

Homework #4

CSE 446/546: Machine Learning

Madeline Brown

Due: **Wednesday** Dec 4, 2024 11:59pm

Points A: 96; B: 25

Outside references:

math 318 lecture notes (I am a TA for that class and we cover SVD and PCA so I used this as a cross-reference)

<https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.linalg.svds.html>

perplexity.ai for syntax and definition help

Please review all homework guidance posted on the website before submitting to Gradescope. Reminders:

- All code must be written in Python and all written work must be typeset (e.g. \LaTeX).
- Make sure to read the “What to Submit” section following each question and include all items.
- Please provide succinct answers and supporting reasoning for each question. Similarly, when discussing experimental results, concisely create tables and/or figures when appropriate to organize the experimental results. All explanations, tables, and figures for any particular part of a question must be grouped together.
- For every problem involving generating plots, please include the plots as part of your PDF submission.
- When submitting to Gradescope, please link each question from the homework in Gradescope to the location of its answer in your homework PDF. Failure to do so may result in deductions of up to 10% of the value of each question not properly linked. For instructions, see https://www.gradescope.com/get_started#student-submission.

Important: By turning in this assignment (and all that follow), you acknowledge that you have read and understood the collaboration policy with humans and AI assistants alike: <https://courses.cs.washington.edu/courses/cse446/24au/assignments/>. Any questions about the policy should be raised at least 24 hours before the assignment is due. There are no warnings or second chances. If we suspect you have violated the collaboration policy, we will report it to the college of engineering who will complete an investigation. Not adhering to these reminders may result in point deductions.

Conceptual Questions

A1. The answers to these questions should be answerable without referring to external materials. Briefly justify your answers with a few words.

- a. [2 points] True or False: Given a data matrix $X \in \mathbb{R}^{n \times d}$ where d is much smaller than n and $k = \text{rank}(X)$, if we project our data onto a k -dimensional subspace using PCA, our projection will have zero reconstruction error (in other words, we find a perfect representation of our data, with no information loss).

True, since X is rank k , the k -dimensional subspace that PCA projects onto will just be the column space of X (since this has dimension k) and thus have zero information loss as the projection will be equal to the original.

- b. [2 points] True or False: Suppose that an $n \times n$ matrix X has a singular value decomposition of USV^T , where S is a diagonal $n \times n$ matrix. Then, the rows of V are equal to the eigenvectors of $X^T X$.

False. The columns of V are equal to the eigenvectors of $X^T X$, as this is the definition of right singular vectors. The eigenvectors of $X^T X$ are the rows of V^T . This is my understanding from my experience with Math 318, notes attached:

Lemma 8.2.1. Suppose $A \in \mathbb{R}^{m \times n}$ and $\text{rank}(A) = r$.

1. The **singular values** of A are the positive square roots of the r positive eigenvalues of $A^T A$ (and AA^T).
2. The **right singular vectors** of A are an orthonormal basis of **eigenvectors of $A^T A$** . The first r of them form an orthonormal basis of $\text{Row}(A)$ and the last $n - r$ of them form an orthonormal basis of $\text{Null}(A)$.
3. The **left singular vectors** of A are an orthonormal basis of eigenvectors of AA^T . The first r of them form an orthonormal basis of $\text{Col}(A)$ and the last $m - r$ of them form an orthonormal basis of $\text{Null}(A^T)$.

- c. [2 points] True or False: choosing k to minimize the k -means objective (see Equation (4) below) is a good way to find meaningful clusters.

False. If we are just minimizing the distance from the mean of every point in a cluster, then we could just make one cluster for each point and this would give the best minimum value of the objective, but the most meaningless clusters, since we would have no clusters at all.

$$\min_{\pi_1, \dots, \pi_k} \sum_{i=1}^k \sum_{j \in \pi_i} \|\mathbf{x}_j - \mu_i\|_2^2, \quad \mu_i = \frac{1}{|\pi_i|} \sum_{j \in \pi_i} \mathbf{x}_j$$

- d. [2 points] True or False: The singular value decomposition of a matrix is unique.

False. If X has SVD USV^T , we can create a new SVD, for example, by taking a factor of -1 out of some columns in V and putting this factor of -1 onto the corresponding columns of U (multiply these columns by -1). The resulting matrices would still satisfy the definition of SVD, as the columns of U and V would still form a basis for correct spaces, but this would not be equivalent to the old SVD.

- e. [2 points] True or False: The rank of a square matrix equals the number of its unique nonzero eigenvalues.

False. The matrix $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ has rank 2 (two pivots) but only has the eigenvalues 1 and 0. There is only one unique nonzero eigenvalue, so this shows that not every square matrix has rank equal to unique nonzero eigenvalues.

What to Submit:

- **Parts a-e:** 1-2 sentence explanation containing your answer.

Think before you train

A2. **The first part of this problem (part a)** explores how you would apply machine learning theory and techniques to a real-world problem. There is one scenario detailing a setting, a dataset, and a specific result we hope to achieve. Your job is to describe how you would handle the scenario with the tools we've learned in this class. Your response should include:

- (1) any pre-processing steps you would take (e.g. any data processing),
- (2) the specific machine learning pipeline you would use (i.e., algorithms and techniques learned in this class),
- (3) how your setup acknowledges the constraints and achieves the desired result.

You should also aim to leverage some of the theory we have covered in this class. Some things to consider may be: the nature of the data (i.e., *How hard is it to learn? Do we need more data? Are the data sources good?*), the effectiveness of the pipeline (i.e., *How strong is the model when properly trained and tuned?*), and the time needed to effectively perform the pipeline.

a. *[10 points]* **Scenario: Disease Susceptibility Predictor**

- Setting: You are tasked by a research institute to create an algorithm that learns the factors that contribute most to acquiring a specific disease.
- Dataset: A rich dataset of personal demographic information, location information, risk factors, and whether a person has the disease or not.
- Result: The company wants a system that can determine how susceptible someone is to this disease when they enter in their own personal information. The pipeline should take limited amount of personal data from a new user and infer more detailed metrics about the person.

To create an algorithm that learns factors which contribute the most to acquiring a specific disease, I will utilize the rich dataset provided and make a model using logistic regression, since we are hoping to predict the binary outcome of if someone has the disease or not. Since we are wanting to find the factors with the greatest contribution, we will also do some regularization, likely lasso or ridge, in order to shrink the coefficients of the less-important factors down to zero, leaving us with only the most important ones. While creating the model, we would also like to use some form of validation to test and fine-tune our parameters, which we can do since we have such a rich dataset (it would be harder with less data). We can section off about 10% of the data for validation (choosing parameters) and 10% for testing and evaluating the model after we use validation to find the best performing parameters. After the model is created and coefficients are chosen for each contributing factor, we can determine how likely someone is to have the disease by taking their personal information and applying the coefficients from the model onto the personal data.

The second part of this problem (parts b, c) focuses on exploring possible shortcomings of machine learning models, and what real-world implications might follow from ignoring these issues.

- b. *[5 points]* Briefly describe (1) some potential shortcomings of your training process from the disease susceptibility predictor scenario above that may result in your algorithm having different accuracy on different populations, and (2) how you may modify your procedure to address these shortcomings.

It is possible that the data being used to train my model does not equally represent the same diverse population of people that the model is likely to be used on, for example, lots of data from one geographic location, but very little from another. This may mean that the model would be less accurate on demographics that are not well represented in the data, meaning inaccurate disease predictions.

In these cases, I might want to apply more regularization in order to further minimize the effect of certain variables (the ones that don't completely reflect the overarching population at hand.) Another option would be, after identifying which variables lack diversity, to create separate models while controlling for said variables. For example, if one of the variables was "state" and there were 100 people from Oregon but only 50 people from Washington, I might separate Oregon and Washington in the data before doing the logistic regression. If the two models give similar results, I would go ahead and combine them again, but if they are extremely different, I would proceed with caution, taking careful note of the size of the coefficient(s) for location. If it is large, I might look at the outlier locations more closely.

- c. *[5 points]* Recall in Homework 2 we trained models to predict crime rates using various features. It is important to note that **datasets describing crime have many shortcomings in describing the entire landscape of illegal behavior in a city, and that these shortcomings often fall disproportionately on minority communities**. Some of these shortcomings include that crimes are reported at different rates in different neighborhoods, that police respond differently to the same crime reported or observed in different neighborhoods, and that police spend more time patrolling in some neighborhoods than others. What real-world implications might follow from ignoring these issues?

If these issues in the dataset are ignored, our model may end up perpetuating biases that the data reflects, or otherwise causing harm to the people impacted by the decisions being made by the model. In the case of the crime-predicting model, as described, the data might show that there is more police activity in areas where crimes are reported more often, but the police response may not accurately reflect the level of crime in that neighborhood. If we have a model that predicts where crime happens based on this information, it will mainly say that there will be crime in those areas where crime has been reported before. If there is a place where crime happens a lot but goes unreported, the model will not tell decision makers to focus on that area. If the police are so biased as to focus on a specific area for their patrolling, resulting in more crime reported for that area (because there are more officers there to record it or perpetuate escalating situations, regardless if a crime was really taking place) the model might tell decision makers that more police need to go there, which would just result in more crimes reported in that area, and the cycle continues.

What to Submit:

- **For part (a):** One clearly-written short paragraph (approximately 4-7 sentences).
- **For part (b):** Clearly-written and well-thought-out answers addressing (1) and (2) (as described in the problem). Two short paragraphs or one medium paragraph suffice.
- **For part (c):** One clearly-written short paragraph on real-world implications that may follow from ignoring dataset issues.

Image Classification on CIFAR-10

A3. In this problem we will explore different deep learning architectures for image classification on the CIFAR-10 dataset. Make sure that you are familiar with `torch.Tensors`, two-dimensional convolutions (`nn.Conv2d`) and fully-connected layers (`nn.Linear`), ReLU non-linearities (`F.relu`), pooling (`nn.MaxPool2d`), and tensor reshaping (`view`).

Hint: For loops are costly. Can you vectorize it or use Numpy operations to make it faster in some ways?

A few preliminaries:

- Make sure to read the “Tips for HW4” EdStem post for additional tips about training your models.
- Each network f maps an image $x^{\text{in}} \in \mathbb{R}^{32 \times 32 \times 3}$ (3 channels for RGB) to an output $f(x^{\text{in}}) = x^{\text{out}} \in \mathbb{R}^{10}$. The class label is predicted as $\arg \max_{i=0,1,\dots,9} x_i^{\text{out}}$. An error occurs if the predicted label differs from the true label for a given image.
- The network is trained via multiclass cross-entropy loss.
- Create a validation dataset by appropriately partitioning the train dataset. *Hint:* look at the documentation for `torch.utils.data.random_split`. Make sure to tune hyperparameters like network architecture and step size on the validation dataset. Do **NOT** validate your hyperparameters on the test dataset.
- At the end of each epoch (one pass over the training data), compute and print the training and validation classification accuracy.
- While one could train a network for hundreds of epochs to reach convergence and maximize accuracy, this can be prohibitively time-consuming, so feel free to train for just a dozen or so epochs.

