教師情報に含まれるノイズに堅牢な弱教師あり領域分割手法

下田 $1^{1,a}$ 柳井 啓司¹

1. 概要

領域分割においては,学習画像における物体のカテゴリ ごとにピクセル単位の領域の教師情報が必要であり,教師 情報の付与は大きなコスト,時間を要する。一般に,高度 な教師情報を必要とする手法を完全教師あり学習領域分 割,画像における物体のカテゴリ情報のみから学習する手 法を弱教師あり学習領域分割と呼ぶ。弱教師あり学習によ る領域分割が可能となれば,大幅な学習データを収集する ためのコストの削減が可能である。本研究においては,弱 教師あり領域分割の精度向上を目指す。近年の弱教師あり 領域分割においては,ノイズを含む教師情報を用いて領域 分割モデルの学習を行うのが主流になっている。そこで, 本研究ではノイズに対して堅牢な弱教師あり領域分割手法 の提案を行う。

2. 導入

近年,深層学習の進歩とともに画像の領域分割制度は飛 躍的に向上している。しかしながら,領域分割の学習に必 要な教師情報のアノテーションコストの高さなどから実用 面において課題が残っている。領域分割の教師情報のアノ テーションコストを下げる手段の一つとして,弱教師あり 学習による領域分割が注目を集めている。領域分割におけ る弱教師あり学習では,教師情報として画像のクラスラベ ルのみを用いる。クラスラベルのアノテーションは領域の アノテーションと比較して大幅に安価であり,弱教師あり 学習による領域分割が可能であれば,大きなコスト削減が 期待できる。また,Web 画像における検索クエリをクラス ラベルとみなせば,クラスラベル付きの画像は自動で収集 することも可能である。

近年の弱教師あり領域分割は二つの段階にわけて行われ ている。まず,学習画像の大まかな領域分割を行い,その 後にその大まかな領域分割結果を活用しより詳細な領域分 割ネットワークを学習する。弱教師あり領域分割において, 大まかな領域を推定する方法としては,画像の識別結果の 可視化や前景領域の抽出があるが,近年はこれら両方を用 いる手法が広く使われている。特に,近年はクラス分類の

^{a)} shimoda-k@mm.inf.uec.ac.jp

可視化方法として Class Activation Mapping(CAM) が高 い精度と汎用性から広く用いられている。Global Average Pooling(GAP)を用いてクラス分類ネットワークを学習す ることで,物体の大まかな可視化を行うことができる。詳 細な領域分割ネットワークの学習においては,完全教師あ り領域分割において高い精度を達成しているネットワー クが用いられる。また,さらなる精度の改善方法として, CRF による色特徴を用いた改善方法,領域の類似度を学 習する改善方法などがある。

現在の弱教師あり領域分割と完全教師あり領域分割の違 いには,教師情報にノイズが含まれているという点がある。 弱教師あり領域分割の精度改善においては,教師情報にお けるノイズの減少,ノイズに堅牢な学習方法の2点が重要 となる。本手法ではノイズに堅牢な学習に着目した。ノイ ズに堅牢な学習手法としてはアンサンブルが有効である ことが広く知られている。本手法においては,異なる学習 データのモデルの出力によるアンサンプルによりノイズに 堅牢な学習を目指した。これにより学習画像の不足が生じ るが,Web 画像を用いることでこれを補った。また,Web 画像と対象のデータセットのドメインの違いから生じる変 化に対応するために,広域の認識結果と局所領域の認識結 果との一貫性からの学習を行った。以下に本研究のコント リビューションを示す。

- 弱教師あり領域分割における異なる学習データのモデルの出力によるアンサンブル手法の提案。
- Web 画像を用いた弱教師あり領域分割の学習データ不足の改善。
- 広域の認識結果と局所領域の認識結果との一貫性からの学習手法の提案。
- 3. 関連研究
- 3.1 クラス分類器の可視化手法

認識結果の可視化においては,画像におけるクラス分類 に寄与した領域を推定する。クラス分類に寄与した領域と 領域分割における対象領域との間には相関があり,認識結 果の可視化は弱教師あり領域分割の手法として活用できる。 Zeiler ら [16] は畳み込み演算と同じ学習パラメータによる 逆畳み込み演算と逆 Pooling により,出力を入力空間に戻 した際に,物体の位置に対応するピクセルが強く応答する

