
Instance-sensitive Image Segmentation with Mask-RCNN and CRF

Haodong Zhou

Dept. of Computer Science
haodoz1@uci.edu

Tzu-Chi Lin

Dept. of Computer Science
tzucl2@uci.edu

Abstract

1 Image understanding problems such as object detection and semantic segmentation
2 have made breakthroughs in recent years due to the development of deep learning
3 techniques. We aim to provide a method to tackle with an instance segmentation
4 problem that can be done in pixel level. Our model used Mask-RCNN incorporates
5 with a conditional random field to further refine the result to the pixel level[4].

6 1 Introduction

7 Instance segmentation is one of a fundamental problem in computer vision that combines object
8 detection and semantic segmentation. Semantic segmentation labels every pixel to its object class,
9 without distinguishing different objects. Object detection does localize the object in a coarse,
10 bounding-box level. These two problems are well-studied topics in scene understanding and have
11 recently risen due to deep learning. The rising of autonomous driving and robotics industry also
12 made this problem more important, since the accurate recognition is one of the main tasks in these
13 two industries.

14 There are two different approaches to instance segmentation. First, based on object detection pipelines,
15 where objects are first localized to different boxes, then applies the semantic segmentation method to
16 further refine the object in each box. Another result is based on segment-based but rather box-based
17 method. However, these methods do not consider the whole image, but rather independent proposals.
18 Moreover, since they process numerous proposal independently, it is hard to produce segmentation
19 maps of the image.

20 We provide a model which is based on the fully convolution neural network, and corporate the
21 conditional random field to refine the result. Fully convolution neural network has been proved
22 an efficient model to segmentation problem. We use CRF as a post-processing method, which has
23 successfully pursued a joint learning method with CNN. We aim to use mean-field inference to
24 convert the system into an end-to-end network, so we can use the convolutional layer as features to
25 approximate CRF mean field inference.

26 2 Related work

27 Convolution Networks are driving advanced in segmentation. Traditional convolution networks solve
28 for whole-image tasks like classification, and for local tasks like bounding box object detection.
29 Because of the lack of the position information of the last layer in the network, it's hard to make a
30 pixel level prediction of an image.

31 The fully convolutional network (FCN) trains an end-to-end, pixels-to-pixels network, which does
32 a pixelwise prediction from supervised pre-trainning [4] . To improve the result of the convolution
33 network, a conditional random field (CRF) is used to improve the marginal performance by Chen, L.
34 C. et al [1]. Chen, L. C. [2, 3] claims the performance of semantic segmentation can be improved



Figure 1: Workflow of Instance Segmentation System

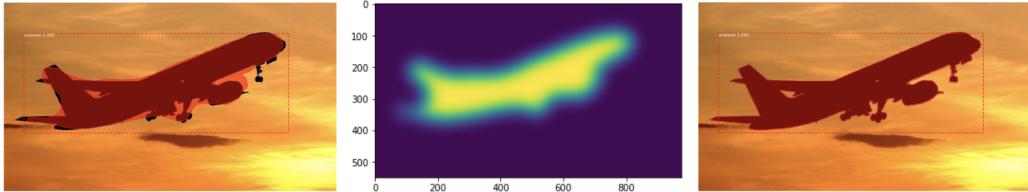


Figure 2: Left: result from Mask-RCNN. Middle: result generated by smoothing the mask, as the input to the CRF, which helps to improve the border performance. Right: result from the CRF refinement.

35 by using the atrous spatial pyramid pooling (ASPP) and using batch normalization within each
 36 ASPP. Also, an Encoder-Decoder with atrous convolution model was introduced to improve the
 37 performance.

38 However, these works mainly focus on image semantic segmentation. For instance-sensitive seg-
 39 mentation, Dai, J. et al [5] proposed a method that uses different score map to distinguish different
 40 instance. Li, Y. et al [6] helps to improve the performance by training a region proposal network
 41 (RPN) in the model, which proposes region-of-interests on the score map for joint object segmentation
 42 and detection.

