

Joint Modeling of Image and Label Statistics for Enhancing Model Generalizability of Medical Image Segmentation*

Shangqi Gao¹, Hangqi Zhou¹, Yibo Gao¹, and Xiahai Zhuang^{*1}

School of Data Science, Fudan University, Shanghai 200433, China
www.sdspeople.fudan.edu.cn/zhuangxiahai/

Abstract. Although supervised deep-learning has achieved promising performance in medical image segmentation, many methods cannot generalize well on unseen data, limiting their real-world applicability. To address this problem, we propose a deep learning-based Bayesian framework, which jointly models image and label statistics, utilizing the domain-irrelevant contour of a medical image for segmentation. Specifically, we first decompose an image into components of contour and basis. Then, we model the expected label as a variable only related to the contour. Finally, we develop a variational Bayesian framework to infer the posterior distributions of these variables, including the contour, the basis, and the label. The framework is implemented with neural networks, thus is referred to as deep Bayesian segmentation. Results on the task of cross-sequence cardiac MRI segmentation show that our method set a new state of the art for model generalizability. Particularly, the BayeSeg model trained with LGE MRI generalized well on T2 images and outperformed other models with great margins, *i.e.*, over 0.47 in terms of average Dice. Our code is available at <https://zmiclab.github.io/projects.html>.

Keywords: Bayesian segmentation · Image decomposition · Model generalizability · Deep learning

1 Introduction

Medical image segmentation is a task of assigning specific class for each anatomical structure. Thanks to the advance of deep learning, learning-based methods achieve promising performance in medical image segmentation [1,2,3]. However, many methods require a large number of images with manual labels for supervised learning [4,5], which limits their applications. For cardiac magnetic resonance (CMR) image segmentation, repeatedly labeling multi-sequence CMR image requires more labor of experts, and therefore is expensive [9,10]. Besides,

* Xiahai Zhuang is the corresponding author. This work was funded by the National Natural Science Foundation of China (grant no. 61971142, 62111530195 and 62011540404) and the development fund for Shanghai talents (no. 2020015).

the models trained at one site often cannot perform well at the other site [6]. Therefore, exploring segmentation methods with better generalizability is attractive and challenging.

Much effort has been made to train an end-to-end network by supervised learning. U-Net is one of the widely used networks, since it is more suitable for image segmentation [1]. Training deep neural networks in a supervised way often requires a lot of labeled data [4], but manual labeling of medical image requires professional knowledge and is very expensive. However, small training dataset can result in the problems of over-fitting and overconfidence, which will mislead the clinical diagnosis [6,14]. To solve the problems, Kohl et al. [7] proposed a probabilistic U-Net (PU-Net) for segmentation of ambiguous images by learning the distribution of segmentation. Their results showed that PU-Net could produce the possible segmentation results as well as the frequencies of occurring. Recently, Isensee et al. [3] developed a self-configuring method, i.e., nnU-Net, for learning-based medical image segmentation. This model could automatically configure its preprocessing, network architecture, training, and postprocessing, and achieves state-of-the-art performance on many tasks. Nevertheless, current learning-based methods deliver unsatisfactory performance when applied to unseen tasks [8], and improving generalizability of deep learning models is very challenging.

In this work, we propose a new Bayesian segmentation (BayeSeg) framework to promote model generalizability by joint modeling of image and label statistics. To the best of our knowledge, this is the first attempt of combining image decomposition, image segmentation, and deep learning. Concretely, we first decompose an image into two parts. One is the contour of this image, and the other is the basis approximating its intensity. Both the contour and basis are unknown, and we assign hierarchical Bayesian priors to model their statistics. After that, since the contour of an image is more likely to be sequence-independent, site-independent, and even modality-independent, we try to generate a label from the contour by explicitly modeling of label statistics. Finally, given an image, we build neural networks to infer the posterior distributions of the contour, basis, and label. Being different from many deep learning models that try to learn a deterministic segmentation from a given image, BayeSeg is aimed to learn the distribution of segmentation.

Our contributions are summarized as follows:

- We propose a new Bayesian segmentation framework, i.e., BayeSeg, by joint modeling of image and label statistics. Concretely, we decompose an image into the contour and basis, and assign hierarchical Bayesian priors to model the statistics of the contour, basis, and expected label.
- We solve the model by developing a variational Bayesian approach of computing the posterior distributions of the contour, basis, and label, and build a deep learning architecture of implementing the approach.
- We validate BayeSeg on the tasks of cross-sequence segmentation and cross-site segmentation, and show the superior generalizability of BayeSeg for unseen tasks.