¹ 電気通信大学

ことを示した。Simonyan ら [12] は, Zeilar らと類似した 手法で特定のクラスについての信号を逆伝搬させること で, CNN の認識結果に対するクラス応答を可視化させた。 CNN の可視化結果から弱教師あり領域分割が可能である ことを示した。派生手法に Guided Backpropagation [13] がある。

3.2 CNN の Activation を活用した弱教師あり領域分割

Oquab ら [5], Pathak ら [8] は FCN の最終層に Global Pooling(GP)を用いることで,出力マップをクラス分類の CNN と同じ次元に変換し,画像ラベルのみを用いて FCN を学習させた。GP により学習させた FCN は大まかな物 体の位置を推定することが可能であり,逆伝搬による可視 化とは異なる形で弱教師あり領域分割を実現した。その 後, Pinheiro ら [9] や Zhou ら [17] により, Global Pooling の派生手法についても研究がなされている。また, Chen ら [6], Pathak ら [7] は Global Pooling を用いずに,弱教 師ありで FCN の出力を直接学習させた。FCN の出力の各 ピクセルが画像ラベルに含まれていないクラスを出力して いる場合には修正を加えるような形で, FCN の出力と画 像ラベルから動的に領域の教師情報を生成し, これを用 いて FCN の学習を行った。一方で、[11] の手法では FCN を GP で学習した認識結果を可視化した。また, Simonyan ら [12] による Backward の可視化手法をマルチクラスの画 像に適応させた。においては,シングルクラスの場合に有 効であり,マルチクラスの場合に極端に精度が下がるとい う問題があったが,各クラスの逆伝搬値について差分をと ることでマルチクラスにも応用可能であることを示した。 類似研究として, Jianming ら [3] による研究がある。

3.3 領域分割結果の再学習手法

Weiら[15]は低次特徴量による物体顕著性マップを用 いて学習画像の領域分割を行い,その領域分割結果を領 域の教師情報として再学習を行った。Wei らの手法は単純 ながら既存の弱教師あり領域分割の精度を大きく上回っ た。Weiらの手法は単純ながら既存の弱教師あり領域分割 の精度を大きく上回り,弱教師あり領域分割の研究に大 きな変化が生まれた。特に,その後の弱教師あり領域分割 において,学習画像の領域分割手法(カテゴリは既知で ある画像の領域分割手法),ノイズを含む領域の教師情報 について堅牢な学習手法の重要性が増した。Chen ら [6], Pathak ら [7] は画像ラベルを活用して教師情報を動的に生 成していたが,その方向性が強くなったという見方もでき る。Kolesnikovら[4]は, Global Pooling [17] による学習 画像の大まかな位置推定結果について再学習を行った。ま た, Kolesnikov らは多くの領域の評価をスキップすること によるノイズに堅牢な学習, KL-Divergence による CRF による平滑化結果との近似などを用いて高精度な弱教師

あり領域分割を達成した。Saleh ら [10] は特徴マップにつ いて CRF を適用した結果から再学習を行い高精度を達成 した。Tokmakov ら [14] は動画から得られるモーションの Segmentation と Gaussian Mixture Model における動画の フレームの Segmentation 結果を用いて領域分割結果の再 学習を行った。本手法においては、物体の領域と相関のあ る情報の活用ではなく、一貫性による教師なしのアプロー チに取り組んだ。

4. 提案手法

4.1 Class Activation Mapping

本手法においては,クラス分類の可視化方法として Class Activation Mapping(CAM)用いる。Class Activation Mapping においては,Global Average Pooling(GAP) を用いてクラス分類ネットワークを学習することで,物体 の大まかな可視化を行うことができる。CAM においては, クラス分類モデルの出力マップをGAP によりベクトルに 変換し,これを全結合層により分類する識別器を学習する。 テスト時に,識別器の出力に対してGAP を行う前に全結 合層の重みをかけることで,出力マップにおける各クラス の分布マップを得ることができる。本手法では,この出力 マップにおける平均の2倍値を閾値として,各クラスの出 力マスクとした。

