# MODEL DOCTOR FOR DIAGNOSING AND TREATING SEGMENTATION ERROR

*Zhijie Jia[†], Lin Chen[†], Kaiwen Hu[†], Lechao Cheng[††*], Zunlei Feng[†], Mingli Song[†]*

[†] Zhejiang University
[††] Zhejiang Lab

## ABSTRACT

Despite the remarkable progress in semantic segmentation tasks with the advancement of deep neural networks, existing U-shaped hierarchical typical segmentation networks still suffer from local misclassification of categories and inaccurate target boundaries. In an effort to alleviate this issue, we propose a Model Doctor for semantic segmentation problems. The Model Doctor is designed to diagnose the aforementioned problems in existing pre-trained models and treat them without introducing additional data, with the goal of refining the parameters to achieve better performance. Extensive experiments on several benchmark datasets demonstrate the effectiveness of our method. Code is available at `https://github.com/zhijiejia/SegDoctor`.

*Index Terms*— Semantic segmentation, Model treatment.

## 1. INTRODUCTION

Image segmentation [1, 2, 3] is a crucial task in the computer vision field, with a wide range of applications [4, 5], including scene understanding, video surveillance, medical image analysis, robotic perception, and so on.

However, the current mainstream semantic segmentation techniques focus on the structural design of deep convolutional neural networks, but ignore the treatment and utilization of existing semantic segmentation models. In addition, the black box [6] structure of deep neural networks also contributes to the lack of ability to analyze problems from segmentation results, making it challenging to target errors and fine-tune the semantic segmentation model. There are currently model-interpretable methods that can assist in better understanding and analyzing models. However, much of the focus has been on visualizing model prediction results through techniques such as Class Activation Mapping (CAM) [7], Grad-CAM [8], and Grad-CAM++ [9]. Through these methods, the patterns that the model prioritizes and the areas of input that the model pays more attention to can be identified. Additionally, some works utilize the interpretable random forests algorithm to dissect deep neural networks [10], and decouple deep neural models, which facilitates rapid identification of the source and location of model

*Corresponding author.



**Fig. 1**. Feature analysis of semantic segmentation model. The content of the red box represents the category error, and the content of the yellow box represents the boundary error.

errors. Nevertheless, these techniques cannot be applied directly and automatically to model treatment.

In the preliminary experiments, we find that errors in semantic segmentation models can generally be divided into two types: semantic category errors and regional boundary errors. Semantic category errors arise from the inclusion of feature errors in deep semantic features, resulting in category classification errors for certain regions. On the other hand, region boundary errors occur due to the lack of fine edge detail features in shallow texture features, resulting in lost boundary information.

In this paper, we introduce a Model Doctor to amend semantic category errors and regional boundary errors, respectively. As shown in Fig. 1, we apply semantic category treatment to deep semantic features extracted by deep neural networks to bridge the gap within classes in deep features and force intra-class features to converge to the category center. For regional boundary treatment, we constrain shallow texture features at various levels to enhance internal feature constraints on objects and preserve more edge detail features. Exhaustive experiments demonstrate that incorporating the proposed method with several semantic segmentation models leads to improved performance on commonly used datasets. Our contributions can be summarized as follows:

- We present a Model Doctor for diagnostic treatment segmentation models, which can be plugged into existing convolutional segmentation models.

- Semantic category treating strategy and region boundary treating strategy are designed to address semantic category errors and region boundary errors, respectively.

- Extensive experiments showcase that the proposed semantic segmentation model treating method can effectively boost the performance of existing semantic segmentation models.