For parts (a) and (b), apply a hyperparameter tuning method (e.g. random search, grid search, etc.) using the validation set, report the hyperparameter configurations you evaluated and the best set of hyperparameters from this set, and plot the training and validation classification accuracy as a function of epochs. Produce a separate line or plot for each hyperparameter configuration evaluated (top 3 configurations is sufficient to keep the plots clean). Finally, evaluate your best set of hyperparameters on the test data and report the test accuracy.

Note 1: Please refer to the provided notebook with starter code for this problem, on the course website. That notebook provides a complete end-to-end example of loading data, training a model using a simple network with a fully-connected output and no hidden layers (this is equivalent to logistic regression), and performing evaluation using canonical Pytorch. We recommend using this as a template for your implementations of the models below.

Note 2: If you are attempting this problem and do not have access to GPU we highly recommend using Google Colab. The provided notebook includes instructions on how to use GPU in Google Colab.

Here are the network architectures you will construct and compare.

- a. **[18 points] Fully-connected output, 1 fully-connected hidden layer:** this network has one hidden layer denoted as $x^{\text{hidden}} \in \mathbb{R}^M$ where M will be a hyperparameter you choose (M could be in the hundreds). The nonlinearity applied to the hidden layer will be the **relu** ($\text{relu}(x) = \max\{0, x\}$). This network can be written as

$$x^{\text{out}} = W_2 \text{relu}(W_1(x^{\text{in}}) + b_1) + b_2$$

where $W_1 \in \mathbb{R}^{M \times 3072}$, $b_1 \in \mathbb{R}^M$, $W_2 \in \mathbb{R}^{10 \times M}$, $b_2 \in \mathbb{R}^{10}$. Tune the different hyperparameters and train for a sufficient number of epochs to achieve a *validation accuracy* of at least 50%. Provide the hyperparameter configuration used to achieve this performance.

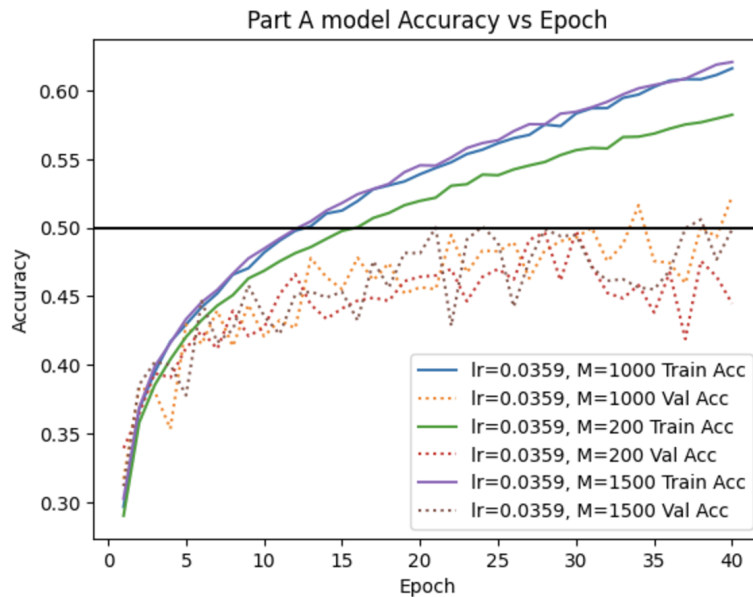
Hyperparameter tuning method used with validation set: grid search (iterating through every combination of lr and M).

Hyperparameter configurations evaluated (learning rates and Ms): Ms = [200, 500, 1000, 1500], learning rates approximately 9.99999e-06, 2.782559e-05, 7.7426e-05, 0.000215, 0.000599, 0.001668, 0.00464, 0.01292, 0.03594, and 0.10000.

Best hyperparameters from above set (top 3, based on validation accuracy.) Note that these models were only trained up to five epochs for this searching phase:

```
learning rate: 0.03593813627958298
M: 1000
Validation accuracy: 0.4330
learning rate: 0.03593813627958298
M: 200
Validation accuracy: 0.4168
learning rate: 0.03593813627958298
M: 1500
Validation accuracy: 0.4168
best validation accuracy= 0.4330078125 for lr=0.03593813627958298, M=1000
```

Plot of training and validation classification accuracy as a function of epochs (all three best models on the plot, for six lines total):



(next page)

Test accuracy of best sets of hyperparameters on the test data (best three models):

```
Test Accuracy for M=1000, lr=0.03593813627958298: 0.5306566455696202
Test Accuracy for M=200, lr=0.03593813627958298: 0.4606408227848101
Test Accuracy for M=1500, lr=0.03593813627958298: 0.5029667721518988
```

- b. [18 points] **Convolutional layer with max-pool and fully-connected output:** for a convolutional layer W_1 with filters of size $k \times k \times 3$, and M filters (reasonable choices are $M = 100$, $k = 5$), we have that $\text{Conv2d}(x^{\text{in}}, W_1) \in \mathbb{R}^{(33-k) \times (33-k) \times M}$.

- Each convolution will have its own offset applied to each of the output pixels of the convolution; we denote this as $\text{Conv2d}(x^{\text{in}}, W) + b_1$ where b_1 is parameterized in \mathbb{R}^M . Apply a **relu** activation to the result of the convolutional layer.
- Next, use a max-pool of size $N \times N$ (a reasonable choice is $N = 14$ to pool to 2×2 with $k = 5$) we have that $\text{MaxPool}(\text{relu}(\text{Conv2d}(x^{\text{in}}, W_1) + b_1)) \in \mathbb{R}^{\lfloor \frac{33-k}{N} \rfloor \times \lfloor \frac{33-k}{N} \rfloor \times M}$.
- We will then apply a fully-connected layer to the output to get a final network given as

$$x^{\text{output}} = W_2(\text{MaxPool}(\text{relu}(\text{Conv2d}(x^{\text{input}}, W_1) + b_1))) + b_2$$

where $W_2 \in \mathbb{R}^{10 \times M(\lfloor \frac{33-k}{N} \rfloor)^2}$, $b_2 \in \mathbb{R}^{10}$.

The parameters M, k, N (in addition to the step size and momentum) are all hyperparameters, but you can choose a reasonable value. Tune the different hyperparameters (number of convolutional filters, filter sizes, dimensionality of the fully-connected layers, step size, etc.) and train for a sufficient number of epochs to achieve a *validation accuracy* of at least 65%. Provide the hyperparameter configuration used to achieve this performance.

The number of hyperparameters to tune, combined with the slow training times, will hopefully give you a taste of how difficult it is to construct networks with good generalization performance. State-of-the-art networks can have dozens of layers, each with their own hyperparameters to tune. Additional hyperparameters you are welcome to play with, if you are so inclined, include: changing the activation function, replace max-pool with average-pool, adding more convolutional or fully connected layers, and experimenting with batch normalization or dropout.

Hyperparameter tuning method used with validation set: grid search (iterating through every combination of lr, M, k, and N), which may have been a mistake given the number of parameters to check.

Hyperparameter configurations evaluated (learning rates and Ms): Ms = [200, 500, 1500], ks = [3,4,5,6,7], Ns = [2,4,8,16,26], learning rates approximately 9.9999e-05, 0.00022, 0.00046, 0.0010, 0.00215, 0.0046, 0.00999, 0.0215, 0.0464, and 0.1000.

Best hyperparameters from above set (top 3, based on validation accuracy.) Note that these models were only trained up to three epochs for this searching phase, and on the 1/10 training data set:

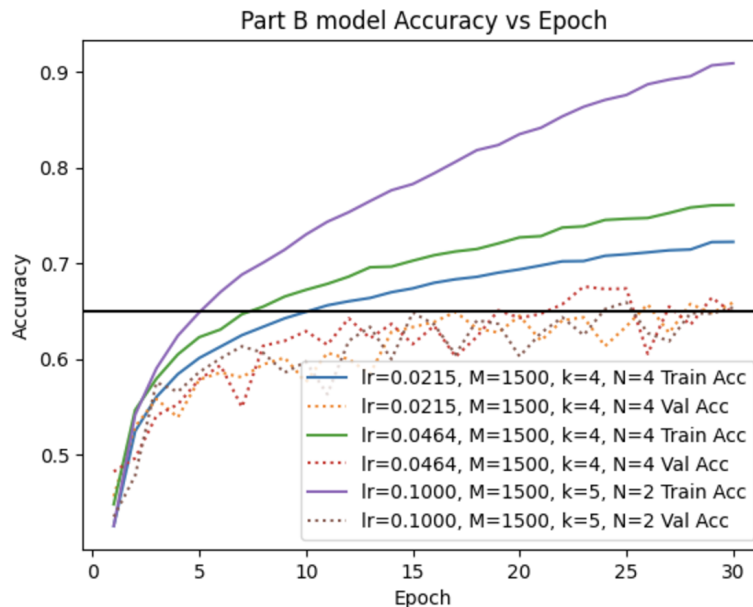
```
learning rate: 0.02154434658586979
M: 1500
k: 4
N: 4
Validation accuracy: 0.4141

learning rate: 0.04641588777303696
M: 1500
k: 4
N: 4
Validation accuracy: 0.4004

learning rate: 0.10000000149011612
M: 1500
k: 5
N: 2
Validation accuracy: 0.3984

best validation accuracy= 0.4140625 for lr=0.02154434658586979, M=1500, k=4, N=4
```


Plot of training and validation classification accuracy as a function of epochs (all three best models on the plot, for six lines total):



Test accuracy of best sets of hyperparameters on the test data (best three models):

```
Test Accuracy for lr=0.0215, M=1500, k=4, N=4: 0.6661392405063291
Test Accuracy for lr=0.0464, M=1500, k=4, N=4: 0.6597112341772152
Test Accuracy for lr=0.1000, M=1500, k=5, N=2: 0.6600079113924051
```

What to Submit:

- **For parts (a)-(b):** A single plot of the training and validation accuracy for the top 3 hyperparameter configurations you evaluated (x-axis is training epoch; y-axis is accuracy; this plot should contain 6 lines total). If it took fewer than 3 hyperparameter configurations to pass the performance threshold, plot all hyperparameter configurations you evaluated. A horizontal line should be plotted at the targeted threshold (50% or 65%). Validation lines should be dotted, and training lines should be solid.
- **For parts (a)-(b):** List the hyperparameter values you searched over and your search method (random, grid, etc.). Provide the values of best performing hyperparameters, and accuracy of best models on test data.
- **For parts (a)-(b):** Code. You should convert your code (the .ipynb notebook) into a Python (.py) file, rename it to `hw4-a3.py`, and submit it to the corresponding Gradescope submission. To download the file from Google Colab, you can go to File ↴ Download ↴ Download as .py.