43 3 Proposed approach

44 3.1 Pipeline

45 Figure 1 and 2 shows a simple pipeline of our model. We aim to combine Mask-RCNN and CRF to
 46 get a refined result. Given an image as input, firstly we use Mask-RCNN to generate a bounding box
 47 and a mask for the image. Considering that the border of the mask is improvable, we use a Gaussian
 48 kernel to smooth the mask, which makes the border of the mask to be the state of "improvable" while
 49 keeps the center of the mask a high probability value.

50 The probability value generated by the smoothing kernel is then consider the unary potential of the
 51 CRF model. We further define the pairwise potential function as describe in section 3.3 and use mean
 52 field approximation to update maximum a posteriori probability (MAP) iteratively as described in
 53 section 3.4.

54 3.2 Mask-RCNN

55 Mask-RCNN is a kind of network that achieve three functions: it can generate a bounding box,
 56 showing the position of the object, generate a class number and a possibility value indicating the type
 57 of the object and the probability, and the mask of the object. The Mask-rcnn shows a good result in
 58 the instance segmentation task.

59 The Faster-RCNN [8] is a network that can do multi-object detection and classification. The basic
 60 work flow of the Faster-RCNN is going through a CNN, using a RPN branch to generate ROIs(Region
 61 of Interest), sending the ROIs into a ROI Pooling layer and a classifier to get the bounding box and a
 62 classification result.

63 The Mask-RCNN [7] is based on the Faster-RCNN. It uses a FCN branch to generate the mask of the
 64 instance. Also, it uses the ROI Align layer instead of the ROI Pooling, which solves the problem
 65 caused by same feature size generated by different size ROIs. This will cause a low precision of the

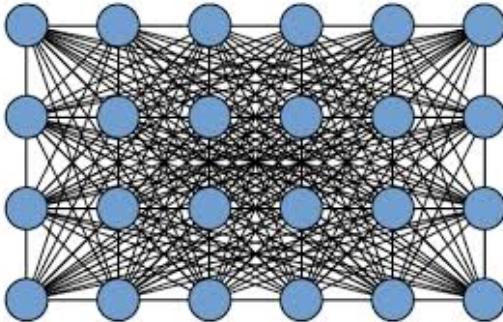


Figure 3: Example of fully connected CRF

66 mask. Mask-RCNN is a instance sensitive segmentation, and it has the state-of-the-art result, which
67 is used as a baseline in our project.

68 3.3 Fully Connected Conditional Random Field Model

69 We use the fully connected CRF model to further refine the result, which is shown in Figure 3. We
70 can see each pixel is represented by each node in CRF model. Consider a random field \mathbf{X} defined
71 over a set of variables $\{X_1, \dots, X_n\}$ which ranges over possible pixel-level image labeling. Suppose
72 there are k label, then the domain of each variable in \mathbf{X} is a set of labels $L = \{l_1, \dots, l_k\}$. Consider
73 another random field \mathbf{I} defined over $\{I_1, \dots, I_n\}$, which indicates the color vector of the pixel.

74 Now we have two conditional random fields to consider, joint these two CRFs to get a new CRF
75 (\mathbf{X}, \mathbf{I}), which can be represented by Gibbs distribution $P(X|I) = \frac{1}{Z(I)} \exp(-\sum \phi_c(X_c|I))$. Also,
76 the Gibbs energy can be expressed as $E(x|I) = \sum \phi_c(x_c|I)$, which as our unary potential function
77 in fully connected CRF. The maximum a-posterior (MAP) is $x^* = \text{argmax}_{x \in L^N} P(x|I)$, which is
78 the distribution we want to compute.