## 2. RELATED WORK

Due to the complexity and ambiguity of deep neural networks, humans cannot give exact explanations for their behavior. At present, the interpretability methods of deep models are mainly divided into two categories [11]: Post-hoc interpretability analysis method and Ad-hoc interpretable modeling method. Post-hoc interpretability analysis method is an interpretable analysis of deep models that have been trained; Ad-hoc interpretable modeling method mainly builds deep models into interpretable models to ensure that the inferences of the models are interpretable. Post-hoc interpretability analysis method mainly include seven categories of techniques, such as feature analysis [12, 13, 14], model checking [6, 15], salient expression [16, 17], surrogate modeling [18], advanced mathematics analysis [19], case interpretation [20], and text interpretation [21]. Ad-hoc interpretable modeling method mainly includes two types of methods: interpretable representation [22] and model improvement [23]. However, the above methods mainly focus on model interpretation and cannot achieve automatic diagnosis and optimization of model defects. Recently, Feng et al. [24] proposed model doctor for the optimization of classified convolutional neural networks, but due to the difference between the segmentation model and the classification model architecture, this method cannot be applied to the semantic segmentation model.

## 3. METHOD

In this paper, we present a novel model therapeutic approach for semantic segmentation models, designed to address the inadequacies in semantic category classification and boundary refinement of these models.

### 3.1. Segmentation Error Diagnosis

In the preliminary experiments, we find that semantic segmentation models are prone to regional boundary problems and category classification problems, and different model problems are related to different feature errors.

#### 3.1.1. Semantic category error

The semantic segmentation model is typically composed of an encoder and a decoder, where the encoder is responsible for extracting image features and the decoder is responsible for restoring image edge details. Given an input image $I$, the output feature map of the last layer of the encoder is $M^e$, computed as $M^e = Encoder(I)$, where the shape of $M^e$ is $(N, C, H, W)$, where $N$ is bach size, $C$ is the number of channel and $(H, W)$ is the feature map size, and the vectors of each $(1, C, 1, 1)$ in $M^e$ correspond to a patch in the original image. The deep features extracted by the encoder $M^e$, possess a wealth of deep semantic information and semantic category information. The widening gap between the deep feature vectors signifies that the semantic category information of the corresponding patches is no longer equivalent, leading to subsequent classification errors.

#### 3.1.2. Regional boundary error

The extracted image features $\{M_1^e, M_2^e, M_3^e, ..., M_l^e, ..., M_L^e\}$ of the encoder exhibit distinct attributes at various depths, where $L$ is the maximum layer number in the encoder. While shallow image features $M_l^e$ are rich in edge detail information, they lack semantic intricacies; Conversely, deep image features $M_l^e$ are abundant in semantic information but deficient in edge detail. The extensive semantic features of $M_l^e$ enable the model to perform efficient class classification, whereas the edge details present in $M_l^e$ aid in partial reconstruction of the object's edge details by the decoder.

Hence, during the decoding phase, the shallow and deep feature maps $\{M_l^e\}_{l=1}^L$ are concatenated and processed by a convolutional function $\mathcal{F}_{conv}$, to produce the feature map $M_{i+1}^d$ of $i$-th layer as follows:

$$M_i^d = \mathbf{Concat}(M_l^e, M_{i-1}^d), l \in \{1, 2, 3, ..., L\}, \quad (1)$$

where the initial input feature map $M_1^d$ of encoder is the output feature map $M_L^e$ of the encoder. However, if the shallow feature $M_i^d$ of the decoder contains errors in shallow detail information, the model will miss crucial detail information during the upsampling process, rendering it insensitive to the object's edge area and incapable of producing fine-grained edge details of the object.

### 3.2. Segmentation Error Treatment

In light of the aforementioned observations, we have developed a segmentation error diagnosis and treating method that encompasses both semantic category correction and regional boundary rectification, aiming at addressing the classification and boundary errors of the semantic segmentation model.

#### 3.2.1. Treating Category error

Consequently, in order to mitigate the impact of semantic category errors on deep features, we devise a category constraint technique for treating semantic category error. It constrains

**Fig. 2**. The framework of the proposed method, which is comprised of two parts: the semantic category treatment applied to the deep features, and the regional boundary treatment applied to the shallow features.

the deep features of the model by minimizing inter-class variations and maximizing intra-class similarity. To achieve this, the cluster center $C_k$ for the $k$-th class in cluster $D_k$ is computed to represent the central tendency of features within each class and provides a basis for comparison with other feature vectors. The cluster center $C_k$ is calculated as follows:

$$\arg\min_{C_k} \sum_{R_k \in D_k} ||R_k - C_k||^2, \qquad (2)$$

where $R_k$ is the feature representation in cluster $D_k$.

