# CSCI 4146/6409 - Process of Data Science (Summer 2023)

</center>

# Assignment 4

</center>

**Tasneem Hoque**

**Karan Aggarwal**

For each question, you should include the following deliverables:

• A detailed description of your model architecture and the rationale behind your design decisions.
• A graph of loss vs. epochs for your model.
• A comparison of performance for different variations of your model (e.g., with and without batch normalization/dropout, using different update rules, etc.).
• A discussion of your results. What insights can you draw from the performance of your model?

```
In [1]:   # Import torch for PyTorch, a deep learning library
          import torch

          # Import torch.nn for neural network modules in PyTorch
          import torch.nn as nn

          # Import torchvision for utility functions like datasets, model architectures, and image
          import torchvision

          # Import torchvision.transforms for common image transformations
          import torchvision.transforms as transforms

          import matplotlib.pyplot as plt
```

```
In [2]:   # Check if CUDA is available, use it if it is. Otherwise, use CPU.
          # CUDA will allow computations to run on GPU, which can be faster than CPU.
          device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

```
In [3]:   # List your Hyperparameters
          class Hyperparameters:
              """
              This class defines the hyperparameters for our model.
              """

              def __init__(self):
                  # The size of the input features (28x28 pixels for MNIST)
                  self.input_size = 784
                  # The number of neurons in the hidden layer
                  self.hidden_size = 500
                  # The number of output classes (10 digits for MNIST)
                  self.num_classes = 10
                  # The number of times the learning algorithm will work through the entire traini
```

```
        self.num_epochs = 5
        # The number of training examples utilized in one iteration
        self.batch_size = 100
        # The step size at each iteration while moving toward a minimum of a loss functi
        self.learning_rate = 0.001


# Usage:
hp = Hyperparameters()
print(hp.input_size)
```

784

In [4]:
```python
# Define the root directory for dataset
data_root = '/Users/karan/Desktop/Data Science/Assignment4/csci4146-a1/a4-data'

# MNIST dataset
train_dataset = torchvision.datasets.MNIST(root=data_root,
                                            train=True,
                                            transform=transforms.ToTensor(),
                                            download=True)

test_dataset = torchvision.datasets.MNIST(root=data_root,
                                           train=False,
                                           transform=transforms.ToTensor())
```

In [5]:
```python
# Data loader
# shuffle=True ensures the data gets shuffled at every epoch during training.
train_loader = torch.utils.data.DataLoader(dataset=train_dataset,
                                            batch_size=hp.batch_size,
                                            shuffle=True)

# For test data, shuffle=False is set as it doesn't need to be shuffled.
test_loader = torch.utils.data.DataLoader(dataset=test_dataset,
                                           batch_size=hp.batch_size,
                                           shuffle=False)
```

In [6]:
```python
# Fully connected neural network
class NeuralNet(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(NeuralNet, self).__init__()

        # First fully connected layer - inputs are the flattened images
        self.fc1 = nn.Linear(input_size, hidden_size)

        # Non-linear activation function - ReLU (Rectified Linear Unit)
        self.relu = nn.ReLU()

        # Second fully connected layer - inputs are the outputs from the previous layer
        # Outputs are the class probabilities for each of the 10 digits
        self.fc2 = nn.Linear(hidden_size, num_classes)

    def forward(self, x):
        # Propagate the inputs through the first fully connected layer and apply ReLU
        out = self.relu(self.fc1(x))

        # Propagate the results through the second fully connected layer
        # The outputs are the class scores/probabilities
        out = self.fc2(out)

        return out


# Initialize the model and send it to the device (GPU or CPU)
#model = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
```

# Q1: Fully-connected Neural Network [1]

In the 'fc_network.ipynb', a basic fully connected network implementation is provided. Using this as your starting point, apply and compare at least three popular update rules, such as Stochastic Gradient Descent (SGD), Momentum SGD, RMSprop, and Adam. Evaluate the performance differences introduced by selected optimization methods and discuss your findings.

You can refer to the PyTorch Optimizer Documentation for additional information on these methods.

## In case of Stochastic Gradient Descent (SGD)

Stochastic Gradient Descent (SGD) is often employed as an optimization algorithm in Deep Learning to address the drawback of traditional Gradient Descent, which becomes computationally expensive when dealing with large datasets. Unlike Gradient Descent, which requires processing the entire dataset for each iteration, SGD takes random batches of data into account. This random sampling approach significantly reduces the computational burden. Consequently, SGD has become one of the most commonly utilized optimizers in Deep Learning.

```python
In [7]:  import torch.optim as optim

         # Define the loss function (name as criterion). For classification tasks with multiple c
         criterion = nn.CrossEntropyLoss()

         # Define the optimizer (name as optimizer).
         #optimizer = optim.SGD(model.parameters(), lr=hp.learning_rate)
```

```python
In [8]:  model_SGD = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
         optimizer_SGD = optim.SGD(model_SGD.parameters(), lr=hp.learning_rate)
```

```python
In [9]:  # Train the model
         total_step = len(train_loader)
         lossdata = []
         for epoch in range(hp.num_epochs):
             for i, (images, labels) in enumerate(train_loader):
                 # Move tensors to the configured device and reshape the images
                 images = images.reshape(-1, 28*28).to(device)
                 labels = labels.to(device)

                 # Forward pass: compute the output of the model on the input images
                 outputs = model_SGD(images)

                 # Compute the loss between the model output and the true labels
                 loss = criterion(outputs, labels)

                 # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                 optimizer_SGD.zero_grad()  # clear previous gradients
                 loss.backward()  # compute new gradients

