Eduardo Morales, Hugo Jair Escalante

INAOE

Self-Supervised Learning

- Supervised learning requires a large number of labeled examples
- Semi-supervised learning can use some labeled examples with a larger set of unlabeled examples
- Usupervised learning tries to find structures in data without labels
- Self-supervised learning (SSL) receives unlabeled data which is annotated by the model
- It employs a surrogate objective to learn useful representations from unlabeled data

- Constrastive learning (CL) is a sub-family of SSL that enhances the performance on classification tasks by contrasting examples
- Learn an embedding space where similar sample pairs (same class) are close while dissimilar sample pairs (different classes) are far away
- The mechanism used is to move closer points of the same distribution (class) and pull apart points belonging to different distributions in an embedding space

- The general idea is that although we may not know information of a class, we can still distinguish between classes (e.g., elephants and tigers)
- Procedure: Select a data sampe, called *anchor*, select another data point belonging to the same distribution, called *positive* and select another data point belonging to a different distribution, called *negative*
- Minimize the distance between the anchor and the positive example and maximize the distance between the anchor and the positive in an embedding space



Negative: I-

 It is common to use one image as anchor and employ data augmentation, such as jittering, rotation, flipping, croping, add noise, and other random affine transformations as positives



- The distance function used in an embedding space can be Eucliden, Mahalanobis, Cosine, ...
- Most of the research work differs in terms of the loss function and number of used examples

Max margin constrastive loss:

$$\mathcal{L}_{contrastive}(x_i, x_j, \theta) = \mathbb{1}[y_i = y_j] \cdot \| f_{\theta}(x_i) - f_{\theta}(x_j) \|_2^2 + \mathbb{1}[y_i \neq y_j] \cdot \max(0, \epsilon - \| f_{\theta}(x_i) - f_{\theta}(x_j) \|_2^2)$$

where x_i and x_j are the samples with labels y_i and y_j , θ are the parameters of the embedding network, and ϵ defines the lowerbound distance between samples of different classes

 Triplet Loss: Use positive and negative samples to compute the loss

$$\mathcal{L}_{triplet}(x_a, x_+, x_-, \theta) = \sum_{\forall x} \max(0, \| f_{\theta}(x_a) - f_{\theta}(x_+) \|_2^2 - \| f_{\theta}(x_a) - f_{\theta}(x_-) \|_2^2 + \epsilon)$$

N-pair Loss: Is an extension of the triplet loss to multiple negative samples

$$\begin{split} \mathcal{L}_{N-\textit{pair}}(x^a, x^+, \{x_i^-\}_{i=1}^{N-1}, \theta) &= \\ &= \log\left(1 + \sum_{i=1}^{N-1} e^{f_{\theta}^T(x^a) \cdot f_{\theta}(x_i^-) - f_{\theta}^T(x^a) \cdot f_{\theta}(x^+)}\right) \\ &= -\log\left(\frac{e^{f_{\theta}^T(x^a) \cdot f_{\theta}(x^+)}}{e^{f_{\theta}^T(x^a) \cdot f_{\theta}(x^+)} + \sum_{i=1}^{N-1} e^{f_{\theta}^T(x^a) \cdot f_{\theta}(x_i^-)}}\right) \end{split}$$

which translates into a softmax loss function

- InfoNCE (Noise-Contrastive Estimation): Identifies a positive sample among a set of unrelated noise samples
- Given a context vector (c), the positive sample should be drawn from the conditional distribution (p(x|c)) while the N - 1 negative samples drawn from (p(s))
- The probability of detecting *x*⁺ correctly is:

$$\begin{split} \mathsf{P}(x^+|X,c) &= \frac{\mathsf{p}(x^+)\prod_{i=1,...,N; i\neq +}\mathsf{p}(x_i)}{\sum_{j=1}^{N}[\mathsf{p}(x_j|c)\prod_{i=1,...,N; i\neq +}\mathsf{p}(x_i)]} \\ &= \frac{\frac{\mathsf{p}(x^+|c)}{\mathsf{p}(x^+)}}{\sum_{j=i}^{N}\frac{\mathsf{p}(x_j|c)}{\mathsf{p}(x_j)}} = \frac{f(x^+,c)}{\sum_{j=i}^{N}f(x_j,c)} \end{split}$$