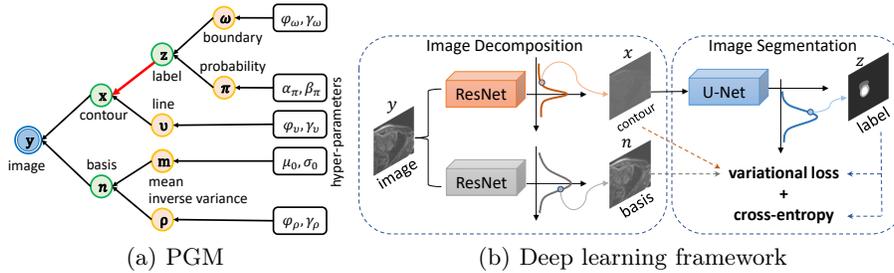


Fig. 1. The framework of Bayesian segmentation (BayeSeg). (a) shows the probabilistic graphical model (PGM) of BayeSeg. Here, the blue circle denotes an observed image, green and orange circles denote unknown variables, and white rectangles denote hyper-parameters. (b) presents the deep learning framework of BayeSeg. Given an image, we first use ResNets to infer the posterior distributions of the contour and basis, and obtain their random samples. Then, we use U-Net to infer the posterior distribution of label, and the resulting random sample is a segmentation. Please refer to Section 2.3 for the details of network architecture and training strategy.

2 Methodology

We propose a Bayesian segmentation (BayeSeg) framework to improve the generalizability of deep learning models. Many learning-based methods trained on one sequence MR images cannot generalize well to the other sequence or site data [6]. To solve the challenge, we propose the BayeSeg mainly consisting of two parts, i.e., (1) statistical modeling of image and label as shown in Fig. 1 (a), and (2) statistical inference of image and label as shown in Fig. 1 (b). At the first stage, we build a probabilistic graphical model (PGM) for the modeling of image and label. That is we decompose an image into its contour and basis, and only the contour is related to an expected label. At the second stage, we first build two residual networks (ResNets) to infer the posterior distributions of the contour and basis, respectively. Then, we build a U-Net to estimate the posterior distribution of the label. An intuitive understanding of “contour” and “basis” is shape and appearance. Since shape is domain-irrelevant, the model predicting a label from the contour will have better generalizability.

Fig. 1 shows the framework of statistical modeling and inference of the proposed BayeSeg. For the statistical modeling as shown in Fig. 1 (a), we first decompose an image \mathbf{y} into its basis \mathbf{n} and contour \mathbf{x} . The former is a Gaussian variable depending on the mean \mathbf{m} and the inverse variance $\boldsymbol{\rho}$. The latter depends on the expected label \mathbf{z} and the line \mathbf{v} for detecting the edges of contour. Similarly, the label \mathbf{z} depends on the segmentation boundary $\boldsymbol{\omega}$ and the segmentation probability $\boldsymbol{\pi}$ of all classes. Finally, Gamma priors are assigned to $\boldsymbol{\rho}$, \mathbf{v} , and $\boldsymbol{\omega}$, a Beta prior is assigned to $\boldsymbol{\pi}$, and a Gaussian prior is assigned to \mathbf{m} . Fig. 1 (b) shows the deep learning framework of inferring related variables. The outputs of ResNets and U-Net will be jointly used to compute a variational loss, which is the key of improving model generalizability.

2.1 Statistical modeling of image and label

This section shows the statistical modeling of image and label. Given an image sampled from the variable $\mathbf{y} \in \mathbb{R}^{d_y}$, where d_y denotes the dimension of \mathbf{y} , we decompose \mathbf{y} into the sum of a contour \mathbf{x} and a basis \mathbf{n} , i.e., $\mathbf{y} = \mathbf{x} + \mathbf{n}$. Then, the basis \mathbf{n} is modeled by a normal distribution with a mean $\mathbf{m} \in \mathbb{R}^{d_y}$ and a covariance $\text{diag}(\boldsymbol{\rho})^{-1} \in \mathbb{R}^{d_y \times d_y}$. Moreover, the contour \mathbf{x} is modeled by a simultaneous autoregressive model (SAR) [11] depending on the expected label $\mathbf{z} \in \mathbb{R}^{d_y \times K}$, the line $\mathbf{v} \in \mathbb{R}^{d_y}$ indicating edges of the contour, and the matrix $\mathbf{D}_x \in \mathbb{R}^{d_y \times d_y}$ describing a neighboring system of the contour, where K is the number of classes for segmentation. Finally, the observation likelihood of an image \mathbf{y} can be expressed as

$$p(\mathbf{y}|\mathbf{x}, \mathbf{m}, \boldsymbol{\rho}) = \mathcal{N}(\mathbf{x} + \mathbf{m}, \text{diag}(\boldsymbol{\rho})^{-1}). \quad (1)$$

The basis \mathbf{n} of an image is determined by a normal distribution, that is, $p(\mathbf{n}|\mathbf{m}, \boldsymbol{\rho}) = \mathcal{N}(\mathbf{n}|\mathbf{m}, \text{diag}(\boldsymbol{\rho})^{-1})$. Specifically, we assign a Gaussian prior to \mathbf{m} , i.e., $p(\mathbf{m}|\boldsymbol{\mu}_0, \sigma_0) = \mathcal{N}(\mathbf{m}|\boldsymbol{\mu}_0, \sigma_0^{-1}\mathbf{I})$, and a Gamma prior to $\boldsymbol{\rho}$, namely, $p(\boldsymbol{\rho}|\phi_\rho, \gamma_\rho) = \prod_{i=1}^{d_y} \mathcal{G}(\rho_i|\phi_{\rho i}, \gamma_{\rho i})$. Here, \mathbf{I} denotes an identity matrix, $\mu_0, \sigma_0, \phi_\rho$, and γ_ρ are predefined hyper-parameters, and $\mathcal{G}(\cdot, \cdot)$ represents the Gamma distribution.

