

Supplementary Material of GATCluster: Self-Supervised Gaussian-Attention Network for Image Clustering

Chuang Niu¹, Jun Zhang², Ge Wang³, and Jimin Liang¹

¹ School of Electronic Engineering, Xidian University, Xi'an, Shaanxi 710071, China

² Tencent AI Lab, Shenzhen, Guangdong 518057, China

³ Rensselaire Polytechnic Institute, Troy, NY 12180, US

1 Visualization mapping function

For visualizing the clustering results in Figure 3, the mapping function maps the label feature $\mathbf{l}_i \in R^k$ to a two-dimensional vector $\mathbf{v}_i \in R^2$. Then, the two-dimensional vectors can be drawn on the figure and rendered by the ground-truth labels. Specifically, the mapping function is defined as

$$\mathbf{v}_i = \left[\sum_{h=i}^k l_{ih} \sin\left(\frac{2\pi h}{k}\right), \sum_{h=i}^k l_{ih} \cos\left(\frac{2\pi h}{k}\right) \right], \quad (\text{A1})$$

where $\|\mathbf{l}_i\|_1 = 1$ due to the probability constraint of label features. For the ImageNet-10, the cluster number k is set to 10.

2 Analysis on trivial solutions of DAC [1]

In this section, we give a formal analysis of why DAC [1] is prone to trivial solutions. Formally, we have the following conclusions:

If the optimal value of Eq. (1) is attained, for $\forall i, j, \mathbf{l}_i = \mathbf{l}_j \Leftrightarrow r_{ij} = 1$ and $r_{ij} = 0 \Rightarrow \mathbf{l}_i \neq \mathbf{l}_j$, but $|\{\mathbf{l}_i\}_{i=1}^N| = k$ cannot be guaranteed so that trivial solutions will be obtained.

Proof. If the optimal value of Eq. (1) is attained, for $\forall i, j$, we have:

$$\mathbf{l}_i \cdot \mathbf{l}_j = \begin{cases} 1, & \text{if } r_{ij} = 1, \\ 0, & \text{if } r_{ij} = 0, \end{cases} \quad (\text{A2})$$

For $\forall i, \|\mathbf{l}_i\|_2 = 1$ and $l_{ih} \geq 0$ ($h = 1, \dots, k$) are satisfied, where $\|\cdot\|_2$ implies L_2 -norm of a vector.

First, we verify that for $\forall i, j, \mathbf{l}_i = \mathbf{l}_j \Leftrightarrow r_{ij} = 1$ and $r_{ij} = 0 \Rightarrow \mathbf{l}_i \neq \mathbf{l}_j$ are satisfied. For arbitrary $\mathbf{l}_i, \mathbf{l}_j$, if $r_{ij} = 1$, we have $\mathbf{l}_i \cdot \mathbf{l}_j = 1$ and $\|\mathbf{l}_i\|_2 = \|\mathbf{l}_j\|_2 = 1$,

and we can obtain the following equation:

$$\begin{aligned}
\|\mathbf{l}_i - \mathbf{l}_j\|_2^2 &= \sum_{h=1}^k (l_{ih} - l_{jh})^2 \\
&= \sum_{h=1}^k (l_{ih} + l_{jh} - 2l_{ih}l_{jh}) \\
&= \sum_{h=1}^k l_{ih}^2 + \sum_{h=1}^k l_{jh}^2 - 2 \sum_{h=1}^k l_{ih}l_{jh} \\
&= \|\mathbf{l}_i\|_2^2 + \|\mathbf{l}_j\|_2^2 - 2\mathbf{l}_i \cdot \mathbf{l}_j \\
&= 1 + 1 - 2 \cdot 1 \\
&= 0.
\end{aligned} \tag{A3}$$