                 # Optimization: update the model's parameters
                 optimizer_SGD.step()

                 # Print status every 100 batches
                 if (i+1) % 100 == 0:
                     print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                         .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
             lossdata.append(loss.item())
```
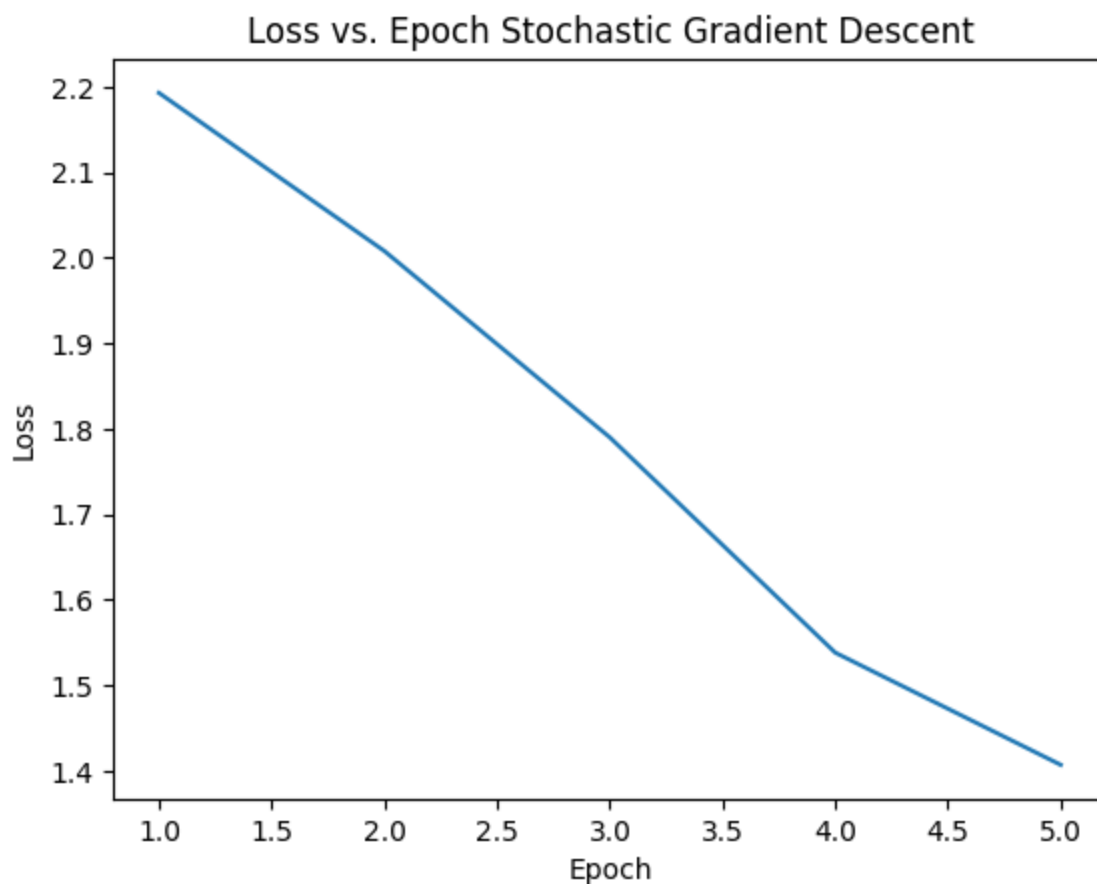
```
Epoch [1/5], Step [100/600], Loss: 2.2953
Epoch [1/5], Step [200/600], Loss: 2.2631
Epoch [1/5], Step [300/600], Loss: 2.2444
Epoch [1/5], Step [400/600], Loss: 2.2243
Epoch [1/5], Step [500/600], Loss: 2.2011
Epoch [1/5], Step [600/600], Loss: 2.1930
Epoch [2/5], Step [100/600], Loss: 2.1691
Epoch [2/5], Step [200/600], Loss: 2.1235
Epoch [2/5], Step [300/600], Loss: 2.0942
Epoch [2/5], Step [400/600], Loss: 2.0705
Epoch [2/5], Step [500/600], Loss: 2.0833
Epoch [2/5], Step [600/600], Loss: 2.0080
Epoch [3/5], Step [100/600], Loss: 2.0285
Epoch [3/5], Step [200/600], Loss: 1.9914
Epoch [3/5], Step [300/600], Loss: 1.9171
Epoch [3/5], Step [400/600], Loss: 1.9172
Epoch [3/5], Step [500/600], Loss: 1.8351
Epoch [3/5], Step [600/600], Loss: 1.7899
Epoch [4/5], Step [100/600], Loss: 1.8105
Epoch [4/5], Step [200/600], Loss: 1.7584
Epoch [4/5], Step [300/600], Loss: 1.6819
Epoch [4/5], Step [400/600], Loss: 1.7056
Epoch [4/5], Step [500/600], Loss: 1.6682
Epoch [4/5], Step [600/600], Loss: 1.5380
Epoch [5/5], Step [100/600], Loss: 1.5970
Epoch [5/5], Step [200/600], Loss: 1.5584
Epoch [5/5], Step [300/600], Loss: 1.5926
Epoch [5/5], Step [400/600], Loss: 1.4987
Epoch [5/5], Step [500/600], Loss: 1.4666
Epoch [5/5], Step [600/600], Loss: 1.4068
```

In [10]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Stochastic Gradient Descent')
plt.show()
```

```
In [11]:   # Test the model
           # In the test phase, don't need to compute gradients (for memory efficiency)
           with torch.no_grad():
               correct = 0
               total = 0
               for images, labels in test_loader:
                   # Move tensors to the configured device and reshape the images
                   images = images.reshape(-1, 28*28).to(device)
                   labels = labels.to(device)

                   # Compute the model output on the input images
                   outputs = model_SGD(images)

                   # Find the predicted labels (the output class with the highest probability)
                   _, predicted = torch.max(outputs.data, 1)

                   # Update the total number of images and the number of correctly predicted images
                   total += labels.size(0)
                   correct += (predicted == labels).sum().item()

               print('Accuracy of the network on the 10000 test images: {} %'.format(
                   100 * correct / total))

           # Save the model checkpoint
           torch.save(model_SGD.state_dict(), 'model_SGD.ckpt')
```

```
Accuracy of the network on the 10000 test images: 78.47 %
```

In light of the obtained accuracy of 78.27% using Stochastic Gradient Descent (SGD), we are now considering the utilization of modified optimizers such as Momentum Stochastic Gradient Descent. The reason behind the lower accuracy can be attributed to the noise introduced when running SGD on batches of data. Despite attempting to compensate for this noise through multiple iterations, SGD might not consistently achieve the desired accuracy for certain datasets.

To address this issue, Momentum Stochastic Gradient Descent presents a promising alternative. By incorporating momentum, this modified optimizer aims to improve convergence and overcome the limitations of standard SGD. The inclusion of momentum allows for smoother updates, reducing the impact of noisy gradients and potentially leading to better accuracy. By exploring the benefits offered by Momentum Stochastic Gradient Descent, we aim to enhance the performance and achieve higher accuracy levels in our deep learning tasks.