 InfoNCE optimizes the negative log probability of classifying the positive sample correctly:

$$egin{aligned} \mathcal{L}_{\mathit{InfoNCE}} &= -\mathbb{E}_X \cdot \log\left[rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight] \ f_k(x_{t+k}, c_t) \propto rac{oldsymbol{p}(x_{t+k} | c_t)}{oldsymbol{p}(x_{t+k})} \end{aligned}$$

In practice, most CL methods use:

$$\mathcal{L}(f) = -\mathbb{E}_{x,x^+,\{x_i^-\}} \left[-\ln \frac{e^{f^T(x)f(x^+)/\tau}}{e^{f^T(x)f(x^+)/\tau} + \sum_{i=1}^{N-1} e^{f^T(x)f(x_i^-)/\tau}} \right]$$

- Self-supervised constractive learning has no labels and generates positive examples using data augmentation techniques
- A clear problem is that if we have two different instances of the same class, it will generate positive examples for one instance (the anchor) and will consider the other instance as a negative example (two different dog breeds)
- Supervised contrastive learning has information of the instances of the same class

Supervised Contrastive Learning



SimCLR

Components:

- Several data augmentation techniques to transform any given example (image) into two correlated views of the same example
- An encoder that uses ResNet50 to obtain a vector representation (*f*(*x*))
- A small neural network (*g*(·)) that maps the vector representations into a common latent space where the contranstive loss is applied
- Loss function:

$$\mathcal{L}_{i,j} = -\log \frac{\exp(\textit{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\textit{sim}(z_i, z_k)/\tau)}$$

• They found that constrastive learning benefits from larger batch sizes and longer training times

SimCLR



Algorithm SimCLR

Algorithm 1 SimCLR's main learning algorithm. **input:** batch size N, constant τ , structure of f, g, \mathcal{T} . for sampled minibatch $\{\boldsymbol{x}_k\}_{k=1}^N$ do for all $k \in \{1, ..., N\}$ do draw two augmentation functions $t \sim T$, $t' \sim T$ # the first augmentation $\tilde{\boldsymbol{x}}_{2k-1} = t(\boldsymbol{x}_k)$ $h_{2k-1} = f(\tilde{x}_{2k-1})$ # representation $z_{2k-1} = q(h_{2k-1})$ # projection # the second augmentation $\tilde{\boldsymbol{x}}_{2k} = t'(\boldsymbol{x}_k)$ $\boldsymbol{h}_{2k} = f(\tilde{\boldsymbol{x}}_{2k})$ # representation $\boldsymbol{z}_{2k} = \boldsymbol{q}(\boldsymbol{h}_{2k})$ # projection end for for all $i \in \{1, \ldots, 2N\}$ and $j \in \{1, \ldots, 2N\}$ do $s_{i,i} = \mathbf{z}_i^{\top} \mathbf{z}_i / (\|\mathbf{z}_i\| \|\mathbf{z}_i\|)$ # pairwise similarity end for define $\ell(i, j)$ as $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{\{k \neq i\}} \exp(s_{i,k}/\tau)}$ $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^{N} \left[\ell(2k-1,2k) + \ell(2k,2k-1) \right]$ update networks f and q to minimize \mathcal{L} end for **return** encoder network $f(\cdot)$, and throw away $q(\cdot)$