Thus, $\mathbf{l}_i = \mathbf{l}_j$ is satisfied if $r_{ij} = 1$. Similarly, if $r_{ij} = 0$, $\mathbf{l}_i \neq \mathbf{l}_j$ is always satisfied. That is,

$$\begin{aligned}
r_{ij} = 1 &\Rightarrow \mathbf{l}_i = \mathbf{l}_j, \\
r_{ij} = 0 &\Rightarrow \mathbf{l}_i \neq \mathbf{l}_j,
\end{aligned} \tag{A4}$$

Furthermore, due to $\forall i, \|\mathbf{l}_i\|_2 = 1$, we have the following proposition that if $\mathbf{l}_i = \mathbf{l}_j$ is satisfied, then

$$r_{ij} = \mathbf{l}_i \cdot \mathbf{l}_j = \mathbf{l}_i \cdot \mathbf{l}_i = 1. \tag{A5}$$

According to Eq. (A4) and Eq. (A5), we have:

$$r_{ij} = 1 \Leftrightarrow \mathbf{l}_i = \mathbf{l}_j. \tag{A6}$$

Eq. (A6) means that $\mathbf{l}_i = \mathbf{l}_j$ if and only if $\mathbf{x}_i, \mathbf{x}_j$ belong to the same clusters. And $r_{ij} = 0 \Rightarrow \mathbf{l}_i \neq \mathbf{l}_j$ implies that $\mathbf{l}_i \neq \mathbf{l}_j$ if $\mathbf{x}_i, \mathbf{x}_j$ belong to different clusters.

However, we assume that $\forall i \mathbf{l}_i \in E^k$ that denotes the standard basis of the k -dimensional Euclidean space. When $|\{\mathbf{l}_i\}_{i=1}^N| < k$, all above equations are true. For example, if $|\{\mathbf{l}_i\}_{i=1}^N| = 1$ meaning that all examples belong to the same cluster, while it is also the optimal solution of Eq. (1). Specifically, we assume that $\forall i \mathbf{l}_i = [1, 0, \dots, 0] \in R^k$, then the estimated $r_{ij} = 1$ being always true and it is easy to verify that this solution is an optimal one of Eq. (1). \square

3 Proof of Label Feature Theorem

In this section, we give formal proof of our Label Feature Theorem:

Label Feature Theorem. *If the optimal value of Eq. (2) is attained, for $\forall i, j, \mathbf{l}_i \in E^k, \mathbf{l}_i \neq \mathbf{l}_j \Leftrightarrow r_{ij} = 0, \mathbf{l}_i = \mathbf{l}_j \Leftrightarrow r_{ij} = 1$, and $|\{\mathbf{l}_i\}_{i=1}^N| = k$, where $|\cdot|$ denotes the cardinality of a set.*

Proof. If the optimal value of Eq. (2) is attained, we have:

1) For $\forall i, j$,

$$\mathbf{y}_i \cdot \mathbf{y}_j = \begin{cases} 1, & \text{if } r_{ij} = 1, \\ 0, & \text{if } r_{ij} = 0, \end{cases} \quad (\text{A7})$$

where for $\forall i$, $\mathbf{y}_i = \frac{\mathbf{l}_i}{\|\mathbf{l}_i\|_2}$, $\|\mathbf{l}_i\|_1 = 1$, $0 \leq l_{ih} \leq 1$ ($h = 1, \dots, k$) are satisfied, and $\|\cdot\|_2$ implies L_2 -norm of a vector so that $\|\mathbf{y}_i\|_2 = 1$.

2) For $\forall i$, the value of $\mathbf{l}_i \cdot \mathbf{l}_i$ is maximized.

First, we verify that for $\forall i, j$, $\mathbf{y}_i = \mathbf{y}_j \Leftrightarrow r_{ij} = 1$ is satisfied.