Homework 4: Image Classification on CIFAR-10

Information before starting

In this problem, we will explore different deep learning architectures for image classification on the CIFAR-10 dataset. Make sure that you are familiar with torch `Tensor`'s, two-dimensional convolutions (`nn.Conv2d`) and fully-connected layers (`nn.Linear`), ReLU non-linearities (`F.relu`), pooling (`nn.MaxPool2d`), and tensor reshaping (`view`). ****Make sure to read through all instructions in both this notebook and in the PDF while completing this problem!****

Copying this Colab Notebook to your Google Drive

Since the course staff is the author of this notebook, you cannot make any lasting changes to it. You should make a copy of it to your Google Drive by clicking **File -> Save a Copy in Drive**.

Problem Introduction

You've already had some practice using the PyTorch library in HW3, but this problem dives into training more complex deep learning models.

The specific task we are trying to solve in this problem is image classification. We're using a common dataset called CIFAR-10 which has 60,000 images separated into 10 classes:

- * airplane
- * automobile
- * bird
- * cat
- * deer
- * dog
- * frog
- * horse
- * ship
- * truck

We've provided an end-to-end example of loading data, training a model, and performing evaluation. We recommend using this code as a template for your implementations of the more complex models. Feel free to modify or reuse any of the functions we provide.

****Unlike other coding problems in the past, this one does not include an autograded component.****

Enabling GPU

We are using Google Colab because it has free GPU runtimes available. GPUs can accelerate training times for this problem by 10-100x when compared to using CPU. To use the GPU runtime on Colab, make sure to **enable** the runtime by going to **Runtime -> Change runtime type -> Select T4 GPU under "Hardware accelerator"**.

Note that GPU runtimes are *limited* on Colab. We recommend limiting your training to short-running jobs (under 15 minutes each) and spread your work over time, if possible. Colab *will* limit your usage of GPU time, so plan ahead and be prepared to take breaks during training. If you have used up your quota for GPU, check back in a day or so to be able to enable GPU again.

Your code will still run on CPU, so if you are just starting to implement your code or have been GPU limited by Colab, you can still make changes and run your code - it will just be quite a bit slower. You can also choose to download your notebook and run locally if you have a personal GPU or have a faster CPU than the one Colab provides. If you choose to do this, you may need to install the packages that this notebook depends on to your ``cse446`` conda environment or to another Python environment of your choice.

To check if you have enabled GPU, run the following cell. If ``device`` is ``cuda``, it means that GPU has been enabled successfully.

```
"""
```

```

import torch

DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
print(DEVICE) # this should print out CUDA

"""### Submitting your assignment

Once you are done with the problem, make sure to put all of your necessary figures into your
PDF submission. Then, download this notebook as a Python file (`.py`) by going to **File ->
Download -> Download `.py`**. Rename this file as `hw4-a3.py` and upload to the Gradescope
submission for HW4 code.

## End-to-end Example

### Background and Setup

1. We first import all of the dependencies required for this problem:
"""

# Commented out IPython magic to ensure Python compatibility.
import torch
from torch import nn
import numpy as np

from typing import Tuple, Union, List, Callable
from torch.optim import SGD
import torchvision
from torch.utils.data import DataLoader, TensorDataset, random_split
import matplotlib.pyplot as plt
from tqdm.notebook import tqdm

# %matplotlib inline

"""2. And check if we are using GPU, if it is available. (Make sure to set your runtime to
enable GPU!)"""

DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
print(DEVICE) # this should print out CUDA

"""*To use the GPU, you will need to send both the model and the data to a device; this
transfers the model from its default location on CPU to the GPU.*

*Note that torch operations on Tensors will fail if they are not located on the same device.
Here's a small example of how to send the model and data to your device:*

```python
model = model.to(DEVICE) # Sending a model to GPU

for x, y in tqdm(data_loader):
 x, y = x.to(DEVICE), y.to(DEVICE)
```

*When reading tensors you may need to send them back to cpu, you can do so with `x =
x.cpu()`*

```

3. Now, let's load the CIFAR-10 data. We can take advantage of public datasets available through PyTorch torchvision!

```
"""
```

```
train_dataset = torchvision.datasets.CIFAR10("./data", train=True, download=True,
transform=torchvision.transforms.ToTensor())
test_dataset = torchvision.datasets.CIFAR10("./data", train=False, download=True,
transform=torchvision.transforms.ToTensor())
```

```
"""4. Finally, like we did in HW3, we'll use the PyTorch `DataLoader` to wrap our datasets.
You've already seen that `DataLoader`s handle batching, shuffling, and iterating over data,
and this is really useful for this problem as well!
```

```
*Since training on all of the 50,000 training samples can be prohibitively expensive, we
define a flag called* `SAMPLE_DATA` *that controls if we should make the dataset smaller for
faster training time. When* `SAMPLE_DATA=true`, *we'll only use 10% of our training data
when training and performing our hyperparameter searches.* **Make sure that you've set
`SAMPLE_DATA=false` when you want to perform your final training loops for submission!**
"""
```

```
SAMPLE_DATA = True # set this to True if you want to speed up training when searching for
hyperparameters!
```

```
batch_size = 128
```

```
if SAMPLE_DATA:
    train_dataset, _ = random_split(train_dataset, [int(0.1 * len(train_dataset)), int(0.9 *
len(train_dataset))]) # get 10% of train dataset and "throw away" the other 90%
```

```
train_dataset, val_dataset = random_split(train_dataset, [int(0.9 * len(train_dataset)),
int( 0.1 * len(train_dataset))])
```

```
# Create separate dataloaders for the train, test, and validation set
```

```
train_loader = DataLoader(
    train_dataset,
    batch_size=batch_size,
    shuffle=True
)
```

```
val_loader = DataLoader(
    val_dataset,
    batch_size=batch_size,
    shuffle=True
)
```

```
test_loader = DataLoader(
    test_dataset,
    batch_size=batch_size,
    shuffle=True
)
```

```
"""Now, we're ready to train!
```

```
### Logistic Regression Example
```

Let's first take a look at our data to get an understanding of what we are doing. As a reminder, CIFAR-10 is a dataset containing images split into 10 classes.

```
"""
```

```
imgs, labels = next(iter(train_loader))
print(f"A single batch of images has shape: {imgs.size()}")
example_image, example_label = imgs[0], labels[0]
c, w, h = example_image.size()
print(f"A single RGB image has {c} channels, width {w}, and height {h}.")

# This is one way to flatten our images
batch_flat_view = imgs.view(-1, c * w * h)
print(f"Size of a batch of images flattened with view: {batch_flat_view.size()}")

# This is another equivalent way
batch_flat_flatten = imgs.flatten(1)
print(f"Size of a batch of images flattened with flatten: {batch_flat_flatten.size()}")

# The new dimension is just the product of the ones we flattened
d = example_image.flatten().size()[0]
print(c * w * h == d)

# View the image
t = torchvision.transforms.ToPILImage()
plt.imshow(t(example_image))

# These are what the class labels in CIFAR-10 represent. For more information,
# visit https://www.cs.toronto.edu/~kriz/cifar.html
classes = ["airplane", "automobile", "bird", "cat", "deer", "dog", "frog",
           "horse", "ship", "truck"]
print(f"This image is labeled as class {classes[example_label]}")
```

```
"""In this problem, we will attempt to predict what class an image is labeled as.
```

```
1. First, let's create our model. Note: for a linear model we could flatten the data before
passing it into the model, but that is not the case for convolutional neural networks.
"""
```

```
def linear_model() -> nn.Module:
    """Instantiate a linear model and send it to device."""
    model = nn.Sequential(
        nn.Flatten(),
        nn.Linear(d, 10)
    )
    return model.to(DEVICE)

"""2. Let's define a method to train this model using SGD as our optimizer."""

def train(
    model: nn.Module, optimizer: SGD,
    train_loader: DataLoader, val_loader: DataLoader,
    epochs: int = 20
) -> Tuple[List[float], List[float], List[float], List[float]]:
    """
    Trains a model for the specified number of epochs using the loaders.
```

Returns:

Lists of training loss, training accuracy, validation loss, validation accuracy for each epoch.

"""

```
loss = nn.CrossEntropyLoss()
train_losses = []
train_accuracies = []
val_losses = []
val_accuracies = []
for e in tqdm(range(epochs)):
    model.train()
    train_loss = 0.0
    train_acc = 0.0

    # Main training loop; iterate over train_loader. The loop
    # terminates when the train loader finishes iterating, which is one epoch.
    for (x_batch, labels) in train_loader:
        x_batch, labels = x_batch.to(DEVICE), labels.to(DEVICE)
        optimizer.zero_grad()
        labels_pred = model(x_batch)
        batch_loss = loss(labels_pred, labels)
        train_loss = train_loss + batch_loss.item()

        labels_pred_max = torch.argmax(labels_pred, 1)
        batch_acc = torch.sum(labels_pred_max == labels)
        train_acc = train_acc + batch_acc.item()

        batch_loss.backward()
        optimizer.step()
    train_losses.append(train_loss / len(train_loader))
    train_accuracies.append(train_acc / (batch_size * len(train_loader)))

    # Validation loop; use .no_grad() context manager to save memory.
    model.eval()
    val_loss = 0.0
    val_acc = 0.0

    with torch.no_grad():
        for (v_batch, labels) in val_loader:
            v_batch, labels = v_batch.to(DEVICE), labels.to(DEVICE)
            labels_pred = model(v_batch)
            v_batch_loss = loss(labels_pred, labels)
            val_loss = val_loss + v_batch_loss.item()

            v_pred_max = torch.argmax(labels_pred, 1)
            batch_acc = torch.sum(v_pred_max == labels)
            val_acc = val_acc + batch_acc.item()
        val_losses.append(val_loss / len(val_loader))
        val_accuracies.append(val_acc / (batch_size * len(val_loader)))

return train_losses, train_accuracies, val_losses, val_accuracies
```

"""3. Now, let's define our hyperparameter search. For this problem, we will be using SGD. The two hyperparameters for our linear model trained with SGD are the learning rate and momentum. Only learning rate will be searched for in this example, but you will be tuning multiple hyperparameters. **Feel free to experiment with hyperparameters and how you search. We recommend implementing random search!**

