

Bootstrap Your Own Latent (BYOL) A New Approach to Self-Supervised Learning

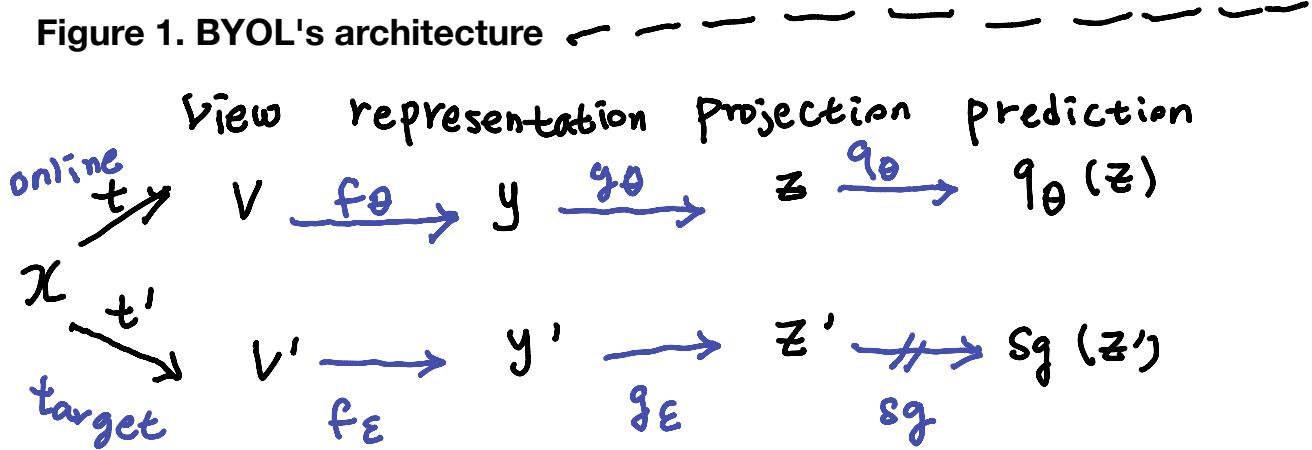
0. Abstract

- A new approach to self-supervised image representation learning
- online and target networks interact and learn from each other
- from an augmented view of an image, we train the online network to predict the target network representation of the same image under different augmented view
- update the target network with a slow-moving average of the online network
- achieves a new state of the art without negative pairs.

1. Introduction

- how to learn good image representations?
- state-of-the-art contrastive methods : by reducing distance between positive pairs, and increasing the distance between negative pairs.
- BYOL achieves higher performance without using negative pairs.
- iteratively bootstraps the output of a network to serve as targets for an enhanced representation
- more robust to the choice of image augmentations than contrastive methods
- start from an augmented view of an image
- train its online network to predict the target network's representation of another augmented view of the same image
- use slow-moving average of the online network as the target network

Figure 1. BYOL's architecture



minimize similarity loss $q_{\theta}(z)$ and $sg(z')$

θ trained weights

ϵ exponential moving average of θ

sg stop-gradient

everything but $f\theta$ is discarded at the end of training

y is used as image representation



3. Method

- Is negative examples indispensable to prevent collapsing?
- from a given representation (target), we can train a new, potentially enhanced representation (online), by predicting the target representation
- use subsequent online networks as new target networks for further training
- use a slowly moving exponential average of the online networks as the target network instead of fixed checkpoints.
- In deep RL, target networks stabilize the bootstrapping updates provided by the Bellman equation

3.1 Description of BYOL

- goal: learn a representation y
- which can then be used for downstream tasks
- online network: defined by a set of weights θ

encoder $f\theta$

projector $g\theta$

predictor $q\theta$

- target network: defined by a set of weights ϵ
- provides the regression targets to train the online network

\mathcal{E} : exponential moving average of θ

γ : target decay rate $\in [0, 1]$

after each training step, we perform the update

$$\mathcal{E} \leftarrow \gamma \mathcal{E} + (1 - \gamma) \theta$$

procedure of BYOL

① T, T' : Two distributions of image augmentations

$$v = t(x) \rightarrow y = f_\theta(v) \rightarrow z_\theta = g_\theta(y) \rightarrow q_\theta(z_\theta)$$

$$v' = t'(x) \rightarrow y' = f_\epsilon(v') \rightarrow z'_\epsilon = g_\epsilon(y')$$