For arbitrary $\mathbf{y}_i, \mathbf{y}_j$, if $r_{ij} = 1$, we have $\mathbf{y}_i \cdot \mathbf{y}_j = 1$ and $\|\mathbf{y}_i\|_2 = \|\mathbf{y}_j\|_2 = 1$, and we can obtain the following equation:

$$\begin{aligned} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 &= \sum_{h=1}^k (y_{ih} - y_{jh})^2 \\ &= \sum_{h=1}^k (y_{ih} + y_{jh} - 2y_{ih}y_{jh}) \\ &= \sum_{h=1}^k y_{ih}^2 + \sum_{h=1}^k y_{jh}^2 - 2 \sum_{h=1}^k y_{ih}y_{jh} \\ &= \|\mathbf{y}_i\|_2^2 + \|\mathbf{y}_j\|_2^2 - 2\mathbf{y}_i \cdot \mathbf{y}_j \\ &= 1 + 1 - 2 \cdot 1 \\ &= 0. \end{aligned} \quad (\text{A8})$$

Thus, $\mathbf{y}_i = \mathbf{y}_j$ is satisfied if $r_{ij} = 1$. That is,

$$r_{ij} = 1 \Rightarrow \mathbf{l}_i = \mathbf{l}_j \quad (\text{A9})$$

Furthermore, due to $\forall i$, $\|\mathbf{y}_i\|_2 = 1$, we have the following proposition that if $\mathbf{y}_i = \mathbf{y}_j$ is satisfied, then

$$r_{ij} = \mathbf{y}_i \cdot \mathbf{y}_j = \mathbf{y}_i \cdot \mathbf{y}_i = 1. \quad (\text{A10})$$

According to Eq. (A9) and Eq. (A10), we have

$$r_{ij} = 1 \Leftrightarrow \mathbf{y}_i = \mathbf{y}_j \quad (\text{A11})$$

In the probability constraint, $\forall i, j, \sum_h l_{ih} = \sum_h l_{jh} = 1$ is always true, then we have:

$$\begin{aligned}
\mathbf{y}_i = \mathbf{y}_j &\Leftrightarrow \forall h, \frac{l_{ih}}{\|\mathbf{l}_i\|_2} = \frac{l_{jh}}{\|\mathbf{l}_j\|_2} \\
&\Leftrightarrow \forall h, l_{ih} = l_{jh} \frac{\|\mathbf{l}_i\|_2}{\|\mathbf{l}_j\|_2} \\
&\Leftrightarrow \sum_h l_{jh} \frac{\|\mathbf{l}_i\|_2}{\|\mathbf{l}_j\|_2} = \sum_h l_{jh} \\
&\Leftrightarrow \left(\frac{\|\mathbf{l}_i\|_2}{\|\mathbf{l}_j\|_2} - 1 \right) \sum_h l_{jh} = 0 \\
&\Leftrightarrow \frac{\|\mathbf{l}_i\|_2}{\|\mathbf{l}_j\|_2} = 1 \text{ (with } \frac{l_{ih}}{\|\mathbf{l}_i\|_2} = \frac{l_{jh}}{\|\mathbf{l}_j\|_2} \text{)} \\
&\Leftrightarrow \forall h, l_{ih} = l_{jh} \\
&\Leftrightarrow \mathbf{l}_i = \mathbf{l}_j
\end{aligned} \tag{A12}$$

According to Eq. (A11) and Eq. (A12), we have $r_{ij} = 1 \Leftrightarrow \mathbf{l}_i = \mathbf{l}_j$, meaning that $\mathbf{l}_i = \mathbf{l}_j$ if and only if $\mathbf{x}_i, \mathbf{x}_j$ belong to the same cluster.

Second, we verify that $\forall i, \mathbf{l}_i \in E^k$ that denotes the standard basis of the k -dimensional Euclidean space, and $r_{ij} = 0 \Leftrightarrow \mathbf{l}_i \neq \mathbf{l}_j$.