79 The Gibbs energy in the fully connected CRF model can be represented as

$$80 E(x) = \sum_i \psi_u(x_i) + \sum_{i \neq j} \psi_p(x_i, x_j)$$

81 The unary potential $\psi_u(x_i)$ is computed independently for each pixel by mask-RCNN. Since the
82 unary potential is produced independently, we could consider pairwise potential $\psi_p(x_i, x_j)$ to make
83 the image more consistent. The pairwise potential has the form

$$84 \psi_p(x_i, x_j) = \mu(x_i, x_j) \sum_{m=1}^K w^{(m)} k^{(m)}(f_i, f_j)$$

85 where μ is the label compatibility function, $k^{(m)}(f_i, f_j)$ is a Gaussian kernel and $w^{(m)}$ is its weight
86 of linear combination. We could express the kernel function as two Gaussian kernel, considered in
87 both color vectors and position

$$88 k(f_i, f_j) = w^{(1)} \exp\left(-\frac{|p_i - p_j|^2}{2\theta_\alpha^2} - \frac{|I_i - I_j|^2}{2\theta_\beta^2}\right) + w^{(2)} \exp\left(-\frac{|p_i - p_j|^2}{2\theta_\gamma^2}\right)$$

89 where \mathbf{I} indicates color vector and \mathbf{p} is the position of the pixel. The first Gaussian kernel is called
90 *appearance kernel*, which is inspired by the nature of the object - nearby pixels with similar color are
91 likely to be the same object. Later kernel is called *smoothness kernel*, which aims to remove isolated
92 region.

93 3.4 Mean Field Approximation

94 The algorithm to compute MAP is based on the mean field approximation to the CRF distribution.
95 We could see this approximation as a message passing algorithm and run it iteratively to get the
96 approximation inference.

97 Assume the KL-divergence we want to compute is $D(Q||P)$. We can derive this Q as

98
$$Q_i(x_i = l) = \frac{1}{z_i} \exp\{-\psi_\mu(x_i) - \sum_{l' \in L} \mu(l, l') \sum_{m=1}^K w^{(m)} \sum_{j \neq i} k^{(m)}(f_i, f_j) Q_j(l')\}$$

99 A detailed derivation is given in Appendix. We can decompose this formula into three steps. First, We
100 can see $\sum_{m=1}^K w^{(m)} \sum_{j \neq i} k^{(m)}(f_i, f_j) Q_j(l')$ as a message passing from all X_i to X_j . $\sum_{l' \in L} \mu(l, l')$
101 is a compatibility transformation. And we update $\exp\{-\psi_\mu(x_i) - Q'\}$ locally at last step.

102 For this message passing step, it takes quadratic time. We can improve this step to linear time by
103 the observation this message passing step could be represented as a Gaussian filter in feature space
104 [9]. We then only need to do a sampling and performs low-pass filter on this sample, which can be
105 performed in O(N) time.

106 4 Implementation Details

107 Our project is based on the tensorflow. The Mask-RCNN part we used the implementation of the
108 matterport, and the model is a pre-trained model ¹, which already performs a good result. The output
109 we get from the Mask-RCNN is a group of bounding box with a binary mask inside.

110 Then we apply CRF to masks one by one. We use a adaptive size Gaussian kernel to smooth the
111 mask. Then the result will be sent into the CRF function.

112 We adaptive our code to denseCRF code ². Since we already have the unary potential from Mask-
113 RCNN, we only have to create two Gaussian kernels for pairwise potential. There are some parameters
114 we could tune on. For smoothness kernel parameter θ_γ , it does not affect the result significantly.
115 We just set $\theta_\gamma = 1$ and found it work well. Unfortunately, the appearance kernel parameters θ_α
116 and θ_β cannot be computed effectively by gradient decent, since their gradient involves a sum of
117 non-Gaussian kernels, which are not amenable to the same acceleration techniques. We simply use a
118 grid search to find θ_α and θ_β .

119 After setting all the parameters in kernels, we need to run inference algorithm to compute MAP
120 and assign each pixel by MAP. These steps could all be done by denseCRF API. Since we already
121 found out 10 iterations may converge for KL-divergence as described in section 5.1, we set all other
122 experiments to run 10 times iterations.