Note: We ask you to plot the accuracies for the 3 best models for each structure, so you will need to return multiple sets of hyperparameters for the homework.

```
def parameter_search(train_loader: DataLoader,
                    val_loader: DataLoader,
                    model_fn: Callable[[], nn.Module]) -> float:
    """
    Parameter search for our linear model using SGD.

    Args:
    train_loader: the train dataloader.
    val_loader: the validation dataloader.
    model_fn: a function that, when called, returns a torch.nn.Module.

    Returns:
    The learning rate with the least validation loss.
    NOTE: you may need to modify this function to search over and return
    other parameters beyond learning rate.
    """
    num_iter = 10
    best_loss = torch.tensor(np.inf)
    best_lr = 0.0

    lrs = torch.linspace(10 ** (-6), 10 ** (-1), num_iter)

    for lr in lrs:
        print(f"trying learning rate {lr}")
        model = model_fn()
        optim = SGD(model.parameters(), lr)
        train_loss, train_acc, val_loss, val_acc = train(
            model,
            optim,
            train_loader,
            val_loader,
            epochs=20
        )

        if min(val_loss) < best_loss:
            best_loss = min(val_loss)
            best_lr = lr

    return best_lr

"""4. Now that we have everything, we can train and evaluate our model."""

best_lr = parameter_search(train_loader, val_loader, linear_model)

model = linear_model()
```

```

optimizer = SGD(model.parameters(), best_lr)

# We are using 20 epochs for this example. You may have to use more.
train_loss, train_accuracy, val_loss, val_accuracy = train(
    model, optimizer, train_loader, val_loader, 20)

"""5. We can also plot the training and validation accuracy for each epoch."""

epochs = range(1, 21)
plt.plot(epochs, train_accuracy, label="Train Accuracy")
plt.plot(epochs, val_accuracy, label="Validation Accuracy")
plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.legend()
plt.title("Logistic Regression Accuracy for CIFAR-10 vs Epoch")
plt.show()

"""The last thing we have to do is evaluate our model on the testing data."""

def evaluate(
    model: nn.Module, loader: DataLoader
) -> Tuple[float, float]:
    """Computes test loss and accuracy of model on loader."""
    loss = nn.CrossEntropyLoss()
    model.eval()
    test_loss = 0.0
    test_acc = 0.0
    with torch.no_grad():
        for (batch, labels) in loader:
            batch, labels = batch.to(DEVICE), labels.to(DEVICE)
            y_batch_pred = model(batch)
            batch_loss = loss(y_batch_pred, labels)
            test_loss = test_loss + batch_loss.item()

            pred_max = torch.argmax(y_batch_pred, 1)
            batch_acc = torch.sum(pred_max == labels)
            test_acc = test_acc + batch_acc.item()
    test_loss = test_loss / len(loader)
    test_acc = test_acc / (batch_size * len(loader))
    return test_loss, test_acc

test_loss, test_acc = evaluate(model, test_loader)
print(f"Test Accuracy: {test_acc}")

```

"""The rest is yours to code. You can structure the code any way you would like.

We do advise making using code cells and functions (train, search, predict etc.) for each subproblem, since they will make your code easier to debug.

Also note that several of the functions above can be reused for the various different models you will implement for this problem; i.e., you won't need to write a new `evaluate()`.

Your Turn!

The rest is yours to code. You are welcome to structure the code any way you would like.

We do advise making using code cells and functions (train, search, predict etc.) for each subproblem, since they will make your code easier to debug.

Also note that several of the functions above can be reused for the various different models you will implement for this problem; i.e., you won't need to write a new `evaluate()`. Before you reuse functions though, make sure they are compatible with what the assignment is asking for.

Submitting Code

And as a last reminder, once you are done with the problem, make sure to put all of your necessary figures into your PDF submission. Then, download this notebook as a Python file (`.py`) by going to **File -> Download -> Download `.py`**. Rename this file as `hw4-a3.py` and upload to the Gradescope submission for HW4 code.

```
preamble:
"""
```

```
# Commented out IPython magic to ensure Python compatibility.
```

```
import torch
from torch import nn
import numpy as np
```

```
from typing import Tuple, Union, List, Callable
from torch.optim import SGD
import torchvision
from torch.utils.data import DataLoader, TensorDataset, random_split
import matplotlib.pyplot as plt
from tqdm.notebook import tqdm
```

```
# %matplotlib inline
```

```
"""load data:"""
```

```
train_dataset = torchvision.datasets.CIFAR10("./data", train=True, download=True,
transform=torchvision.transforms.ToTensor())
test_dataset = torchvision.datasets.CIFAR10("./data", train=False, download=True,
transform=torchvision.transforms.ToTensor())
```

```
"""use less data for debugging:"""
```

```
SAMPLE_DATA = False # set this to True if you want to speed up training when searching for
hyperparameters!
```

```
"""separate data sets:"""
```

```
batch_size = 128
```

```
if SAMPLE_DATA:
    train_dataset, _ = random_split(train_dataset, [int(0.1 * len(train_dataset)), int(0.9 *
len(train_dataset))]) # get 10% of train dataset and "throw away" the other 90%
```

```
train_dataset, val_dataset = random_split(train_dataset, [int(0.9 * len(train_dataset)),
int( 0.1 * len(train_dataset))])
```

```

# Create separate dataloaders for the train, test, and validation set
train_loader = DataLoader(
    train_dataset,
    batch_size=batch_size,
    shuffle=True
)

val_loader = DataLoader(
    val_dataset,
    batch_size=batch_size,
    shuffle=True
)

test_loader = DataLoader(
    test_dataset,
    batch_size=batch_size,
    shuffle=True
)

"""define models:"""

#part A model
def A_model(M) -> nn.Module: #M is hyperparameter
    """Fully-connected output, 1 fully-connected hidden layer"""
    model = nn.Sequential(
        nn.Flatten(),
        nn.Linear(3072, M),
        nn.ReLU(),
        nn.Linear(M,10)
    )
    return model.to(DEVICE)

"""train function (not modified):"""

def train(
    model: nn.Module, optimizer: SGD,
    train_loader: DataLoader, val_loader: DataLoader,
    epochs: int = 20
) -> Tuple[List[float], List[float], List[float], List[float]]:
    """
    Trains a model for the specified number of epochs using the loaders.

    Returns:
    Lists of training loss, training accuracy, validation loss, validation accuracy for each
    epoch.
    """

    loss = nn.CrossEntropyLoss()
    train_losses = []
    train_accuracies = []
    val_losses = []
    val_accuracies = []
    for e in tqdm(range(epochs)):
        model.train()

```

```

train_loss = 0.0
train_acc = 0.0

# Main training loop; iterate over train_loader. The loop
# terminates when the train loader finishes iterating, which is one epoch.
for (x_batch, labels) in train_loader:
    x_batch, labels = x_batch.to(DEVICE), labels.to(DEVICE)
    optimizer.zero_grad()
    labels_pred = model(x_batch)
    batch_loss = loss(labels_pred, labels)
    train_loss = train_loss + batch_loss.item()

    labels_pred_max = torch.argmax(labels_pred, 1)
    batch_acc = torch.sum(labels_pred_max == labels)
    train_acc = train_acc + batch_acc.item()

    batch_loss.backward()
    optimizer.step()
train_losses.append(train_loss / len(train_loader))
train accuracies.append(train_acc / (batch_size * len(train_loader)))

# Validation loop; use .no_grad() context manager to save memory.
model.eval()
val_loss = 0.0
val_acc = 0.0

with torch.no_grad():
    for (v_batch, labels) in val_loader:
        v_batch, labels = v_batch.to(DEVICE), labels.to(DEVICE)
        labels_pred = model(v_batch)
        v_batch_loss = loss(labels_pred, labels)
        val_loss = val_loss + v_batch_loss.item()

        v_pred_max = torch.argmax(labels_pred, 1)
        batch_acc = torch.sum(v_pred_max == labels)
        val_acc = val_acc + batch_acc.item()
    val_losses.append(val_loss / len(val_loader))
    val accuracies.append(val_acc / (batch_size * len(val_loader)))

return train_losses, train accuracies, val_losses, val accuracies

"""parameter search (modified to include M, and to return top three model parameter
sets):"""

def parameter_search(train_loader: DataLoader,
                    val_loader: DataLoader,
                    model_fn: Callable[[int], nn.Module]) -> Tuple[float, int]:
    """
    Parameter search for our linear model.

    Args:
    train_loader: the train dataloader.
    val_loader: the validation dataloader.
    model_fn: a function that, when called, returns a torch.nn.Module.

```

```

Returns:
The learning rate with the least validation loss.
NOTE: you may need to modify this function to search over and return
other parameters beyond learning rate.
"""
num_iter = 10
best_loss = torch.tensor(np.inf)
best_lr = 0.0
best_M = 0 # modified
best_val_accuracy = 0.0 # modified

Ms = [200, 500, 1000, 1500] # modified
#lrs = torch.linspace(10 ** (-6), 10 ** (-1), num_iter)
lrs = torch.logspace(-5, -1, num_iter)

# keep track of all params tried to return the top three
all_params = []

for M in Ms: # modified
    for lr in lrs:
        print(f"trying learning rate {lr}, M={M}") # modified
        model = model_fn(M) # modified

        optim = SGD(model.parameters(), lr)

        train_loss, train_acc, val_loss, val_acc = train(
            model,
            optim,
            train_loader,
            val_loader,
            epochs=5 # number of epochs changed from 20
        )

        if max(val_acc) > best_val_accuracy: # modified
            best_val_accuracy = max(val_acc) # modified
            best_lr = lr
            best_M = M # modified

        all_params.append((lr, M, max(val_acc)))

    # if validation accuracy is at least 50%, stop
    # if best_val_accuracy > 0.5: # modified
    #     break

# print top three parameter combinations
all_params.sort(key=lambda x: x[2], reverse=True) # sort by validation accuracy
for i, (lr, M, val_acc) in enumerate(all_params[:3], 1): # only look at top three
    results
    print(f"learning rate: {lr}")
    print(f"M: {M}")
    print(f"Validation accuracy: {val_acc:.4f}")

print(f"best validation accuracy= {best_val_accuracy} for lr={best_lr}, M={best_M}")
return best_lr, best_M

```

```

"""train and evaluate model"""

best_lr, best_M = parameter_search(train_loader, val_loader, A_model)

model = A_model(best_M)
optimizer = SGD(model.parameters(), best_lr)

# We are using 20 epochs for this example. You may have to use more.
train_loss, train_accuracy, val_loss, val_accuracy = train(
    model, optimizer, train_loader, val_loader, 40)

#second best model
model2 = A_model(200)
optimizer2 = SGD(model2.parameters(), 0.03593813627958298)
train_loss2, train_accuracy2, val_loss2, val_accuracy2 = train(
    model2, optimizer2, train_loader, val_loader, 40)

#third best model
model3 = A_model(1500)
optimizer3= SGD(model3.parameters(), 0.03593813627958298)
train_loss3, train_accuracy3, val_loss3, val_accuracy3 = train(
    model3, optimizer3, train_loader, val_loader, 40)

"""plot training and validation for each epoch:"""

epochs = range(1, 41)

#best model
plt.plot(epochs, train_accuracy, label="lr=0.0359, M=1000 Train Acc")
plt.plot(epochs, val_accuracy, label="lr=0.0359, M=1000 Val Acc", ls=":")

#second best model
plt.plot(epochs, train_accuracy2, label="lr=0.0359, M=200 Train Acc")
plt.plot(epochs, val_accuracy2, label="lr=0.0359, M=200 Val Acc", ls=":")

#third best model
plt.plot(epochs, train_accuracy3, label="lr=0.0359, M=1500 Train Acc")
plt.plot(epochs, val_accuracy3, label="lr=0.0359, M=1500 Val Acc", ls=":")

plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.axhline(y=0.5, c='black')
plt.legend()
plt.title("Part A model Accuracy vs Epoch")
plt.show()

"""evaluate (do not need to modify):"""

def evaluate(
    model: nn.Module, loader: DataLoader
) -> Tuple[float, float]:
    """Computes test loss and accuracy of model on loader."""
    loss = nn.CrossEntropyLoss()
    model.eval()
    test_loss = 0.0

