

Lecture 19: Spectral Representation Learning

Lecturer: Bo Dai Scribes: Sandilya Sai Garimella & Vidhya Kewale

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 [4]

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 [3]

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

1. conditional probability in the same modality

$$p(x'|x) = p(x') \exp(\varphi(x)^T \varphi(x'))$$

2. cross-modal conditional probabilities

$$p(y|x) = p(y) \exp(\varphi(x)^T \nu(y))$$

$$p(x|y) = p(x) \exp(\varphi(x)^T \nu(y))$$

19.2 SimCLR

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

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

We formulate the loss function using the follows:

$$\begin{aligned} \max_{\varphi} \sum_{i=1}^n \left[\varphi(x_i)^\top \varphi(x'_i) - \log \sum_{j=1}^k \exp(\varphi(x_i)^\top \varphi(x'_j)) \right] \\ \max_{\varphi} f(\theta); \quad D_{\theta} l(\theta) = \sum_{i=1}^n \varphi(x_i)^\top \varphi(x'_i) \cdot (\nabla_{\theta} \varphi(x_i) + \nabla_{\theta} \varphi(x'_i)) \\ \sum_{i=1}^n \left(\frac{\sum_{r=1}^k \exp(\varphi(x_i)^\top \varphi(x'_r)) \cdot (\nabla_{\theta} \varphi(x_i) + \nabla_{\theta} \varphi(x'_r))}{\sum_{r=1}^k \exp(\varphi(x_r)^\top \varphi(x_i))} \right) = O(nk) \end{aligned}$$

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

The explanation with data is we have:

$$\begin{aligned} \{x_i\}_{i=1}^B \\ \sim \{x'_i\}_{i=1}^B \end{aligned}$$

Therefore the computation cost will be:

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

To circumvent this quadratic computation cost, we can use a binary-based NCE instead of a ranking-based NCE. With this, instead of $O(nk)$, we can get $O(2B) \sim O(B)$. This was presented in the paper cited here [5].

Coming back to this expression to derive spectral learning and Bootstrap your own latent (BYOL) [2]:

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

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

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

The L2 loss function is now defined as:

$$\begin{aligned} l_2 \int \left\| p(x' | x) - p(x') \varphi(x')^\top \varphi(x) \right\|^2 dx dx' \\ = \int p(x' | x)^2 dx dx' - 2 \int p(x' | x) p(x') \varphi(x')^\top \varphi(x) dx dx' \end{aligned}$$

$$\begin{aligned}
p(x' | x) p(x) &= p(x') p(x) p(x')^\top p(x) \text{ from} \\
p(x' | x) &= p(x') \varphi(x')^\top \varphi(x) \\
\int \left\| \frac{p(x', x)}{\sqrt{p(x)} \sqrt{p(x')}} \sqrt{p(x')} \sqrt{p(x)^2} \varphi(x')^\top \varphi(x) \right\|^2 dx dx' \\
&= \int \left(\frac{p(x', x)}{\sqrt{p(x)} \sqrt{p(x')}} \right)^2 dx dx' - 2 \int (p(x', x) p(x')^T \varphi(x)) dx dx' + \\
&\quad \int p(x') p(x) (\varphi(x')^T \varphi(x))^2 dx dx'
\end{aligned}$$

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:

$$= -2E_{p(x, x')} [\varphi(x')^\top \varphi(x)] + E_{p(x p(x))} [\varphi(x')^\top \varphi(x)]^2.$$

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

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

but we write this as

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

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:

$$\min_{\varphi, \nu} \int \left\| \left(\rho(x', x) \frac{\varphi(x')^T \varphi(x)}{\sqrt{\rho(x') \rho(x)}} \sqrt{\rho(x')} \sqrt{\rho(x)} \right) \right\|^2 dx' dx$$

Alternative Optimization

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

$$(\text{min problem above}) \propto 2E_{p(x', x)} [\nu(x)^T \varphi(x)] - E_{p(x', x)} [\varphi(x')^T \varphi(x) \varphi(x)^\top \varphi(x')]$$

With the above expanded, we can do separate sampling.

$$\begin{aligned}
\Lambda_t &= E_{p(x)} [\nu_\Psi(x) \nu_\Psi(x)^T] \\
-2E_{p(x, x')} [\varphi(x') \nu(x)^T] &+ E_{p(x)} \varphi(x)^T \Lambda_t \varphi(x)
\end{aligned}$$

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:

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

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

$$E_{p(x, x')} [\Psi(x) \Psi(x')^T]$$

$$E_{p(x)} [\Psi(x) T \Psi(x)] = I_{d \times d}$$

Penalty method:

$$\max_{\Psi} E_{p(x, x')} [\Psi(x) \Psi(x)^T] - \lambda \text{trace}(E_{p(x)} (\Psi(x) \Psi(x)^T) - I)^2$$

The above is a variant of VICreg [1].

References

- [1] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [2] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- [3] Alexey Kovalev. Clip from openai: what is it and how you can try it out yourself, 2021.
- [4] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training, 2021.
- [5] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023.