

EECS 182 Deep Neural Networks
Spring 2023 Anant Sahai

Homework 7

This homework is due on Friday, March 17, 2022, at 10:59PM.

1. Auto-encoder : Learning without Labels

So far, with supervised learning algorithms we have used labelled datasets $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ to learn a mapping f_θ from input x to labels y . In this problem, we now consider algorithms for *unsupervised learning*, where we are given only inputs x , but no labels y . At a high-level, unsupervised learning algorithms leverage some structure in the dataset to define proxy tasks, and learn *representations* that are useful for solving downstream tasks.

Autoencoders present a family of such algorithms, where we consider the problem of learning a function $f_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^k$ from input x to a *intermediate representation* z of x (usually $k \ll m$). For autoencoders, we use reconstructing the original signal as a proxy task, i.e. representations are more likely to be useful for downstream tasks if they are low-dimensional but encode sufficient information to approximately reconstruct the input. Broadly, autoencoder architectures have two modules:

- An encoder $f_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^k$ that maps input x to a representation z .
- A decoder $g_\phi : \mathbb{R}^k \rightarrow \mathbb{R}^m$ that maps representation z to a reconstruction \hat{x} of x .

In such architectures, the parameters (θ, ϕ) are learnable, and trained with the learning objective of minimizing the reconstruction error $\ell(\hat{x}_i, x_i) = \|x - \hat{x}\|_2^2$ using gradient descent.

$$\theta_{\text{enc}}, \phi_{\text{dec}} = \underset{\Theta, \Phi}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \ell(\hat{x}_i, x_i)$$

$$\hat{x} = g_\phi(f_\theta(x))$$

Note that above optimization problem does not require labels \mathbf{y} . In practice however, we would want to use the *pretrained* models to solve the *downstream* task at hand, e.g. classifying MNIST digits. *Linear Probing* is one such approach, where we fix the encoder weights, and learn a simple linear classifier over the features $z = \text{encoder}(x)$.

(a) Designing AutoEncoders

Please follow the instructions in [this notebook](#). You will train autoencoders, denoising autoencoders, and masked autoencoders on a synthetic dataset and the MNIST dataset. Once you finished with the notebook,

- Download `submission_log.json` and submit it to “Homework 7 (Code)” in Gradescope.
- Answer the following questions in your submission of the written assignment:
 - (i) **Show your visualization** of the vanilla autoencoder with different latent representation sizes.
 - (ii) Based on your previous visualizations, answer this question: **How does changing the latent representation size of the autoencoder affect the model’s performance in terms of reconstruction accuracy and linear probe accuracy? Why?**

(b) PCA & AutoEncoders

In the case where the encoder f_θ, g_ϕ are linear functions, the model is termed as a *linear autoencoder*. In particular, assume that we have data $x_i \in \mathbb{R}^m$ and consider two weight matrices: an encoder $W_1 \in \mathbb{R}^{k \times m}$ and decoder $W_2 \in \mathbb{R}^{m \times k}$ (with $k < m$). Then, a linear autoencoder learns a low-dimensional embedding of the data $\mathbf{X} \in \mathbb{R}^{m \times n}$ (which we assume is zero-centered without loss of generality) by minimizing the objective,

$$\mathcal{L}(W_1, W_2; \mathbf{X}) = \frac{1}{n} \|\mathbf{X} - W_2 W_1 \mathbf{X}\|_F^2 \quad (1)$$

We will assume $\sigma_1^2 > \dots > \sigma_k^2 > 0$ are the k largest eigenvalues of $\frac{1}{n} \mathbf{X} \mathbf{X}^\top$. The assumption that the $\sigma_1, \dots, \sigma_k$ are positive and distinct ensures identifiability of the principal components, and is common in this setting. Therefore, the top- k eigenvalues of \mathbf{X} are $S = \text{diag}(\sigma_1, \dots, \sigma_k)$, with corresponding eigenvectors are the columns of $\mathbf{U}_k \in \mathbb{R}^{m \times k}$. A well-established result from (Baldi & Hornik, 1989) shows that principal components are the unique optimal solution to linear autoencoders (up to sign changes to the projection directions). In this subpart, we take some steps towards proving this result.

