

# HOMWORK SET 4

CSCI 5980/8980 Think Deep Learning (Fall 2020)

**Due** 11:59 pm, Dec 12 2020

**Instruction** Typesetting your submission in L<sup>A</sup>T<sub>E</sub>X is now **optional**, and but you need to submit it as a single PDF file in Canvas. No late submission will be accepted. For each problem, you should acknowledge your collaborators if any. For problems containing multiple subproblems, there are often close logic connections between the subproblems. So always remember to build on previous ones, rather than work from scratch.

**Notation** We will use small letters (e.g.,  $u$ ) for scalars, small boldface letters (e.g.,  $\mathbf{a}$ ) for vectors, and capital boldface letters (e.g.,  $\mathbf{A}$ ) for matrices. For a matrix  $\mathbf{A}$ ,  $\mathbf{a}^i$  (supscripting) means its  $i$ -th row as a *row vector*, and  $\mathbf{a}_j$  (subscripting) means the  $j$ -th column as a column vector, and  $a_{ij}$  means its  $(i, j)$ -th element.  $\mathbb{R}$  is the set of real numbers.  $\mathbb{R}^n$  is the space of  $n$ -dimensional real vectors, and similarly  $\mathbb{R}^{m \times n}$  is the space of  $m \times n$  real matrices. The dotted equal sign  $\doteq$  means defining.

## Problem 1 (Autoencoders, deep factorization, and deep sparse coding; 6/12)

- (a) The geometric view of PCA says that PCA tries to fit a subspace to a collection of data points. From a modeling perspective, this means PCA makes sense only when the data points lie near a subspace. For example, if  $\mathbf{x}_i = \mathbf{B}\mathbf{z}_i + \varepsilon_i$  for all  $i$ , where  $\mathbf{B}$  is a basis for a subspace and  $\varepsilon_i$ 's represent small-magnitude noise, trying to find a basis for the subspace spanned by  $\mathbf{B}$  may make a lot of sense. What happens when a small fraction of the data points deviate significantly from the subspace?

To investigate this, let's generate an orthonormal subspace basis  $\mathbf{B} \in \mathbb{R}^{200 \times 20}$  and 98 random data points on the subspace as  $\mathbf{B}\mathbf{z}_i$ 's where each  $\mathbf{z}_i \in \mathbb{R}^{20}$  is iid Gaussian. So this portion of data is perfectly clean. Now let's generate 2 points that are iid Gaussian in  $\mathbb{R}^{200}$ —these two points will be far from the subspace  $\mathbf{B}$  almost surely and they are "outliers". Now we have 100 data points and we collect them into a data matrix  $\mathbf{X} \in \mathbb{R}^{200 \times 100}$ . Please normalize the 100 data points so that they all have unit  $\ell_2$  norm now. *Also, do NOT perform centering to the points for the following steps.*

- Perform PCA via SVD or eigen-decomposition on  $\mathbf{X}$  (again, no centering step) and numerically compare the subspace spanned by the top 20 singular vectors with  $\mathbf{B}$ . Remember for two subspaces spanned by two bases  $\mathbf{B}_1$  and  $\mathbf{B}_2$ , their distance can be measured by  $\|\mathbf{B}_1\mathbf{B}_1^\dagger - \mathbf{B}_2\mathbf{B}_2^\dagger\|_F$ . What do you observe? (1/12)
- Specializing the factorization formulation for PCA for our setting is

$$\min_{\mathbf{A} \in \mathbb{R}^{200 \times 20}, \mathbf{Z} \in \mathbb{R}^{20 \times 100}} \sum_{i=1}^{100} \|\mathbf{x}_i - \mathbf{A}\mathbf{z}_i\|_2^2. \quad (1)$$

It's equivalent to the above PCA we have done of course. Now let's consider a slight modified version:

$$\min_{\mathbf{A} \in \mathbb{R}^{200 \times 20}, \mathbf{Z} \in \mathbb{R}^{20 \times 100}} \sum_{i=1}^{100} \|\mathbf{x}_i - \mathbf{A}\mathbf{z}_i\|_2, \quad (2)$$

i.e., sum of the  $\ell_2$  norm, but norm squared. Numerically solve this optimization problem (choose whatever methods you're comfortable with, and auto-differentiation is also