NNCLR

- Use positive examples (nearest neighbors) from other instances in the dataset as positives rather than augmenting the same image
- The NNCLR loss function:

$$\mathcal{L}_{i}^{\textit{NNCLR}} = -\log rac{\exp(\textit{NN}(z_{i}, \textit{Q}) \cdot z_{i}^{+} / au)}{\sum_{k=1}^{\textit{N}} \exp(\textit{NN}(z_{i}, \textit{Q}) \cdot z_{k}^{+} / au)}$$

where NN(z, Q) is the nearest neighbor operator defined as:

$$\mathit{NN}(z, Q) = rgmin_{q \in Q} \parallel z - q \parallel_2$$

• It minimizes the average loss over all the elements in the mini-batch: $\mathcal{L}^{NNCLR} = \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i}^{NNCLR}$

NNCLR







ORE

- Open-World Object Detection, identifies new objects without supervision in an incremental way
- It uses contrastive clustering in the latent space
- For each known class *i* ∈ *K*, it stores a prototype vector *p_i*, let *f_c* ∈ *R^d* be a feature vector for an object of class *c*, the contrastive loss is:

$$\mathcal{L}_{cont}(f_c) = \sum_{i=0}^{C} I(f_c, p_i)$$

where

$$I(f_c, p_i) = \left\{ egin{array}{ll} D(f_c, p_i) & i = c \ \max\{0, riangle - D(f_c, p_i)\} & ext{otherwise} \end{array}
ight.$$

 Where *D* is any distance function and △ defines how close a similar and dissimilar item can be

(INAOE)



- The prototypes are created by the mean of the feature vectors of each class
- When new prototypes are added, the prototypes are updated with a momentum parameter
- They use Faster R-CNN (just the region proposal network) to generate bounding boxes with a corresponding objectness score
- Then select the top-k region proposals, based on their objectness scores, as unknown objects

ORE

• Transforms the classifiction head into an energy function:

$$E(f;g) = -T \log \sum_{i=1}^{C} \exp(\frac{g_i(f)}{T})$$

- And models the energy distribution of the known and unknown energy values as Weibull distributions ψ_{kn}(f) and ψ_{unkn}(f)
- The predictions are labelled as unknowns if ψ_{kn}(f) < ψ_{unkn}(f)
- They store a balanced set of examplars for each class for finetuning the model

ORE



МоСо

- Some researchers have stored, what is caled, a *memory bank* which consists on the embedding of the data for unsupervised learning with contrastive loss
- Momentum Contrast (MoCo) maintains a dictionary as a queue of data samples, where the samples in the dictionary are progressively replaced, the current mini-batch is added and the oldest mini-batch is removed
- Denoting the parameters of f_k (key) and θ_k and those of f_q (query) and θ_q , θ_k is updated as; $\theta_k \leftarrow m\theta_k + (1 - m)\theta_q$
- Only the parameters of θ_q are updated by backpropagation

MoCo

They use the InfoNCE contrastive loss function:

$$\mathcal{L}_{m{q}} = -\log rac{\exp(m{q}\cdotm{k}_+/ au)}{\sum_{i=o}^{K}\exp(m{q}\cdotm{k}_i/ au)}$$

the sum is over one positive and K negatives

- Can be used for different pretext tasks: images, patches, set of patches
- The networks f_q and f_k can be identical, partially shared or different

МоСо



PCL

- In instance contrastive learning, they learn an embedding with a constractive loss function (e.g., InfoNCE) that includes one positive and N negative samples
- Prototypical Contrastive Learning (PCL) uses prototypes and bridges contrastive learning with clustering
- Perform *k*-means on the features *v*[']_i = *f*_θ['](*x*_i) given the momentum encoder to obtain *k* clusters
- Prototype c_i is the centroid of the *i*-th cluster
- Uses a similar loss as InfoNCE, considering *r* negative prototypes, and clustering the samples *M* times with different number of clusters: $K = \{k_m\}_{m=1}^{M}$