The contour \mathbf{x} of an image is determined by a SAR mainly depending on the expected label \mathbf{z} and the line \mathbf{v} ,

$$p(\mathbf{x}|\mathbf{z}, \mathbf{v}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\mathbf{0}, [\mathbf{D}_x^T \text{diag}(\mathbf{z}_k \mathbf{v}) \mathbf{D}_x]^{-1}), \quad (2)$$

where, \mathbf{z}_k denotes the segmentation of the k -th class, and $\mathbf{D}_x = \mathbf{I} - \mathbf{B}_x$ is non-singular. \mathbf{B}_x describes a neighboring system of each pixel. For examples, if the values of \mathbf{B}_x for the nearest four pixels equals to 0.25 while others are zeros, then \mathbf{D}_x is aimed to compute the average difference of each pixel with its four neighbors. The line \mathbf{v} can indicate the edges of the contour, and it is assigned a Gamma prior, i.e., $p(\mathbf{v}|\phi_{vi}, \gamma_{vi}) = \prod_{i=1}^{d_y} \mathcal{G}(v_i|\phi_{vi}, \gamma_{vi})$.

The label \mathbf{z} is modeled by another SAR depending on the segmentation boundary $\boldsymbol{\omega} \in \mathbb{R}^{d_y \times K}$ and the segmentation probability $\boldsymbol{\pi} \in \mathbb{R}^K$ of all classes, namely,

$$p(\mathbf{z}|\boldsymbol{\pi}, \boldsymbol{\omega}) = \prod_{k=1}^K \mathcal{N}(\mathbf{z}|\mathbf{0}, [-\ln(1 - \pi_k) \mathbf{D}_z^T \text{diag}(\boldsymbol{\omega}_k) \mathbf{D}_z]^{-1}), \quad (3)$$

where, the definition of \mathbf{D}_z is the same as \mathbf{D}_x in Eq. (2); $\boldsymbol{\omega}_k$ can indicate the boundary of the k -th segmentation \mathbf{z}_k ; and π_k denotes the probability of a pixel belonging to the k -th class. Finally, we assign Gamma prior to $\boldsymbol{\omega}$, i.e., $p(\boldsymbol{\omega}) = \prod_{i=1}^{d_y} \prod_{k=1}^K \mathcal{G}(\omega_{ki}|\phi_{\omega ki}, \gamma_{\omega ki})$, and give Beta prior to $\boldsymbol{\pi}$, namely, $p(\boldsymbol{\pi}) = \prod_{k=1}^K \mathcal{B}(\pi_k|\alpha_{\pi k}, \beta_{\pi k})$. The details of Gaussian distribution, Gamma distribution, and Beta distribution are provided in Appendix A.

2.2 Variational inference of image and label

This section shows a variational method of inferring the contour, basis and label given an image \mathbf{y} by maximum a *posteriori* (MAP) estimation. Let $\boldsymbol{\psi} =$

$\{\mathbf{m}, \boldsymbol{\rho}, \mathbf{x}, \mathbf{v}, \mathbf{z}, \boldsymbol{\omega}, \boldsymbol{\pi}\}$ denote the set of all variables to infer, then our aim is to infer the posterior distribution $p(\boldsymbol{\psi}|\mathbf{y})$. Since direct computation is intractable, we use variational Bayesian (VB) method [12] to solve the problem. Concretely, we approximate the posterior distribution $p(\boldsymbol{\psi}|\mathbf{y})$ via a variational distribution $q(\boldsymbol{\psi})$ by assuming the variables in $\boldsymbol{\psi}$ are independent, namely,

$$q(\boldsymbol{\psi}) = q(\mathbf{m})q(\boldsymbol{\rho})q(\mathbf{x})q(\mathbf{v})q(\mathbf{z})q(\boldsymbol{\omega})q(\boldsymbol{\pi}). \quad (4)$$

After that, we minimize the KL divergence between $q(\boldsymbol{\psi})$ and $p(\boldsymbol{\psi}|\mathbf{y})$, and which results in our final variational loss as follows,

$$\min_{q(\boldsymbol{\psi})} \mathcal{L}_{var} = \text{KL}(q(\boldsymbol{\psi})||p(\boldsymbol{\psi})) - \mathbb{E}[\ln p(\mathbf{y}|\boldsymbol{\psi})] \quad (5)$$

The details of further unfolding the variational loss are provided in Appendix B.