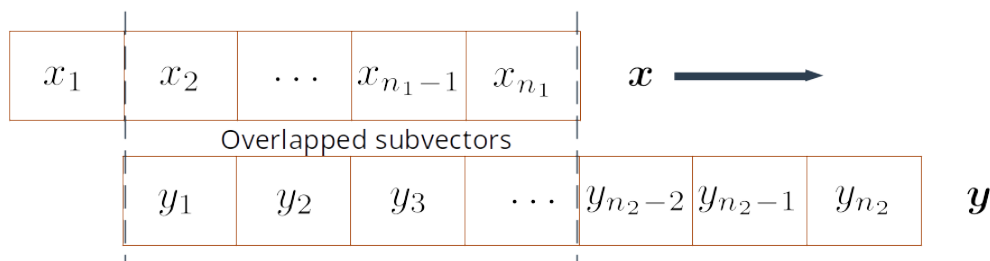
allowed: remember that PyTorch or TensorFlow can also be used to solve generic unconstrained optimization problems, not necessarily deep learning problems). Compare the subspace obtained here with  $B$ . Do you get a better estimate than the plain PCA above? (2/12)

(b) Read the Science paper *Reducing the Dimensionality of Data with Neural Networks* (<https://science.sciencemag.org/content/313/5786/504>).

- Reproduce the first two rows of Fig 2(B), i.e., PCA and autoencoder on MNIST. You should use exactly the same architecture as provided in Fig 1 (right). The original paper uses layer-wise pretraining for initialization and conjugate-gradient for training. You probably don't need these; instead, you can choose modern initialization and optimization methods in PyTorch or TensorFlow. If you are unsure how to use PyTorch to build a neural network and perform training, this tutorial will be helpful: [https://pytorch.org/tutorials/beginner/blitz/neural\\_networks\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/neural_networks_tutorial.html) (2/12)
- Reproduce the plot in Fig 3, i.e., PCA and autoencoder for visualization in the two-dimensional space. The autoencoder architecture is described in the caption of Fig 3. Again, you don't need to use the original initialization and training methods. (1/12)

**Problem 2 (Correlation and template matching; 6/12)** The word “convolutional” in convolutional neural networks is a misnomer. Cross-correlation, which is a close relative of convolution and commonly used in signal processing, is actually used. In this problem, we explore some basic properties and applications of the correlation operation. We use the standard notation  $\star$  to denote cross-correlation, as against  $*$  which is often used to denote convolution.

(a) For two vectors  $\mathbf{x} \in \mathbb{R}^{n_1}$ ,  $\mathbf{y} \in \mathbb{R}^{n_2}$ , the cross-correlation  $\mathbf{x} \star \mathbf{y}$  is obtained as follows:



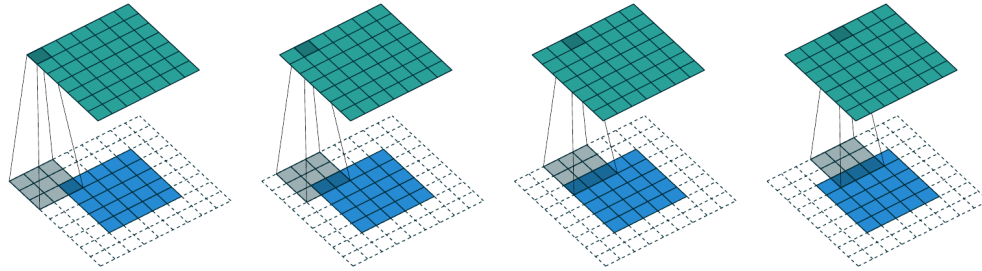
We fix the position of  $\mathbf{y}$ , and shift  $\mathbf{x}$  to the left until  $\mathbf{x}$  and  $\mathbf{y}$  only have one overlapped element spatially, i.e.,  $x_{n_1}$  with  $y_1$  — that's the starting point. We calculate the inner product of two overlapped subvectors—in the beginning only two scalars actually. Then we repeatedly do this: shift  $\mathbf{x}$  to the right by one element and calculate the corresponding inner product of the two overlapped subvectors (i.e., think of a sliding window). We end the process until  $\mathbf{x}$  and  $\mathbf{y}$  overlap only at one element, i.e.,  $x_1$  with  $y_{n_2}$ . The cross-correlation  $\mathbf{x} \star \mathbf{y}$  is basically the vector that collects all the inner product values we have obtained in the left-to-right order. It is easy to see that  $\mathbf{x} \star \mathbf{y} \in \mathbb{R}^{n_1+n_2-1}$ .