- (i) Write out the first order optimality conditions that the minima of Eq. 1 would satisfy.
- (ii) Show that the principal components \mathbf{U}_k satisfy the optimality conditions outlined in (i).

2. Self-supervised Linear Autoencoders

We consider linear models consisting of two weight matrices: an encoder $W_1 \in \mathbb{R}^{k \times m}$ and decoder $W_2 \in \mathbb{R}^{m \times k}$ (assume $1 < k < m$). The traditional autoencoder model learns a low-dimensional embedding of the n points of training data $\mathbf{X} \in \mathbb{R}^{m \times n}$ by minimizing the objective,

$$\mathcal{L}(W_1, W_2; \mathbf{X}) = \frac{1}{n} \|\mathbf{X} - W_2 W_1 \mathbf{X}\|_F^2 \quad (2)$$

We will assume $\sigma_1^2 > \dots > \sigma_k^2 > \sigma_{k+1}^2 \geq 0$ are the $k + 1$ largest eigenvalues of $\frac{1}{n} \mathbf{X} \mathbf{X}^\top$. The assumption that the $\sigma_1, \dots, \sigma_k$ are positive and distinct ensures identifiability of the principal components.

Consider an ℓ_2 -regularized linear autoencoder where the objective is:

$$\mathcal{L}_\lambda(W_1, W_2; \mathbf{X}) = \frac{1}{n} \|\mathbf{X} - W_2 W_1 \mathbf{X}\|_F^2 + \lambda \|W_1\|_F^2 + \lambda \|W_2\|_F^2. \quad (3)$$

where $\|\cdot\|_F^2$ represents the Frobenius norm squared of the matrix (i.e. sum of squares of the entries).

- (a) You want to use SGD-style training in PyTorch (involving the training points one at a time) and self-supervision to find W_1 and W_2 which optimize (3) by treating the problem as a neural net being trained in a supervised fashion. **Answer the following questions and briefly explain your choice:**

- (i) How many linear layers do you need?

- ☐ 0
☐ 1
☐ 2
☐ 3

- (ii) What is the loss function that you will be using?

- ☐ nn.L1Loss
☐ nn.MSELoss
☐ nn.CrossEntropyLoss

(iii) Which of the following would you need to optimize (3) exactly as it is written? (Select all that are needed)

- ☐ Weight Decay
- ☐ Dropout
- ☐ Layer Norm
- ☐ Batch Norm
- ☐ SGD optimizer

(b) Do you think that the solution to (3) when we use a small nonzero λ has an inductive bias towards finding a W_2 matrix with approximately orthonormal columns? Argue why or why not?

(Hint: Think about the SVDs of $W_1 = U_1 \Sigma_1 V_1^\top$ and $W_2 = U_2 \Sigma_2 V_2^\top$. You can assume that if a $k \times m$ or $m \times k$ matrix has all k of its nonzero singular values being 1, then it must have orthonormal rows or columns. Remember that the Frobenius norm squared of a matrix is just the sum of the squares of its singular values. Further think about the minimizer of $\frac{1}{\sigma^2} + \sigma^2$. Is it unique?)

3. Justifying Scaled-Dot Product Attention

Suppose $q, k \in \mathbb{R}^d$ are two random vectors with $q, k \sim N(\mu, \sigma^2 I)$, where $\mu \in \mathbb{R}^d$ and $\sigma \in \mathbb{R}^+$. In other words, each component q_i of q is drawn from a normal distribution with mean μ and standard deviation σ , and the same is true for k .

- (a) Define $\mathbb{E}[q^T k]$ in terms of μ, σ and d .
- (b) Considering a practical case where $\mu = 0$ and $\sigma = 1$, define $\text{Var}(q^T k)$ in terms of d .
- (c) Continue to use the setting in part (b), where $\mu = 0$ and $\sigma = 1$. Let s be the scaling factor on the dot product. Suppose we want $\mathbb{E}[\frac{q^T k}{s}]$ to be 0, and $\text{Var}(\frac{q^T k}{s})$ to be $\sigma = 1$. What should s be in terms of d ?