In the context of deep features, the image feature for a given class $k$ is denoted as $R_k$. To alleviate semantic errors and improve the model's classification accuracy, a feature distance constraint is imposed to force the intra-class image features to gravitate towards the centroid of the class cluster $C_k$, which can mitigate intra-class feature divergence. The feature distance penalty $\zeta_{sim}$ is calculated as follows:

$$\zeta_{sim} = 1 - \mathcal{D}(C_k, R_k), \mathcal{D}(C_k, R_k) = \frac{C_k \cdot R_k}{||C_k|| \times ||R_k||}, \quad (3)$$

where '·' denotes vector multiplication, $\mathcal{D}(C_k, R_k)$ represents the feature distance between the feature representation $R_k$ and the cluster center $C_k$.

### 3.2.2. Treating boundary error

In accordance with the information presented in Section 3.1.2, if the shallow image features contain erroneous texture feature information, this can result in inaccuracies in the decoder's fine edge reconstruction. To address this, superpixel technology is incorporated as superpixel branch, which is a coarse segmentation method that helps preserve edge details and enforce consistency within shallow image features. The SpixelFCN algorithm proposed in [25] is a noteworthy implementation of superpixel segmentation that leverages a fully convolutional network to achieve rapid and remarkable results. In this work, we drew inspiration from SpixelFCN to devise the superpixel branch, aiming to preserve the shallow

texture features. The superpixel branch is assembled by a block consisting of three conv-bn-relu layers, which performs the upsampling operation and generates the link probability connecting the pixel to the neighboring superpixels.

For shallow feature map $M_l^e$, the superpixel branch $\mathcal{F}_{sup}$ predicts the probability of $p$ being associated with surrounding superpixels as follows:

$$p = \sigma(\mathcal{F}_{sup}(M_l^e)), \qquad (4)$$

where $\sigma(\cdot)$ represents the sigmoid function. Then the reconstruction of pixel feature $\mathbf{f}'(\cdot)$ and pixel coordinates $\mathbf{v}'$ are calculated as follows:

$$\mathbf{v}' = \sum_{s \in \mathcal{N}_\mathbf{v}} \frac{\sum_{\mathbf{v}:s \in \mathcal{N}_\mathbf{v}} \mathbf{v} \cdot p}{\sum_{\mathbf{v}:s \in \mathcal{N}_\mathbf{v}} p} \cdot p, \qquad (5)$$

$$\mathbf{f}'(\mathbf{v}) = \sum_{s \in \mathcal{N}_\mathbf{v}} \frac{\sum_{\mathbf{v}:s \in \mathcal{N}_\mathbf{v}} \mathbf{f}(\mathbf{v}) \cdot p}{\sum_{\mathbf{v}:s \in \mathcal{N}_\mathbf{v}} p} \cdot p, \qquad (6)$$

where $\mathbf{v} = [x, y]^T$ denotes the original pixel's position, and $\mathcal{N}_\mathbf{v}$ is the set of surrounding superpixels of $\mathbf{v}$. The penalty function of the superpixel branch is divided into two parts: feature constraint and coordinate constraint, which is specified as follows:

$$\zeta_{sp} = \sum_\mathbf{v} CE(\mathbf{f}(\mathbf{v}), \mathbf{f}'(\mathbf{v})) + \frac{m}{s}||\mathbf{v} - \mathbf{v}'||_2, \qquad (7)$$

where $\mathbf{f}(\cdot)$ represents one-hot encoding vector of semantic label. $s$ denotes the superpixel sampling interval, and $m$ is a weight-balancing term, and $CE(\cdot, \cdot)$ denotes Cross-Entropy.

Overall, with the shallow feature map $M_l^e$ as input, the first terms of $\zeta_{sp}$ encourages the trained superpixel branch $\mathcal{F}_{sup}$ to group pixels with similar category property, and the second term enforces the superpixels to be spatially compact.