Due to $\forall i, h, 0 \leq l_{ih} \leq 1$, then $0 \leq l_{ih}^2 \leq l_{ih} \Rightarrow \mathbf{l}_i \cdot \mathbf{l}_i = \sum_h l_{ih}^2 \leq \sum_h l_{ih} = 1$. Therefore, if the optimal solution of Eq. (2) is attained, the value of $\mathbf{l}_i \cdot \mathbf{l}_i$ is maximized, i.e., $\mathbf{l}_i \cdot \mathbf{l}_i = 1$. Then, we verify the proposition that $\mathbf{l}_i \cdot \mathbf{l}_i = 1$ is satisfied if and only if \mathbf{l}_i is the one-hot vector, i.e., $\forall i, \exists m, l_{im} = 1, \forall h \neq m, l_{ih} = 0$.

Considering that $\forall i, \sum_h l_{ih} = 1$ and $\forall i, h, 0 \leq l_{ih} \leq 1$, its opposite proposition is that $\forall i, h, 0 \leq l_{ih} < 1$. We assume the opposite proposition is correct, then we have $\forall i, h, l_{ih}^2 < l_{ih} \Rightarrow \forall i, \mathbf{l}_i \cdot \mathbf{l}_i = \sum_h l_{ih}^2 < \sum_h l_{ih} = 1$. It is contradictory to that $\mathbf{l}_i \cdot \mathbf{l}_i = 1$, so the original proposition is true.

As $\forall i, \mathbf{l}_i$ is the one-hot vector, we can easily verify that $r_{ij} = 0 \Rightarrow \mathbf{l}_i \neq \mathbf{l}_j$ and $\mathbf{l}_i \neq \mathbf{l}_j \Rightarrow r_{ij} = 0$, thus $r_{ij} = 0 \Leftrightarrow \mathbf{l}_i \neq \mathbf{l}_j$.

Finally, we verify the proposition that $|\{\mathbf{l}_i\}_{i=1}^N| = k$ is satisfied.

We have the nonempty cluster constraint that $\forall h, p_h = \frac{1}{N} \sum_{i=1}^N l_{ih} > 0$, so that $\forall h, \exists i, l_{ih} > 0$. As we have already proved that $\forall i, \mathbf{l}_i \in E^k$, therefore, $|\{\mathbf{l}_i\}_{i=1}^N| = |E^k| = k$. □

4 Architecture details for all experiments

In this section, we describe the architecture details for each experiment. The proposed GATCluster consists of three components, i.e., the image feature module, the label feature module, and the attention module. The architectures of the image feature module for different experiments are summarized in Table A4. All of them are implemented by the VGG-style convolutional neural networks with

Table A1. Architecture of label feature module.

STL10/ImageNet-10/ImageNet-Dog/Cifar10/Cifar100-20
<p style="text-align: center;"> $1 \times 1-1-0-64$ Conv BN ReLU GlobalPooling Linear(k, k) ReLU Linear(k, k) ReLU Linear(k, k) Softmax </p>

Table A2. Architecture of attention module.

STL10/ImageNet-10/ImageNet-Dog/Cifar10/Cifar100-20
<p style="text-align: center;"> Linear(H*W, 3) ReLU Gaussian-kernel Channel-Multiplication GlobalPooling Linear(k, k) ReLU Linear(k, k) ReLU Linear(k, k) Softmax </p>

batch normalization. On different datasets, the architectures slightly differ from each other in the layer number, kernel size of the first layer and padding value. To evaluate the effectiveness of attention map resolution, we simply stacked extra convolutional layers to the last block and the feature map resolution is reduced due to the valid convolution, i.e., 3×3 kernel size with zero-padding. Similarly, to further increase the input image size (i.e., 160×160 and 192×192) and keep the attention-map-size unchanged (i.e., 6×6), we also stacked extra convolutional layers to the last block, which is similar to that of ImageNet-10-128. Without too much redundancy, the architecture details for the large input-size of 160×160 and 192×192 are not shown in the Table A4. For the label feature module and attention module, all experiments share the same network configuration, as shown in Table A1 and A2.