② l₂-normalize

$$\overline{q_\theta}(z_\theta) = q_\theta(z_\theta) / \|q_\theta(z_\theta)\|_2$$

$$\overline{z'_\epsilon} = z'_\epsilon / \|z'_\epsilon\|_2$$

③ define mean squared error between $\overline{q_\theta}(z_\theta)$ and $\overline{z'_\epsilon}$

$$L_\theta^{\text{BYOL}} = \|\overline{q_\theta}(z_\theta) - \overline{z'_\epsilon}\|_2^2 = 2 - 2 \frac{\langle \overline{q_\theta}(z_\theta), \overline{z'_\epsilon} \rangle}{\|\overline{q_\theta}(z_\theta)\|_2 \cdot \|\overline{z'_\epsilon}\|_2}$$

④ symmetrize the loss L_θ^{BYOL} by separately feeding v' to online network and v to the target network to compute $\tilde{L}_\theta^{\text{BYOL}}$

⑤ at each training step, perform stochastic optimization step to minimize $L_\theta^{\text{BYOL}} + \tilde{L}_\theta^{\text{BYOL}}$ respect to θ

⑥ only keep the encoder f_θ

Appendix A. Algorithm

Input:

D, T, T' set of images and distributions of transformations

$\theta, f_\theta, g_\theta, q_\theta$ initial online parameters, encoder, projector, predictor

$\varepsilon, f_\varepsilon, g_\varepsilon$ initial target parameters, target encoder, target projector

optimizer updates online parameters using the loss gradient

K, N total number of optimization steps and batch size

$\{\eta_k\}_{k=1}^K$ target network update schedule, learning rate schedule

for $k=1$ to K do

$B \leftarrow \{x_i \sim D\}_{i=1}^N$ sample a batch of N images

for $x_i \in B$ do

$t \sim T, t' \sim T'$ sample image transformations

$z_1 \leftarrow g_\theta(f_\theta(t(x_i))), z_2 \leftarrow g_\theta(f_\theta(t'(x_i)))$ compute projections

$z'_1 \leftarrow g_\varepsilon(f_\varepsilon(t'(x_i))), z'_2 \leftarrow g_\varepsilon(f_\varepsilon(t(x_i)))$

end

$$l_i \leftarrow -2 \cdot \left(\frac{\langle q_\theta(z_1), z'_1 \rangle}{\|q_\theta(z_1)\|_2 \cdot \|z'_1\|_2} + \frac{\langle q_\theta(z_2), z'_2 \rangle}{\|q_\theta(z_2)\|_2 \cdot \|z'_2\|_2} \right)$$

$\delta_\theta \leftarrow \frac{1}{N} \sum_{i=1}^N \partial_\theta l_i$ compute the total loss gradient

$\theta \leftarrow \text{optimizer}(\theta, \delta_\theta, \eta_k)$ update online parameters

$\varepsilon \leftarrow \gamma_k \varepsilon + (1 - \gamma_k) \theta$ update target parameters

output: encoder f_θ

3.2 Implementation details

1) Image augmentations

- use the same set of image augmentations as SimCLR
- resized to 224 x 224

2) Architecture

- base parametric encoder: convolutional residual network with 50 layers and post-activation (ResNet-50 (1x) v1)
- representation corresponds to the output of the final average pooling layer, which has feature dimension of 2048
- projector, predictor: MLP (linear layer with output size 4096 followed by batch normalization, ReLU, and a final linear layer with output dimension 256)

3) Optimization

- LARS optimizer with a cosine learning rate schedule over 1000 epochs
- LearningRate = $0.2 \times \text{BatchSize} / 256$
- a global weight decay parameter of 1.5×10^{-6}
- the exponential moving average parameter starts from 0.996 and is increased to one during training
- batch size 4096

4. Experimental evaluation

- BYOL's representation after self-supervised pretraining on the training set of the ImageNet ILSVRC-2012 dataset

1) Linear evaluation on ImageNet

- evaluate BYOL's representation by training a linear classifier on top of the frozen representation
- 74.3% top 1 accuracy, 91.6% top 5 accuracy

2) Semi-supervised training on ImageNet

- fine-tune BYOL's representation on a classification task with a small subset of ImageNet's train set, using label information
- consistently outperforms previous approaches

3) Transfer to other classification tasks

- evaluate representation on other classification dataset to assess whether the features learned on ImageNet are generic
- perform linear evaluation and fine-tuning on the set of classification task with dataset CIFAR, SUN397, VOC2007 etc.
- outperforms SimCLR on all benchmarks

4) Transfer to other vision tasks

- evaluate representation on different tasks (semantic segmentation, object detection, depth estimation)
- assess whether BYOL's representation generalizes beyond classification tasks
- VOC2012 semantic segmentation task: classify each pixel in the image
- object detection using a Raster R-CNN architecture: fine tuen on trainval2007, report results on test2007
- depth estimation on NYU v2 dataset: depth map of a scene is estimated given a single RGB image
- Depth prediction measures how well a network represents geometry, and how well that information call be localized to pixel accuracy