2.3 Neural networks and training strategy

This section shows the network architecture of achieving the variational inference and the training strategy for image segmentation. As Fig. 1 (b) shows, *at the decomposition stage*, we adopt two ResNets [13] to separately infer the variational posteriors of the contour \mathbf{x} and basis \mathbf{n} , i.e., $q(\mathbf{x})$ and $q(\mathbf{n})$, respectively. The ResNet of inferring the contour consists of 10 residual blocks, and each block has a structure of ‘‘Conv + ReLU + Conv’’. The output of this ResNet has two channels. One is the element-wise mean of the contour, and the other is its element-wise variance. The contour \mathbf{x} in the figure denotes a random sample from $q(\mathbf{x})$. The ResNet of inferring the basis consists of 6 residual blocks, and each block has a structure of ‘‘Conv + BN + ReLU + Conv + BN’’. Similarly, this ResNet will output the mean and variance of the basis, and the basis \mathbf{n} is randomly sampled from its variational posterior distribution. *At the segmentation stage*, we adopt a U-Net [1] to infer the variational posterior of the label \mathbf{z} , i.e., $q(\mathbf{z})$. The output of this U-Net has $2K$ channels. The first K channels denote the element-wise mean of the label, and the left channels represent its element-wise variance. The label \mathbf{z} in Fig. 1 (b) is a random sample from the resulting posterior distribution, and it will be taken as a stochastic segmentation for training.

BayeSeg is trained in an end-to-end manner by balancing between cross-entropy and the variational loss in (5). For convenience, the cross-entropy between a stochastic segmentation and the provided manual segmentation is notated as \mathcal{L}_{ce} . Then, our total loss of training BayeSeg is given by,

$$\min_{q(\boldsymbol{\psi})} \mathcal{L}_{ce} + \lambda \mathcal{L}_{var}, \quad (6)$$

where, the balancing weight λ is set to 100 in our experiments. Besides, other hyper-parameters in Fig. 1 (a) is summarized as follows. Each element of $\boldsymbol{\phi}$. of related variables is set to 2; $\alpha_\pi = 2$ and $\beta_\pi = 2$ for the segmentation probability $\boldsymbol{\pi}$; $\boldsymbol{\mu}_0 = \mathbf{0}$ and $\sigma_0 = 1$ for the mean of the basis; The elements of γ_ρ , γ_v , and γ_ω

are set to 10^{-6} , 10^{-8} , and 10^{-4} , respectively. Note that BayeSeg could properly decompose an image into the contour and basis due to the priors with respect to \mathbf{v} and $\boldsymbol{\rho}$. That is the contour and basis are adaptively balanced after selecting proper γ_ρ and γ_v .

3 Experiments

3.1 Tasks and datasets

We used the LGE of MSCMRseg [9] to train models. To validate the performance on the task of cross-sequence segmentation, we tested models using LGE and T2 of MSCMRseg. To validate the performance on the task of cross-site segmentation, we tested models using LGE of MSCMRseg and ACDC [14].

MSCMRseg [10,9] was provided by MICCAI’19 Multi-sequence Cardiac MR Segmentation Challenge. This dataset consists of 45 multi-sequence CMR images, including LGE, C0, and T2. Each case comes from the same patient who underwent cardiomyopathy. The manual segmentation results of left ventricle (LV), right ventricle (RV), and myocardium (Myo) for all images are available. In this study, we randomly split the 45 cases into three sets consisting of 25, 5, and 15 cases, respectively. Then, we only used the 25 LGE images for training, and the 5 LGE images for validation. Finally, we tested models on the 15 multi-sequence cases to show the performance of cross-sequence segmentation.

ACDC [14] was provided by MICCAI’17 Automatic Cardiac Diagnosis Challenge. This dataset consists of shot-axis cardiac cine-MRIs of 100 patients for training, and of 50 patients for test. Only the manual segmentation results of training data are provided for LV, RV, and Myo during the end-diastolic (ED) and end-systolic (ES) phases. In our study, we tested models using the 100 training images to show the performance of cross-site segmentation.

BayeSeg was implemented by Pytorch, and trained by Adam optimizer with the initial learning rate to be 10^{-4} . The learning rate was dropped every 500 epochs by a factor of 0.1, and the training was stopped when up to 2000 epochs. At the test stage, we took the mean of \mathbf{z} as the final segmentation label, since the variational posterior $q(\mathbf{z})$ is a Gaussian distribution whose mode is its mean. All experiments were run on a TITAN RTX GPU with 24G memory.

3.2 Cross-sequence segmentation

To study the performance of BayeSeg on the task cross-sequence segmentation, we trained four models using the 25 LGE images of MSCMRseg. Concretely, we trained a U-Net, which had the same architecture as the U-Net in Fig. 1 (b), by minimizing the cross-entropy. Then, we trained PU-Net on the same dataset using its public code and default settings. Moreover, we trained a baseline, which shares the same architecture as BayeSeg, without using the variational loss in (6). Finally, we trained the BayeSeg by minimizing the total loss in (6). For fair comparisons, all models were trained using consistent data augmentation,

Table 1. Evaluation on the task of cross-sequence cardiac segmentation. *Note that all models were only trained on LGE of MSCMRseg, but tested on LGE and T2.* Here, G denotes the drop of average dice, and it is used to measure model generalizability.