```

```

test_acc = 0.0
with torch.no_grad():
    for (batch, labels) in loader:
        batch, labels = batch.to(DEVICE), labels.to(DEVICE)
        y_batch_pred = model(batch)
        batch_loss = loss(y_batch_pred, labels)
        test_loss = test_loss + batch_loss.item()

        pred_max = torch.argmax(y_batch_pred, 1)
        batch_acc = torch.sum(pred_max == labels)
        test_acc = test_acc + batch_acc.item()
    test_loss = test_loss / len(loader)
    test_acc = test_acc / (batch_size * len(loader))
    return test_loss, test_acc

test_loss, test_acc = evaluate(model, test_loader)
print(f"Test Accuracy for M=1000, lr=0.03593813627958298: {test_acc}")

#test accuracy of all three top models
test_loss, test_acc = evaluate(model, test_loader)
test_loss2, test_acc2 = evaluate(model2, test_loader)
test_loss3, test_acc3 = evaluate(model3, test_loader)
print(f"Test Accuracy for M=1000, lr=0.03593813627958298: {test_acc}")
print(f"Test Accuracy for M=200, lr=0.03593813627958298: {test_acc2}")
print(f"Test Accuracy for M=1500, lr=0.03593813627958298: {test_acc3}")

"""# Part b:"""

#part b model
def B_model(M, k , N) -> nn.Module: #M is hyperparameter
    """Convolutional layer with max-pool and fully-connected output"""
    model = nn.Sequential(
        nn.Conv2d(3, M, k),
        nn.ReLU(),
        nn.MaxPool2d(N),
        nn.Flatten(),
        nn.Linear(M * ((33-k)//N) * ((33-k)//N), 10)
    )
    return model.to(DEVICE)

def parameter_search_B(train_loader: DataLoader,
                        val_loader: DataLoader,
                        model_fn: Callable[[int], nn.Module]) -> Tuple[float, int, int, int]:
    """
    Parameter search for our linear model.

    Args:
    train_loader: the train dataloader.
    val_loader: the validation dataloader.
    model_fn: a function that, when called, returns a torch.nn.Module.

    Returns:
    The learning rate with the least validation loss.
    NOTE: you may need to modify this function to search over and return
    other parameters beyond learning rate.

```

```

"""
num_iter = 10
best_loss = torch.tensor(np.inf)
best_lr = 0.0
best_M = 0 # modified
best_k = 0 # modified
best_N = 0 # modified
best_val_accuracy = 0.0 # modified

Ms = [200, 500, 1500] # modified
#lrs = torch.linspace(10 ** (-6), 10 ** (-1), num_iter)
lrs = torch.logspace(-4, -1, num_iter)
ks = [3,4,5,6,7]
Ns = [2,4,8,16,26]

# keep track of all params tried to return the top three
all_params = []
for M in Ms:
    for lr in lrs:
        for k in ks:
            for N in Ns:
                print(f"trying learning rate {lr}, M={M}, k={k}, N={N}") # modified
                model = model_fn(M, k, N) # modified

                optim = SGD(model.parameters(), lr)

                train_loss, train_acc, val_loss, val_acc = train(
                    model,
                    optim,
                    train_loader,
                    val_loader,
                    epochs=3 # number of epochs changed from 20
                )

                if max(val_acc) > best_val_accuracy: # modified
                    best_val_accuracy = max(val_acc) # modified
                    best_lr = lr
                    best_M = M # modified
                    best_k = k
                    best_N = N

                all_params.append((lr, M, k, N, max(val_acc)))

                # if validation accuracy is at least 65%, stop
                # if best_val_accuracy > 0.65: # modified
                #     break

# print top three parameter combinations
all_params.sort(key=lambda x: x[4], reverse=True) # sort by validation accuracy
for i, (lr, M, k, N, val_acc) in enumerate(all_params[:3], 1): # only look at top three
    results
    print(f"learning rate: {lr}")
    print(f"M: {M}")
    print(f"k: {k}")
    print(f"N: {N}")

```

```

        print(f"Validation accuracy: {val_acc:.4f}")
        print("-----")

        print(f"best validation accuracy= {best_val_accuracy} for lr={best_lr}, M={best_M},
k={best_k}, N={best_N}")
        return best_lr, best_M, best_k, best_N

best_lr, best_M, best_k, best_N = parameter_search_B(train_loader, val_loader, B_model)

model_B = B_model(best_M, best_k, best_N)
optimizer_B = SGD(model_B.parameters(), best_lr)

train_loss_B, train_accuracy_B, val_loss_B, val_accuracy_B = train(
    model_B, optimizer_B, train_loader, val_loader, 30) #change number of epochs?

#second best model
model2B = B_model(1500,4,4)
optimizer2B = SGD(model2B.parameters(), 0.04641588777303696)
train_loss2B, train_accuracy2B, val_loss2B, val_accuracy2B = train(
    model2B, optimizer2B, train_loader, val_loader, 30)

#third best model
model3B = B_model(1500,5,2)
optimizer3B = SGD(model3B.parameters(), 0.10000000149011612)
train_loss3B, train_accuracy3B, val_loss3B, val_accuracy3B = train(
    model3B, optimizer3B, train_loader, val_loader, 30)

epochs = range(1, 31)
plt.plot(epochs, train_accuracy_B, label="Train Accuracy")
plt.plot(epochs, val_accuracy_B, label="Validation Accuracy", ls=":")
plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.axhline(y=0.65, c='black')
plt.legend()
plt.title("Part B model Accuracy vs Epoch")
plt.show()

#plot of all three best models for part b
epochs = range(1, 31)

#best model
plt.plot(epochs, train_accuracy_B, label="lr=0.0215, M=1500, k=4, N=4 Train Acc")
plt.plot(epochs, val_accuracy_B, label="lr=0.0215, M=1500, k=4, N=4 Val Acc", ls=":")

#second best model
plt.plot(epochs, train_accuracy2B, label="lr=0.0464, M=1500, k=4, N=4 Train Acc")
plt.plot(epochs, val_accuracy2B, label="lr=0.0464, M=1500, k=4, N=4 Val Acc", ls=":")

#third best model
plt.plot(epochs, train_accuracy3B, label="lr=0.1000, M=1500, k=5, N=2 Train Acc")
plt.plot(epochs, val_accuracy3B, label="lr=0.1000, M=1500, k=5, N=2 Val Acc", ls=":")

plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.axhline(y=0.65, c='black')

```



```

plt.legend()
plt.title("Part B model Accuracy vs Epoch")
plt.show()

#Part B test accuracy of all three top models
test_lossB, test_accB = evaluate(model_B, test_loader)
test_loss2B, test_acc2B = evaluate(model2B, test_loader)
test_loss3B, test_acc3B = evaluate(model3B, test_loader)
print(f"Test Accuracy for lr=0.0215, M=1500, k=4, N=4: {test_accB}")
print(f"Test Accuracy for lr=0.0464, M=1500, k=4, N=4: {test_acc2B}")
print(f"Test Accuracy for lr=0.1000, M=1500, k=5, N=2: {test_acc3B}")

```

Matrix Completion and Recommendation System

A4.

Note: Please refer to the provided notebook with starter code for this problem, on the course website. The notebook provides a template for each part of the problem, and includes code to help load the data properly. We recommend creating a copy of the starter notebook and completing the assignment by filling out the template.

Hint: For loops are costly. Can you vectorize it or use Numpy operations to make it faster in some ways?

You will build a personalized movie recommendation system. We will use the 100K MovieLens dataset available at <https://grouplens.org/datasets/movielens/100k/>. There are $m = 1682$ movies and $n = 943$ users. Each user rated at least 20 movies, but some watched many more. The total dataset contains 100,000 total ratings from all users. The goal is to recommend movies the users haven't seen.

Consider a matrix $R \in \mathbb{R}^{m \times n}$ where the entry $R_{i,j} \in \{1, \dots, 5\}$ represents the j th user's rating on movie i . A higher value represents that the user is more satisfied with the movie.

We may think of our historical data as some observed entries of this matrix while many remain unknown, and we wish to estimate the unknown ratings that each user would assign to each movie. We could use these ratings to recommend the “best” movies for each user.

The dataset contains user and movie metadata which we will ignore. We strictly use the ratings contained in the `u.data` file. Use this data file and the following python code to construct a training and test set:

```
import csv
import numpy as np
data = []
with open('u.data') as csvfile:
    spamreader = csv.reader(csvfile, delimiter='\t')
    for row in spamreader:
        data.append([int(row[0])-1, int(row[1])-1, int(row[2])])
data = np.array(data)

num_observations = len(data)    # num_observations = 100,000
num_users = max(data[:,0])+1   # num_users = 943, indexed 0,...,942
num_items = max(data[:,1])+1   # num_items = 1682 indexed 0,...,1681

np.random.seed(1)
num_train = int(0.8*num_observations)
perm = np.random.permutation(data.shape[0])
train = data[perm[0:num_train],:]
test = data[perm[num_train:],:]
```

The arrays `train` and `test` contain R_{train} and R_{test} , respectively. Each line takes the form “j, i, s”, where j is the user index, i is the movie index, and s is the user's score.

Using `train`, you will train a model that can predict $\hat{R} \in \mathbb{R}^{m \times n}$, how every user would rate every movie. You will evaluate your model based on the average squared error on `test`:

$$\mathcal{E}_{\text{test}}(\hat{R}) = \frac{1}{|\text{test}|} \sum_{(i,j,R_{i,j}) \in \text{test}} (\hat{R}_{i,j} - R_{i,j})^2.$$

Low-rank matrix factorization is a baseline method for personalized recommendation. It learns a vector representation $u_i \in \mathbb{R}^d$ for each movie and a vector representation $v_j \in \mathbb{R}^d$ for each user, such that the inner product

$\langle u_i, v_j \rangle$ approximates the rating $R_{i,j}$. You will build a simple latent factor model.