123 5 Experiments

124 5.1 Convergence of mean field approximation

125 We test the convergence of the mean field approximation first by analyzing KL-divergence between Q
126 and P. Figure 4 shows the KL-divergence over iterations of the inference algorithm. We could see that
127 it converges after 10 iterations. Hence, we set 10 as the number of iterations in all other experiments.

128 5.2 Grid search for parameters

129 Figure 6 and 7 shows a simple result for grid search of parameters θ_α and θ_β . We can see that for
130 $\theta_\alpha = 60$ and $\theta_\beta = 10$ is good for our result. We then set $\theta_\alpha = 60$ and $\theta_\beta = 10$ to the following
131 experiment.

132 5.3 Evaluation method

133 The performance is evaluated by intersection-over-union (IOU) between the predicted instances and
134 the ground-truth instances. We report the standard COCO metrics including AP (averaged over IoU
135 thresholds), AP_{50} , AP_{75} , and AP_S , AP_M , AP_L (AP at different scales), which is the same method
136 with the Mask-RCNN

¹github.com/matterport/Mask_RCNN/

²github.com/lucasb-eyer/pydensecrf/

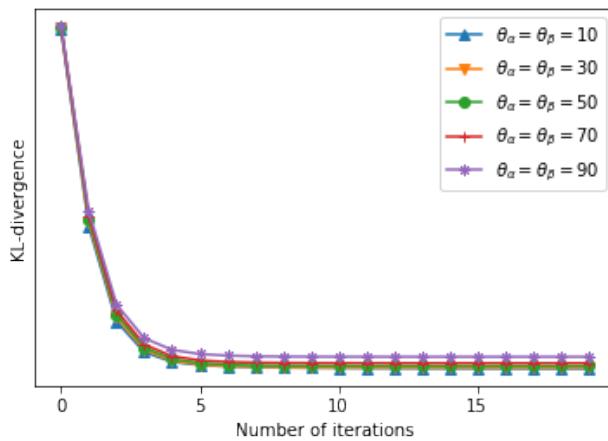


Figure 4: Convergence of KL-divergence

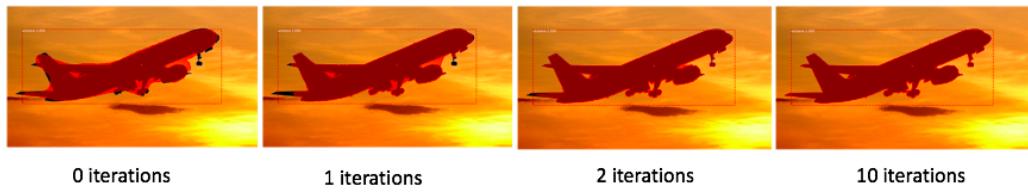


Figure 5: Iteration result of the image

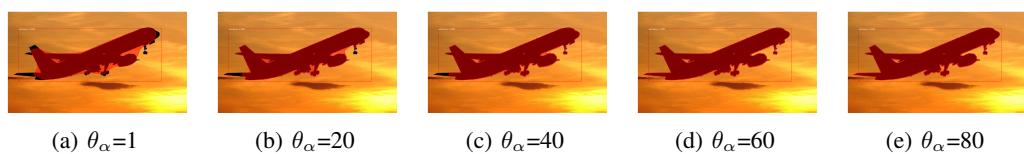


Figure 6: Grid search for θ_α

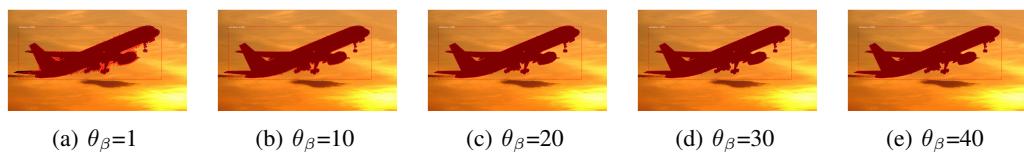


Figure 7: Grid search for θ_β

Table 1: Instance Segmentation Result on COCO2014 Dataset

| method | AP | AP_{50} | AP_{75} | AP_S | AP_M | AP_L |
|---------------------|------|-----------|-----------|--------|--------|--------|
| Mask R-CNN | 35.7 | 58.0 | 37.8 | 15.5 | 38.1 | 52.4 |
| Mask R-CNN with CRF | 35.2 | 58.7 | 36.4 | 17.3 | 36.5 | 49.8 |

137 5.4 Result and Analysis

138 The result based on the evaluation method is in the Table 1. Our result overall is similar with the
 139 result of the Mask R-CNN. Parameters is pretty hard to automatically adapt for each mask, and the
 140 current result is based on the parameters which makes small objects performs better than original