### 3.3. Overview

Finally, the total loss function adopted is the original Cross-Entropy loss $\zeta_{ce}$ combined with the category error loss $\zeta_{sim}$

**Table 1**. The performance on different models and datasets.

| Dataset → | VOC 2012 | | Cityscapes | |
|---|---|---|---|---|
| Method ↓ | Origin | +Treatment | Origin | +Treatment |
| FPN | 61.7 | 62.5(**+0.8**) | 66.5 | 67.9(**+1.4**) |
| UNet | 54.0 | 55.2(**+1.2**) | 69.5 | 70.1(**+0.6**) |
| CCNet | 57.1 | 58.7(**+1.6**) | 70.8 | 72.0(**+1.2**) |
| PSPNet | 68.1 | 69.0(**+0.9**) | 72.8 | 74.1(**+1.3**) |
| Deeplab v3+ | 67.3 | 68.4(**+1.1**) | 74.2 | 74.9(**+0.7**) |

**Table 2**. The ablation study on different treating strategies.

| Method | mIoU |
|---|---|
| UNet | 54.0 |
| + Treating category | 54.4 (**+0.4**) |
| + Treating boundary | 54.8 (**+0.7**) |
| + Treating category & boundary | 55.2 (**+1.2**) |

and the boundary error loss $\zeta_{sp}$ as follows:

$$Loss = \zeta_{ce} + \alpha\zeta_{sim} + \beta\zeta_{sp}, \quad (8)$$

where $\alpha$ and $\beta$ denote the balance parameters.

## 4. EXPERIMENT

### 4.1. Dataset and Experiment setting

**Dataset.** Our experimental evaluation is performed on two publicly available datasets, namely the PASCAL VOC 2012 [1] dataset and Cityscapes dataset [2]. The PASCAL VOC 2012 dataset, a semantic segmentation dataset with 20 categories, comprises 10,582 images in its training set and 1,449 images in its validation set. The Cityscapes dataset, a driving dataset for panoramic segmentation with 19 categories, comprises 2,975 images in the training set and 500 images in the validation set.

**Experiment Setting.** During the training of models, we randomly crop images to $512 \times 512$ (VOC) and $512 \times 1024$ (Cityscapes) and utilize horizontal and vertical flipping augmentations. The batch size is set to 8 for all datasets, and the optimization is performed using Stochastic Gradient Descent (SGD). The initial learning rate is set at 0.01 and the cosine annealing rate decay policy is employed. The balance parameters are set as follows: $\alpha = 1$ and $\beta = 0.01$. The performance of the semantic segmentation is reported using the mean Intersection over Union (mIoU) metric.

### 4.2. Compatibility with Existing Segmentation Models

In the experiment, we adopt some mainstream segmentation network to verify the effectiveness of the proposed method. Results in Table 1 demonstrate that the proposed approach is



| Input | Original | + Treatment | GT |

**Fig. 3**. Visual results on PASCAL VOC 2012 dataset.

able to enhance the performance of different models on the PASCAL VOC 2012 dataset and the Cityscapes dataset.

### 4.3. Visual Results

We demonstrate the efficacy of the proposed method by incorporating it into the UNet network on the VOC 2012 dataset, resulting in improved semantic segmentation performance. As depicted in Fig. 3, our method produces more accurate and nuanced structures, as evidenced by several visualizations from the VOC 2012 validation set.

### 4.4. Ablation Study

In this section, we conduct the ablation study on two treatment strategies. The ablation study experiment is conducted with UNet on the VOC 2012 dataset. As shown in Table 2, the semantic category treatment strategy and regional boundary treatment strategy both effectively enhance the performance of the segmentation model.

## 5. CONCLUSION

In this paper, a new method called Model Docter is introduced to address semantic category errors and regional boundary errors in semantic segmentation. Semantic category treatment is applied to deep semantic features extracted by deep neural networks to reduce gaps within classes and correct misclassifications. Regional boundary treatment is imposed on shallow texture features to enhance internal feature constraints and preserve edge detail features. The proposed approach has been tested on several datasets and models and can be combined with other models for further refinement.

