

# Lab 4: Data Imputation using an Autoencoder

In this lab, you will build and train an autoencoder to impute (or "fill in") missing data.

We will be using the Adult Data Set provided by the UCI Machine Learning Repository [1], available at <https://archive.ics.uci.edu/ml/datasets/adult>. The data set contains census record files of adults, including their age, marital status, the type of work they do, and other features.

Normally, people use this data set to build a supervised classification model to classify whether a person is a high income earner. We will not use the dataset for this original intended purpose.

Instead, we will perform the task of imputing (or "filling in") missing values in the dataset. For example, we may be missing one person's marital status, and another person's age, and a third person's level of education. Our model will predict the missing features based on the information that we do have about each person.

We will use a variation of a denoising autoencoder to solve this data imputation problem. Our autoencoder will be trained using inputs that have one categorical feature artificially removed, and the goal of the autoencoder is to correctly reconstruct all features, including the one removed from the input.

In the process, you are expected to learn to:

1. Clean and process continuous and categorical data for machine learning.
2. Implement an autoencoder that takes continuous and categorical (one-hot) inputs.
3. Tune the hyperparameters of an autoencoder.
4. Use baseline models to help interpret model performance.

[1] Dua, D. and Karra Taniskidou, E. (2017). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.

## What to submit

Submit a PDF file containing all your code, outputs, and write-up. You can produce a PDF of your Google Colab file by going to File > Print and then save as PDF. The Colab instructions have more information.

Do not submit any other files produced by your code.

Include a link to your colab file in your submission.

## Colab Link

Include a link to your Colab file here. If you would like the TA to look at your Colab file in case your solutions are cut off, **please make sure that your Colab file is publicly accessible at the time of submission.**

Colab Link: <https://drive.google.com/file/d/1b4K9SDfVINWxGjm-So3toDwTYMevbNcS/view?usp=sharing>

```
In [ ]: import csv
import numpy as np
import random
```

```
import torch
import torch.utils.data
```

## Part 0

We will be using a package called `pandas` for this assignment.

If you are using Colab, `pandas` should already be available. If you are using your own computer, installation instructions for `pandas` are available here: <https://pandas.pydata.org/pandas-docs/stable/install.html>

```
In [ ]: import pandas as pd
```

## Part 1. Data Cleaning [15 pt]

The adult.data file is available at <https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data>

The function `pd.read_csv` loads the adult.data file into a pandas dataframe. You can read about the pandas documentation for `pd.read_csv` at [https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read\\_csv.html](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_csv.html)

```
In [ ]: header = ['age', 'work', 'fnlwgt', 'edu', 'yrelu', 'marriage', 'occupation',
                  'relationship', 'race', 'sex', 'capgain', 'caploss', 'workhr', 'country']
df = pd.read_csv(
    "https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data",
    names=header,
    index_col=False)
```

```
C:\Users\Admin\AppData\Local\Temp\ipykernel_23840\1831985018.py:3: ParserWarning: Length of head
er or names does not match length of data. This leads to a loss of data with index_col=False.
df = pd.read_csv(
```

```
In [ ]: df.shape # there are 32561 rows (records) in the data frame, and 14 columns (features)
```

```
Out[ ]: (32561, 14)
```

### Part (a) Continuous Features [3 pt]

For each of the columns `["age", "yrelu", "capgain", "caploss", "workhr"]`, report the minimum, maximum, and average value across the dataset.

Then, normalize each of the features `["age", "yrelu", "capgain", "caploss", "workhr"]` so that their values are always between 0 and 1. Make sure that you are actually modifying the dataframe `df`.

Like numpy arrays and torch tensors, pandas data frames can be sliced. For example, we can display the first 3 rows of the data frame (3 records) below.

```
In [ ]: df[:3] # show the first 3 records
```

Out [ ]:	age	work	fnlwgt	edu	yrelu	marriage	occupation	relationship	race	sex	capgain	caploss	work
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	

Alternatively, we can slice based on column names, for example `df["race"]` , `df["hr"]` , or even index multiple columns like below.

```
In [ ]: subdf = df[["age", "yrelu", "capgain", "caploss", "workhr"]]
subdf
```

Out [ ]:	age	yrelu	capgain	caploss	workhr
0	39	13	2174	0	40
1	50	13	0	0	13
2	38	9	0	0	40
3	53	7	0	0	40
4	28	13	0	0	40
...	...	...	...	...	...
32556	27	12	0	0	38
32557	40	9	0	0	40
32558	58	9	0	0	40
32559	22	9	0	0	20
32560	52	9	15024	0	40