**Question:** Calculate  $[3, 2, 1] \star [4, 6, 3, 9]$ . (0.5/12) For general  $\mathbf{x}, \mathbf{y}$ , is it true that  $\mathbf{x} \star \mathbf{y} = \mathbf{y} \star \mathbf{x}$ ? If not, what relationship between  $\mathbf{x} \star \mathbf{y}$  and  $\mathbf{y} \star \mathbf{x}$  do you observe? (0.5/12)

(b) In convolutional neural networks, we have building blocks of the form  $\mathbf{w} \star \mathbf{x}$ , where  $\mathbf{w}$  represents a group of learnable weights, often called *filter* following the signal processing

convention. For simplicity, let's assume  $w \in \mathbb{R}^3$  and  $x \in \mathbb{R}^4$ . Show that  $w \star x$  can be written equivalently as  $C_w x$  for a certain matrix  $C_w \in \mathbb{R}^{6 \times 4}$  and write down  $C_w$  explicitly in terms of elements of  $w$ . (1/12)

- (c) To apply reverse-mode auto differentiation, we need to specify  $\frac{\partial}{\partial w} (w \star x)$ , i.e., the associated Jacobian. Assume again  $w \in \mathbb{R}^3$  and  $x \in \mathbb{R}^4$ , can you derive the analytic form of the Jacobian? (1/12; Hint: is it possible to write  $w \star x$  as  $C_x w$  for a certain  $C_x$ ?)
- (d) The 2D cross-correlation is a natural generalization of the 1D cross-correlation to matrices.



(image credit: <https://arxiv.org/abs/1603.07285>; check out [https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic) to see the dynamic demonstration under the **Full padding, no strides** setting. Note that padding zeros, as indicated as the additional dotted boxes, is equivalent to ignoring the out-of-boundary elements.)

Compared to the 1D version, now we start from the top-left corner and end at the bottom-right corner. We scan row by row and inner products are now taken between the overlapped submatrices. All the inner product values are naturally organized into a matrix. In the pictorial illustration above, we are considering  $X \star Y$ , where  $X \in \mathbb{R}^{3 \times 3}$  is the gray matrix, and  $Y \in \mathbb{R}^{5 \times 5}$  is the blue matrix, the resulting green matrix  $X \star Y \in \mathbb{R}^{7 \times 7}$ , where  $7 = 3 + 5 - 1$ .

**Question:** In Numpy, implement a 2D cross-correlation function. The function should take in two general matrices  $Z_1 \in \mathbb{R}^{n_1 \times n_2}$  and  $Z_2 \in \mathbb{R}^{m_1 \times m_2}$  and return the resulting cross-correlation matrix. To debug your implementation, please generate a couple of random cases and benchmark against the Scipy built-in function `scipy.signal.correlate2d` (remember to set `mode = 'full'`) <https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.correlate2d.html>. (1/12)

- (e) Most basic image processing algorithms are implemented as cross-correlation of a small filter  $X$  with the image of interest  $Y$ . Check out the examples at the bottom of the page <https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.convolve2d.html>. Use their image ascent, let's test your implementation of 2D cross-correlation. Try two filters

$$X_1 = \begin{bmatrix} -3 & 0 & 3 \\ -10 & 0 & 10 \\ -3 & 0 & 3 \end{bmatrix} \quad \text{and} \quad X_2 = X_1^T.$$

Let's say they generate two resulting matrices  $G_1$  and  $G_2$ . Calculate  $\sqrt{G_1^2 + G_2^2}$ , where the operations are pointwise. Display your result (i.e., `imshow` as in the online example). Does your result look alike the gradient magnitude plot, except for the image boundaries? (1/12)

- (f) Another way of thinking about cross-correlation is template matching. Imagine that  $X$  is a 2D pattern of interest. During the cross-correlation process, the inner product measures the agreement of the 2D pattern and local patches in  $Y$ . If the value is relatively large, very likely we find a match. After we finish the cross-correlation calculation, we can spot the locations of the largest values in the cross-correlation matrix as candidate matching locations. Study the example here <https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.correlate2d.html> and compare the performance of your implementation with that of the example using `scipy.signal.correlate2d`. (1/12)