4.2 異なる学習データのモデルの出力によるアンサンブル

セクション 4.1 で得られた大まかな物体の領域を用い て,領域分割モデルを学習する。しかし,この際に用いる 教師情報は多くのノイズを含んでおり,そのままこれを 教師情報として学習しても十分な精度を得ることは難し い。そこで,本研究では異なる学習データのモデルの出力 によるアンサンブルによりノイズに堅牢な学習を目指す。 図1に提案手法の概要を示す。全体の学習データをSと する。このとき,この学習データを n 個にわけた場合に おけるサブセットを $S = (S_0, S_1, ..., S_n)$ とする。この各サ ブセットにつき,それぞれ領域分割用の子モデルを用意す る。各サブセットを教師情報として学習するそれぞれの 領域分割モデルのパラメータを $(heta_0, heta_1,..., heta_n)$,また,そ れぞれのモデルから得られる出力から生成したマスクを (m₀, m₁, ..., m_n) とする。本手法においてはこの異なるパ ラメータから得られたマスクについてアンサンブルを行う。 データ全体に対しての親モデルのパラメータを θ , 出力を $p(\theta)$ とする。各マスク m_i についての segmentaton loss を $s_{loss_i} = seg_{loss}(p(\theta), m_i)$ とする。ただし, クラスラベル のセットをC, yは確率分布pにおけるクラスである。こ のとき,生成マスクのアンサンブルにより得られるロスは 以下の式になる。

$$s_loss = \frac{1}{n} \sum_{i}^{n} s_loss_i \tag{1}$$

このアンサンブルにより得られるロスの最適化によりノイズに堅牢な学習が期待できる。



図 1 異なる学習データのモデルの出力によるアンサンプルについ ての概要。

4.3 Web 画像による学習データの補間

学習データをサブセットに分ける場合,各モデルの学習 データ数が 1/n になってしまうという問題がある。そこ で本研究では Web 画像を用いて全体の学習データを増や した。各クエリで収集した Web 画像についてクエリをラ ベルとしてこれを弱教師ありデータとして扱った。Bing API, Flickr APIを用いて, Web 画像を収集した。本研究 では Pascal VOC dataset における 20 クラスの物体につい てそれぞれ約 5000 枚ずつ,計10 万枚の画像を収集した。

4.4 広域領域の認識結果と局所領域の認識結果の一貫性 の学習

Web 画像は弱教師あり画像として扱うことができるが, Pascal VOC dataset などの実世界画像におけるデータセッ トとは異なる部分がある。特に,Web 画像は検索クエリを ラベルとして扱っているために,各画像につきラベルは一 つである。一方で,Pascal VOC のような実世界画像デー タのデータセットにおいては複数ラベルも含まれている。 また,物体のサイズも大きく異なっている傾向がある。こ のようなドメインの違いにより生まれる変化により,Web 画像を用いることで直接的に画像枚数の上昇を精度向上に つなげることは難しい。そこで,本研究ではWeb 画像と Pascal VOC データにおけるドメインの違いによる精度の 低下を,広域領域の認識結果と局所領域の認識結果の一貫 性の学習により抑制した。本手法の概要を図2に示した。

各入力画像を x, ランダムにクロップを行いスケールを 変化させた画像を x_s とする。本研究ではこれらの入力に おける出力で共通している領域の出力が同一になるよう に学習を行う。本手法における広域領域の認識結果と局 所領域の認識結果の一貫性の学習の概要を図に示した。xについての出力を p(x), x_s についての出力を $p(x_s)$ とす る。p(x) を $p(x_s)$ の出力について対応をとったマスクを $m(x|x_s)$, $p(x_s)$ を p(x) の出力について対応をとったマス クを $m(x_s|x)$ とする。このとき,一貫性のロス *c*-loss を以 下の式で定義する。

$$c \text{-loss} = \text{seg-loss}(p(x), m(x_s|x)) + \text{seg-loss}(p(x_s), m(x|x_s))$$
(2)

ただし,対応のない領域については Ignore label として評価を行わない。また, Pascal VOC の画像における物体のサイズが小さいことを考慮し,Web 画像のスケール変化はパディングによる縮小化,Pascal VOC 画像のスケール変化はクロップによる拡大を行う。これにより,物体のスケールの違い,複数ラベル物体の画像などドメインの違いによる過学習を抑制することが期待できる。