32561 rows × 5 columns

Numpy works nicely with pandas, like below:

```
In [ ]: np.sum(subdf["caploss"])
```

Out [ ]: 2842700

Just like numpy arrays, you can modify entire columns of data rather than one scalar element at a time. For example, the code

```
df["age"] = df["age"] + 1
```

would increment everyone's age by 1.

```
In [ ]: columns = ["age", "yrelu", "capgain", "caploss", "workhr"]
result = df[columns].describe().loc[['min', 'max', 'mean']].transpose()
```

```
result.columns = ['Minimum', 'Maximum', 'Average']
print(result)
```

	Minimum	Maximum	Average
age	17.0	90.0	38.581647
yrelu	1.0	16.0	10.080679
capgain	0.0	99999.0	1077.648844
caploss	0.0	4356.0	87.303830
workhr	1.0	99.0	40.437456

```
In [ ]: from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df.loc[:, columns] = scaler.fit_transform(df[columns])
df[:3]
```

```
Out [ ]:
```

	age	work	fnlwgt	edu	yrelu	marriage	occupation	relationship	race	sex	capgain	caplos
0	0.301370	State-gov	77516	Bachelors	0.800000	Never-married	Adm-clerical	Not-in-family	White	Male	0.02174	0.
1	0.452055	Self-emp-not-inc	83311	Bachelors	0.800000	Married-civ-spouse	Exec-managerial	Husband	White	Male	0.00000	0.
2	0.287671	Private	215646	HS-grad	0.533333	Divorced	Handlers-cleaners	Not-in-family	White	Male	0.00000	0.

## Part (b) Categorical Features [1 pt]

What percentage of people in our data set are male? Note that the data labels all have an unfortunate space in the beginning, e.g. " Male" instead of "Male".

What percentage of people in our data set are female?

```
In [ ]: # hint: you can do something like this in pandas
sum(df["sex"] == " Male")
```

```
Out [ ]: 21790
```

```
In [ ]: num_rows = df.shape[0]
print("Male percentage: ", sum(df["sex"] == " Male")/num_rows*100)
print("Female percentage: ", sum(df["sex"] == " Female")/num_rows*100)
```

```
Male percentage: 66.92054912318419
Female percentage: 33.07945087681583
```

## Part (c) [2 pt]

Before proceeding, we will modify our data frame in a couple more ways:

1. We will restrict ourselves to using a subset of the features (to simplify our autoencoder)
2. We will remove any records (rows) already containing missing values, and store them in a second dataframe. We will only use records without missing values to train our autoencoder.

Both of these steps are done for you, below.

How many records contained missing features? What percentage of records were removed?

```
In [ ]: contcols = ["age", "yrelu", "capgain", "caploss", "workhr"]
catcols = ["work", "marriage", "occupation", "edu", "relationship", "sex"]
features = contcols + catcols
df = df[features]
```

```
In [ ]: missing = pd.concat([df[c] == " ?" for c in catcols], axis=1).any(axis=1)
df_with_missing = df[missing]
df_not_missing = df[~missing]
```

```
In [ ]: print(df_with_missing.shape[0], " records contained missing features")
print(df_with_missing.shape[0]/num_rows*100, " percent of records were removed")
```

1843 records contained missing features  
5.660145572924664 percent of records were removed

## Part (d) One-Hot Encoding [1 pt]

What are all the possible values of the feature "work" in `df_not_missing`? You may find the Python function `set` useful.

```
In [ ]: unique_work_values = set(df_not_missing["work"])
unique_work_values
```

```
Out[ ]: {' Federal-gov',
' Local-gov',
' Private',
' Self-emp-inc',
' Self-emp-not-inc',
' State-gov',
' Without-pay'}
```

We will be using a one-hot encoding to represent each of the categorical variables. Our autoencoder will be trained using these one-hot encodings.

We will use the pandas function `get_dummies` to produce one-hot encodings for all of the categorical variables in `df_not_missing`.

```
In [ ]: data = pd.get_dummies(df_not_missing)
```

```
In [ ]: data[:3]
```

```
Out[ ]:
```

	age	yrelu	capgain	caploss	workhr	work_ Federal-gov	work_ Local-gov	work_ Private	work_ Self-emp-inc	work_ Self-emp-not-inc	...	edu_ Prof-school	edu_ Some-college	re
0	0.301370	0.800000	0.02174	0.0	0.397959	0	0	0	0	0	...	0	0	
1	0.452055	0.800000	0.00000	0.0	0.122449	0	0	0	0	1	...	0	0	
2	0.287671	0.533333	0.00000	0.0	0.397959	0	0	1	0	0	...	0	0	

3 rows × 57 columns

## Part (e) One-Hot Encoding [2 pt]

The dataframe `data` contains the cleaned and normalized data that we will use to train our denoising autoencoder.