4. Argmax Attention

Recall from lecture that we can think about attention as being *queryable softmax pooling*. In this problem, we ask you to consider a hypothetical argmax version of attention where it returns exactly the value corresponding to the key that is most similar to the query, where similarity is measured using the traditional inner-product.

- (a) Perform **argmax attention** with the following keys and values:

Keys:

$$\left\{ \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}, \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$

Corresponding Values:

$$\left\{ \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \right\}$$

using the following query:

$$\mathbf{q} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

What would be the output of the attention layer for this query?

Hint: For example, $\text{argmax}([1, 3, 2]) = [0, 1, 0]$

- (b) Note that instead of using *softmax* we used *argmax* to generate outputs from the attention layer. **How does this design choice affect our ability to usefully train models involving attention?**

(Hint: think about how the gradients flow through the network in the backward pass. Can we learn to improve our queries or keys during the training process?)

5. Kernelized Linear Attention

The softmax attention is widely adopted in transformers (Luong et al., 2015; Vaswani et al., 2017), however the $\mathcal{O}(N^2)$ (N stands for the sequence length) complexity in memory and computation often makes it less desirable for processing long document like a book or a passage, where the N could be beyond thousands. There is a large body of the research studying how to resolve this ¹.

Under this context, this question presents a formulation of attention via the lens of the kernel. A large portion of the context is adopted from Tsai et al. (2019). In particular, attention can be seen as applying a kernel over the inputs with the kernel scores being the similarities between inputs. This formulation sheds light on individual components of the transformer’s attention, and helps introduce some alternative attention mechanisms that replaces the “softmax” with linearized kernel functions, thus reducing the $\mathcal{O}(N^2)$ complexity in memory and computation.

We first review the building block in the transformer. Let $x \in \mathbb{R}^{N \times F}$ denote a sequence of N feature vectors of dimensions F . A transformer Vaswani et al. (2017) is a function $T : \mathbb{R}^{N \times F} \rightarrow \mathbb{R}^{N \times F}$ defined by the composition of L transformer layers $T_1(\cdot), \dots, T_L(\cdot)$ as follows,

$$T_l(x) = f_l(A_l(x) + x). \quad (4)$$

The function $f_l(\cdot)$ transforms each feature independently of the others and is usually implemented with a small two-layer feedforward network. $A_l(\cdot)$ is the self attention function and is the only part of the transformer that acts across sequences.

We now focus on the the self attention module which involves softmax. The self attention function $A_l(\cdot)$ computes, for every position, a weighted average of the feature representations of all other positions with a weight proportional to a similarity score between the representations. Formally, the input sequence x is projected by three matrices $W_Q \in \mathbb{R}^{F \times D}$, $W_K \in \mathbb{R}^{F \times D}$ and $W_V \in \mathbb{R}^{F \times M}$ to corresponding representations Q , K and V . The output for all positions, $A_l(x) = V'$, is computed as follows,

$$\begin{aligned} Q &= xW_Q, K = xW_K, V = xW_V, \\ A_l(x) &= V' = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V. \end{aligned} \quad (5)$$

Note that in the previous equation, the softmax function is applied rowwise to QK^T . Following common terminology, the Q , K and V are referred to as the “queries”, “keys” and “values” respectively.

¹<https://huggingface.co/blog/long-range-transformers>

Equation 5 implements a specific form of self-attention called softmax attention where the similarity score is the exponential of the dot product between a query and a key. Given that subscripting a matrix with i returns the i -th row as a vector, we can write a generalized attention equation for any similarity function as follows,

$$V'_i = \frac{\sum_{j=1}^N \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^N \text{sim}(Q_i, K_j)}. \quad (6)$$

Equation 6 is equivalent to equation 5 if we substitute the similarity function with $\text{sim}_{\text{softmax}}(q, k) = \exp(\frac{q^T k}{\sqrt{D}})$. This can lead to

$$V'_i = \frac{\sum_{j=1}^N \exp(\frac{Q_i^T K_j}{\sqrt{D}}) V_j}{\sum_{j=1}^N \exp(\frac{Q_i^T K_j}{\sqrt{D}})}. \quad (7)$$

For computing the resulting self-attended feature $A_l(x) = V'$, we need to compute all V'_i $i \in 1, \dots, N$ in equation 7.