# 6. REFERENCES

[1] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman, "The pascal visual object classes challenge: A retrospective," *IJCV*, vol. 111, pp. 98–136, 2015.

[2] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele, "The cityscapes dataset for semantic urban scene understanding," in *CVPR*, 2016, pp. 3213–3223.

[3] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba, "Scene parsing through ade20k dataset," in *CVPR*, 2017, pp. 633–641.

[4] Shervin Minaee, Yuri Y Boykov, Fatih Porikli, Antonio J Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos, "Image segmentation using deep learning: A survey," *IEEE T-PAMI*, 2021.

[5] King-Sun Fu and JK Mui, "A survey on image segmentation," *PR*, vol. 13, no. 1, pp. 3–16, 1981.

[6] Pang Wei Koh and Percy Liang, "Understanding blackbox predictions via influence functions," in *ICML*. PMLR, 2017, pp. 1885–1894.

[7] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba, "Learning deep features for discriminative localization," in *CVPR*, 2016, pp. 2921–2929.

[8] Ramprasaath R Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra, "Grad-cam: Why did you say that?," *arXiv preprint arXiv:1611.07450*, 2016.

[9] Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian, "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks," in *WACV*. IEEE, 2018, pp. 839–847.

[10] Jie Lei, Zhe Wang, Zunlei Feng, Mingli Song, and Jiajun Bu, "Understanding the prediction process of deep networks by forests," in *BigMM*. IEEE, 2018, pp. 1–7.

[11] Feng-Lei Fan, Jinjun Xiong, Mengzhou Li, and Ge Wang, "On interpretability of artificial neural networks: A survey," *IEEE Transactions on Radiation and Plasma Medical Sciences*, vol. 5, no. 6, pp. 741–760, 2021.

[12] Alexey Dosovitskiy and Thomas Brox, "Inverting visual representations with convolutional networks," in *CVPR*, 2016, pp. 4829–4837.

[13] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent, "Visualizing higher-layer features of a deep network," *University of Montreal*, vol. 1341, no. 3, pp. 1, 2009.

[14] Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft, "Convergent learning: Do different neural networks learn the same representations?," *arXiv preprint arXiv:1511.07543*, 2015.

[15] Aayush Bansal, Ali Farhadi, and Devi Parikh, "Towards transparent systems: Semantic characterization of failure modes," in *ECCV*. Springer, 2014, pp. 366–381.

[16] Jerome H Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of statistics*, pp. 1189–1232, 2001.

[17] Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin, "Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation," *JCGS*, vol. 24, no. 1, pp. 44–65, 2015.

[18] R Krishnan, G Sivakumar, and P Bhattacharya, "Extracting decision trees from trained neural networks," *PR*, vol. 32, no. 12, 1999.

[19] Yiping Lu, Aoxiao Zhong, Quanzheng Li, and Bin Dong, "Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations," in *ICML*. PMLR, 2018, pp. 3276–3285.

[20] Janet L Kolodner, "An introduction to case-based reasoning," *Artificial intelligence review*, vol. 6, no. 1, pp. 3–34, 1992.

[21] Andrej Karpathy and Li Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *CVPR*, 2015, pp. 3128–3137.

[22] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," *Advances in neural information processing systems*, vol. 29, 2016.

[23] Lingyang Chu, Xia Hu, Juhua Hu, Lanjun Wang, and Jian Pei, "Exact and consistent interpretation for piecewise linear neural networks: A closed form solution," in *SIGKDD*, 2018, pp. 1244–1253.

[24] Zunlei Feng, Jiacong Hu, Sai Wu, Xiaotian Yu, Jie Song, and Mingli Song, "Model doctor: A simple gradient aggregation strategy for diagnosing and treating cnn classifiers," in *AAAI*, 2022, vol. 36, pp. 616–624.

[25] Fengting Yang, Qian Sun, Hailin Jin, and Zihan Zhou, "Superpixel segmentation with fully convolutional networks," in *CVPR*, 2020, pp. 13964–13973.