How many **columns** (features) are in the dataframe `data` ?

Briefly explain where that number come from.

```
In [ ]: print("Num columns in data: ", data.shape[1])
```

```
Num columns in data: 57
```

57 columns include 5 columns of numerical data (age, yredu, capgain, caploss, workhr) and, for each of the categorical columns, all possible unique values, one-hot encoded. For example, there are 7 possible values for "work", so 7 columns will be added. We can confirm that this is true:

```
In [ ]: print("Num columns in data: ", len(contcols) + df_not_missing[catcols].nunique().sum())
```

```
Num columns in data: 57
```

## Part (f) One-Hot Conversion [3 pt]

We will convert the pandas data frame `data` into numpy, so that it can be further converted into a PyTorch tensor. However, in doing so, we lose the column label information that a panda data frame automatically stores.

Complete the function `get_categorical_value` that will return the named value of a feature given a one-hot embedding. You may find the global variables `cat_index` and `cat_values` useful. (Display them and figure out what they are first.)

We will need this function in the next part of the lab to interpret our autoencoder outputs. So, the input to our function `get_categorical_values` might not actually be "one-hot" -- the input may instead contain real-valued predictions from our neural network.

```
In [ ]: datanp = data.values.astype(np.float32)
```

```
In [ ]: cat_index = {} # Mapping of feature -> start index of feature in a record
cat_values = {} # Mapping of feature -> list of categorical values the feature can take
```

```
# build up the cat_index and cat_values dictionary
for i, header in enumerate(data.keys()):
    if "_" in header: # categorical header
        feature, value = header.split()
        feature = feature[:-1] # remove the last char; it is always an underscore
        if feature not in cat_index:
            cat_index[feature] = i
            cat_values[feature] = [value]
        else:
            cat_values[feature].append(value)
```

```
def get_onehot(record, feature):
    """
    Return the portion of `record` that is the one-hot encoding
    of `feature`. For example, since the feature "work" is stored
    in the indices [5:12] in each record, calling `get_range(record, "work")`
    is equivalent to accessing `record[5:12]`.
```

Args:

- record: a numpy array representing one record, formatted

```

        the same way as a row in `data.np`
        - feature: a string, should be an element of `catcols`
    """
    start_index = cat_index[feature]
    stop_index = cat_index[feature] + len(cat_values[feature])
    return record[start_index:stop_index]

def get_categorical_value(onehot, feature):
    """
    Return the categorical value name of a feature given
    a one-hot vector representing the feature.

    Args:
        - onehot: a numpy array one-hot representation of the feature
        - feature: a string, should be an element of `catcols`

    Examples:

    >>> get_categorical_value(np.array([0., 0., 0., 0., 0., 1., 0.]), "work")
    'State-gov'
    >>> get_categorical_value(np.array([0.1, 0., 1.1, 0.2, 0., 1., 0.]), "work")
    'Private'
    """
    # <----- TODO: WRITE YOUR CODE HERE ----->
    # You may find the variables `cat_index` and `cat_values`
    # (created above) useful.
    return cat_values[feature][np.argmax(onehot)]

# Testing
get_categorical_value(np.array([0.1, 0., 1.1, 0.2, 0., 1., 0.]), "work")

```

Out[ ]: 'Private'

In [ ]: *# more useful code, used during training, that depends on the function  
# you write above*

```

def get_feature(record, feature):
    """
    Return the categorical feature value of a record
    """
    onehot = get_onehot(record, feature)
    return get_categorical_value(onehot, feature)

def get_features(record):
    """
    Return a dictionary of all categorical feature values of a record
    """
    return { f: get_feature(record, f) for f in catcols }

```

## Part (g) Train/Test Split [3 pt]

Randomly split the data into approximately 70% training, 15% validation and 15% test.

Report the number of items in your training, validation, and test set.

In [ ]: *# set the numpy seed for reproducibility  
# <https://docs.scipy.org/doc/numpy/reference/generated/numpy.random.seed.html>  
np.random.seed(50)*

*# todo*

```

np.random.shuffle(data)
train_index = int(data.shape[0] * 0.7)

```

```

val_index = int(datanp.shape[0] * 0.85)

train_data = datanp[:train_index]
val_data = datanp[train_index:val_index]
test_data = datanp[val_index:]

print(f"Training set shape: {train_data.shape}")
print(f"Validation set shape: {val_data.shape}")
print(f"Test set shape: {test_data.shape}")

```

Training set shape: (21502, 57)  
 Validation set shape: (4608, 57)  
 Test set shape: (4608, 57)

## Part 2. Model Setup [5 pt]

### Part (a) [4 pt]

Design a fully-connected autoencoder by modifying the `encoder` and `decoder` below.