- (a) **Identify the conditions that needs to be met by the sim function to ensure that V_i in Equation 6 remains finite (the denominator never reaches zero).**
- (b) The definition of attention in equation 6 is generic and can be used to define several other attention implementations.
 - (i) One potential attention variant is the “polynomial kernel attention”, where the similarity function as $\text{sim}(q, k)$ is measured by polynomial kernel \mathcal{K} ². **Considering a special case for a “quadratic kernel attention” that the degree of “polynomial kernel attention” is set to be 2, derive the $\text{sim}(q, k)$ for “quadratic kernel attention”. (NOTE: any constant factor is set to be 1.)**
 - (ii) One benefit of using kernelized attention is that we can represent a kernel using a feature map $\phi(\cdot)$ ³. **Derive the corresponding feature map $\phi(\cdot)$ for the quadratic kernel.**
 - (iii) **Considering a general kernel attention, where the kernel can be represented using feature map that $\mathcal{K}(q, k) = (\phi(q)^T \phi(k))$, rewrite kernel attention of equation 6 with feature map $\phi(\cdot)$.**
- (c) We can rewrite the softmax attention in terms of equation 6 as equation 7. **For all the V'_i ($i \in \{1, \dots, N\}$), derive the time complexity (asymptotic computational cost) and space complexity (asymptotic memory requirement) of the above softmax attention in terms of sequence length N , D and M .**
NOTE: for memory requirement, we need to store any intermediate results for backpropagation, including all Q, K, V
- (d) Assume we have a kernel \mathcal{K} as the similarity function and the kernel can be represented with a feature map $\phi(\cdot)$, we can rewrite equation 6 with $\text{sim}(x, y) = \mathcal{K}(x, y) = (\phi(Q_i)^T \phi(K_j))$ in part (b). We can then further simplify it by making use of the associative property of matrix multiplication to

$$V'_i = \frac{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j) V_j^T}{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j)}. \quad (8)$$

²https://en.wikipedia.org/wiki/Polynomial_kernel

³https://en.wikipedia.org/wiki/Kernel_method

Note that the feature map $\phi(\cdot)$ is applied row-wise to the matrices Q and K .

Considering using a linearized polynomial kernel $\phi(x)$ of degree 2, and assume $M \approx D$, derive the computational cost and memory requirement of this kernel attention as in (8).

6. Debugging DNNs (Optional)

- (a) Your friends want to train a classifier for a new app they're designing. They implement a deep convolutional network and train models with two configurations: a 20 layer model and a 56 layer model. However, they observe the following training curves and are surprised that the 20-layer network has better training as well as test error.

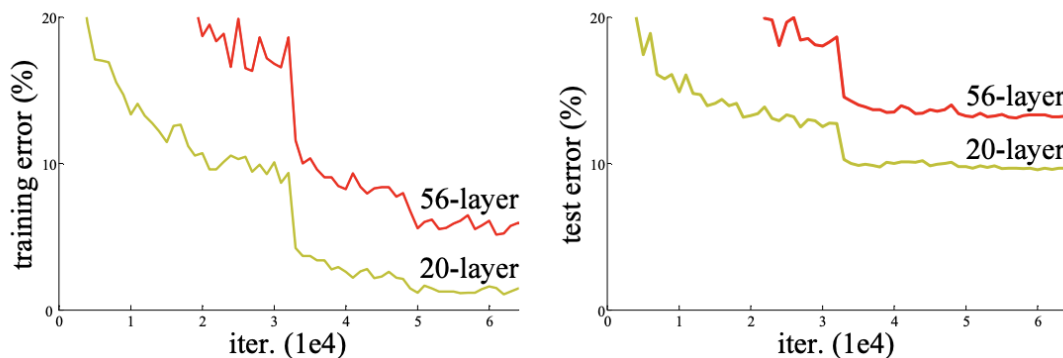


Figure 1: Training deep networks on CIFAR10

What are the potential reasons for this observation? Are there changes to the architecture design that could help mitigate the problem?