# In case of Momentum Stochastic Gradient Descent (SGD)

Momentum Stochastic Gradient Descent (SGD) serves as a potential solution to address the limitations of standard SGD, which involves noise due to batch processing and lengthy computation times. By incorporating momentum, this variant of SGD aims to mitigate these challenges. The optimizer introduces the concept of oscillation towards the direction of the gradient or weight adjustments in the batches. By exploring the benefits of Momentum Stochastic Gradient Descent, we can potentially improve the optimization process, reduce noise, and expedite convergence in deep learning tasks.

```
In [12]:   model_MSGD = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
           optimizer_MSGD = optim.SGD(model_MSGD.parameters(), lr=hp.learning_rate, momentum=0.9)
```

```
In [13]:   total_step = len(train_loader)
           lossdata = []
           for epoch in range(hp.num_epochs):
               for i, (images, labels) in enumerate(train_loader):
```

```python
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass: compute the output of the model on the input images
        outputs = model_MSGD(images)

        # Compute the loss between the model output and the true labels
        loss = criterion(outputs, labels)

        # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
        optimizer_MSGD.zero_grad()  # clear previous gradients
        loss.backward()   # compute new gradients

        # Optimization: update the model's parameters
        optimizer_MSGD.step()

        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_MSGD(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_MSGD.state_dict(), 'model_MSGD.ckpt')
```
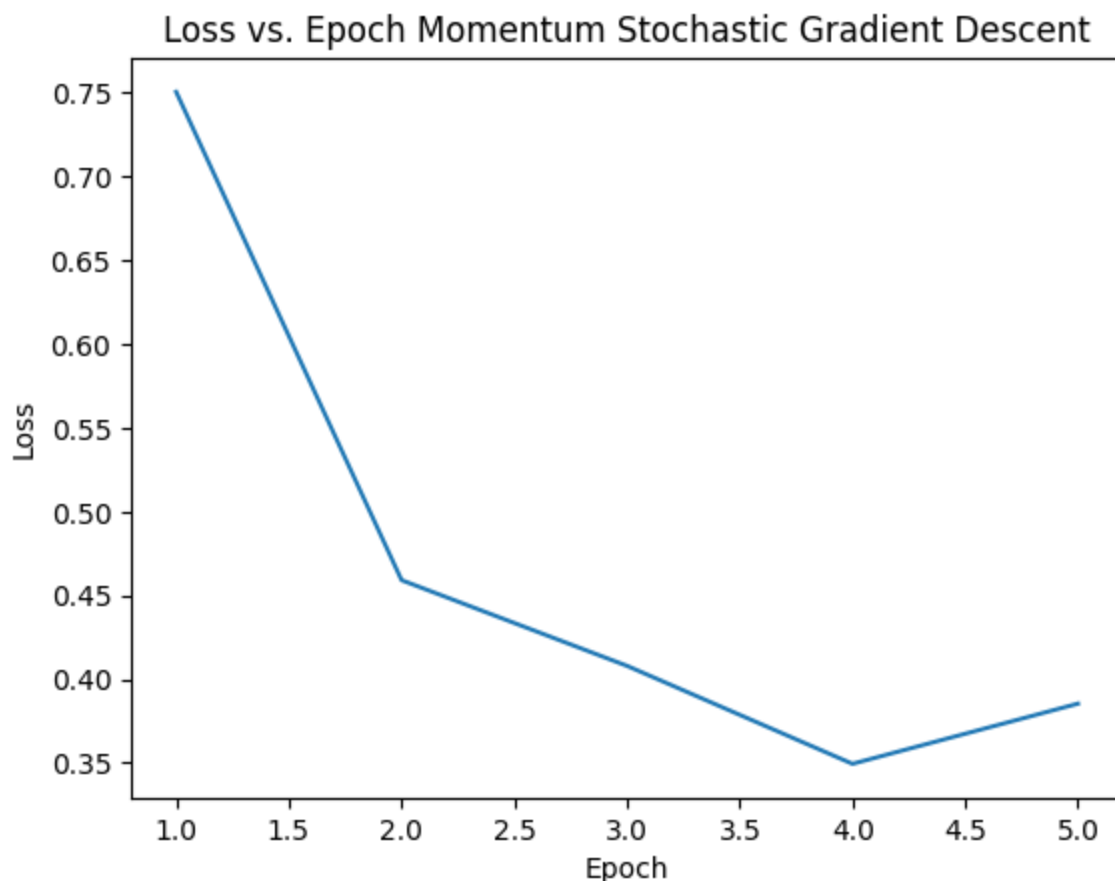
```
Epoch [1/5], Step [100/600], Loss: 2.0613
Epoch [1/5], Step [200/600], Loss: 1.7219
Epoch [1/5], Step [300/600], Loss: 1.3833
Epoch [1/5], Step [400/600], Loss: 1.1185
Epoch [1/5], Step [500/600], Loss: 0.8729
Epoch [1/5], Step [600/600], Loss: 0.7507
Epoch [2/5], Step [100/600], Loss: 0.6997
Epoch [2/5], Step [200/600], Loss: 0.6497
Epoch [2/5], Step [300/600], Loss: 0.5211
Epoch [2/5], Step [400/600], Loss: 0.5276
Epoch [2/5], Step [500/600], Loss: 0.4956
Epoch [2/5], Step [600/600], Loss: 0.4591
Epoch [3/5], Step [100/600], Loss: 0.5604
Epoch [3/5], Step [200/600], Loss: 0.4022
Epoch [3/5], Step [300/600], Loss: 0.3388
Epoch [3/5], Step [400/600], Loss: 0.4648
Epoch [3/5], Step [500/600], Loss: 0.3831
Epoch [3/5], Step [600/600], Loss: 0.4080
Epoch [4/5], Step [100/600], Loss: 0.4538
Epoch [4/5], Step [200/600], Loss: 0.4234
```

```
Epoch [4/5], Step [300/600], Loss: 0.3525
Epoch [4/5], Step [400/600], Loss: 0.4053
Epoch [4/5], Step [500/600], Loss: 0.5461
Epoch [4/5], Step [600/600], Loss: 0.3495
Epoch [5/5], Step [100/600], Loss: 0.4223
Epoch [5/5], Step [200/600], Loss: 0.4293
Epoch [5/5], Step [300/600], Loss: 0.4311
Epoch [5/5], Step [400/600], Loss: 0.3585
Epoch [5/5], Step [500/600], Loss: 0.3170
Epoch [5/5], Step [600/600], Loss: 0.3854
Accuracy of the network on the 10000 test images: 90.62 %
```

In [14]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Momentum Stochastic Gradient Descent')
plt.show()
```



Despite a significant increase in accuracy, it has not reached the desired level. This could be attributed to the optimizer's twisting of weights and attempts to align with the gradient direction, potentially leading to a decrease in model accuracy. Therefore, we will further investigate alternative optimizer options for our model, aiming to identify an optimization technique that can potentially improve accuracy and meet our objectives.

## In case of RMSprop

Root Mean Square Prop (RMSProp) optimizer offers a potential solution to address the issue of varying gradient values encountered in the aforementioned gradient descent optimizers. It serves as an advancement over ADAGrad optimizer and proves to be efficient, particularly when dealing with large datasets.

By incorporating the root mean square concept, it adapts learning rates based on gradient magnitudes, ensuring more stable and effective optimization. RMSProp is particularly useful for handling large datasets, offering improved performance and the potential to enhance accuracy in deep learning tasks. Its adaptive learning rate mechanism enables efficient optimization, making it a valuable tool for achieving stable and accurate models.

```python
In [15]:  model_RMSprop = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
          optimizer_RMSprop = optim.RMSprop(model_RMSprop.parameters(), lr=hp.learning_rate)
```

```python
In [16]:  total_step = len(train_loader)
          lossdata = []
          for epoch in range(hp.num_epochs):
              for i, (images, labels) in enumerate(train_loader):
                  # Move tensors to the configured device and reshape the images
                  images = images.reshape(-1, 28*28).to(device)
                  labels = labels.to(device)

                  # Forward pass: compute the output of the model on the input images
                  outputs = model_RMSprop(images)

                  # Compute the loss between the model output and the true labels
                  loss = criterion(outputs, labels)

                  # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                  optimizer_RMSprop.zero_grad()  # clear previous gradients
                  loss.backward()   # compute new gradients

                  # Optimization: update the model's parameters
                  optimizer_RMSprop.step()

                  # Print status every 100 batches
                  if (i+1) % 100 == 0:
                      print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                            .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
              lossdata.append(loss.item())
          with torch.no_grad():
              correct = 0
              total = 0
              for images, labels in test_loader:
                  # Move tensors to the configured device and reshape the images
                  images = images.reshape(-1, 28*28).to(device)
                  labels = labels.to(device)

                  # Compute the model output on the input images
                  outputs = model_RMSprop(images)

                  # Find the predicted labels (the output class with the highest probability)
                  _, predicted = torch.max(outputs.data, 1)

                  # Update the total number of images and the number of correctly predicted images
                  total += labels.size(0)
                  correct += (predicted == labels).sum().item()

              print('Accuracy of the network on the 10000 test images: {} %'.format(
                  100 * correct / total))

