

Computer Vision

Multi-label Image Classification

Group 25

Riajul Islam, Andreas Calonijs Kreth, Christine Midtgaard

Abstract—Here is a summary of the project and the conclusions.

Index Terms—Multi-label learning, deep learning, computer vision, multi-label classification, deep learning for MLC.

I. INTRODUCTION

Multi-label classification is the supervised learning problem where an instance may be associated with multiple labels. Image classification is a computer vision task that requires assigning a label or multiple labels to an image. Single-label classification, or multi-class classification, refers to the problem where an image contains only one object to be identified. However, natural images usually contain multiple objects or concepts, highlighting the importance of multi-label classification [1]. In this project we investigate two methods aimed to solve two different problems within multi-label learning: Query2Label targets the challenges of imbalance and object localization, whereas MLSPL addresses the challenge of training an effective multi-label classifier from minimal supervision.

This project focuses on testing existing solutions to two common issues that arise when training a network that assigns multiple labels to an input image: (i) the general multi-label learning issue accurately identifying multiple objects in an image under class imbalance, representing the complexity of real-world images, and (ii) the scenario where the training data underwent sparse supervision, which is often the case for data.

II. PROBLEM OVERVIEW

A. Related Work

Multi-label learning is a well studied problem within computer vision [2].

a) *Loss Functions*:

b) *PU learning*: Learning from positive and unlabeled data: a survey by Bekker and Davis [3].

Learning to Classify Texts Using Positive and Unlabeled Data by Li and Liu [4].

c) *Partially Observed Labels*:

d) *Heuristics*: Heuristics can be used to reduce the required annotation effort [34, 18], but this runs the risk of increasing error in the labels [2].

III. BACKGROUND

Contrary to binary or multi-class classification where each instance is associated with only one label, multi-label classification allows assigning multiple labels to a single input

[2]. Given K categories, an input image $x \in \mathcal{X}$ is associated with a binary vector of labels $y = [y_1, \dots, y_K]$ from the label space $\mathcal{Y} = \{0, 1\}^K$, where $y_k = 1$ represents that k is present in x and $y_k = 0$ otherwise. The goal is to create a model that outputs the probability of the presence of a category $p = [p_1, \dots, p_K]$ [2], [5].

A. Challenges in Multi-Label Learning

a) *Label Correlation*:

b) *Class Imbalance*:

c) *Sparse Supervision*:

B. Convolutional Neural Networks (CNNs)

In the field of computer vision there exists different deep learning methods that extract features from input images, and the two main neural networks are convolutional neural networks and vision transformers.

C. Transformer Architectures

Transformers, introduced by Vaswani et al. [6], are a type of neural network architecture initially developed to model long-range dependencies in sequence data. They achieve this using attention and self-attention mechanisms, which enables the model to have a long-term memory of inputs, and dynamically weigh the importance of each element in the input sequence. Originally, transformer-based models were mainly developed for natural language processing tasks, but the introduction of the Vision Transformer (ViT) by Dosovitskiy et al. [7], transformers have been widely used for computer vision tasks.

a) *Architecture*: The original transformer consists of an encoder-decoder structure. Both encoder and decoder are made of a stack of identical layers, where encoder layers contain a multi-head self-attention mechanism followed by a position-wise feed-forward network. In addition, decoder layers consist of an encoder-decoder attention layer.

The transformer encoder is of interest, as the authors of [5] use them to extract features and automatically learn label embedding, described in section TODO.

b) *Attention*: The core principle of the transformer architecture is the attention mechanism, which allows the model to attend to all positions in an input sequence when processing each element. Thus allowing the model to weigh the relevance of different positions in a sequence. Given a query and a set of key-value pairs, the attention function computes a weighted sum of the values, where the weights are determined by similarities between the query and the keys.

c) *Self-Attention*: The transformer model uses self-attention by relating every element in the input sequence to every other element. The attention function is a function that maps a query and a set of key-value pairs to an output. The scaled dot-product self-attention is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where Q, K, V , are the query, key, and value vectors, and d_k is the dimension of the queries and key vectors and serves as a scaling factor [6].

d) *Multi-Head Attention*:

e) *Vision Transformers (ViTs)*: The CvT by Wu et al. [8].

D. Loss Functions

a) *Binary Cross-Entropy*:

b) *Focal Loss*:

c) *Assymmetric Loss*:

d) *Online Estimation of Unobserved Labels (ROLE)*:

$$\mathcal{L}_{ROLE}(\mathbf{F}_B, \tilde{\mathbf{Y}}_B) = \frac{\mathcal{L}'(\mathbf{F}_B|\tilde{\mathbf{Y}}) + \mathcal{L}'(\tilde{\mathbf{Y}}|\mathbf{F}_B)}{2} \quad (2)$$

for a batch B , where $\tilde{\mathbf{Y}} \in [0, 1]^{N \times L}$ represent the estimated labels, and $\mathbf{F} \in [0, 1]^{N \times L}$ is the matrix of classifier predictions.

The goal is to train the label estimator $g(\cdot; \phi)$ and the image classifier $f(\cdot; \theta)$

IV. METHOD

In real-world datasets, obtaining full label annotations is practically impossible. This project focuses on finding solutions to the multi-label learning problems:

- The challenge of multi-label learning in scenarios where only a single positive label per image is available during training (MLSPL).
- Label imbalance (Query2Label).
- Feature localization (Query2Label).

