks (right). Map Data © 2018 Google, DigitalGlobe.

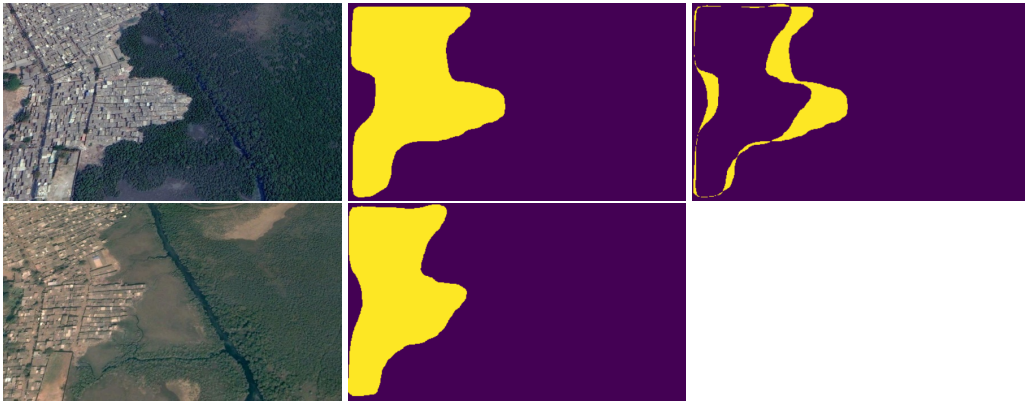


Figure 3: 35.25% change between 2018 (top left) and 2005 (bottom left). Predicted masks (middle) and change subtraction map (top right). Map Data © 2004 Google, Copernicus
Map Data © 2018 Google, DigitalGlobe.

Table 1: Results on test set and and slum-wise analysis

Slums	100 m		1000 m	
	IoU	AP ₅₀	IoU	AP ₅₀
Test Set	0.86	80.2	0.73	38.3
Govandi	0.89	60.3	0.80	59.2
Bhandup	0.88	95.2	0.78	75.9
Dharavi	0.90	75.4	0.67	15.5

5 Conclusion

In this work, we present an instance segmentation based approach to address the problem of slum mapping and monitoring. We show that our method achieves strong performance on these tasks. We hope that our work will be useful to organizations and other entities in their slum improvement and rehabilitation initiatives.

References

- [1] UN-Habitat, *Informal Settlements*, p. 1–8. UN-Habitat: New York, NY, USA, 2015.

- [2] M. Kuffer, K. Pfeffer, R. Sliuzas, and I. Baud, "Extraction of slum areas from vhr imagery using glcm variance," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 9, pp. 1830–1840, 2016.
- [3] S. Eckert, "Urban expansion and its impact on urban agriculture – remote sensing based change analysis of kizinga and mzingu valley – dar es salaam , tanzania," 2011.
- [4] S. Kabir, D.-C. He, M. A. Sanusi, and W. M. A. W. Hussina, "Texture analysis of ikonos satellite imagery for urban land use and land cover classification," *The Imaging Science Journal*, vol. 58, no. 3, pp. 163–170, 2010.
- [5] J. Bhavika, "62% of mumbai lives in slums: Census," *Hindustan Times*, 2010.
- [6] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask r-cnn," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [7] M. Kuffer, K. Pfeffer, and R. Sliuzas, "Slums from space - 15 years of slum mapping using remote sensing," *Remote Sensing*, vol. 8, p. 455, 2016.
- [8] D. Kohli, P. Warwadekar, N. Kerle, R. Sliuzas, and A. Stein, "Transferability of object-oriented image analysis methods for slum identification," *Remote Sensing*, vol. 5, pp. 4209–4228, 2013.
- [9] P. Hofmann, T. Blaschke, and J. Strobl, "Quantifying the robustness of fuzzy rule sets in object-based image analysis," *International Journal of Remote Sensing*, vol. 32, no. 22, pp. 7359–7381, 2011.
- [10] M. Fauvel, Y. Tarabalka, J. A. Benediktsson, J. Chanussot, and J. C. Tilton, "Advances in spectral-spatial classification of hyperspectral images," *Proceedings of the IEEE*, vol. 101, pp. 652–675, March 2013.
- [11] F. Dell'Acqua, M. Stasolla, and P. Gamba, "Unstructured human settlement mapping with sar sensors," *2006 IEEE International Symposium on Geoscience and Remote Sensing*, pp. 3619–3622, 2006.
- [12] X. Huang, H. Liu, and L. Zhang, "Spatiotemporal detection and analysis of urban villages in mega city regions of china using high-resolution remotely sensed imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, pp. 3639–3657, July 2015.
- [13] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.
- [14] L. Bruzzone and F. Bovolo, "A novel framework for the design of change-detection systems for very-high-resolution remote sensing images," *Proceedings of the IEEE*, vol. 101, pp. 609–630, 2013.
- [15] P. C. C. Author, I. Jonckheere, K. Nackaerts, B. Muys, and E. Lambin, "Review article digital change detection methods in ecosystem monitoring: a review," *International Journal of Remote Sensing*, vol. 25, no. 9, pp. 1565–1596, 2004.
- [16] X. Wang, S. Liu, P. Du, H. Liang, J. Xia, and Y. Li, "Object-based change detection in urban areas from high spatial resolution images based on multiple features and ensemble learning," *Remote Sensing*, vol. 10, no. 2, 2018.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [18] T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 936–944, 2017.
- [19] W. Abdulla, "Mask r-cnn for object detection and instance segmentation on keras and tensorflow." https://github.com/matterport/Mask_RCNN, 2017.
- [20] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014.