In this work, we mainly focus on the learning framework and the corresponding theory of the deep unsupervised clustering. We did not study the network architectures comprehensively. To further improve the clustering performance, we believe it is a promising way to search an effective architecture for deep clustering, especially, in an automatic manner.

5 Definitions of evaluation metrics

In this section, we introduce the definition of the evaluation metrics as following.

Accuracy (ACC). Given the predicted label q_i and the ground-truth label r_i of each sample, a mapping function O is first obtained by Hungarian algorithm to find the correspondence between q_i and r_i . If $r_i = O(q_i)$, the sample is correctly predicted, or otherwise, the sample is falsely predicted. Then, the ACC is calculated as:

$$ACC = \frac{\sum_{i=1}^N \delta(r_i, O(q_i))}{N} \quad (\text{A13})$$

where

$$\delta(a, b) = \begin{cases} 1, & a = b, \\ 0, & \text{otherwise,} \end{cases} \quad (\text{A14})$$

and N is the total number of samples.

Normalized Mutual Information (NMI). Given the predicted partition $U = \{U_1, \dots, U_R\}$ and the ground-truth partition $V = \{V_1, \dots, V_C\}$, the NMI

is defined by:

$$NMI = \frac{\sum_{i=1}^R \sum_{j=1}^C |U_i \cap V_j| \log \frac{N|U_i \cap V_j|}{|U_i||V_j|}}{\sqrt{(\sum_{i=1}^R |U_i| \log \frac{|U_i|}{N})(\sum_{j=1}^C |V_j| \log \frac{|V_j|}{N})}} \quad (\text{A15})$$

where N is the total number of samples, R and C mean the partition number, $|U_i|$ presents the sample number of the subset U_i .

Adjusted Rand Index (ARI). Given the predicted partition $U = \{U_1, \dots, U_R\}$ and the ground-truth partition $V = \{V_1, \dots, V_C\}$, the information on class overlap between the two partitions U and V can be obtained as shown in the Table A3, where n_{ij} denotes the number of sample of objects that are common to

Table A3. Contingency table.

Class	V_1	V_2	\dots	V_C	Sums
U_1	n_{11}	n_{12}	\dots	n_{1C}	$n_{1\cdot}$
U_2	n_{21}	n_{22}	\dots	n_{2C}	$n_{2\cdot}$
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
\cdot	\cdot	\cdot	\cdot	\cdot	\cdot
U_R	n_{R1}	n_{R2}	\dots	n_{RC}	$n_{R\cdot}$
Sums	$n_{\cdot 1}$	$n_{\cdot 2}$	\dots	$n_{\cdot C}$	$n_{\cdot\cdot} = N$

subset U_i and V_j , $n_{i\cdot}$ and $n_{\cdot j}$ denotes the number of samples in the subset U_i and U_j respectively, R and C mean the partition number, and N is the total number of samples. Then, the ARI is defined as:

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{n_{i\cdot}}{2} \sum_j \binom{n_{\cdot j}}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{n_{i\cdot}}{2} + \sum_j \binom{n_{\cdot j}}{2} \right] - \left[\sum_i \binom{n_{i\cdot}}{2} \sum_j \binom{n_{\cdot j}}{2} \right] / \binom{n}{2}}, \quad (\text{A16})$$

where the $\binom{n}{2}$ is the binomial coefficient.

6 More implementation details

The proposed GATCluster was implemented based on Pytorch. By adjusting the sub-batch size m_1 , our algorithm can fit for different GPUs with both large and relative small memory. For example, we conducted the GATCluster on several types of devices, including the GTX 1080 with 8G memory, the TITAN Xp and TITAN V with 12G memory, and the Tesla V100 with 32G memory. In this work, the models were trained and tested on a single GPU.

References

1. Chang, J., Wang, L., Meng, G., Xiang, S., Pan, C.: Deep adaptive image clustering. In: ICCV. pp. 5880–5888 (Oct 2017). <https://doi.org/10.1109/ICCV.2017.626>