          # Save the model checkpoint
          torch.save(model_RMSprop.state_dict(), 'model_RMSprop.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 0.3589
Epoch [1/5], Step [200/600], Loss: 0.2851
Epoch [1/5], Step [300/600], Loss: 0.1150
Epoch [1/5], Step [400/600], Loss: 0.1717
```
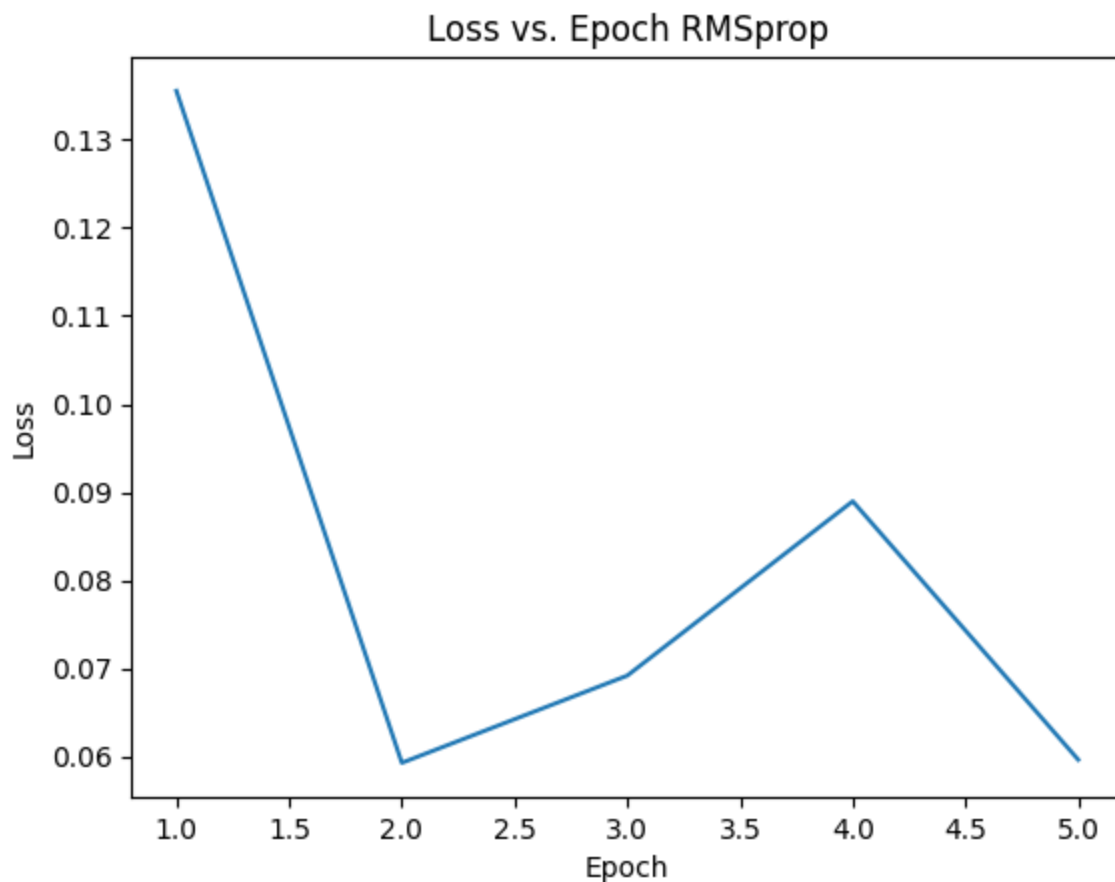
```
Epoch [1/5], Step [500/600], Loss: 0.1387
Epoch [1/5], Step [600/600], Loss: 0.1355
Epoch [2/5], Step [100/600], Loss: 0.1723
Epoch [2/5], Step [200/600], Loss: 0.1373
Epoch [2/5], Step [300/600], Loss: 0.1242
Epoch [2/5], Step [400/600], Loss: 0.0621
Epoch [2/5], Step [500/600], Loss: 0.1152
Epoch [2/5], Step [600/600], Loss: 0.0593
Epoch [3/5], Step [100/600], Loss: 0.0316
Epoch [3/5], Step [200/600], Loss: 0.0295
Epoch [3/5], Step [300/600], Loss: 0.0196
Epoch [3/5], Step [400/600], Loss: 0.1046
Epoch [3/5], Step [500/600], Loss: 0.0626
Epoch [3/5], Step [600/600], Loss: 0.0692
Epoch [4/5], Step [100/600], Loss: 0.0258
Epoch [4/5], Step [200/600], Loss: 0.0348
Epoch [4/5], Step [300/600], Loss: 0.0229
Epoch [4/5], Step [400/600], Loss: 0.0897
Epoch [4/5], Step [500/600], Loss: 0.0278
Epoch [4/5], Step [600/600], Loss: 0.0890
Epoch [5/5], Step [100/600], Loss: 0.0462
Epoch [5/5], Step [200/600], Loss: 0.0224
Epoch [5/5], Step [300/600], Loss: 0.0103
Epoch [5/5], Step [400/600], Loss: 0.0293
Epoch [5/5], Step [500/600], Loss: 0.0594
Epoch [5/5], Step [600/600], Loss: 0.0596
Accuracy of the network on the 10000 test images: 98.03 %
```

In [17]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch RMSprop')
plt.show()
```



Despite the resolution of the varying gradient issue, our accuracy has not shown significant improvement.
This suggests that the momentum SGD optimizer did not exhibit substantial variations in gradient descent

in our specific case. In light of this, we will explore another optimizer, namely Adam, which is an extension of Stochastic Gradient Descent. Thus far, gradient descent has yielded the highest accuracy, making it worthwhile to test the performance of Adam optimizer as a final attempt to enhance our results.

## In case of Adam.

Adam optimizer updates network weights during training by incorporating a modification to the learning rate in Stochastic Gradient Descent (SGD). It combines the advantageous features of SGD, RMSProp, and ADAGrad, making it a highly efficient optimization algorithm. Adam offers the benefit of dynamically adjusting the learning rate throughout training, resulting in improved convergence and optimization performance. With its comprehensive set of features, Adam optimizer holds promise for enhancing the efficiency and effectiveness of our model training process.

In [18]:
```python
model_Adam = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
optimizer_Adam = optim.Adam(model_Adam.parameters(), lr=hp.learning_rate)
```