 141 Mask-RCNN.

142 Some other results with specified tuned parameters are showed in the Figure 8, 9 and 10 . We can
 143 find that after 10 iterations of CRF, the border of the mask is improved. However, to achieve these
 144 result, we increase the weight of the color, which causes that the part of the bird which has the similar
 145 color with the background color to has poor performance.

146 6 Future Work

147 To further improve the performance, we may consider to put CRF model to end-to-end model so that
 148 we could jointly optimize the parameters of the CNN and the CRF, such as CRF as RNN model[10].

149 7 Appendix

150 7.1 Detailed Derivation of Q

151 We give a detailed derivation of Q in this section.

152 First we have Gibbs distribution of denseCRF

$$153 P(X) = \frac{1}{Z} P'(X) = \frac{1}{Z} \exp(\sum_i \psi_u(x_i) + \sum_{i \neq j} \psi_p(x_i, x_j))$$

154 We can derive KL-divergence as following

$$\begin{aligned} 155 D(Q||P) &= \sum_x Q(x) \log\left(\frac{Q(x)}{P(x)}\right) = -\sum_x Q(x) \log P(x) + \sum Q(x) \log Q(x) \\ 156 &= -E_{X \in Q}[\log P(X)] + E_{X \in Q}[\log Q(x)] \\ 157 &= -E_{X \in Q}[\log P'(X)] + E_{X \in Q}[\log Z] + \sum_i E_{X_i \in Q_i}[\log Q_i(X_i)] \\ 158 &= -E_{X \in Q}[\log P'(X)] + \log Z + \sum_i E_{X_i \in Q_i}[\log Q_i(X_i)] \end{aligned}$$

159 Since our goal is to compute Q and $\log Z$ do not have Q, we ignore this term.

160 At the same time, Q have to satisfy

$$161 \sum_{x_i} Q_i x_i = 1$$

162 By Lagrange Multiplier, we can get

$$163 L(Q_i) = -E_{X \in Q}[\log P'(X)] + \sum_i E_{X_i \in Q_i}[\log Q_i(X_i)] + \lambda(\sum_{x_i} Q_i(x_i) - 1)$$

164 By partial differential we can get

$$165 \frac{\partial L(Q_i)}{\partial Q_i(x_i)} = -E_{X \in Q}[\log P'(X|x_i)] - \log Q_i(x_i) - 1 + \lambda$$

166 Set this to 0, we can get

$$167 Q_i(x_i) = \exp(\lambda - 1) \exp(-E_{X \in Q}[\log P'(X|x_i)])$$

168 We can see $\exp(\lambda - 1)$ as a constant number, so Q becomes

169 $Q_i(x_i) = \frac{1}{Z_i} \exp(-E_{\bar{X} \in Q_i} [\log P'(X|x_i)])$

170 where Z is a normalization term.