A. Overview of Approach

A brief description of what we did to solve the multi-label classification problem.

B. Dataset

The MS-COCO 2014 [9] dataset is used as a benchmark for evaluation both Query2Label and MLSPL. MS-COCO (Microsoft Common Objects in Context) is a large-scale dataset commonly used for object detection, segmentation, and multi-label image classification. COCO consists of 82,081 training images and 80 classes, and a validation set of 40,137 images.

C. Query2Label: A Simple Transformer Way to Multi-Label Classification

Query2Label: A Simple Transformer Way to Multi-Label Classification (Query2Label) by Liu et al. [5] is a two-stage framework for multi-label classification. It uses transformer decoders to extract features with multi-head attentions focusing on different parts of an object category and learn label embeddings from data automatically.

D. Multi-Label Learning from Single Positive Labels

Multi-Label Learning from Single Positive Labels (MLSPL) by Cole et al. [2]. How it uses weak supervision and contrastive learning.

E. Experiments

A description of our setup versus the author's setups.

All experiments were run on a single NVIDIA GeForce RTX 3060 GPU (12 GB VRAM), with Python 3.11.8, NVIDIA driver 550.90.07 and CUDA 12.4.

a) *Query2Label*: All Query2Label experiments were conducted on environment versions `cuda==12.4`, `torch==2.1.0+cu121`, `torchvision==0.16.0+cu121`, `python==3.11.8`. We evaluated the Query2Label model using the best performing backbone model, the CvT-w24 backbone, with a 384 input resolution, pretrained on the ImageNet-22k dataset, as described by Liu et al. [5]. For our experiments, we used a pretrained model checkpoint released by the authors, which was trained on the MS-COCO 2024 dataset. The model was tested on the MS-COCO 2024 validation set with a batchsize of 16 and otherwise no fine-tuning or modification.

b) *Multi-Label Learning from Single Positive Labels*: All MLSPL experiments were conducted on environment versions `cuda==12.4`, `torch==2.2.1+cu121`, `torchvision==0.17.1+cu121`, `python==3.11.8`.

c) *Data preparation*: The dataset preparation for the MS-COCO dataset in MLSPL relies on converting the standard multi-label MS-COCO dataset to a single positive label format, simulating a weakly supervised setting. Cole et al. does this beginning with the fully annotated multi-label image dataset and corrupt it by discarding annotations, simulating single positive training data by randomly selecting one positive label for each training example [2].

To convert into a single positive label format, the instructions provided by Cole et al. are followed: firstly, both images and annotations for training and validation are downloaded. Secondly, pre-extracted features for COCO, provided by the authors are downloaded. Lastly, the script `format_coco.py` is used to produce uniformly formatted image lists and labels.

d) Hyperparameters:

F. Evaluation Metrics

To assess model performance, we adopt mean Average Precision (mAP), a standard metric widely reported in multi-label classification tasks as it is used to analyze the performance object detection and segmentation. Both Query2Label and MLSPL report results in terms of mAP.

V. RESULTS AND DISCUSSION

This section compares results from the experiments to the those of the papers. Discuss why they might not be the same, or describe the similarities. Discuss whether the methods are able to solve the problems.

TABLE I
COMPARISON OF MAP RESULTS BETWEEN OUR EXPERIMENTS AND REPORTED MAP RESULTS ON THE MS-COCO 2014 DATASET.

Method	Backbone	Resolution	mAP(Ours)	mAP (Paper)
Q2L-R101-448	ResNet-101	448×448	84.9	84.9
Q2L-R101-576	ResNet-101	576×576	86.5	86.5
Q2L-SwinL	Swin-L(22k)	384×384	90.5	90.5
Q2L-CvT	CvT-w24(22k)	384×384	91.3	91.3

A. Results from Query2Label

The results from our experiments are compared to those of the paper in table I.

B. Results from Multi-Label Learning from Single Positive Labels

VI. CONCLUSION

REFERENCES

- [1] T. Ridnik, G. Sharir, A. Ben-Cohen, E. Ben-Baruch, and A. Noy, “ML-decoder: Scalable and versatile classification head,” 2021. [Online]. Available: <https://arxiv.org/abs/2111.12933>
- [2] E. Cole, O. M. Aodha, T. Lorieul, P. Perona, D. Morris, and N. Jojic, “Multi-label learning from single positive labels,” 2021. [Online]. Available: <https://arxiv.org/abs/2106.09708>
- [3] J. Bekker and J. Davis, “Learning from positive and unlabeled data: a survey,” *Machine Learning*, vol. 109, no. 4, p. 719–760, Apr. 2020. [Online]. Available: <http://dx.doi.org/10.1007/s10994-020-05877-5>
- [4] X. Li and B. Liu, “Learning to classify texts using positive and unlabeled data,” 01 2003, pp. 587–594.
- [5] S. Liu, L. Zhang, X. Yang, H. Su, and J. Zhu, “Query2label: A simple transformer way to multi-label classification,” 2021. [Online]. Available: <https://arxiv.org/abs/2107.10834>
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” 2023. [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [7] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” 2021. [Online]. Available: <https://arxiv.org/abs/2010.11929>
- [8] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, “Cvt: Introducing convolutions to vision transformers,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.15808>
- [9] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, “Microsoft coco: Common objects in context,” 2015. [Online]. Available: <https://arxiv.org/abs/1405.0312>