The input to this autoencoder will be the features of the `data`, with one categorical feature recorded as "missing". The output of the autoencoder should be the reconstruction of the same features, but with the missing value filled in.

**Note:** Do not reduce the dimensionality of the input too much! The output of your embedding is expected to contain information about ~11 features.

```

In [ ]: from torch import nn

class AutoEncoder(nn.Module):
    def __init__(self):
        super(AutoEncoder, self).__init__()
        self.name = "AutoEncoder"
        self.encoder = nn.Sequential(
            nn.Linear(57, 42), # TODO -- FILL OUT THE CODE HERE!
            nn.Linear(42, 27),
            nn.Linear(27, 12)
        )
        self.decoder = nn.Sequential(
            nn.Linear(12, 27), # TODO -- FILL OUT THE CODE HERE!
            nn.Linear(27, 42),
            nn.Linear(42, 57),
            nn.Sigmoid() # get to the range (0, 1)
        )

    def forward(self, x):
        x = self.encoder(x)
        x = self.decoder(x)
        return x

```

### Part (b) [1 pt]

Explain why there is a sigmoid activation in the last step of the decoder.

(**Note:** the values inside the data frame `data` and the training code in Part 3 might be helpful.)

The values inside the input data frame `data` are all within the range [0,1] so the reconstructed output should also be in the same range. This is guaranteed by the Sigmoid function



## Part 3. Training [18]

### Part (a) [6 pt]

We will train our autoencoder in the following way:

- In each iteration, we will hide one of the categorical features using the `zero_out_random_features` function
- We will pass the data with one missing feature through the autoencoder, and obtain a reconstruction
- We will check how close the reconstruction is compared to the original data -- including the value of the missing feature

Complete the code to train the autoencoder, and plot the training and validation loss every few iterations. You may also want to plot training and validation "accuracy" every few iterations, as we will define in part (b). You may also want to checkpoint your model every few iterations or epochs.