図 2 広域領域の認識結果と局所領域の認識結果の一貫性の学習に ついての概要。

5. 実験

領域分割のベンチマークのデータセットとして PASCAL VOC 2012 [1] を用いた。PASCAL VOC 2012 における 領域分割ベンチマークでは、20 の異なるクラスと Back Ground クラスを含む 21 のクラスの領域分割の精度を比較 する。また、1464 枚の train 画像、1449 枚の validation 画 像、1456 枚の test 画像があるが、近年は [2] により提供さ れた 10582 枚の train_aug 画像を用いるのが一般的となっ ている。本実験においてもこれを学習画像として用いた。 また追加画像として Web 画像 10 万枚を追加し、これも学 習画像として用いた。本手法における領域分割の出力結果 の例を図 3 に示す。

6. 結論

本手法においては,弱教師あり領域分割における異なる 学習データのモデルの出力によるアンサンブル手法の提 案を行った。上記の手法により学習データの不足が生じた が,これをWeb画像を用いて補完した。また,広域の認識 結果と局所領域の認識結果との一貫性を用いることでWeb 画像とPascal VOC dataのドメインの違いから生じる精 度低下を抑えた。今後の課題としては異なる学習データの モデルの出力によるアンサンブルを行うためのモデルの省 メモリ化,サリエンシーなどを活用したSeed 値の改善な どが考えられる。



図 3 Pascal VOC validation set における領域分割結果の例。13段目については成功例、4 段目については失敗例を示した。

謝辞

本研究は JSPS 科研費 17J10261 の助成を受けたもの です。

参考文献

- Everingham, M., Eslami, S. M. A., Van Gool, L., Williams, C. K. I., Winn, J. and Zisserman, A.: The Pascal Visual Object Classes Challenge: A Retrospective, *International Journal of Computer Vision*, Vol. 111, No. 1, pp. 98–136 (2015).
- [2] Hariharan, B., Arbelaez, P., Bourdev, L., Maji, S. and J., M.: Semantic contours from inverse detectors, *ICCV* (2011).
- [3] Jianming, Z., Zhe, L., Jonathan, B., Xiaouhui, S. and Sclaroff, S.: Top-down Neural Attention by Excitation Backprop, *ECCV* (2016).
- [4] Kolesnikov, A. and H.Lampert, C.: Seed, Expand and Constrain: Three Principles for Weakly-Supervised Image Segmentation, *ECCV* (2016).
- [5] Oquab, M., Bottou, L., Laptev, I. and Sivic, J.: Learning and Transferring Mid-level Image Representations Using Convolutional Neural Networks, *CVPR* (2014).
- [6] Papandreou, G., Chen, L.-C., Murphy, K. and Yuille, A. L.: Weakly-and semi-supervised learning of a dcnn for semantic image segmentation, *ICCV* (2015).
- [7] Pathak, D., Krahenbuhl, P. and Darrell, T.: Constrained convolutional neural networks for weakly supervised segmentation, *ICCV* (2015).
- [8] Pathak, D., Shelhamer, E., Long, J. and Darrell, T.:

Fully convolutional multi-class multiple instance learning, ICLR (2015).

- [9] Pedro, P. and Ronan, C.: From Image-level to Pixel-level Labeling with Convolutional Networks, *CVPR* (2015).
- [10] Saleh, F., Akbarian, M., Salzmann, M., Petersson, L., Gould, S. and M.Alvares, J.: Built-in Foreground/Background Prior for Weakly-Supervised Semantic Segmentation, ECCV (2016).
- [11] Shimoda, W. and Yanai, K.: Distinct Class Saliency Maps for Weakly Supervised Semantic Segmentation, *ECCV* (2016).
- [12] Simonyan, K., Vedaldi, A. and Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps, *ICLR WS* (2014).
- [13] Springenberg, J. T., Dosovitskiy, A., Brox, T. and Riedmiller, M.: Striving for Simplicity: The All Convolutional Net, *ICLR WS* (2015).
- [14] Tokmakov, P., Alahari, K. and Schmid, C.: Weakly-Supervised Semantic Segmentation using Motion Cues, *ECCV* (2016).
- [15] Wei, Y., Liang, X., Chen, Y., Shen, X., Cheng, M., Zhao, Y. and Yan, S.: STC: A Simple to Complex Framework for Weakly-supervised Semantic Segmentation, *ECCV* (2016).
- [16] Zeiler, M. and Fergus, R.: Adaptive deconvolutional networks for mid and high level feature learning, *ICCV* (2011).
- [17] Zhou, B., Khosla, A., Lapedriza, A. and Oliva, A. Torralba, A.: Learning Deep Features for Discriminative Localization, *CVPR* (2016).