

Boosting Semi-Supervised Semantic Segmentation with Probabilistic Representations

Haoyu Xie, Changqi Wang, Mingkai Zheng, Minjing Dong, Shan You, Chong Fu, Chang Xu

School of Computer Science and Engineering, Northeastern University, Shenyang, China

School of Computer Science, Faculty of Engineer, The University of Sydney, Sydney, Australia

SenseTime Research, Beijing, China

Key Laboratory of Intelligent Computing in Medical Image, Ministry of Education, Northeastern University, China

Introduction

TASK:

- Pixel-wise contrastive learning in Computer vision.
- Semi-supervised semantic segmentation.

MOTIVATION:

- Improve the quality of pixel-wise representation considering the probability, allow them to perform better under the inaccurate pseudo labels

CONTRIBUTION:

- Define pixel-wise representations from a new perspective of probability theory.
- Through modelling the mapping from pixels to representations as the probability via multivariate Gaussian distributions, tune the contribution of the ambiguous representations to tolerate the risk of inaccurate pseudo-labels.
- Define prototypes in the form of distributions, which indicates the confidence of a class, while the point prototype cannot.
- Propose a Probabilistic Representation Contrastive Learning (PRCL) framework that improves representation quality by taking its probability into consideration.

Methodology

PROBABILISTIC REPRESENTATION



- We denote the probability of mapping a pixel x_i to a representation z_i as $p(z_i|x_i)$ and define the representation as a random variable following it. For simplicity, we take the form of multivariate Gaussian distribution $\mathcal{N}(\mu, \sigma^2 I)$ as:

$$p(z_i|x_i) = \mathcal{N}(z_i; \mu_i, \sigma_i^2 I). \quad (1)$$

DISTRIBUTION PROTOTYPE

- The prototype is the posterior distribution after the n^{th} observations $\{z_1, z_2, \dots, z_n\}$. Under the assumption that all the observations are conditionally independent, the distribution prototype can be derived as:

$$p(\rho|z_1, z_2, \dots, z_{n+1}) = \alpha \frac{p(\rho|z_{n+1})}{p(\rho)} p(\rho|z_1, z_2, \dots, z_n), \quad (2)$$

where α is a normalization factor. In addition to Equation 1, we can rewrite the prototype as

$$\rho \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}^2 I), \quad (3)$$

where

$$\hat{\mu} = \sum_{i=1}^n \frac{\hat{\sigma}^2}{\sigma_i^2} \mu_i \quad (4)$$

$$\frac{1}{\hat{\sigma}^2} = \sum_{i=1}^n \frac{1}{\sigma_i^2}. \quad (5)$$

SIMILARITY

- We leverage Mutual Likelihood Score (MLS) to measure the similarity between two distributions z_i and z_j , as follows:

$$\begin{aligned} MLS(z_i, z_j) &= \log(p(z_i = z_j)) \\ &= -\frac{1}{2} \sum_{l=1}^D \left(\frac{(\mu_i^{(l)} - \mu_j^{(l)})^2}{\sigma_i^{2(l)} + \sigma_j^{2(l)}} + \log(\sigma_i^{2(l)} + \sigma_j^{2(l)}) \right) \\ &\quad - \frac{D}{2} \log 2\pi, \end{aligned} \quad (6)$$

