

Lecture 19: Representation Learning from Spectral Decomposition view

Lecturer: Bo Dai

Scribes: Vidhya Kewale & Sandilya Sai Garimella

Note: *LaTeX template courtesy of UC Berkeley EECS Department.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

19.1 Recap

SimCLR (**S**imple **C**ontrastive **L**earning of Visual **R**epresentations)

- self-supervised contrastive learning method for visual representations
- works with a single modality: images
- trains the model to bring augmented views of the same image closer together in the embedding space, while pushing apart representations of different images

CLIP (**C**ontrastive **L**anguage-**I**mage **P**re-training)

- extends contrastive learning to multiple modalities using image-text pairs
- learns aligned representations across the visual and textual domains
- trains the model to maximize the cosine similarity between embeddings of matching image-text pairs, while minimizing it for non-matching pairs

Energy-Based Models (EBM) and Noise-Contrastive Estimation (NCE)

1. conditional probability in the same modality

$$p(x' | x) = p(x') \exp(\varphi(x)^\top \varphi(x')) \quad (1.1)$$

2. cross-modal conditional probabilities

$$p(y | x) = p(y) \exp(\varphi(x)^\top \nu(y)) \quad (1.2)$$

$$p(x | y) = p(x) \exp(\varphi(x)^\top \nu(y)) \quad (1.3)$$

19.2 SimCLR

All of them use ranking-based NCE to estimate a special EBM. SimCLR is as follows:

$$p(x' | x) = p(x') \exp(\varphi(x)^\top \varphi(x')) \quad (2.1)$$

$$D = \{(x_i, x'_i, (x_i^{1'}, \dots, x_i^{k'}))\}_{i=1}^n \quad (2.2)$$

We formulate the loss function using the follows:

$$\max_{\varphi_\theta} \sum_{i=1}^n \left[\varphi(x_i)^\top \varphi(x'_i) - \log \sum_{j=1}^k \exp(\varphi(x_i)^\top \varphi(x_i^{j'})) \right] \quad (2.3)$$

$$\max_{\varphi} f(\theta); \quad D_\theta l(\theta) = \sum_{i=1}^n \varphi(x_i)^\top \varphi(x'_i) \cdot (\nabla_\theta \varphi(x_i) + \nabla_\theta \varphi(x'_i)) \quad (2.4)$$

$$= \sum_{i=1}^n \left(\frac{\sum_{j=1}^k \exp(\varphi(x_i)^\top \varphi(x_i^{j'})) \cdot (\nabla_\theta \varphi(x_i) + \nabla_\theta \varphi(x_i^{j'}))}{\sum_{j=1}^k \exp(\varphi(x_i)^\top \varphi(x_i^{j'}))} \right) = \mathcal{O}(nk) \quad (2.5)$$

Computation cost is $\mathcal{O}(nk)$. We typically also use $k=n$, therefore the computation cost becomes $\mathcal{O}(n^2)$.

To further demonstrate with data, assume we have:

$$\{x_i\}_{i=1}^B \sim \{x'_i\}_{i=1}^B \quad (2.6)$$

Therefore, the computation cost will be:

$$i, \{x - i\} \quad (B - 1) \Rightarrow \sim \mathcal{O}(B^2) \quad (2.7)$$

To circumvent this quadratic computation cost, we can use a binary-based NCE instead of a ranking-based NCE. With this, instead of $\mathcal{O}(nk)$, we can get $\mathcal{O}(2B) \sim \mathcal{O}(B)$.