You will implement multiple estimators and use the inner product $\langle u_i, v_j \rangle$ to predict if user j likes movie i in the test data. For simplicity, we will put aside best practices and choose hyperparameters by using those that minimize the test error. You may use fundamental operators from `numpy` or `pytorch` in this problem (`numpy.linalg.lstsq`, `SVD`, `autograd`, etc.) but not any precooked algorithm from a package like `scikit-learn`. If there is a question whether some package is not allowed for use in this problem, it probably is not appropriate.

- a. [5 points] Our first estimator pools all users together and, for each movie, outputs as its prediction the average user rating of that movie in `train`. That is, if $\mu \in \mathbb{R}^m$ is a vector where μ_i is the average rating of the users that rated the i th movie, write this estimator \hat{R} as a rank-one matrix. Compute the estimate \hat{R} . What is $\mathcal{E}_{\text{test}}(\hat{R})$ for this estimate?

Since we are just using the average rating for all users as our predictor, the estimator \hat{R} would be the matrix with all columns equal to μ , or $\hat{R}_{ij} = \mu_i$, for all integers $i \in [1, m]$ and $j \in [1, n]$. We can also write this in the form, $\hat{R} = \mu \mathbf{1}^T$, where $\mathbf{1}$ is the vector of all 1s. Writing it in this form shows that \hat{R} is rank 1.

From the calculation in the code, $\mathcal{E}_{\text{test}}(\hat{R}) \approx 1.064$.

- b. [5 points] Allocate a matrix $\tilde{R}_{i,j} \in \mathbb{R}^{m \times n}$ and set its entries equal to the known values in the training set, and 0 otherwise. Let $\hat{R}^{(d)}$ be the best rank- d approximation (in terms of squared error) approximation to \tilde{R} . This is equivalent to computing the singular value decomposition (SVD) and using the top d singular values. This learns a lower-dimensional vector representation for users and movies, assuming that each user would give a rating of 0 to any movie they have not reviewed.
- For each $d = 1, 2, 5, 10, 20, 50$, compute the estimator $\hat{R}^{(d)}$. We recommend using an efficient solver such as `scipy.sparse.linalg.svds`.
 - Plot the average squared error of predictions on the training set and test set on a single plot, as a function of d .

Note that, in most applications, we would not actually allocate a full $m \times n$ matrix. We do so here only because our data is relatively small and it is instructive.

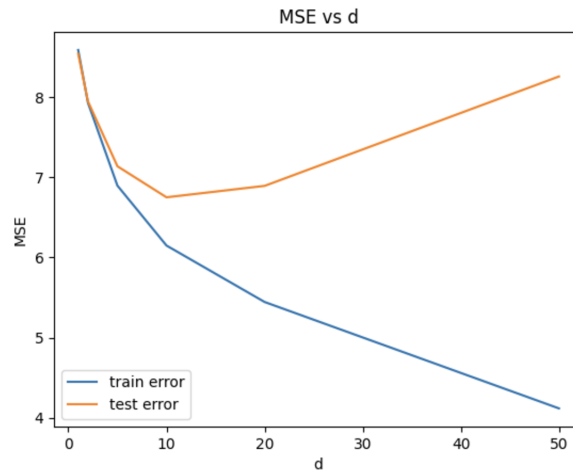


Figure 1: Plot of MSE vs d for rank- d approximations of the ranking matrix using SVD (part b). The MSE is really high but hopefully based on the Ed discussion this is expected because we are setting all missing values to 0.

- c. [10 points] Replacing all missing values by a constant may impose strong and potentially incorrect assumptions on the unobserved entries of R . A more reasonable choice is to minimize the MSE (mean squared error) only on rated movies. Define a loss function:

$$\mathcal{L}(\{u_i\}_{i=1}^m, \{v_j\}_{j=1}^n) := \sum_{(i,j,R_{i,j}) \in \text{train}} (\langle u_i, v_j \rangle - R_{i,j})^2 + \lambda \sum_{i=1}^m \|u_i\|_2^2 + \lambda \sum_{j=1}^n \|v_j\|_2^2 \quad (1)$$

where $\lambda > 0$ is the regularization coefficient. We will implement algorithms to learn vector representations by minimizing (1). Note: we define the loss function here as the sum of squared errors; be careful to calculate and plot the mean squared error for your results.

Since this is a non-convex optimization problem, the initial starting point and hyperparameters may affect the quality of \hat{R} . You may need to tune λ and σ to optimize the loss you see.

- *Alternating minimization:* First, randomly initialize $\{u_i\}$ and $\{v_j\}$. Then, alternate between (1) minimizing the loss function with respect to $\{u_i\}$ by treating $\{v_j\}$ as fixed; and (2) minimizing the loss function with respect to $\{v_j\}$ by treating $\{u_i\}$ as fixed. Repeat (1) and (2) until both $\{u_i\}$ and $\{v_j\}$ converge. Note that when one of $\{u_i\}$ or $\{v_j\}$ is given, minimizing the loss function with respect to the other part has a closed-form solution. Indeed, it can be shown that when minimizing with respect to a *single* u_i (with $\{v_j\}$ fixed), the gradient is given by:

$$\nabla_{u_i} L(\{u_i\}_{i=1}^m, \{v_j\}_{j=1}^n) = 2 \left(\sum_{j \in r(i)} v_j v_j^T + \lambda I \right) u_i - 2 \sum_{j \in r(i)} R_{i,j} v_j \quad (2)$$

where here $r(i)$ is a shorthand for the set of users who have reviewed movie i in the training set, or more formally, $r(i) = \{j : (j, i, R_{i,j}) \in \text{train}\}$. Setting the overall gradient to be equal to 0 gives us that

$$\arg \min_{u_i} L(\{u_i\}_{i=1}^m, \{v_j\}_{j=1}^n) = \left(\sum_{j \in r(i)} v_j v_j^T + \lambda I \right)^{-1} \left(\sum_{j \in r(i)} R_{i,j} v_j \right) \quad (3)$$

Note that this update rule is for a single vector u_i , whereas you should update all of the $\{u_i\}_{i=1}^m$ in one round. When it comes to the alternate step which involves fixing $\{u_i\}$ and minimizing $\{v_j\}$, an analogous calculation will give you a very similar update rule.

- Try $d \in \{1, 2, 5, 10, 20, 50\}$ and plot the mean squared error of train and test as a function of d .

Some hints:

- Common choices for initializing the vectors $\{u_i\}_{i=1}^m, \{v_j\}_{j=1}^n$ include: entries drawn from `np.random.rand()` scaled by some scale factor $\sigma > 0$ (σ is an additional hyperparameter), or using one of the solutions from part b or c.
- The only $m \times n$ matrix you need to allocate is probably for \tilde{R} .
- It is **crucial** that the squared error part of the loss is only defined w.r.t. $R_{i,j}$ that actually exist in the training set. Consider implementing some type of data structures that allow you to keep track of $r(i)$ as well as the reverse mapping $r^{-1}(j)$ from movies to relevant users.

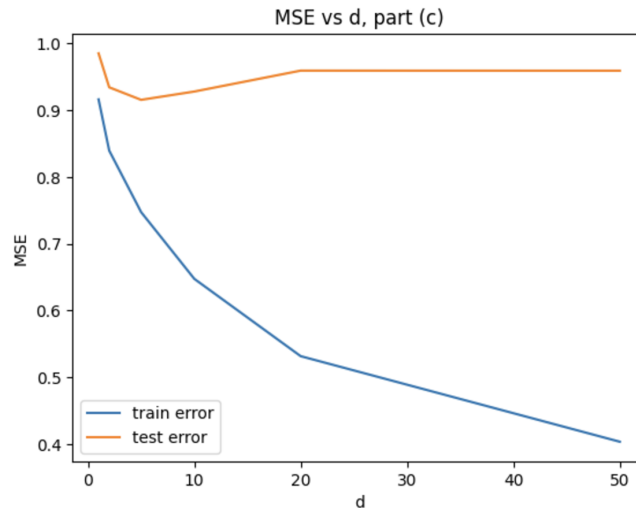


Figure 2: Plot of MSE vs d (part c).

What to Submit:

- **For part a:** A mathematical expression for \hat{R} . Value for $\mathcal{E}_{test}(\hat{R})$.
- **For part b:** Plot of MSE on training and test set vs. d .
- **For part c:** Plot of MSE on training and test set vs. d .
- **For parts a-c:** Code. You should convert your code (the .ipynb notebook) into a Python (.py) file, rename it to `hw4-a4.py`, and submit it to the corresponding Gradescope submission. To download the file from Google Colab, you can go to File ↴ Download ↴ Download as .py.

Homework 4: Matrix Completion and Recommendation System

Information before starting

In this problem, we will be building a personalized movie recommendation system! To make these recommendations, we'll build on what we've learned in lecture about SVD and what we've practiced so far with Python and Python packages such as NumPy and PyTorch.

Copying this Colab Notebook to your Google Drive

Since the course staff is the author of this notebook, you cannot make any lasting changes to it. You should make a copy of it to your Google Drive by clicking **File -> Save a Copy in Drive**.

Problem Introduction

We will use the 100K MovieLens dataset available at <https://grouplens.org/datasets/movielens/100k/> to estimate unknown user ratings given their previous ratings. Run the code block below to download the dataset.

```
"""
```

```
# @title loading dataset
!rm -rf ml-100k*
!wget https://files.grouplens.org/datasets/movielens/ml-100k.zip
!unzip ml-100k.zip
```

```
!mv ml-100k/u.data .
```

```
"""### Compute
```

This problem should not require using GPU. Since Google Colab will limit your GPU usage, we recommend saving your GPU quota for HW4 A3 and making sure that your runtime is set to CPU by going to **Runtime -> Change runtime type -> Select CPU** under "Hardware accelerator".

```
### Submitting your assignment
```

Once you are done with the problem, make sure to put all of your necessary figures into your PDF submission. Then, download this notebook as a Python file (`.py`) by going to **File -> Download -> Download `.py`**. Rename this file as `hw4-a4.py` and upload to the Gradescope submission for HW4 code.

```
## Code: Setup
```

Let's start by importing the packages that we'll need to complete this problem.

```
import csv
import numpy as np
from scipy.sparse.linalg import svds
import matplotlib.pyplot as plt
import torch
```

"""Now, let's load the 100K MovieLens data. If you have downloaded the `u.data` file and uploaded to the "Files" tab, the following code block will construct training and test sets for you. There are $m = 1682$ movies and $n = 943$ users in the dataset, and each user has rated at least 20 movies. The total dataset has 100,000 total ratings from all users, and our goal will be to estimate the unknown ratings that each user would assign to each movie. These ratings can then be used to recommend the "best" movies for each user!"""

```
data = []
with open('u.data') as csvfile:
    spamreader = csv.reader(csvfile, delimiter='\t')
    for row in spamreader:
        data.append([int(row[0])-1, int(row[1])-1, int(row[2])])
data = np.array(data)

num_observations = len(data) # num_observations = 100,000
num_users = max(data[:,0])+1 # num_users = 943, indexed 0,...,942
num_items = max(data[:,1])+1 # num_items = 1682 indexed 0,...,1681

np.random.seed(1)
num_train = int(0.8*num_observations)
perm = np.random.permutation(data.shape[0])
train = data[perm[0:num_train],:]
test = data[perm[num_train:,:],:]

print(f"Successfully loaded 100K MovieLens dataset with",
      f"{len(train)} training samples and {len(test)} test samples")
```

"""For this problem, we will consider a matrix $R \in \mathbb{R}^{m \times n}$ where the entry $R_{i,j} \in \{1, \dots, 5\}$ represents the j th user's rating on movie i . A higher value represents that the user is more satisfied with the movie.