Use `nn.MSELoss()` as your loss function. (Side note: you might recognize that this loss function is not ideal for this problem, but we will use it anyway.)

```
In [ ]: def zero_out_feature(records, feature):
        """ Set the feature missing in records, by setting the appropriate
        columns of records to 0
        """
        start_index = cat_index[feature]
        stop_index = cat_index[feature] + len(cat_values[feature])
        records[:, start_index:stop_index] = 0
        return records

def zero_out_random_feature(records):
    """ Set one random feature missing in records, by setting the
    appropriate columns of records to 0
    """
    return zero_out_feature(records, random.choice(catcols))

def get_model_name(name, learning_rate, epoch):
    """ Generate a name for the model consisting of all the hyperparameter values

    Args:
        config: Configuration object containing the hyperparameters
    Returns:
        path: A string with the hyperparameter name and value concatenated
    """
    path = "C:/Users/Admin/Downloads/lab4points/model_{0}_lr{1}_epoch{2}".format(name,
                                                                                 learning_rate,
                                                                                 epoch)
    return path

def train(model, train_loader, valid_loader, num_epochs=5, learning_rate=1e-4):
    """ Training loop. You should update this."""
    torch.manual_seed(42)
    criterion = nn.MSELoss()
    optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)

    #####
    # Set up some numpy arrays to store the training/test loss/error
    train_acc = np.zeros(num_epochs)
    train_loss = np.zeros(num_epochs)
    val_acc = np.zeros(num_epochs)
```

```
val_loss = np.zeros(num_epochs)
```

```
#####
```

```
# Train the network
```

```
# Loop over the data iterator and sample a new batch of training data
```

```
# Get the output from the network, and optimize our loss function.
```

```
for epoch in range(num_epochs):
```

```
    total_train_err = 0.0
```

```
    total_train_loss = 0.0
```

```
    i = 0
```

```
    for data in train_loader:
```

```
        datam = zero_out_random_feature(data.clone()) # zero out one categorical feature
```

```
        if torch.cuda.is_available():
```

```
            datam = datam.cuda()
```

```
            data = data.cuda()
```

```
        recon = model(datam)
```

```
        loss = criterion(recon, data)
```

```
        loss.backward()
```

```
        optimizer.step()
```

```
        optimizer.zero_grad()
```

```
        total_train_loss += loss.item()
```

```
        i += 1
```

```
train_acc[epoch] = get_accuracy(model, train_loader)
```

```
train_loss[epoch] = float(total_train_loss) / (i+1)
```

```
val_acc[epoch] = get_accuracy(model, valid_loader)
```

```
val_loss[epoch] = get_loss(model, valid_loader, criterion)
```

```
print(("Epoch {}: Train accuracy: {}, Train loss: {} | "+  
      "Validation accuracy: {}, Validation loss: {}").format(  
    epoch + 1,
```

```
    train_acc[epoch],
```

```
    train_loss[epoch],
```

```
    val_acc[epoch],
```

```
    val_loss[epoch]))
```

```
# Save the current model (checkpoint) to a file  
model_path = get_model_name(model.name, learning_rate, epoch)
```

```
torch.save(model.state_dict(), model_path)
```

```
print('Finished Training')
```

```
# Write the train/test loss/err into CSV file for plotting later
```

```
epochs = np.arange(1, num_epochs + 1)
```

```
np.savetxt("{}_train_acc.csv".format(model_path), train_acc)
```

```
np.savetxt("{}_train_loss.csv".format(model_path), train_loss)
```

```
np.savetxt("{}_val_acc.csv".format(model_path), val_acc)
```

```
np.savetxt("{}_val_loss.csv".format(model_path), val_loss)
```

```
def plot_training_curve(path):
```

```
    """ Plots the training curve for a model run, given the csv files  
    containing the train/validation accuracy/loss.
```

```
    Args:
```

```
        path: The base path of the csv files produced during training
```

```
    """
```

```
    import matplotlib.pyplot as plt
```

```
    train_acc = np.loadtxt("{}_train_acc.csv".format(path))
```

```
    val_acc = np.loadtxt("{}_val_acc.csv".format(path))
```

```
    train_loss = np.loadtxt("{}_train_loss.csv".format(path))
```

```
    val_loss = np.loadtxt("{}_val_loss.csv".format(path))
```

```
    plt.title("Train vs Validation Accuracy")
```

```
    n = len(train_acc) # number of epochs
```

```
    plt.plot(range(1,n+1), train_acc, label="Train")
```

```
    plt.plot(range(1,n+1), val_acc, label="Validation")
```

```
    plt.xlabel("Epoch")
```

```
    plt.ylabel("Accuracy")
```

```

plt.legend(loc='best')
plt.show()
plt.title("Train vs Validation Loss")
plt.plot(range(1,n+1), train_loss, label="Train")
plt.plot(range(1,n+1), val_loss, label="Validation")
plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.legend(loc='best')
plt.show()

def get_data_loader(batch_size):
    train_loader = torch.utils.data.DataLoader(train_data, batch_size=batch_size, num_workers=1)
    val_loader = torch.utils.data.DataLoader(val_data, batch_size=batch_size, num_workers=1)
    test_loader = torch.utils.data.DataLoader(test_data, batch_size=batch_size, num_workers=1)
    return train_loader, val_loader, test_loader

```

## Part (b) [3 pt]

While plotting training and validation loss is valuable, loss values are harder to compare than accuracy percentages. It would be nice to have a measure of "accuracy" in this problem.

Since we will only be imputing missing categorical values, we will define an accuracy measure. For each record and for each categorical feature, we determine whether the model can predict the categorical feature given all the other features of the record.

A function `get_accuracy` is written for you. It is up to you to figure out how to use the function. **You don't need to submit anything in this part.** To earn the marks, correctly plot the training and validation accuracy every few iterations as part of your training curve.

```

In [ ]: def get_accuracy(model, data_loader):
    """Return the "accuracy" of the autoencoder model across a data set.
    That is, for each record and for each categorical feature,
    we determine whether the model can successfully predict the value
    of the categorical feature given all the other features of the
    record. The returned "accuracy" measure is the percentage of times
    that our model is successful.

    Args:
        - model: the autoencoder model, an instance of nn.Module
        - data_loader: an instance of torch.utils.data.DataLoader

    Example (to illustrate how get_accuracy is intended to be called.
        Depending on your variable naming this code might require
        modification.)

    >>> model = AutoEncoder()
    >>> vdl = torch.utils.data.DataLoader(data_valid, batch_size=256, shuffle=True)
    >>> get_accuracy(model, vdl)
    """
    total = 0
    acc = 0
    for col in catcols:
        for item in data_loader: # minibatches
            inp = item.detach().numpy()
            if torch.cuda.is_available():
                item = item.cuda()
            out = model(zero_out_feature(item.clone(), col)).detach()
            out = out.cpu().numpy()
            for i in range(out.shape[0]): # record in minibatch
                acc += int(get_feature(out[i], col) == get_feature(inp[i], col))

```

```

        total += 1
    return acc / total

def get_loss(model, data_loader, criterion):
    tot_loss = 0
    i = 0
    for data in data_loader:
        datam = zero_out_random_feature(data.clone())
        if torch.cuda.is_available():
            datam = datam.cuda()
            data = data.cuda()
        recon = model(datam)
        loss = criterion(recon, data)
        tot_loss += loss.item()
        i += 1

    return tot_loss/(i + 1)

```

## Part (c) [4 pt]

Run your updated training code, using reasonable initial hyperparameters.

