
[Appendix]

Unsupervised Representation Learning via Neural Activation Coding

Yookoon Park¹ Sangho Lee² Gunhee Kim² David M. Blei¹

1. Experimental Details for VAE Pretraining

We use diagonal Gaussians for both the variational posterior and the generative distribution of VAEs. For the encoder, we attach a linear output layer on ResNet-18 to predict the mean and the variances of a Gaussian distribution. The decoder takes a similar architecture with transposed convolution layers. We jointly train the encoder and the decoder on CIFAR-10 for 100 epochs with a batch size of 128. We use the Adam optimizer (Kingma & Ba, 2015) with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e - 8$ and no weight decay. The global learning rate is set to $5e-4$. When finetuning, we apply a smaller learning rate of $5e-5$ to the pretrained ResNet encoder, while keeping the learning rate high for the linear output layer.

2. MI and Average Hamming Distance

We argue that maximizing the mutual information (MI) $I(\mathbf{X}, \tilde{\mathbf{C}})$ over a noisy communication channel learns the codewords that have high Hamming distance to each other. We here show a relationship between the mutual information and the average Hamming distance between the codewords. Specifically, the mutual information *lower-bounds* the average Hamming distance. Recall that the Hamming distance between two codewords $\mathbf{c}_i, \mathbf{c}_j \in \{-1, 1\}^D$ is

$$d_H(\mathbf{c}_i, \mathbf{c}_j) = D - \frac{\mathbf{c}_i \cdot \mathbf{c}_j}{2}. \quad (1)$$

The average Hamming distance is defined as

$$\overline{d}_H = \frac{1}{N(N-1)} \sum_{j \neq i} d_H(\mathbf{c}_i, \mathbf{c}_j). \quad (2)$$

Let $\tilde{\mathbf{c}}$ be the noisy message transmitted through a binary symmetric channel with flip probability p . Then

$$\mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i}[\tilde{\mathbf{c}} \cdot \mathbf{c}_j] = (1 - 2p)(\mathbf{c}_i \cdot \mathbf{c}_j). \quad (3)$$

¹Computer Science Department, Columbia University, New York, USA ²Department of Computer Science and Engineering, Seoul National University, Seoul, South Korea. Correspondence to: David M. Blei <david.blei@columbia.edu>.

Finally,

$$I(\mathbf{X}, \tilde{\mathbf{C}}) \quad (4)$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[\log \frac{\exp((\tilde{\mathbf{c}} \cdot \mathbf{c}_i) \frac{1}{2} \log \frac{1-p}{p})}{\frac{1}{N} \sum_{j=1}^N \exp((\tilde{\mathbf{c}} \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p})} \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[(\tilde{\mathbf{c}} \cdot \mathbf{c}_i) \frac{1}{2} \log \frac{1-p}{p} \right] \right. \quad (5)$$

$$\left. - \mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[\log \frac{1}{N} \sum_{j=1}^N \exp((\tilde{\mathbf{c}} \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p}) \right] \right)$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\frac{D(1-2p)}{2} \log \frac{1-p}{p} \right. \quad (6)$$

$$\left. - \mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[\log \frac{1}{N} \sum_{j=1}^N \exp((\tilde{\mathbf{c}} \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p}) \right] \right)$$

$$\leq \frac{1}{N} \sum_{i=1}^N \left(\text{Const} - \mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[\frac{1}{N} \sum_{j=1}^N (\tilde{\mathbf{c}} \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p} \right] \right)$$

(Jensen's inequality)

$$= \frac{1}{N} \sum_{i=1}^N \left(\text{Const} + \sum_{j \neq i} (1-2p)(\mathbf{c}_i \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p} \right)$$

$$= \frac{N-1}{N} (1-2p) \log \left(\frac{1-p}{p} \right) \overline{d}_H + \text{Const}. \quad (7)$$

$$\leq \frac{1}{N} \sum_{i=1}^N \left(\text{Const} - \mathbb{E}_{\tilde{\mathbf{c}}|\mathbf{c}_i} \left[\frac{1}{N} \sum_{j=1}^N (\tilde{\mathbf{c}} \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p} \right] \right)$$

(Jensen's inequality)

$$= \frac{1}{N} \sum_{i=1}^N \left(\text{Const} + \sum_{j \neq i} (1-2p)(\mathbf{c}_i \cdot \mathbf{c}_j) \frac{1}{2} \log \frac{1-p}{p} \right)$$

$$= \frac{N-1}{N} (1-2p) \log \left(\frac{1-p}{p} \right) \overline{d}_H + \text{Const}. \quad (8)$$

Therefore, maximizing the mutual information objective increases the average Hamming distance between the codewords.

References

- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- Kingma, D. and Ba, J. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015.