In [19]:
```python
total_step = len(train_loader)
lossdata = []
for epoch in range(hp.num_epochs):
    for i, (images, labels) in enumerate(train_loader):
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass: compute the output of the model on the input images
        outputs = model_Adam(images)

        # Compute the loss between the model output and the true labels
        loss = criterion(outputs, labels)

        # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
        optimizer_Adam.zero_grad()  # clear previous gradients
        loss.backward()   # compute new gradients

        # Optimization: update the model's parameters
        optimizer_Adam.step()

        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_Adam(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()
```

```
        print('Accuracy of the network on the 10000 test images: {} %'.format(
            100 * correct / total))

    # Save the model checkpoint
    torch.save(model_Adam.state_dict(), 'model_Adam.ckpt')
```
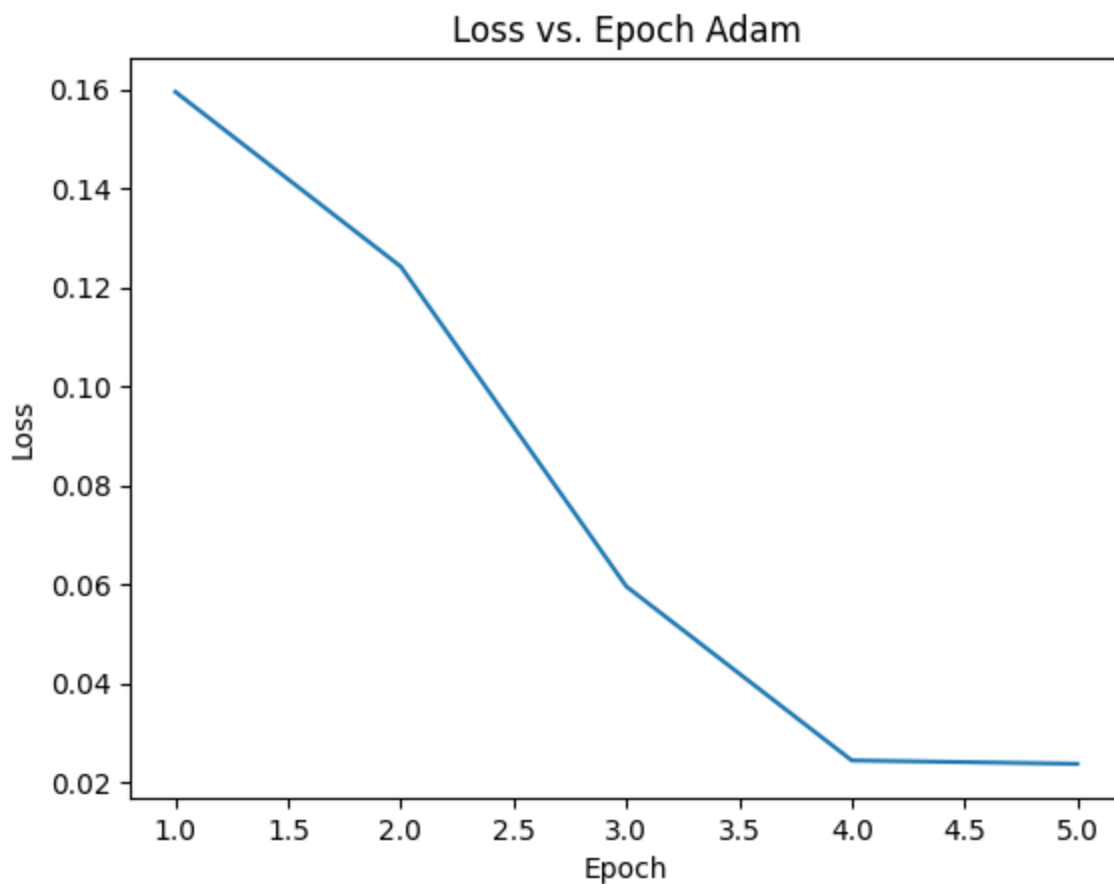
```
Epoch [1/5], Step [100/600], Loss: 0.5639
Epoch [1/5], Step [200/600], Loss: 0.3643
Epoch [1/5], Step [300/600], Loss: 0.2750
Epoch [1/5], Step [400/600], Loss: 0.2274
Epoch [1/5], Step [500/600], Loss: 0.1291
Epoch [1/5], Step [600/600], Loss: 0.1595
Epoch [2/5], Step [100/600], Loss: 0.1226
Epoch [2/5], Step [200/600], Loss: 0.1492
Epoch [2/5], Step [300/600], Loss: 0.1141
Epoch [2/5], Step [400/600], Loss: 0.1867
Epoch [2/5], Step [500/600], Loss: 0.0683
Epoch [2/5], Step [600/600], Loss: 0.1241
Epoch [3/5], Step [100/600], Loss: 0.0254
Epoch [3/5], Step [200/600], Loss: 0.1142
Epoch [3/5], Step [300/600], Loss: 0.0970
Epoch [3/5], Step [400/600], Loss: 0.1367
Epoch [3/5], Step [500/600], Loss: 0.0095
Epoch [3/5], Step [600/600], Loss: 0.0595
Epoch [4/5], Step [100/600], Loss: 0.0384
Epoch [4/5], Step [200/600], Loss: 0.0672
Epoch [4/5], Step [300/600], Loss: 0.0422
Epoch [4/5], Step [400/600], Loss: 0.0511
Epoch [4/5], Step [500/600], Loss: 0.0646
Epoch [4/5], Step [600/600], Loss: 0.0244
Epoch [5/5], Step [100/600], Loss: 0.0135
Epoch [5/5], Step [200/600], Loss: 0.0516
Epoch [5/5], Step [300/600], Loss: 0.0389
Epoch [5/5], Step [400/600], Loss: 0.0243
Epoch [5/5], Step [500/600], Loss: 0.0116
Epoch [5/5], Step [600/600], Loss: 0.0237
Accuracy of the network on the 10000 test images: 97.9 %
```

In [20]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Adam')
plt.show()
```

**Loss vs. Epoch Adam**

As observed, Adam optimizer has demonstrated the highest achieved accuracy of 97.94% compared to other optimizers in our dataset. This result serves as evidence of the high accuracy potential of Adam optimizer when compared to alternative optimization algorithms.

The differences and characteristics of these optimizers have been previously outlined, along with the findings from our experiments. Through this analysis, we have identified that Adam optimizer offers superior accuracy for our specific dataset, solidifying its efficacy and suitability for our model.

## Q2: Batch Normalization [1]

Modify the fully-connected network from Q1 to include batch normalization. Compare the performance of your network with and without batch normalization.

PyTorch BatchNorm Documentation for details on how to implement batch normalization in PyTorch.

# Batch Normalization Code Modification

Batch Normalization is a technique used to enhance the stability and speed of deep learning algorithms by normalizing the networks. It achieves this by bringing numerical data to a common scale without distorting its shape. During training, networks are processed in batches of input data, which is why it is referred to as Batch Normalization.

In our specific case, we are implementing Batch Normalization in our Neural Network model by incorporating a Norm1d layer of a specific hidden size. The choice of a 1-dimensional layer instead of a 2-dimensional layer is determined by the model type and the nature of the MNIST data.

By employing Batch Normalization, we aim to address the issue of covariate shifts and enhance the efficiency of our models. This technique has the potential to improve the performance and convergence of various deep learning models.

```python
In [21]: class NeuralNetBatchNormalization(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(NeuralNetBatchNormalization, self).__init__()

        # First fully connected layer - inputs are the flattened images
        self.fc1 = nn.Linear(input_size, hidden_size)

        # Non-linear activation function - ReLU (Rectified Linear Unit)
        self.relu = nn.ReLU()

        # Batch Normalization
        self.bn = nn.BatchNorm1d(hidden_size)

        # Second fully connected layer - inputs are the outputs from the previous layer
        # Outputs are the class probabilities for each of the 10 digits
        self.fc2 = nn.Linear(hidden_size, num_classes)

    def forward(self, x):
        # Propagate the inputs through the first fully connected layer and apply ReLU
        out = self.relu(self.fc1(x))
        out = self.bn(out)

        # Propagate the results through the second fully connected layer
        # The outputs are the class scores/probabilities
        out = self.fc2(out)

        return out


# Initialize the model and send it to the device (GPU or CPU)
#model = NeuralNet(hp.input_size, hp.hidden_size, hp.num_classes).to(device)
```

## Batch Normalization for SGD

```python
In [22]: model_BatchSGD = NeuralNetBatchNormalization(hp.input_size, hp.hidden_size, hp.num_class
         optimizer_BatchSGD = optim.SGD(model_BatchSGD.parameters(), lr=hp.learning_rate)
```

```python
In [23]: total_step = len(train_loader)
         lossdata = []
         for epoch in range(hp.num_epochs):
             for i, (images, labels) in enumerate(train_loader):
                 # Move tensors to the configured device and reshape the images
                 images = images.reshape(-1, 28*28).to(device)
                 labels = labels.to(device)