Code: Assignment

The rest is yours to code! We provide some scaffolding for your implementation, but feel free to modify it and implement however you would like to. You may use fundamental operators from `NumPy` and `PyTorch` in this problem, such as `numpy.linalg.lstsq`, SVD, autograd, etc., but you may not use any pre-cooked algorithm from a package like `scikit-learn`.

Part (a)

Our first estimator pools all users together and, for each movie, outputs as its prediction the average user rating of that movie in ``train``. That is, if $\mu \in \mathbb{R}^m$ is a vector where μ_i is the average rating of the users that rated the i -th movie. Write this estimator \widehat{R} as a rank-one matrix.

Compute the estimate \widehat{R} . What is $\mathcal{E}_{\text{test}}(\widehat{R})$ for this estimate?

"""

Your code goes here. You should:

`mu = np.zeros(num_items) #will store movie rating averages`

`#mu = np.bincount(train[:, 1], train[:, 2]) / np.bincount(train[:, 1]) #sum of movie ratings divided by total occurrences of said movie`

1. Compute estimate and

`for i in range(num_items):`

`if len(train[train[:, 1] == i, 2]) == 0: #use train data? train[:, 1]=i is movie i, train[:,2] is where the ratings live`

`mu[i] = 0 #if no ratings, set to 0`

`else:`

`mu[i] = np.mean(train[train[:, 1] == i, 2]) #if there are ratings, take mean of all ratings for that movie`

2. Evaluate test error

`predictions = mu[test[:, 1]] #vector of ratings of movies in test set`

`print(f"test error {np.mean((test[:, 2]-predictions)**2)}") #mean squared difference in ratings vs predicted`

"""### Part (b)

Allocate a matrix $\widetilde{R}_{i,j} \in \mathbb{R}^{m \times n}$ and set its entries equal to the known values in the training set, and 0 otherwise.

Let $\widehat{R}^{(d)}$ be the best rank- d approximation (in terms of squared error) approximation to \widetilde{R} . This is equivalent to computing the singular value decomposition (SVD) and using the top d singular values. This learns a lower-dimensional vector representation for users and movies, assuming that each user would give a rating of 0 to any movie they have not reviewed.

- For each $d = 1, 2, 5, 10, 20, 50$, compute the estimator $\widehat{R}^{(d)}$. We recommend using an efficient solver, such as ```scipy.sparse.linalg.svds```.

```

- Plot the average squared error of predictions on the training set and test set on a single
plot, as a function of  $d$ .
"""

# Your code goes here
# Create the matrix  $R$   $\tilde{R}$ .
r_twiddle = np.zeros((num_items, num_users)) #default to 0
r_twiddle[train[:, 1], train[:, 0]] = train[:, 2] #update corresponding entry if there is a
rating: (movie,user)=rating

# Your code goes here
def construct_estimator(d, r_twiddle):
    #take SVD of  $\tilde{R}$ 
    u, s, vh = svds(r_twiddle, d)

    #why are the singular values given in ascending order???
    #print(s)

    #reverse all three to get them in the correct order (turns out this wasn't my problem :/)
    u = u[:, ::-1]
    s = s[::-1]
    vh = vh[:, ::-1, :]

    return u @ np.diag(s) @ vh
    #raise NotImplementedError("Your code goes here")

def get_error(d, r_twiddle, dataset):
    r_hat = construct_estimator(d, r_twiddle)

    predictions = r_hat[dataset[:, 1], dataset[:, 0]] #(movie, user)
    return np.mean((dataset[:, 2]-predictions)**2) #MSE
    #raise NotImplementedError("Your code goes here")

# Your code goes here
# Evaluate train and test error for:  $d = 1, 2, 5, 10, 20, 50$ .
d= [1,2,5,10,20,50]

#Plot the average squared error of predictions on the training set and test set on a single
plot, as a function of  $d$ .
train_error = []
test_error = []

for i in d:
    train_error.append(get_error(i, r_twiddle, train))
    test_error.append(get_error(i, r_twiddle, test))

# Your code goes here
# Plot both train and test error as a function of  $d$  on the same plot.
plt.plot(d, train_error, label="train error")
plt.plot(d, test_error, label="test error")
plt.xlabel("d")
plt.ylabel("MSE")
plt.title("MSE vs d")
plt.legend()

```



```
plt.show()
```

```
"""### Part (c)
```

Replacing all missing values by a constant may impose strong and potentially incorrect assumptions on the unobserved entries of R . A more reasonable choice is to minimize the mean squared error (MSE) only on rated movies. Define a loss function:

```
$$
```

$$\mathcal{L} \left(\{u_i\}_{i=1}^m, \{v_j\}_{j=1}^n \right) := \sum_{(i, j, R_{i,j}) \in \{\text{rm train}\}} (\langle u_i, v_j \rangle - R_{i,j})^2 + \lambda \sum_{i=1}^m \|u_i\|_2^2 + \lambda \sum_{j=1}^n \|v_j\|_2^2$$

```
$$
```

where $\lambda > 0$ is the regularization coefficient. We will implement algorithms to learn vector representations by minimizing the above loss. You may need to tune λ and σ to optimize the loss.

Implement alternating minimization (as defined in the homework spec) and plot the MSE of ``train`` and ``test`` for $d \in \{1, 2, 5, 10, 20, 50\}$.

Note: we define the loss function here as the sum of squared errors; be careful to calculate and plot the mean squared error for your results

```
"""
```

```
# Your code goes here. You are welcome to change the parameter lists and/or write new
# functions to complete this part of the assignment.
# In particular, you will likely also want to use R twiddle, and you may want to create
# global data structures to store observed entries.
# These global data structures might look like mappings of users to the movies they've
# reviewed, and of movies to the users who have reviewed that movie.
```

```
r_twiddle = np.zeros((num_items, num_users)) #default to 0
r_twiddle[train[:, 1], train[:, 0]] = train[:, 2] #update corresponding entry if there is a
rating: (movie,user)=rating
```

```
#Consider implementing some type of data structures that allow you to keep track of  $r(i)$ 
as well as the reverse mapping  $r^{-1}(j)$  from movies to relevant users.
```

```
#r(i) (i=movie, user)
r_i = {i: train[train[:, 1] == i, 0] for i in range(num_items)}
#r^{-1}(j) (movie, j=user)
r_inv_j = {j: train[train[:, 0] == j, 1] for j in range(num_users)}
```

```
def closed_form_u(d, V, U, l):
```

```
    #fix V, update U
    U_updated = np.zeros_like(U)
```

```
    for j in range(num_users): # for each user j
        movies = r_inv_j[j] #movies rated by user j

        if len(movies) == 0: #if rated no movies we move on
            continue
        else:
            v_j = V[movies] #V vector of rated movies by user j
            ratings = r_twiddle[movies, j] #ratings for those movies
```

```
    #equation (3)
```

```

        A = v_j.T @ v_j + 1 * np.eye(d)
        b = v_j.T @ ratings
        U_updated[j] = np.linalg.solve(A, b)

    return U_updated
    #raise NotImplementedError("Your code goes here")

def closed_form_v(d, V, U, l):
    V_updated = np.zeros_like(V)

    for i in range(num_items):
        users = r_i[i]
        if len(users) == 0:
            continue
        else:
            u_i = U[users]
            ratings = r_twiddle[i, users]

            #equation (3)
            A = u_i.T @ u_i + 1 * np.eye(d)
            b = u_i.T @ ratings
            V_updated[i] = np.linalg.solve(A, b)

    return V_updated
    #raise NotImplementedError("Your code goes here")

def construct_alternating_estimator(d, r_twiddle, l=10.0, delta=1e-1, sigma=0.1, U=None,
V=None):
    #randomly initialize U and V
    if U is None:
        U = sigma * np.random.randn(num_users, d)
    if V is None:
        V = sigma * np.random.randn(num_items, d)

    prev_error = float('inf')

    #alternating
    while True:
        #update U
        U_updated = closed_form_u(d, V, U, l)
        #then V
        V_updated = closed_form_v(d, V, U_updated, l)

        #run until U and V converge (change in error is less than delta)
        #equation (1)
        error = np.sum((r_twiddle - V_updated @ U_updated.T)**2) + 1 *
np.linalg.norm(U_updated)**2 + 1 * np.linalg.norm(V_updated)**2
        if abs(error - prev_error) < delta:
            break

    U, V = U_updated, V_updated

    prev_error = error

```

```

    return V @ U.T
    #raise NotImplementedError("Your code goes here")

# Your code goes here
# Any additional functions that you may write to help implement alternating minimization.

#same as above except with construct_alternating_estimator
def get_error(d, r_twiddle, dataset):
    r_hat = construct_alternating_estimator(d, r_twiddle)

    predictions = r_hat[dataset[:, 1], dataset[:, 0]] #(movie, user)
    return np.mean((dataset[:, 2]-predictions)**2) #MSE
    #raise NotImplementedError("Your code goes here")

# Your code goes here
# Evaluate train and test error for: d = 1, 2, 5, 10, 20, 50.
d= [1,2,5,10,20,50]
#d=[1,2]

#Plot the average squared error of predictions on the training set and test set on a single
plot, as a function of .
train_error = []
test_error = []

for i in d:
    print(f"d={i}")
    train_error.append(get_error(i, r_twiddle, train))
    print(f"train error: {get_error(i, r_twiddle, train)}")
    test_error.append(get_error(i, r_twiddle, test))
    print(f"test error: {get_error(i, r_twiddle, test)}")

# Your code goes here
# Plot both train and test error as a function of d on the same plot.
plt.plot(d, train_error, label="train error")
plt.plot(d, test_error, label="test error")
plt.xlabel("d")
plt.ylabel("MSE")
plt.title("MSE vs d, part (c)")
plt.legend()
plt.show()

```

k -means clustering

A5. Given a dataset $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and an integer $1 \leq k \leq n$, recall the following k -means objective function

$$\min_{\pi_1, \dots, \pi_k} \sum_{i=1}^k \sum_{j \in \pi_i} \|\mathbf{x}_j - \mu_i\|_2^2, \quad \mu_i = \frac{1}{|\pi_i|} \sum_{j \in \pi_i} \mathbf{x}_j. \quad (4)$$

Above, $\{\pi_i\}_{i=1}^k$ is a partition of $\{1, 2, \dots, n\}$. The objective (4) is NP-hard¹ to find a global minimizer of. Nevertheless, Lloyd's algorithm (discussed in lecture) typically works well in practice.²

- [5 points]** Implement Lloyd's algorithm for solving the k -means objective (4). Do not use any off-the-shelf implementations, such as those found in `scikit-learn`.
- [5 points]** Run Lloyd's algorithm on the *training* dataset of MNIST with $k = 10$. Show the image representing the center of each cluster, as a set of k 28×28 images.