Include your training curve in your submission.

```

In [ ]: model = AutoEncoder()
        if torch.cuda.is_available():
            model = model.cuda()
        train_loader, val_loader, test_loader = get_data_loader(64)
        train(model, train_loader, val_loader, learning_rate=0.001, num_epochs=5)

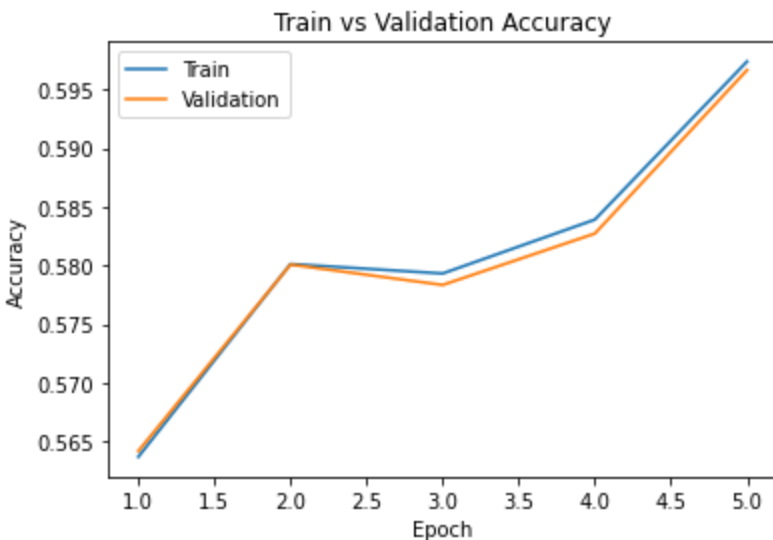
```

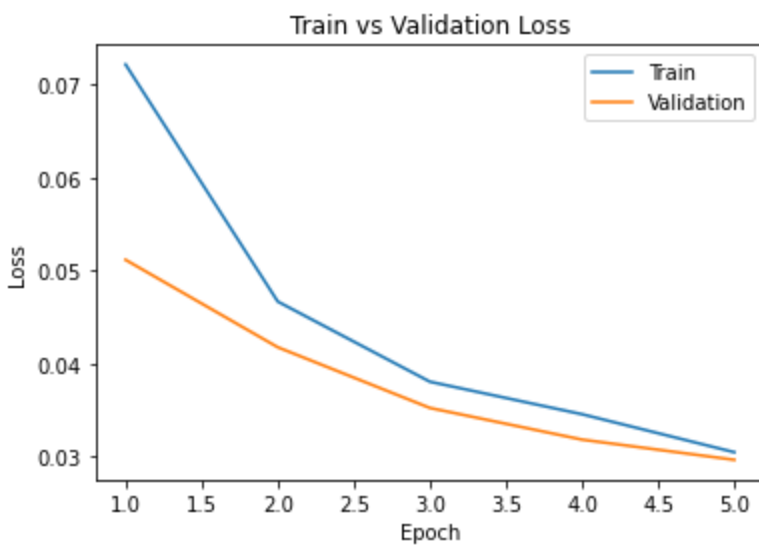
Epoch 1: Train accuracy: 0.5637072520385701, Train loss: 0.0721849154139308 | Validation accuracy: 0.5641999421296297, Validation loss: 0.05116876076957951  
 Epoch 2: Train accuracy: 0.5801243295197346, Train loss: 0.046680705030434914 | Validation accuracy: 0.580078125, Validation loss: 0.041764182587192486  
 Epoch 3: Train accuracy: 0.5793182029578644, Train loss: 0.03805337548366463 | Validation accuracy: 0.5783420138888888, Validation loss: 0.03523205692739519  
 Epoch 4: Train accuracy: 0.5839069233869718, Train loss: 0.03455432992461704 | Validation accuracy: 0.5827184606481481, Validation loss: 0.03182434032939068  
 Epoch 5: Train accuracy: 0.5973707872135925, Train loss: 0.030470345893360742 | Validation accuracy: 0.5966435185185185, Validation loss: 0.02965119993952039  
 Finished Training

```

In [ ]: model_path = get_model_name(model.name, learning_rate=0.001, epoch=4)
        plot_training_curve(model_path)

```





## Part (d) [5 pt]

Tune your hyperparameters, training at least 4 different models (4 sets of hyperparameters).

