

Homework 0
CSC 277 / 477
End-to-end Deep Learning
Fall 2024

John Doe - `jdoe@ur.rochester.edu`

Deadline: See Blackboard

Instructions

Your homework solution must be typed and prepared in \LaTeX . It must be output to PDF format. To use \LaTeX , we suggest using <http://overleaf.com>, which is free.

Your submission must cite any references used (including articles, books, code, websites, and personal communications). All solutions must be written in your own words, and you must program the algorithms yourself. **If you do work with others, you must list the people you worked with.** Submit your solutions as a PDF to Blackboard.

Your programs must be written in Python. The relevant code should be in the PDF you turn in. If a problem involves programming, then the code should be shown as part of the solution. One easy way to do this in \LaTeX is to use the verbatim environment, i.e., `\begin{verbatim} YOUR CODE \end{verbatim}`.

About Homework 0: Homework 0 is intended to review prerequisite skills, help you become familiar with LaTeX, and ensure you have your programming environment prepared. *Later assignments will be more challenging! Do not think all assignments will be this short or require little time to train networks.* Copy and paste this template into an editor, e.g., www.overleaf.com, and then just type the answers in. You can use a math editor to make this easier, e.g., CodeCogs Equation Editor or MathType. You may use the AI (LLM) plugin for Overleaf for help you with \LaTeX formatting.

Problem 1 - Linear Algebra Review #1

Let matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ and matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$, where $n \neq m$.

Part 1 (1 point)

If it is possible to compute the matrix product \mathbf{AB} , give the size of the matrix produced. Otherwise, write, 'Not possible.'

Answer:

Not possible.

Part 2 (1 point)

If it is possible to compute the matrix product \mathbf{BA} , give the size of the matrix produced. Otherwise, write, 'Not possible.'

Answer:

The size of the matrix produced is $n \times m$.

Problem 2 - Linear Algebra Review #2

Let $\mathbf{A} = \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix}$ and $\mathbf{B} = \begin{pmatrix} 1 & 5 & 0 \\ 2 & 10 & 2 \end{pmatrix}$ and $\mathbf{D} = \begin{pmatrix} 1 & 1 & 2 & 3 & 5 & 4 \\ 1 & 2 & 3 & 6 & 6 & 6 \\ 1 & 1 & 2 & 3 & 5 & 4 \\ 1 & 0 & 1 & 0 & 4 & 2 \end{pmatrix}$.

Low rank approximations are used in algorithms for updating large neural networks on modest hardware. In this problem we will assess your understanding of rank in linear algebra.

Part 1 (1 points)

Compute $\text{rank}(\mathbf{A})$.

Answer:

$$\text{rank}(\mathbf{A}) = 1.$$

Part 2 (1 points)

Compute $\text{rank}(\mathbf{B})$.

Answer:

$$\text{rank}(\mathbf{B}) = 2.$$

Part 3 (1 points)

Compute $\text{rank}(\mathbf{D})$.

Answer:

$$\text{rank}(\mathbf{D}) = 2.$$

Part 4 (1 points)

Compute \mathbf{AB} and $\text{rank}(\mathbf{AB})$.

Answer:

$$\mathbf{AB} = \begin{pmatrix} 0 & 0 & 0 \\ 4 & 20 & 10 \end{pmatrix}$$

$$\text{rank}(\mathbf{AB}) = 1.$$

Part 5 - Compression with linearly dependent columns (5 points)

Now we are going to use linearly independent columns of matrix \mathbf{D} to create an compressed version of matrix \mathbf{D} . We will use the $\mathbf{D} = \mathbf{CR}$ decomposition, i.e., rank factorization

(column-row factorization). Do the decomposition and provide $\mathbf{C} \in \mathbb{R}^{4 \times r}$ and $\mathbf{R} \in \mathbb{R}^{r \times 6}$. Calculate the compression ratio by comparing the total number of elements in original matrix and in decomposed matrices. Explain why r is the rank of \mathbf{D} .

You may read about the CR factorization [here](#) and [here](#). You can implement this factorization using singular value decomposition or by using reduced row echelon form.

Answer:

$$\mathbf{C} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1 \\ 1 & 0 \end{pmatrix}$$

NOTE: C consists of the linearly independent columns of D

$$\mathbf{R} = \begin{pmatrix} 1 & 0 & 1 & 2 & 3 & 4 \\ 0 & 1 & 1 & 2 & 3 & 2 \end{pmatrix}$$

NOTE: R is formed by expressing all columns of D as linear combinations of the columns in C

$$\text{Compression ratio: } \frac{24}{4 \times 2 + 2 \times 6} = \frac{24}{20} = 1.2$$

The reason why r is the rank of \mathbf{D} is because there are only 2 linearly independent columns in \mathbf{D} .