171 We put P' to this formula as we define at the beginning.

172 $Q_i(x_i) = \frac{1}{Z_i} \exp(-E_{\bar{X} \in Q_i} [\sum_i \psi_u(x_i) + \sum_{i \neq j} \psi_p(x_i, x_j)|x_i])$

173 Since x_i is known as label l, we now have

174 $Q_i(x_i = l) = \frac{1}{Z_i} \exp(-\psi_\mu(l) - \sum_{j \neq i} E_{\bar{X} \in Q_j} \psi_p(l, X_j))$
 175 $= \frac{1}{Z_i} \exp(-\psi_\mu(l) - \sum_{m=1}^K w^{(m)} \sum_{j \neq i} E_{X \in Q_j} [\mu(l, X_j) k^{(m)}(f_i, f_j)])$
 176 $= \frac{1}{Z_i} \exp(-\psi_\mu(l) - \sum_{j \neq i} \sum_{l' \in L} k^{(m)}(f_i, f_j) \mu(l, l') Q_j(l'))$
 177 $= \frac{1}{z_i} \exp\{-\psi_\mu(x_i) - \sum_{l' \in L} \mu(l, l') \sum_{m=1}^K w^{(m)} \sum_{j \neq i} k^{(m)}(f_i, f_j) Q_j(l')\}$

178 where we can the result of

179 $Q_i(x_i = l) = \frac{1}{z_i} \exp\{-\psi_\mu(x_i) - \sum_{l' \in L} \mu(l, l') \sum_{m=1}^K w^{(m)} \sum_{j \neq i} k^{(m)}(f_i, f_j) Q_j(l')\}$

180 7.2 Contribution of each member

181 Haodong Zhou is mainly responsible for Mask-RCNN, including its description in this paper and
 182 implementation of the code. Also, section 2 is written by him.

183 Tzu-Chi Lin is mainly responsible for CRF model, including its description in this paper and
 184 implementation of the code. Also, section 1 is written by him.

185 We together figure out how to combine Mask-RCNN and CRF together.

186 References

- 187 [1] Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2018) Deeplab: Semantic image
 188 segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, **40**(4), pp.834-848.
 189
- 190 [2] Chen, L.C., Papandreou, G., Schroff, F., & Adam, H. (2017) Rethinking atrous convolution for semantic
 191 image segmentation. *arXiv preprint*, arXiv:1706.05587.
 192
- 193 [3] Chen, L. C., Zhu, Y., Papandreou, G., Schroff, F., & Adam, H. (2017) Encoder-decoder with atrous separable
 193 convolution for semantic image segmentation. *arXiv preprint* arXiv:1802.02611.
 194
- 195 [4] Long, Jonathan, Evan Shelhamer, & Trevor Darrell. (2015) Fully convolutional networks for semantic
 195 segmentation. *Proceedings of the IEEE conference on computer vision and pattern recognition* pp.3431-3440.
 196
- 196 [5] Dai, J., He, K., Li, Y., Ren, S., & Sun, J. (2016) Instance-sensitive fully convolutional networks. *European Conference on Computer Vision* pp. 534-549.
 197
- 198 [6] Li, Y., Qi, H., Dai, J., Ji, X., & Wei, Y. (2016) Fully convolutional instance-aware semantic segmentation.
 199 *arXiv preprint* arXiv:1611.07709.
 200
- 201 [7] He, K., Gkioxari, G., Dollár, P. & Girshick, R., 2017, October. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on* (pp. 2980-2988). IEEE.
 202
- 203 [8] Ren, S., He, K., Girshick, R., & Sun, J., (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. *In Advances in neural information processing systems* (pp. 91-99).
 204
- 205 [9] Krahenbuhl, P. and Koltun, V. (2011) Efficient inference in fully connected crfs with gaussian edge potentials. In NIPS.
 206
- 207 [10] Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V., Su, Z., Du, D., Huang, C., and Torr, P. (2015)
 207 Conditional random fields as recurrent neural networks, in *IEEE ICCV*.



Figure 8: Example With Tuned Parameters



Figure 9: Example With Tuned Parameters



Figure 10: Example With Tuned Parameters