Do not include all your training curves. Instead, explain what hyperparameters you tried, what their effect was, and what your thought process was as you chose the next set of hyperparameters to try.

Epoch size was pretty small, so I increased to 15 and see how model behaves Learning rate: 0.001, Epochs: 15

Training accuracy: 0.5874104734443307, Validation accuracy: 0.5829716435185185

Accuracy did not increase too much, evened out at around 0.59, but fluctuated too much.

Therefore, try decrease learning rate for stable results

Learning rate: 0.0005, Epochs: 10

Training accuracy: 0.6032229560040926, Validation accuracy: 0.5981264467592593

Very little improvement in accuracy, but much more stable results. Try changing batch size

Batch size: 500, Learning rate: 0.0005, Epochs: 10

Training accuracy: 0.5550801475831706, Validation accuracy: 0.5555555555555556

Performance decreased, so try to decrease batch size.

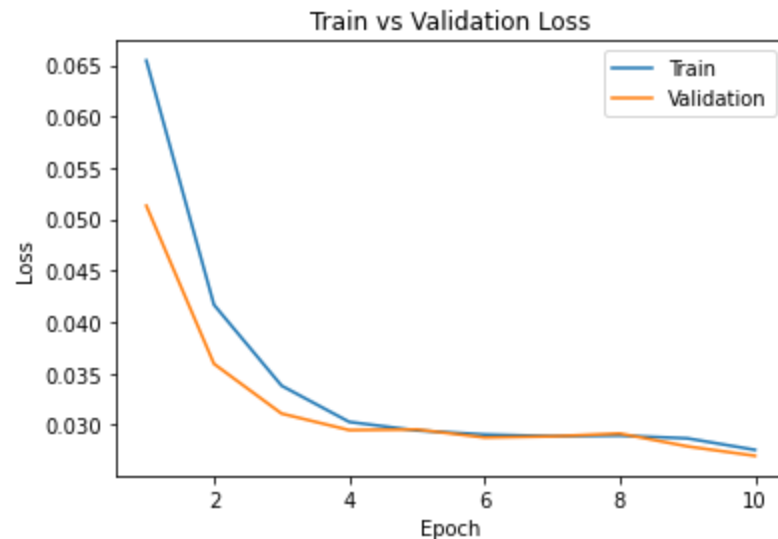
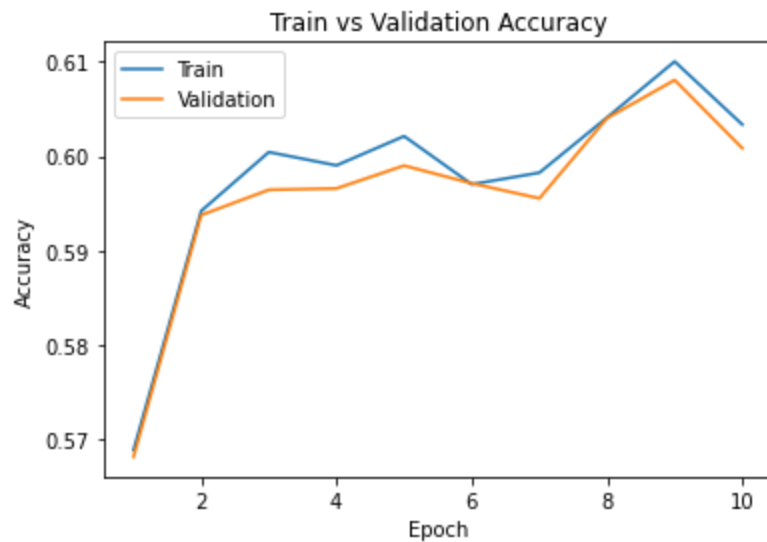
Batch size: 20, Learning rate: 0.0005, Epochs: 10

Training accuracy: 0.6100362756952842, Validation accuracy: 0.6080729166666666

This gives best result so far

```
In [ ]: model = AutoEncoder()
if torch.cuda.is_available():
    model = model.cuda()
train_loader, val_loader, test_loader = get_data_loader(20)
train(model, train_loader, val_loader, learning_rate=0.0005, num_epochs=10)
model_path = get_model_name(model.name, learning_rate=0.0005, epoch=9)
plot_training_curve(model_path)
```