Method	LGE of MSCMRseg (15 samples)				T2 of MSCMRseg (15 samples)				G
	LV	Myo	RV	Avg	LV	Myo	RV	Avg	
U-Net	.855±.045	.727±.064	.733±.097	.772±.093	.203±.183	.095±.093	.055±.062	.118±.139	.654
PU-Net	.898±.027	.768±.056	.729±.089	.798±.096	.279±.162	.166±.122	.195±.130	.213±.147	.585
Baseline	.893±.023	.783±.045	.727±.069	.801±.085	.481±.129	.117±.079	.090±.123	.230±.211	.571
BayeSeg	.887±.028	.774±.048	.763±.060	.808±.073	.846±.119	.731±.117	.528±.206	.701±.202	.107

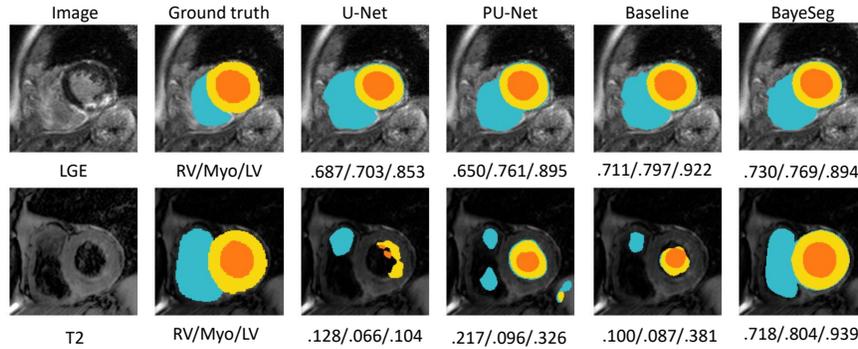


Fig. 2. Visualization of results on the task of cross-sequence segmentation. Here, we choose the median case of BayeSeg according to the dices of 15 LGE images.

including random flip and rotation. In the test stage, we evaluated all models on the 15 multi-sequence images, i.e., LGE and T2, and reported the dices of LV, Myo, and RV as well as average dice.

Table 1 reports the results of compared methods. One can see that PU-Net, Baseline, and BayeSeg achieved comparable performance, when the training and test sequences were consistent. If the training sequence and test sequence are different, the performance of BayeSeg dropped weakly, but that of others decreased dramatically. Besides, BayeSeg greatly outperformed Baseline in this case, which demonstrates that the variational loss, induced by joint modeling of image and label statistics, is the key of improving model generalizability. To qualitatively evaluate all models, we chose the median case of BayeSeg according to the average dices of 15 LGE images, and visualized the segmentation results in Figure 2. This figure shows BayeSeg delivers the best performance in segmenting the unseen T2 sequence, which again confirms the effectiveness of our framework in improving model generalizability.

3.3 Cross-site segmentation

To study the performance of BayeSeg on the task of cross-center segmentation, we tested all models trained in the previous section on ACDC. Table 2 reports the results of these models. This table showed that the performance of all methods dropped when the training and test samples came from two different sites, but

Table 2. Evaluation on the task of cross-site cardiac segmentation. *Note that all models were only trained on LGE of MSCMRseg, but tested on LGE and ACDC.* Here, G denotes the drop of average dice, and it is used to measure model generalizability.

Method	LGE of MSCMRseg (15 samples)				ACDC (100 samples)				G
	LV	Myo	RV	Avg	LV	Myo	RV	Avg	
U-Net	.855±.045	.727±.064	.733±.097	.772±.093	.721±.187	.602±.183	.659±.202	.660±.197	.112
PU-Net	.898±.027	.768±.056	.729±.089	.798±.096	.743±.152	.641±.146	.604±.215	.663±.184	.126
Baseline	.893±.023	.783±.045	.727±.069	.801±.085	.776±.134	.667±.150	.585±.227	.676±.192	.125
BayeSeg	.887±.028	.774±.048	.763±.060	.808±.073	.792±.130	.694±.123	.659±.175	.715±.155	.093

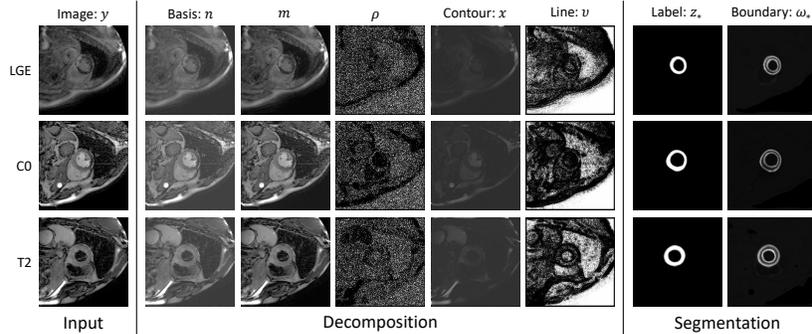


Fig. 3. Visualization of posteriors inferred by BayeSeg. Here, the subscript $*$ denotes the result of myocardium.