Note on Time to Run — The runtime of a good implementation for this problem should be fairly fast (a few minutes); if you find it taking upwards of one hour, please check your implementation! (Hint: **For loops are costly**. Can you vectorize it or use Numpy operations to make it faster in some ways? If not, is looping through data-points or through centers faster?)

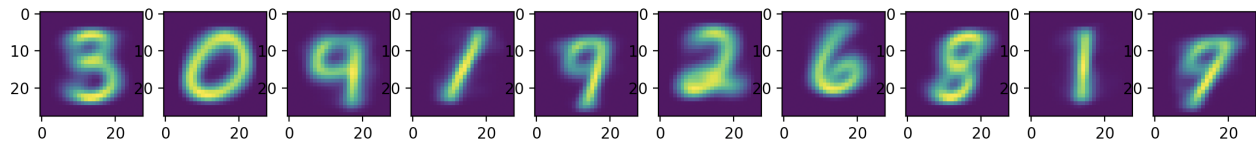


Figure 3: Images representing the ten centers given by Lloyd's algorithm on the MNIST dataset.

why are there so many 9s :/

What to Submit:

- **For part (a):** Nothing required in PDF submission.
- **For part (b):** 10 images of cluster centers.
- **For parts (a)-(b):** Code through corresponding Gradescope coding submission.

¹To be more precise, it is both NP-hard in d when $k = 2$ and k when $d = 2$.

²See the references on the Wikipedia page for k -means and k -means++ for more details.

Random Fourier Features

B1. Kernel methods such as Logistic Regression are considered memory-based learners. Rather than learning a mapping from a set of input features $\mathcal{X} \subset \mathbb{R}^d$ to outputs in \mathcal{Y} , they *remember* all training examples (\mathbf{x}_i, y_i) and learn a corresponding weight for them.

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \omega_i k(\mathbf{x}_i, \mathbf{x})$$

After learning the weight vector $\mathbf{w} = [\mathbf{w}_1, \dots, \mathbf{w}_N]$, we can make prediction on unseen samples using the *kernel function* k between all training samples and \mathbf{x} . Kernel methods are attractive because they rely on the *kernel trick*. Any positive definite function $k(\mathbf{x}, \mathbf{x}')$ with $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ defines a function ψ mapping \mathbb{R}^d to a higher-dimensional space such that the inner product between datapoints can be quickly computed as $\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$. In essence, the kernel trick is an efficient way to learn a linear decision boundary in a higher dimension space than that of \mathcal{X} .

The kernel trick can be prohibitively expensive for large datasets. This is because the memory-based algorithm accesses the data through evaluations of the kernel matrix $k(x, x')$ which grows in proportion to the dataset size N .

Instead of relying on the implicit feature mapping ψ provided by the kernel trick, suppose we can approximate the kernel function k as the inner product of two vectors in \mathbb{R}^D . Mathematically, we would like to find a mapping \mathbf{z} :

$$\mathbf{z} : \mathbb{R}^d \rightarrow \mathbb{R}^D \quad \text{such that} \quad k_p(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle \approx \langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle$$

With this approximation, we no longer require the *kernel trick* to express $\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle$ as $k(\mathbf{x}, \mathbf{x}')$. Rather, we can approximate it by directly computing the tractable inner product $\langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle$.

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \omega_i k(\mathbf{x}_i, \mathbf{x}) = \sum_{i=1}^N \omega_i \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}) \rangle \approx \sum_{i=1}^N \omega_i \langle \mathbf{z}(\mathbf{x}_i), \mathbf{z}(\mathbf{x}) \rangle = \left(\sum_{i=1}^N \omega_i \mathbf{z}(\mathbf{x}_i)^T \right) \mathbf{z}(\mathbf{x}) = \beta^T \mathbf{z}(\mathbf{x})$$

Assuming $\mathbf{z}(\mathbf{x}) = \sigma(M\mathbf{x} + b)$ for some nonlinear function σ , this “approximate” Logistic Regression *can potentially be evaluated much quicker* than the kernel Logistic Regression. To see why, note that the left-hand-side requires evaluating $k(\mathbf{x}_i, \mathbf{x})$ for all $i \in \{1, \dots, N\}$, in general, if ω_i is not sparse. On the other hand, the right-hand-side just requires computing $\mathbf{z}(\mathbf{x}) = \sigma(M\mathbf{x} + b)$ which is dominated by the time to compute a $D \times d$ matrix-vector product, and then inner product with β which is \mathbb{R}^D . Thus, the total computation time for the left-hand-side scales linearly with N , and the right-hand-side scales with just d and D , independent of N ! When training the approximate Logistic Regression we also get similar computational savings if $N \gg \max\{d, D\}$.

- a. **[15 points] Deriving Random Fourier Features:** Bochner’s theorem states that a continuous kernel $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ on \mathbb{R}^d is positive definite if and only if k is the Fourier transform of a non-negative measure. While we won’t delve into the logic of Fourier transforms here, this theorem lets us express the kernel as follows: for any probability distribution $p(\mathbf{w})$ define

$$k_p(\mathbf{x}, \mathbf{x}') := \int_{\mathbb{R}^d} p(\mathbf{w}) e^{i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')} d\mathbf{w} = \mathbb{E}_{\mathbf{w}} \left[e^{i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')} \right]$$

where $i = \sqrt{-1}$, the imaginary unit. While any choice of $p(\mathbf{w})$ induces a valid kernel, in this problem we’ll be focusing on the Gaussian distribution, namely

$$p(\mathbf{w}) = (2\pi\sigma^2)^{-\frac{D}{2}} e^{-\frac{1}{2\sigma^2} \|\mathbf{w}\|_2^2} = (2\pi/\gamma^2)^{-\frac{D}{2}} e^{-\gamma^2 \|\mathbf{w}\|_2^2/2} \quad \text{where } \gamma = \frac{1}{\sigma}$$

In this sub-problem, we'll use this Fourier-transform interpretation of k to derive a randomized mapping $\mathbf{z} : \mathbb{R}^d \rightarrow \mathbb{R}^D$ which is an unbiased estimate of the kernel function i.e.

$$\mathbb{E}_{\mathbf{w}}[\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x}')] = k_p(\mathbf{x}, \mathbf{x}')$$

If $\mathbf{z}(x)^T \mathbf{z}(x')$ serves as a good approximation to the kernel matrix, we can apply the aforementioned approximation algorithm.

- i Use Euler's formula $e^{iy} = \cos(y) + i \sin(y)$ to show that $k_p(\mathbf{x}, \mathbf{x}') = E_{\mathbf{w}} [\cos(\mathbf{w}^T (\mathbf{x} - \mathbf{x}'))]$.

Hint: If both x and A are real, then $A = \int f(x) + ig(x)dx = \int f(x)dx$.

- ii We begin by defining $z_{\mathbf{w}} : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$z_{\mathbf{w}}(\mathbf{x}) = \sqrt{2} \cos(\mathbf{w}^T \mathbf{x} + b) \quad \text{where } \mathbf{w} \sim p(\mathbf{w}), b \sim \text{Uniform}(0, 2\pi)$$

Note that this is not yet the mapping vector \mathbf{z} , but rather a mapping to \mathbb{R} . Use part (i) to show that the expected product of $z_{\mathbf{w}}(\mathbf{x})$ s is an unbiased estimate of the kernel function i.e.

$$E_{\mathbf{w}, b} [z_{\mathbf{w}}(\mathbf{x}) z_{\mathbf{w}}(\mathbf{x}')] = k_p(\mathbf{x}, \mathbf{x}')$$

Hint: For this problem you may use the following identity: $2 \cos(\alpha) \cos(\beta) = \cos(\alpha + \beta) + \cos(\alpha - \beta)$.

- iii Now we're ready to define our random Fourier features $\mathbf{z} : \mathbb{R}^d \rightarrow \mathbb{R}^D$. Let \mathbf{z} be the D -dimensional concatenation of $z_{\mathbf{w}}(\mathbf{x})$ s:

$$\mathbf{z}(\mathbf{x}) = \left[\frac{1}{\sqrt{D}} z_{\mathbf{w}_1}(\mathbf{x}), \frac{1}{\sqrt{D}} z_{\mathbf{w}_2}(\mathbf{x}), \dots, \frac{1}{\sqrt{D}} z_{\mathbf{w}_D}(\mathbf{x}) \right]^T$$

Use parts (i) and (ii) to show that the expected inner product of the mapping \mathbf{z} is an unbiased estimate of the kernel function i.e.

$$E_{\mathbf{w}}[\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x}')] = k_p(\mathbf{x}, \mathbf{x}')$$

- b. **[5 points] Random Fourier Features and the RBF Kernel.** As mentioned in part (b), using different distributions $p(\mathbf{w})$ induces different valid kernels. Using the $p(\mathbf{w})$ given in part (a), show that expected value of the inner product between random Fourier features is the RBF kernel i.e.

$$E_{\mathbf{w}}[\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x}')] = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\gamma^2}\right)$$

Hint: The PDF for a variable $X \in \mathbb{R}^d$ following normal distribution with mean μ and covariance matrix Σ is as follows:

$$P(X = x) = ((2\pi)^d |\Sigma|)^{-1/2} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

where $|\Sigma| = \det(\Sigma)$ denote the determinant of matrix Σ . In addition, if $\Sigma = \text{diag}(\sigma^2, \dots, \sigma^2)$, then $|\Sigma| = \sigma^{2d}$, and $\Sigma^{-1} = \text{diag}(\sigma^{-2}, \dots, \sigma^{-2})$.

- c. **[5 points] Concentration Bounds** In part (a) we derived our function \mathbf{z} which serve as a good approximation to the kernel function. Our results let us get an upper bound our approximation error for the kernel function. Explain why we can apply Hoeffding's inequality to obtain

$$p(|\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x}') - k(\mathbf{x}, \mathbf{x}')| \geq \epsilon) \leq 2 \exp(-D\epsilon^2/8)$$

What to Submit:

- Part a: Separate proofs for subproblems (i), (ii) and (iii)
- Part b: proof
- Part c: proof, 1-2 sentence explanation about which conditions are met that allow us to apply Hoeffding's inequality e.g. "B is bounded by [...]"

Administrative

A6.

- a. *[2 points]* About how many hours did you spend on this homework? There is no right or wrong answer :)

At least 13 hours, not counting most of the runtime for A3 :/