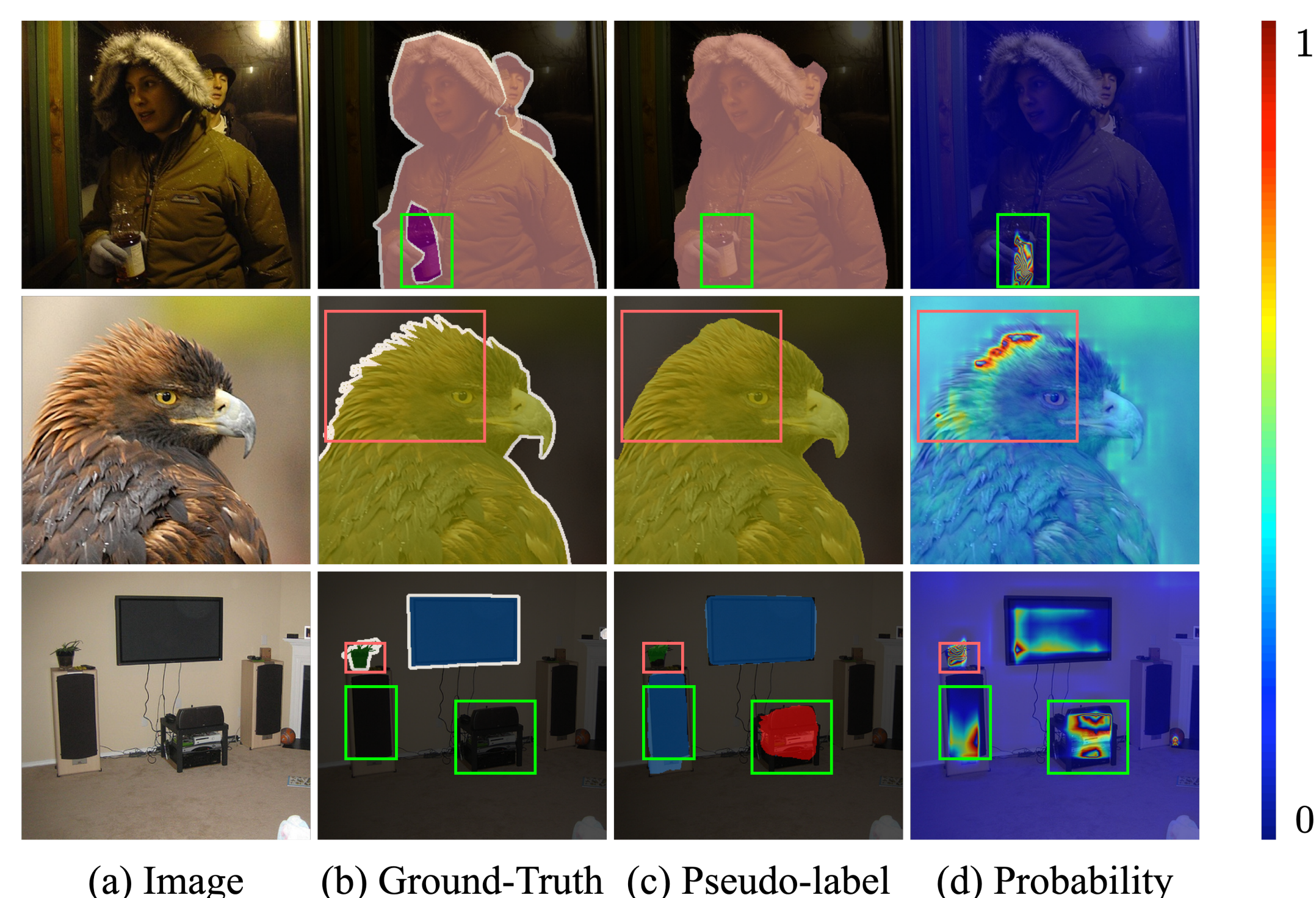
- In the first term, the weight of l_2 distance is small when the σ^2 is large, which indicates that the similarity between z_i and z_j becomes lower due to the low probabilities, even if they are very similar in the view of l_2 distance. The probability has been taken into consideration besides the simple similarity measure for representations learning.

- In the second term, σ^2 is penalized for the low probability representations, which makes all the representations more reliable.

- Besides, σ^2 and μ can interact with each other. The learnable σ^2 is associated with l_2 distance. This means that σ^2 can be learned via the relations among representations. On the other hand, the μ can also be optimized via the σ^2 . This is consistent with intuitive cognition.

Results

In Figure , columns from left to right represent input image, ground-truth, pseudo-label, and probability, respectively. For the fourth column, the red color represents the large σ^2 (low probability). The green boxes mark the mismatches caused by inaccurate pseudo-labels (e.g., person and bottle) and the red boxes mark the fuzzy pixels (e.g., furry edge of the bird). These two cases are highlighted by σ^2 and make low contribution in training process.



(a) Image (b) Ground-Truth (c) Pseudo-label (d) Probability