Epoch 1: Train accuracy: 0.5688308064366105, Train loss: 0.06543973176121823 | Validation accuracy: 0.5681061921296297, Validation loss: 0.05130617460235953  
Epoch 2: Train accuracy: 0.5941850370508169, Train loss: 0.041673589682858524 | Validation accuracy: 0.59375, Validation loss: 0.03592992577187974  
Epoch 3: Train accuracy: 0.6004247666883701, Train loss: 0.03380038219650199 | Validation accuracy: 0.5964265046296297, Validation loss: 0.031088896705929576  
Epoch 4: Train accuracy: 0.5990140452050972, Train loss: 0.030276259910678466 | Validation accuracy: 0.5965711805555556, Validation loss: 0.029482345939774453  
Epoch 5: Train accuracy: 0.6021145319815211, Train loss: 0.02944352250189945 | Validation accuracy: 0.5989945023148148, Validation loss: 0.029543199560380186  
Epoch 6: Train accuracy: 0.5970064800173627, Train loss: 0.029052429373342348 | Validation accuracy: 0.5971137152777778, Validation loss: 0.028760255579354948  
Epoch 7: Train accuracy: 0.5982389235109912, Train loss: 0.028893387768358796 | Validation accuracy: 0.5955222800925926, Validation loss: 0.02887234645735087  
Epoch 8: Train accuracy: 0.6040910923014914, Train loss: 0.02894687116623423 | Validation accuracy: 0.6040219907407407, Validation loss: 0.02915070450800503  
Epoch 9: Train accuracy: 0.6100362756952842, Train loss: 0.028678435877228403 | Validation accuracy: 0.6080729166666666, Validation loss: 0.027898323553582204  
Epoch 10: Train accuracy: 0.6033469754751496, Train loss: 0.02755750241942752 | Validation accuracy: 0.6008391203703703, Validation loss: 0.02697892059122437  
Finished Training



## Part 4. Testing [12 pt]

### Part (a) [2 pt]

Compute and report the test accuracy.

```
In [ ]: print(get_accuracy(model, test_loader))
```

0.5985243055555556

## Part (b) [4 pt]

Based on the test accuracy alone, it is difficult to assess whether our model is actually performing well. We don't know whether a high accuracy is due to the simplicity of the problem, or if a poor accuracy is a result of the inherent difficulty of the problem.

It is therefore very important to be able to compare our model to at least one alternative. In particular, we consider a simple **baseline** model that is not very computationally expensive. Our neural network should at least outperform this baseline model. If our network is not much better than the baseline, then it is not doing well.

For our data imputation problem, consider the following baseline model: to predict a missing feature, the baseline model will look at the **most common value** of the feature in the training set.

For example, if the feature "marriage" is missing, then this model's prediction will be the most common value for "marriage" in the training set, which happens to be "Married-civ-spouse".

What would be the test accuracy of this baseline model?

```
In [ ]: most_common = get_features(train_data.sum(axis = 0))

correct = 0.0
for test in test_data:
    test_values = get_features(test)
    for feature in most_common:
        if most_common[feature] == test_values[feature]:
            correct += 1

accuracy = correct / (len(test_data) * len(catcols))
print("Accuracy: ", accuracy * 100, "%")
```

Accuracy: 45.68504050925926 %

## Part (c) [1 pt]

How does your test accuracy from part (a) compared to your baseline test accuracy in part (b)?

```
In [ ]: # Model gives 59.85% accuracy, larger than baseline which gives 46%
```

## Part (d) [1 pt]

Look at the first item in your test data. Do you think it is reasonable for a human to be able to guess this person's education level based on their other features? Explain.

```
In [ ]: # No, because it's difficult to predict based on just the other features
# It is very hard for human to learn the pattern from the entire dataset and draw correlations.
# Also, human have limited knowledge, which very likely lead to biased results.
```

## Part (e) [2 pt]

What is your model's prediction of this person's education level, given their other features?

```
In [ ]: input_arr = zero_out_feature(torch.tensor(test_data[0]).view(1,57), "edu")
        if torch.cuda.is_available():
            input_arr = input_arr.cuda()
        predict = model(input_arr).detach().cpu().numpy()
        print(get_feature(predict[0], "edu"))
```

Assoc-acdm

```
In [ ]: # Model predicts Assoc-acdm
```

## Part (f) [2 pt]

What is the baseline model's prediction of this person's education level?

```
In [ ]: pred = most_common["edu"]
        print(pred)
```

HS-grad

```
In [ ]: # The baseline model predicts that the person is a high school graduate.
```