BayeSeg delivered the least performance drop. Therefore, BayeSeg generalized well on the unseen samples from the different site, thanks to the joint modeling of image and label statistics.

3.4 Interpretation of joint modeling

In this section we interpreted the joint modeling of image and label statistics. Fig. 3 shows the posteriors inferred by our BayeSeg for three different sequences of MSCMRseg. One can see that, at the decomposition stage, an image was mainly decomposed into its basis and contour. The basis n was modeled as a Gaussian distribution with the mean m and the inverse variance ρ . It was an approximation of the image, and therefore x was left as the contour. To avoid the smoothness of this contour, we assigned the line v to detect its edges. The large values of v indeed showed the smooth areas of contour, while the small values indicate the edges. At the segmentation stage, we choose to segment the contour, since it is more likely to be sequence-independent, site-independent, and even modality-independent. To achieve better segmentation around the boundary of some object, such as myocardium, we assigned the ω to detect the segmentation boundary. This variable successfully indicated the inner and outer boundaries of myocardium, as shown in Fig. 3.

4 Conclusion

In this work, we proposed a new Bayesian segmentation framework by joint modeling of image and label statistics. Concretely, we decomposed an image into its basis and contour, and estimated the segmentation of this image from the more stable contour. Our experiments have shown that the proposed framework could address the problem of over-fitting and greatly improve the generalizability of deep learning models.

A Preliminary

If n is a variable which follows Gaussian distribution, then its probability density function is given by

$$p(n|m, \rho) = \mathcal{N}(n|m, \rho^{-1}) = \frac{1}{\sqrt{2\pi/\rho}} \exp^{-\frac{\rho}{2}(n-m)^2}, \quad (7)$$

If ω is a variable which follows Gamma distribution, then its probability density function is given by

$$p(\omega|\phi, \gamma) = \mathcal{G}(\omega|\phi, \gamma) = \frac{\phi^\gamma}{\Gamma(\gamma)} \omega^{\gamma-1} e^{-\phi\omega}, \quad (8)$$

where, $\Gamma(\cdot)$ denotes the Gamma function.

If π is a variable which follows Beta distribution, then its probability density function is given by

$$p(\pi|\alpha, \beta) = \mathcal{B}(\omega|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \pi^{\alpha-1} (1 - \pi)^{\beta-1}. \quad (9)$$

B Variational Inference

To estimate the variational posteriors, we minimize the KL divergence between $q(\boldsymbol{\psi})$ and $p(\boldsymbol{\psi}|\mathbf{y})$, which results in

$$\operatorname{argmin}_{q(\boldsymbol{\psi})} \text{KL}(q(\boldsymbol{\psi})||p(\boldsymbol{\psi}|\mathbf{y})) = \operatorname{argmin}_{q(\boldsymbol{\psi})} \text{KL}(q(\boldsymbol{\psi})||p(\boldsymbol{\psi})) - \mathbb{E}[\ln p(\mathbf{y}|\boldsymbol{\psi})], \quad (1)$$

Moreover, we covert it to the following problem by reparameterization,

$$\operatorname{argmin}_{q(\boldsymbol{\psi})} \text{KL}(q(\boldsymbol{\psi})||p(\boldsymbol{\psi})) - \mathbb{E}_{q(\boldsymbol{\rho})}[\ln p(\mathbf{y}|\mathbf{x}, \mathbf{m}, \boldsymbol{\rho})]. \quad (2)$$

B.1 Explicit computation of $q(\mathbf{v})$, $q(\boldsymbol{\omega})$, $q(\boldsymbol{\pi})$, $q(\boldsymbol{\rho})$

Minimizing (2) over $q(\mathbf{v})$, $q(\boldsymbol{\omega})$, $q(\boldsymbol{\pi})$, $q(\boldsymbol{\rho})$ successively results in the explicit formulas of computing the parameters of these distributions as follows,

$$\begin{cases} \hat{\alpha}_{vi} = \gamma_{vi} + K/2 \\ \hat{\beta}_{vi} = \frac{1}{2} \sum_{k=1}^K \hat{\boldsymbol{\mu}}_{zki} [(\mathbf{D}_x \hat{\boldsymbol{\mu}}_x)_i^2 + \langle \hat{\boldsymbol{\sigma}}_x, \mathbf{d}_{xi}^2 \rangle] + \phi_{vi} \\ \hat{\boldsymbol{\mu}}_v = \frac{\hat{\boldsymbol{\alpha}}_v}{\hat{\boldsymbol{\beta}}_v} = \frac{2\gamma_v + K}{\sum_{k=1}^K \hat{\boldsymbol{\mu}}_{zk} [(\mathbf{D}_x \hat{\boldsymbol{\mu}}_x)^2 + 2\hat{\boldsymbol{\sigma}}_x]} + 2\phi_v \end{cases},$$