Problem 3 - Softmax Properties

Part 1 (5 points)

Recall the softmax function, which is the most common activation function used for the output of a neural network trained to do classification. In a vectorized form, it is given by

$$\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\sum_{j=1}^K \exp(a_j)},$$

where $\mathbf{a} \in \mathbb{R}^K$ and has the “logits” from the output layer of the network before the softmax activation function. The exp function in the numerator is applied element-wise and a_j denotes the j 'th element of \mathbf{a} .

Show that the softmax function is invariant to constant offsets to its input, i.e.,

$$\text{softmax}(\mathbf{a} + c\mathbf{1}) = \text{softmax}(\mathbf{a}),$$

where $c \in \mathbb{R}$ is some constant and $\mathbf{1}$ denotes a column vector of 1's.

Solution:

Consider adding a constant c to each of the element in vector \mathbf{a} , the softmax function for $\mathbf{a} + c\mathbf{1}$ is:

Proof.

$$\text{softmax}(\mathbf{a} + c\mathbf{1}) = \frac{\exp(\mathbf{a} + c\mathbf{1})}{\sum_{j=1}^K \exp(a_j + c)}$$

Since $\exp(a + c\mathbf{1})$ applies element-wise, we have:

$$\exp(\mathbf{a} + c\mathbf{1}) = \exp(\mathbf{a}) \cdot \exp(c)$$

Thus, we can rewrite the softmax function as:

$$\text{softmax}(\mathbf{a} + c\mathbf{1}) = \frac{\exp(\mathbf{a}) \cdot \exp(c)}{\exp(c) \cdot \sum_{j=1}^K \exp(a_j)} = \frac{\exp(\mathbf{a})}{\sum_{j=1}^K \exp(a_j)}$$

Hence,

$$\text{softmax}(\mathbf{a} + c\mathbf{1}) = \text{softmax}(\mathbf{a})$$

The softmax function is invariant to adding a constant vector to its input. □

Part 2 (3 points)

In practice, why is the observation that the softmax function is invariant to constant offsets to its input important when implementing it in a neural network?

Solution:

The invariance of the softmax function to constant offsets is important in neural networks because:

1. **Numerical Stability:** Subtracting the maximum logit before applying softmax prevents overflow/underflow in the exponentials, ensuring stable computations.
2. **Interpretability:** The relative differences between logits determine the output probabilities, making the absolute scale irrelevant and allowing focus on relative confidence between classes.
3. **Model Flexibility:** This property permits shifts in logits (e.g., due to batch normalization) without affecting the final output, enabling more flexible model design. (i.e. due to the flexibility on “paddings”)

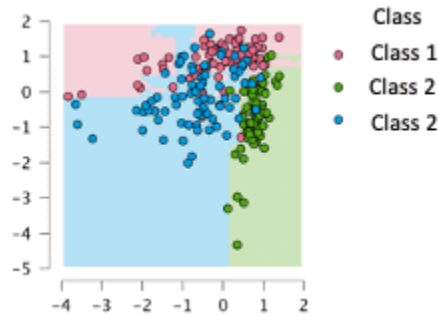


Figure 1: An example decision boundary on a different dataset.

Problem 4 - Feedforward Fully-Connected Networks (10 points)

This class assumes that you already know neural network basics. For this problem you need to use PyTorch and/or PyTorch Lightning. PyTorch Lightning is a wrapper around PyTorch that attempts to simplify a lot of the code, enforcing best practices, and making the code more portable across hardware platforms.

In this problem you will create a neural network that will be trained on the Iris dataset for multi-class classification. We will use a 2-dimensional version of the dataset. Download `iris-train.txt` and `iris-test.txt`. Each row is one data instance. The first column is the label (1, 2 or 3) and the next two columns are features. Build up your data loading process re-arrange the dataset into a suitable format for training.

We will train two neural networks. The first will have only the output layer, i.e., it is described by a 2×3 matrix with a bias and is a linear classifier. The second is a nonlinear classifier that will have one hidden layer with 5 units. Use AdamW as your optimizer and use CrossEntropyLoss. Train the two network for 1000 epochs on the training dataset. Try running your code a few times with different random initializations, since you may hit a poor local optima. You should normalize the data when given to the network by subtracting the mean of the training data. You only need to report on your best run.

Please provide the following:

1. Show the train loss curves in two plots for both networks (2 pts). Label one linear neural network and the other nonlinear neural network.
2. Measure the accuracy of both networks on the testing *and* training datasets (3 pts). Put this in a \LaTeX table.
3. Create decision boundaries for both networks and show them. Make sure to label your plots. Feel free to find helper code online for displaying the decision boundary (5 pts). An example decision boundary is provided in Fig. 1. Show the training data

points in the figure, and color them appropriately for their label.

4. Provide your code for credit.

Answer:

Code for Problem 4:

CODE HERE