Coming back to this expression, to derive spectral learning and Bootstrap your own latent (BYOL):

$$p(x^* | x) = p(x') \exp(\varphi(x)^\top \varphi(x)) \quad (2.8)$$

We remove the exponential because it makes the gradient calculation harder:

$$p(x' | x) = p(x') \varphi(x')^\top \varphi(x) \quad (2.9)$$

The L2 loss function is now defined as:

$$l_2 \int \left\| p(x' | x) - p(x') \varphi(x')^\top \varphi(x) \right\|^2 dx dx' \quad (2.10)$$

$$= \int p(x' | x)^2 dx dx' - 2 \int p(x' | x) p(x') \quad (2.11)$$

$$p(x')^\top \varphi(x) dx dx' \quad (2.12)$$

We know $p(x' | x) p(x) = p(x') p(x) p(x')^\top p(x)$ from $p(x' | x) = p(x') \varphi(x')^\top \varphi(x)$:

$$\int \left\| \frac{p(x', x)}{\sqrt{p(x)} \sqrt{p(x')}} \sqrt{p(x')} \sqrt{p(x)^2} \varphi(x')^\top \varphi(x) \right\|^2 dx dx' \quad (2.13)$$

$$= \int \left(\frac{p(x', x)}{\sqrt{p(x)} \sqrt{p(x')}} \right)^2 dx dx' - 2 \int (p(x', x) p(x')^\top \varphi(x)) dx dx' + \int p(x') p(x) (\varphi(x')^\top \varphi(x))^2 dx dx' \quad (2.14)$$

We observe that the terms in the integrals can be simplified using the definition of expectation; therefore we can apply sampling here. The above simplifies to:

$$= -2 \mathbb{E}_{p(x, x')} [\varphi(x')^\top \varphi(x)] + \mathbb{E}_{p(x) p(x)} [(\varphi(x')^\top \varphi(x))^2]. \quad (2.15)$$

From above, we can see that we sample only once but can use it for computing both expectation terms.

$$p(x', x) = p(x) \varphi(x)^\top p(x') \varphi(x') \quad (2.16)$$

but we write this as

$$p(x', x) = \Psi(x)^\top \Psi(x') \quad (2.17)$$

This is called the Eigen-decomposition spectral perspective of representation.

19.2.1 BYOL w/o ν

The loss function is, using similar reason to above:

$$\min_{\varphi, \nu} \int \left(\frac{\rho(x', x)}{\sqrt{\rho(x')} \sqrt{\rho(x)}} - \sqrt{\rho(x')} \sqrt{\rho(x)} \nu(x')^\top \varphi(x) \right)^2 dx dx' \quad (2.18)$$

Alternative Optimization

Add a constraint such that $\nu = \varphi$.

$$\text{(min problem above)} \propto 2\mathbb{E}_{p(x',x)} [\nu(x)^T \varphi(x)] - \mathbb{E}_{p(x')} [\varphi(x')^T \mathbb{E}_{p(x)} [\varphi(x)\varphi(x)^T] \varphi(x')] \quad (2.19)$$

With the above expanded, we can do separate sampling.

$$\Lambda_t = \mathbb{E}_{p(x)} [\nu_{\Psi}(x)\nu_{\Psi}(x)^T] \quad (2.20)$$

$$-2\mathbb{E}_{p(x,x')} [\varphi(x')\nu(x)^T] + \mathbb{E}_{p(x)} [\varphi(x)^T \Lambda_t \varphi(x)] \quad (2.21)$$

19.3 PCA

Finding the maximal eigenspace while matching the y 's are different.

We have the following, noting that the trace operator is invariant under cyclic permutations:

$$\hat{\rho} = (x, x') \in \mathbb{R}^{n \times n}, n \text{ samples} \quad (3.1)$$

$$\Psi(x) \in \mathbb{R}^{n \times d} \quad (3.2)$$

$$\mathbb{E}_{p(x,x')} [\Psi(x)\Psi(x')^T] \quad (3.3)$$

$$\mathbb{E}_{p(x)} [\Psi(x)^T \Psi(x)] = I_{d \times d} \quad (3.4)$$

Penalty method:

$$\max_{\Psi} \mathbb{E}_{p(x,x')} [\Psi(x)\Psi(x)^T] - \lambda \cdot \text{trace}(\mathbb{E}_{p(x)} (\Psi(x)\Psi(x)^T) - I)^2 \quad (3.5)$$