$$\begin{cases} \hat{\alpha}_{\omega ki} = \gamma_{\omega ki} + 1/2 \\ \hat{\beta}_{\omega ki} = \frac{1}{2} [\Psi(\hat{\alpha}_{\pi k} + \hat{\beta}_{\pi k}) - \Psi(\hat{\beta}_{\pi k})] [(\mathbf{D}_z \hat{\boldsymbol{\mu}}_{zk})_i + \langle \hat{\boldsymbol{\sigma}}_{zk}^2, \mathbf{d}_{zi}^2 \rangle] + \phi_{\omega ki} \\ \hat{\boldsymbol{\mu}}_{\omega k} = \frac{\hat{\boldsymbol{\alpha}}_{\omega k}}{\hat{\boldsymbol{\beta}}_{\omega k}} = \frac{2\gamma_{\omega k} + 1}{[\Psi(\hat{\alpha}_{\pi k} + \hat{\beta}_{\pi k}) - \Psi(\hat{\beta}_{\pi k})] [(\mathbf{D}_z \hat{\boldsymbol{\mu}}_{zk})^2 + 2\hat{\boldsymbol{\sigma}}_{zk}^2] + 2\phi_{\omega k}} \end{cases},$$

$$\begin{cases} \hat{\alpha}_{\pi k} = \alpha_{\pi k} + d_y/2 \\ \hat{\beta}_{\pi k} = \frac{1}{2} \sum_{i=1}^{d_y} \hat{\boldsymbol{\mu}}_{\omega ki} [(\mathbf{D}_z \hat{\boldsymbol{\mu}}_{zk})_i^2 + 2\hat{\boldsymbol{\sigma}}_{zki}^2] + \beta_{\pi k} \end{cases},$$

and $\hat{\boldsymbol{\mu}}_\rho = \hat{\boldsymbol{\alpha}}_\rho / \hat{\boldsymbol{\beta}}_\rho = (2\gamma_\rho + 1) / ([\mathbf{y} - (\mathbf{x} + \mathbf{m})]^2 + 2\phi_\rho)$. Here, $\Psi(\cdot)$ denotes the Digamma function. Finally, the related variational posterior distributions are given by

$$q(\mathbf{v}) = \prod_{i=1}^{d_y} \mathcal{G}(v_i | \hat{\beta}_{vi}, \hat{\alpha}_{vi}) \text{ and } q(\boldsymbol{\omega}) = \prod_{k=1}^K \prod_{i=1}^{d_y} \mathcal{G}(\omega_{ki} | \hat{\beta}_{\omega ki}, \hat{\alpha}_{\omega ki})$$

$$q(\boldsymbol{\pi}) = \prod_{k=1}^K \mathcal{B}(\pi_k | \hat{\alpha}_{\pi k}, \hat{\beta}_{\pi k}) \text{ and } q(\boldsymbol{\rho}) = \prod_{i=1}^{d_y} \mathcal{G}(\rho_i | \hat{\beta}_{\rho i}, \hat{\alpha}_{\rho i})$$

B.2 Variational inference of $q(\mathbf{x})$, $q(\mathbf{z})$ and $q(\mathbf{m})$

Minimizing (2) over $q(\mathbf{x})$, $q(\mathbf{z})$ and $q(\mathbf{m})$ successively results in the losses of further inferring the parameters of these distributions as follows,

$$\mathcal{L}_y = \frac{1}{2} \|\mathbf{y} - (\mathbf{x} + \mathbf{m})\|_{diag(\hat{\boldsymbol{\mu}}_\rho)}^2,$$

where, $\mathbf{x} = \hat{\boldsymbol{\sigma}}_x \odot \boldsymbol{\epsilon} + \hat{\boldsymbol{\mu}}_x$, $\mathbf{m} = \hat{\boldsymbol{\sigma}}_m \odot \boldsymbol{\epsilon} + \hat{\boldsymbol{\mu}}_m$, and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

$$\begin{cases} \mathcal{L}_{\hat{\boldsymbol{\mu}}_z} = \frac{1}{2} \sum_{k=1}^K [\Psi(\hat{\alpha}_{\pi k} + \hat{\beta}_{\pi k}) - \Psi(\hat{\beta}_{\pi k})] \|\mathbf{D}_z \hat{\boldsymbol{\mu}}_{zk}\|_{diag(\hat{\boldsymbol{\mu}}_{\omega k})}^2 \\ \mathcal{L}_{\hat{\boldsymbol{\sigma}}_z} = \frac{1}{2} \sum_{k=1}^K [\Psi(\hat{\alpha}_{\pi k} + \hat{\beta}_{\pi k}) - \Psi(\hat{\beta}_{\pi k})] [\langle 2\hat{\boldsymbol{\mu}}_{\omega k}, \hat{\boldsymbol{\sigma}}_{zk}^2 \rangle - \langle \mathbf{1}, \ln(\hat{\boldsymbol{\sigma}}_{zk}^2) \rangle] \end{cases}$$