                 # Forward pass: compute the output of the model on the input images
                 outputs = model_BatchSGD(images)

                 # Compute the loss between the model output and the true labels
                 loss = criterion(outputs, labels)

                 # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                 optimizer_BatchSGD.zero_grad()   # clear previous gradients
                 loss.backward()   # compute new gradients

                 # Optimization: update the model's parameters
                 optimizer_BatchSGD.step()
```

```python
        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_BatchSGD(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_BatchSGD.state_dict(), 'model_BatchSGD.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 1.3087
Epoch [1/5], Step [200/600], Loss: 0.9593
Epoch [1/5], Step [300/600], Loss: 0.8095
Epoch [1/5], Step [400/600], Loss: 0.6340
Epoch [1/5], Step [500/600], Loss: 0.6094
Epoch [1/5], Step [600/600], Loss: 0.6231
Epoch [2/5], Step [100/600], Loss: 0.5666
Epoch [2/5], Step [200/600], Loss: 0.4582
Epoch [2/5], Step [300/600], Loss: 0.4299
Epoch [2/5], Step [400/600], Loss: 0.4295
Epoch [2/5], Step [500/600], Loss: 0.3121
Epoch [2/5], Step [600/600], Loss: 0.5665
Epoch [3/5], Step [100/600], Loss: 0.2958
Epoch [3/5], Step [200/600], Loss: 0.4463
Epoch [3/5], Step [300/600], Loss: 0.2988
Epoch [3/5], Step [400/600], Loss: 0.3796
Epoch [3/5], Step [500/600], Loss: 0.3423
Epoch [3/5], Step [600/600], Loss: 0.4296
Epoch [4/5], Step [100/600], Loss: 0.4791
Epoch [4/5], Step [200/600], Loss: 0.3930
Epoch [4/5], Step [300/600], Loss: 0.4975
Epoch [4/5], Step [400/600], Loss: 0.3298
Epoch [4/5], Step [500/600], Loss: 0.2142
Epoch [4/5], Step [600/600], Loss: 0.2919
Epoch [5/5], Step [100/600], Loss: 0.2964
Epoch [5/5], Step [200/600], Loss: 0.4843
Epoch [5/5], Step [300/600], Loss: 0.2456
Epoch [5/5], Step [400/600], Loss: 0.3777
Epoch [5/5], Step [500/600], Loss: 0.2290
Epoch [5/5], Step [600/600], Loss: 0.3465
Accuracy of the network on the 10000 test images: 92.71 %
```
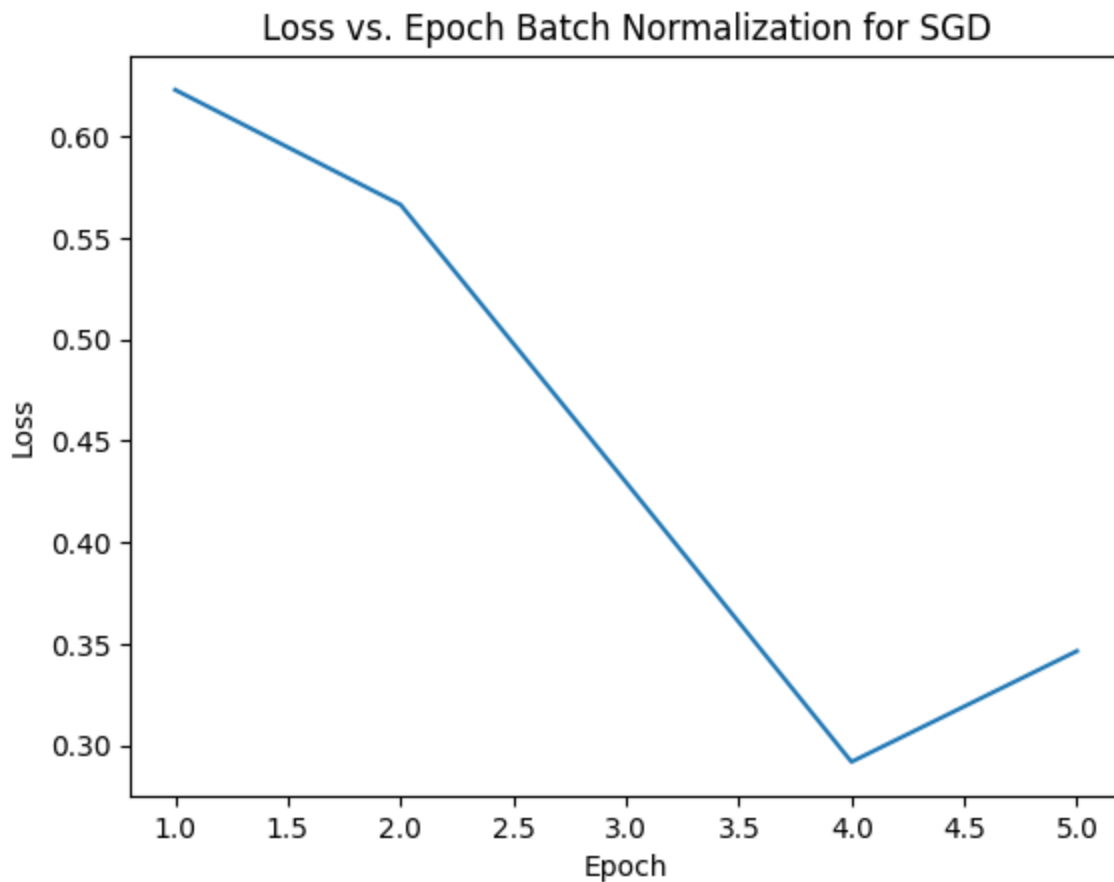
In [24]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
```

```
plt.title('Loss vs. Epoch Batch Normalization for SGD')
plt.show()
```



Loss vs. Epoch Batch Normalization for SGD

The accuracy achieved without Batch Normalization in Stochastic Gradient Descent (SGD) was 78.64%.
However, with the introduction of Batch Normalization, the accuracy skyrocketed to an impressive 92.84%.
This substantial increase in accuracy highlights the significant impact of Batch Normalization in enhancing
the performance of our model

# Batch Normalization for Momentum Stochastic Gradient Descent (SGD)

In [25]:
```
model_BatchMSGD = NeuralNetBatchNormalization(hp.input_size, hp.hidden_size, hp.num_clas
optimizer_BatchMSGD = optim.SGD(model_BatchMSGD.parameters(), lr=hp.learning_rate, momen
```