- (b) You and your teammate want to compare batch normalization and layer normalization for the ImageNet classification problem. You use ResNet-152 as a neural network architecture. The images have input dimension $3 \times 224 \times 224$ (channels, height, width). You want to use a batch size of 1024; however, the GPU memory is so small that you cannot load the model and all 1024 samples at once — you can only fit 32. Your teammate proposes using a gradient-accumulation algorithm:

Gradient accumulation refers to running the forward and backward pass of a model a fixed number of steps (`accumulation_steps`) without updating the model parameters, while aggregating the gradients. Instead, the model parameters are updated every (`accumulation_steps`). This allows us to increase the effective batch size by a factor of `accumulation_steps`.

You implement the algorithm in PyTorch as:

```
model.train()
optimizer.zero_grad()
for i, (inputs, labels) in enumerate(training_set):
    predictions = model(inputs)
    loss = loss_function(predictions, labels)
    loss = loss / accumulation_steps
    loss.backward()
    if (i+1) % accumulation_steps == 0:
        optimizer.step()
        optimizer.zero_grad()
```

Note that the `.backward()` operator in PyTorch implicitly keeps accumulating the gradient for all the parameters, unless zero'd out with an `optimizer.zero_grad()` call.

Before running actual experiments, your friend suggests that you should test whether the gradient accumulation algorithm is implemented correctly. To do so, you collect the output logits (i.e. the outputs of the last layer) from two models — ResNet-152, one with batchnorm and the other with layernorm — using different combinations of batch sizes and the number of accumulation steps that keep the effective batch size to 32.

[Note that the effective batch size is product of the batch size and `accumulation_steps`. In other words, the possible combinations for effective batch-size 32 are:

$(batch_size, accumulation_steps) = (1, 32), (2, 16), (4, 8), (8, 4), (16, 2), (32, 1).$]

Here, the (32,1) combination is the approach without the “gradient accumulation” trick, and we want to see whether the others agree with this.

On running these tests, you observe that one of models: either with batchnorm or with layernorm, doesn't pass the test. **Which one do you expect to not pass the test and why?**

- (c) You are training a CIFAR model and observe that the model is diverging (instead of the training loss decreasing over iterations). **Debug the pseudocode and give a correction that you believe would actually result in reasonable convergence during training.**

(Note: You can assume that the datasets are loaded correctly, model is trained with SGD optimizer with learning rate= 0.001, batchsize= 100)

(HINT: Ideas from the previous part of this question might be relevant.)

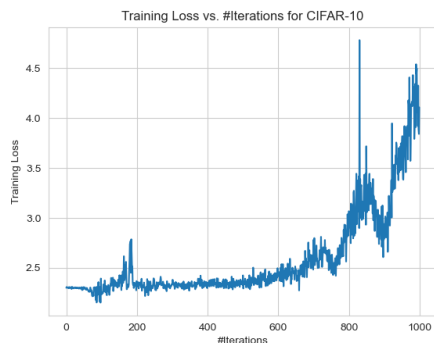


Figure 2: Training loss for CIFAR10

```
model.train()
optimizer.zero_grad()
for (inputs, labels) in training_set:
    predictions = model(inputs)
    loss = loss_fn(predictions, labels)
    loss.backward()
    optimizer.step()
```

7. Homework Process and Study Group

Citing sources and collaborators are an important part of life, including being a student!

We also want to understand what resources you find helpful and how much time homework is taking, so we can change things in the future if possible.

- What sources (if any) did you use as you worked through the homework?**
- If you worked with someone on this homework, who did you work with?**
List names and student ID's. (In case of homework party, you can also just describe the group.)
- Roughly how many total hours did you work on this homework? Write it down here where you'll need to remember it for the self-grade form.**

References

- Baldi, P. and Hornik, K. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Luong, M.-T., Pham, H., and Manning, C. D. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1412–1421, 2015.
- Tsai, Y.-H. H., Bai, S., Yamada, M., Morency, L.-P., and Salakhutdinov, R. Transformer dissection: An unified understanding for transformer’s attention via the lens of kernel. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4344–4353, 2019.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Contributors:

- Kumar Krishna Agrawal.
- Linyuan Gong.
- Sheng Shen.
- Anant Sahai.
- David M. Chan.
- Saagar Sanghavi.
- Shaojie Bai.
- Angelos Katharopoulos.
- Hao Peng.
- Suhong Moon.