PCL

$$\mathcal{L}_{ProtoNCE} = \sum_{i=1}^{n} - \left(\log \frac{e^{(v_i \cdot v_i'/\tau)}}{\sum_{j=0}^{r} e^{(v_i \cdot v_j'/\tau)}} + \frac{1}{M} \sum_{m=1}^{M} \log \frac{e^{(v_i \cdot c_s^m/\phi_s^m)}}{\sum_{j=0}^{r} e^{(v_i \cdot c_j^m/\phi_j^m)}}\right)$$

- · E-step: finds the distribution of prototypes via clustering
- · M-step: optimizes the network via contrastive learning

PCL



CURL

- Contrastive Unsupervised Representation for Reinforcement Learning (CURL) uses contrastive learning to maximize agreement between augmented versions of the same observation (temporally sequential frames)
- Main idea: Extract high-level features from raw images using constractive learning and performs off-policy control on top of the extracted features
- Instance discrimination is performed over frame stacks (instead of a single image)
- Employ random augmentation (cropping) across the batch retaining spatio-temporal information (anchor and positive)

CURL

- The constrastive representation is trained jointly with the RL algorithm
- The similarity measure is: $sim(q, k) = q^T Wk$, where W is a learned parameter matrix
- InfoNCE loss learns encoders f_q and f_k mapping raw anchors (query) x_q and targets (keys) x_k into latents $q = f_q(x_q)$ and $k = f_k(x_k)$ on which to apply the similarity dot product, where it is common to have $f_q = f_k$
- They used momentum contrastive (MoCo) for encoding the keys: $\theta_k = m\theta_k + (1 m)\theta_q$











- Learns from scratch on a database of 400 million (image,text) pairs
- Jointly train an image CNN and a text transformer to predict the caption of an image
- For each batch of N (image,text) pairs, CLIP is trained to predict all N × N pairs learning image and text encoders to maximize cosine similarity of real pairs and minimize cosine similarity of incorrect pairings

CLIP

(1) Contrastive pre-training



(2) Create dataset classifier from label text



- It has been used for audio, text and speech
- There are some works on multimodal domains
- In unsupervised learning it is still an open research question how to distinguish between positive and negative samples

References

- A. Radford, J.W. Kim, C. Hallacy, A. Ramesh, G. Goh. S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, I. Sutskever (2021). Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020v1 [cs.CV]
- L. Weng (2021). Contrastive representation learning. https://lilianweng.github.io/posts/2021-05-31-contrastive/
- P. Kumar, P. Rawat, S. Chauhan (2022). Contrastive self-supervised learning: review, progress, challenges and future research directions. International Journal of Multimedia Information Retrieval 11:461–488
- A. Srinivas, M. Laskin, P. Abbeel (2020). CURL: Contrastive Unsupervised Representations for Reinforcement Learning arXiv:2004.04136v4 [cs.LG]
- A. van den Oord. Y. Li, O. Vinyals (2019). Representation Learning with Contrastive Predictive Coding arXiv:1807.03748v2 [cs.LG]
- K. He, H. Fan, Y. Wu, S. Xie, R. Girshick (2020). Momentum Contrast for Unsupervised Visual Representation Learning CVPR 2020, 9729-9738
- D. Dwibedi, Y. Aytar, J. Tompson, P. Sermanet, A. Zisserman (2021). With a Little Help from My Friends: Nearest-Neighbor Contrastive Learning of Visual Representations. arXiv:2104.14548v2 [cs.CV]
- K.J. Joseph, S. Khan, F.S. Khan, V.N. Balasubramanian (2021). Towards Open World Object Detection arXiv:2103.02603v2 [cs.CV]
- J. Li, P. Zhou, C. Xiong, S.C.H. Hoi (2021). Prototypical Contrastive Learning of Unsupervised Representations. arXiv:2005.04966v5 [cs.CV]
- T. Chen, S. Kornblith, M. Norouzi, G. Hinton (2020). A Simple Framework for Contrastive Learning of Visual Representations. arXiv:2002.05709v3 [cs.LG]