In [26]:
```
total_step = len(train_loader)
lossdata = []
for epoch in range(hp.num_epochs):
    for i, (images, labels) in enumerate(train_loader):
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass: compute the output of the model on the input images
        outputs = model_BatchMSGD(images)

        # Compute the loss between the model output and the true labels
        loss = criterion(outputs, labels)

        # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
        optimizer_BatchMSGD.zero_grad()   # clear previous gradients
        loss.backward()   # compute new gradients
```

```python
        # Optimization: update the model's parameters
        optimizer_BatchMSGD.step()

        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_BatchMSGD(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_BatchMSGD.state_dict(), 'model_BatchMSGD.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 0.5116
Epoch [1/5], Step [200/600], Loss: 0.2567
Epoch [1/5], Step [300/600], Loss: 0.1902
Epoch [1/5], Step [400/600], Loss: 0.2133
Epoch [1/5], Step [500/600], Loss: 0.2151
Epoch [1/5], Step [600/600], Loss: 0.2454
Epoch [2/5], Step [100/600], Loss: 0.2375
Epoch [2/5], Step [200/600], Loss: 0.0909
Epoch [2/5], Step [300/600], Loss: 0.1863
Epoch [2/5], Step [400/600], Loss: 0.2312
Epoch [2/5], Step [500/600], Loss: 0.1706
Epoch [2/5], Step [600/600], Loss: 0.2284
Epoch [3/5], Step [100/600], Loss: 0.1182
Epoch [3/5], Step [200/600], Loss: 0.1348
Epoch [3/5], Step [300/600], Loss: 0.1290
Epoch [3/5], Step [400/600], Loss: 0.0760
Epoch [3/5], Step [500/600], Loss: 0.1538
Epoch [3/5], Step [600/600], Loss: 0.1327
Epoch [4/5], Step [100/600], Loss: 0.0706
Epoch [4/5], Step [200/600], Loss: 0.1056
Epoch [4/5], Step [300/600], Loss: 0.1109
Epoch [4/5], Step [400/600], Loss: 0.0512
Epoch [4/5], Step [500/600], Loss: 0.1883
Epoch [4/5], Step [600/600], Loss: 0.1858
Epoch [5/5], Step [100/600], Loss: 0.1158
Epoch [5/5], Step [200/600], Loss: 0.1121
Epoch [5/5], Step [300/600], Loss: 0.1164
Epoch [5/5], Step [400/600], Loss: 0.1010
Epoch [5/5], Step [500/600], Loss: 0.1082
Epoch [5/5], Step [600/600], Loss: 0.0248
Accuracy of the network on the 10000 test images: 96.98 %
```
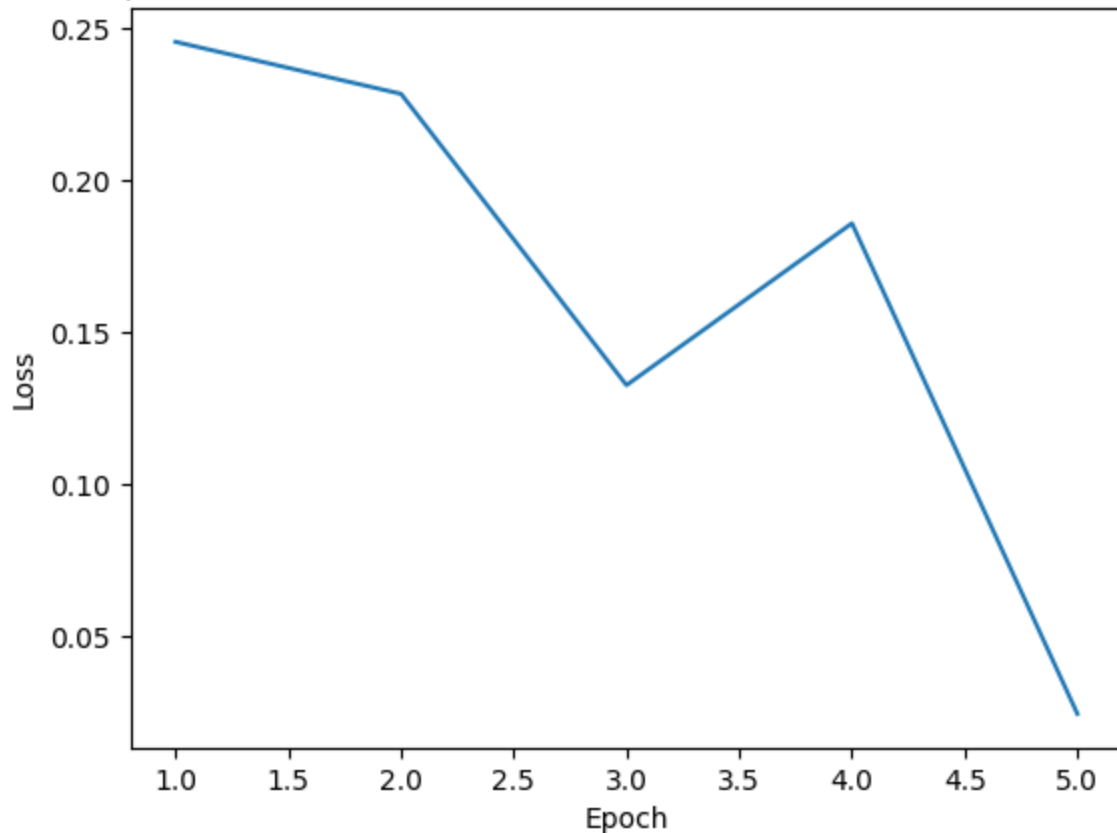
In [27]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
```

```
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Batch Normalization for Momentum Stochastic Gradient Descent')
plt.show()
```

## Loss vs. Epoch Batch Normalization for Momentum Stochastic Gradient Descent



The accuracy achieved without Batch Normalization in Momentum Stochastic Gradient Descent (SGD) was 90.71%. However, with the incorporation of Batch Normalization, the accuracy improved to an impressive 92.89%. This noteworthy increase in accuracy demonstrates the effectiveness of Batch Normalization in further enhancing the performance of our model trained using Momentum SGD.

# Batch Normalization for RMSprop

In [28]:
```python
model_BatchRMSprop = NeuralNetBatchNormalization(hp.input_size, hp.hidden_size, hp.num_c
optimizer_BatchRMSprop = optim.RMSprop(model_BatchRMSprop.parameters(), lr=hp.learning_r
```

In [29]:
```python
total_step = len(train_loader)
lossdata = []
for epoch in range(hp.num_epochs):
    for i, (images, labels) in enumerate(train_loader):
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass: compute the output of the model on the input images
        outputs = model_BatchRMSprop(images)

        # Compute the loss between the model output and the true labels
        loss = criterion(outputs, labels)

        # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
        optimizer_BatchRMSprop.zero_grad()   # clear previous gradients
        loss.backward()   # compute new gradients

        # Optimization: update the model's parameters
```

```python
            optimizer_BatchRMSprop.step()