$$\begin{cases} \mathcal{L}_{\hat{\mu}_x} = \frac{1}{2} \sum_{k=1}^K \|D_x \hat{\mu}_x\|_{diag(\hat{\mu}_{zk} \hat{\mu}_v)}^2 \\ \mathcal{L}_{\hat{\sigma}_x} = \frac{1}{2} \sum_{k=1}^K [\langle 2\hat{\mu}_{zk} \hat{\mu}_v, \hat{\sigma}_x^2 \rangle - \frac{1}{K} \langle \mathbf{1}, \ln(\hat{\sigma}_x^2) \rangle] \end{cases}$$

$$\begin{cases} \mathcal{L}_{\hat{\mu}_m} = \frac{\sigma_0}{2} \|\hat{\mu}_m\|_2^2 \\ \mathcal{L}_{\hat{\sigma}_m} = \frac{1}{2} [\langle \sigma_0 \mathbf{1}, \hat{\sigma}_m^2 \rangle - \langle \mathbf{1}, \ln(\hat{\sigma}_m^2) \rangle] \end{cases}$$

Overall, the final **variational loss** of inferring the basis, contour, and label of an image is summarized as

$$\mathcal{L}_{var} = \mathcal{L}_y + \mathcal{L}_{\hat{\mu}_z} + \mathcal{L}_{\hat{\sigma}_z} + \mathcal{L}_{\hat{\mu}_x} + \mathcal{L}_{\hat{\sigma}_x} + \mathcal{L}_{\hat{\mu}_m} + \mathcal{L}_{\hat{\sigma}_m}.$$

Finally, the related variational posterior distributions are given by

$$\begin{aligned} q(\mathbf{m}) &= \mathcal{N}(\mathbf{m} | \hat{\mu}_m, diag(\hat{\sigma}_m^2)), \\ q(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \hat{\mu}_x, diag(\hat{\sigma}_x^2)), \\ q(\mathbf{z}) &= \prod_{k=1}^K \mathcal{N}(z | \hat{\mu}_{zk}, diag(\hat{\sigma}_{zk}^2)). \end{aligned}$$

References

1. Ronneberger, O., Fischer, P., Brox, T.: U-Net: convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computation and Computer-Assisted Intervention, pp. 234–241. Springer, 2015
2. Cicek, Ö., Abdulkadir, A., Lienkamp S.: 3D U-Net: learning dense volumetric segmentation from sparse annotation. In: International Conference on Medical Image Computation and Computer-Assisted Intervention, pp. 424–432. Springer, 2016
3. Isensee, F., Jaeger, P., Kohl, S., Petersen, J., MaierHein, K.: nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, **18**(2), 203–211 (2021)
4. Zhao, A., Balakrishnan, G., Durand, F., Gutttag, J., Dalca, A.: Data augmentation using learned transformations for one-shot medical image segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8543–8553. IEEE, 2019
5. Chen, H., Dou, Q., Yu, L., Qin, J., Heng, P.: VoxResNet: Deep voxelwise residual networks for brain segmentation from 3D MR images. *NeuroImage*, **170**, 446–455 (2018)
6. Li, Z., Kamnitsas, K., Glocker, B.: Overfitting of neural nets under class imbalance: analysis and improvement for segmentation. In: International Conference on Medical Image Computation and Computer-Assisted Intervention, pp. 402–410. Springer, 2019
7. Kohl, S., Romera-Paredes, B., Meyer, C., Fauw, J., Ledsam, J., Maier-Hein, K., Eslami, S., Rezende, D.: A probabilistic U-net for segmentation of ambiguous images. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems, pp. 6965–6975. ACM, 2018

8. Cheng, O., Chen, C., Li, S., Li, Z., Qin, C., Bai, W., Rueckert, D.: Causality-inspired single-source domain generalization for medical image segmentation. arXiv:2111.12525, 2021
9. Zhuang, X.: Multivariate mixture model for myocardial segmentation combining multi-source images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **41**(12), 2933–2946 (2019)
10. Zhuang, X.: Multivariate mixture model for cardiac segmentation from multi-sequence MRI. In: *International Conference on Medical Image Computation and Computer-Assisted Intervention*, pp. 581–588. Springer, 2016
11. Shekhar, S., Xiong, H.: *Simultaneous Autoregressive Model*. Springer, Boston, MA (2008)
12. Blei, D., Kucukelbir, A., McAuliffe, J.: Variational inference: a review for statisticians. *Journal of the American Statistical Association*, **112**(518), 859–877 (2017)
13. He, K., Zhang, X., Sun, J., Ren, S.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778. IEEE, 2016
14. Berbard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P., Cetin, I. et al.: Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: is the problem solved?. *IEEE Transactions on Medical Imaging*, **37**(11), 2514–2525 (2018)