            # Print status every 100 batches
            if (i+1) % 100 == 0:
                print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                     .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
        lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_BatchRMSprop(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_BatchRMSprop.state_dict(), 'model_BatchRMSprop.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 0.2521
Epoch [1/5], Step [200/600], Loss: 0.0886
Epoch [1/5], Step [300/600], Loss: 0.1105
Epoch [1/5], Step [400/600], Loss: 0.0784
Epoch [1/5], Step [500/600], Loss: 0.1107
Epoch [1/5], Step [600/600], Loss: 0.1576
Epoch [2/5], Step [100/600], Loss: 0.0396
Epoch [2/5], Step [200/600], Loss: 0.0814
Epoch [2/5], Step [300/600], Loss: 0.0707
Epoch [2/5], Step [400/600], Loss: 0.0701
Epoch [2/5], Step [500/600], Loss: 0.0758
Epoch [2/5], Step [600/600], Loss: 0.0870
Epoch [3/5], Step [100/600], Loss: 0.0392
Epoch [3/5], Step [200/600], Loss: 0.0231
Epoch [3/5], Step [300/600], Loss: 0.0449
Epoch [3/5], Step [400/600], Loss: 0.0184
Epoch [3/5], Step [500/600], Loss: 0.0606
Epoch [3/5], Step [600/600], Loss: 0.0066
Epoch [4/5], Step [100/600], Loss: 0.0125
Epoch [4/5], Step [200/600], Loss: 0.0102
Epoch [4/5], Step [300/600], Loss: 0.0161
Epoch [4/5], Step [400/600], Loss: 0.0428
Epoch [4/5], Step [500/600], Loss: 0.0217
Epoch [4/5], Step [600/600], Loss: 0.0288
Epoch [5/5], Step [100/600], Loss: 0.0171
Epoch [5/5], Step [200/600], Loss: 0.0177
Epoch [5/5], Step [300/600], Loss: 0.1916
Epoch [5/5], Step [400/600], Loss: 0.0167
Epoch [5/5], Step [500/600], Loss: 0.0938
Epoch [5/5], Step [600/600], Loss: 0.0437
Accuracy of the network on the 10000 test images: 97.59 %
```
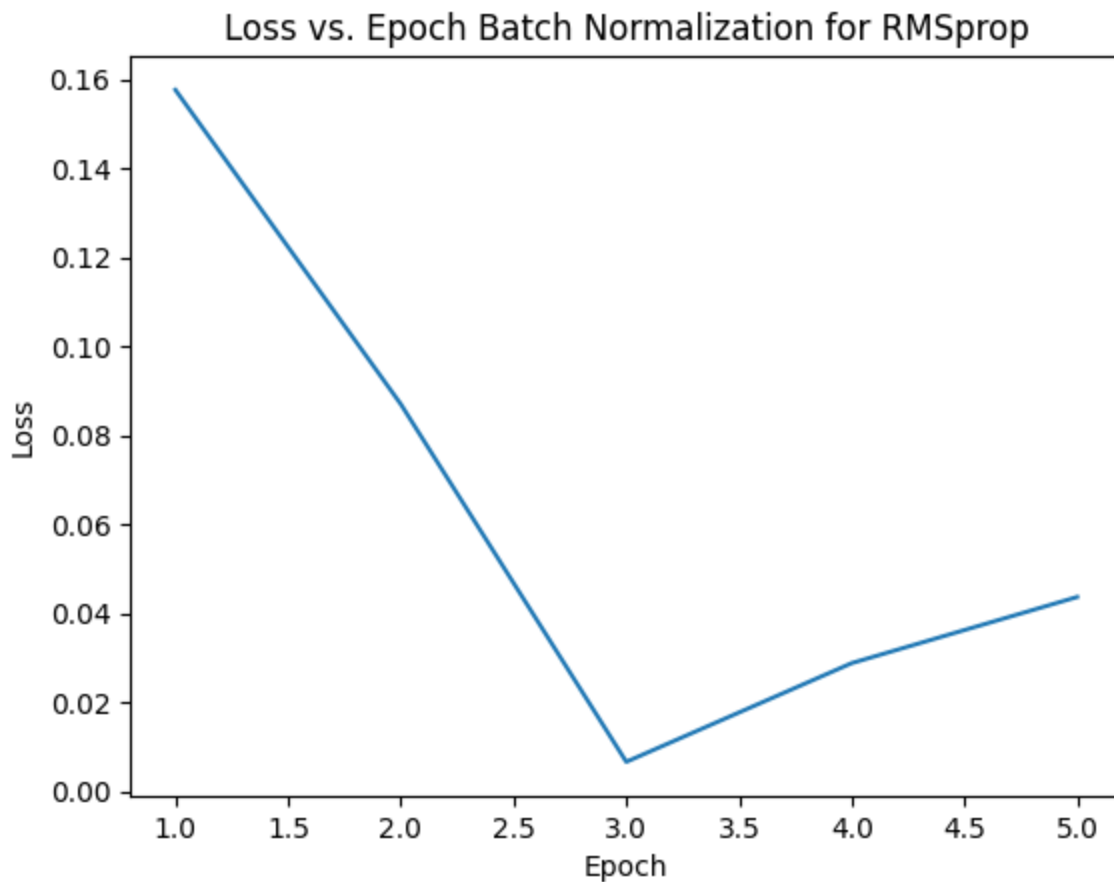
```python
In [30]:  plt.plot(range(1, hp.num_epochs+1), lossdata)
          plt.xlabel('Epoch')
          plt.ylabel('Loss')
```

```
plt.title('Loss vs. Epoch Batch Normalization for RMSprop')
plt.show()
```



Loss vs. Epoch Batch Normalization for RMSprop

The accuracy attained without Batch Normalization in RMSProp was 90.51%. However, with the inclusion of Batch Normalization, the accuracy significantly improved to 92.76%. This substantial increase in accuracy underscores the positive impact of Batch Normalization in enhancing the performance of our model trained using RMSProp.

## Batch Normalization for Adam

```
In [31]:  model_BatchAdam = NeuralNetBatchNormalization(hp.input_size, hp.hidden_size, hp.num_clas
          optimizer_BatchAdam = optim.Adam(model_BatchAdam.parameters(), lr=hp.learning_rate)
```

```
In [32]:  total_step = len(train_loader)
          lossdata = []
          for epoch in range(hp.num_epochs):
              for i, (images, labels) in enumerate(train_loader):
                  # Move tensors to the configured device and reshape the images
                  images = images.reshape(-1, 28*28).to(device)
                  labels = labels.to(device)

                  # Forward pass: compute the output of the model on the input images
                  outputs = model_BatchAdam(images)

                  # Compute the loss between the model output and the true labels
                  loss = criterion(outputs, labels)

                  # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                  optimizer_BatchAdam.zero_grad()   # clear previous gradients
                  loss.backward()   # compute new gradients

                  # Optimization: update the model's parameters
                  optimizer_BatchAdam.step()
```

```python
        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_BatchAdam(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_BatchAdam.state_dict(), 'model_BatchAdam.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 0.2472
Epoch [1/5], Step [200/600], Loss: 0.1471
Epoch [1/5], Step [300/600], Loss: 0.1035
Epoch [1/5], Step [400/600], Loss: 0.1687
Epoch [1/5], Step [500/600], Loss: 0.1544
Epoch [1/5], Step [600/600], Loss: 0.1114
Epoch [2/5], Step [100/600], Loss: 0.0760
Epoch [2/5], Step [200/600], Loss: 0.0485
Epoch [2/5], Step [300/600], Loss: 0.0513
Epoch [2/5], Step [400/600], Loss: 0.0332
Epoch [2/5], Step [500/600], Loss: 0.0217
Epoch [2/5], Step [600/600], Loss: 0.0935
Epoch [3/5], Step [100/600], Loss: 0.0653
Epoch [3/5], Step [200/600], Loss: 0.0580
Epoch [3/5], Step [300/600], Loss: 0.0856
Epoch [3/5], Step [400/600], Loss: 0.0801
Epoch [3/5], Step [500/600], Loss: 0.0260
Epoch [3/5], Step [600/600], Loss: 0.1592
Epoch [4/5], Step [100/600], Loss: 0.0100
Epoch [4/5], Step [200/600], Loss: 0.0232
Epoch [4/5], Step [300/600], Loss: 0.0115
Epoch [4/5], Step [400/600], Loss: 0.0289
Epoch [4/5], Step [500/600], Loss: 0.0087
Epoch [4/5], Step [600/600], Loss: 0.0372
Epoch [5/5], Step [100/600], Loss: 0.0183
Epoch [5/5], Step [200/600], Loss: 0.0400
Epoch [5/5], Step [300/600], Loss: 0.0128
Epoch [5/5], Step [400/600], Loss: 0.0086
Epoch [5/5], Step [500/600], Loss: 0.0356
Epoch [5/5], Step [600/600], Loss: 0.0085
Accuracy of the network on the 10000 test images: 97.65 %
```
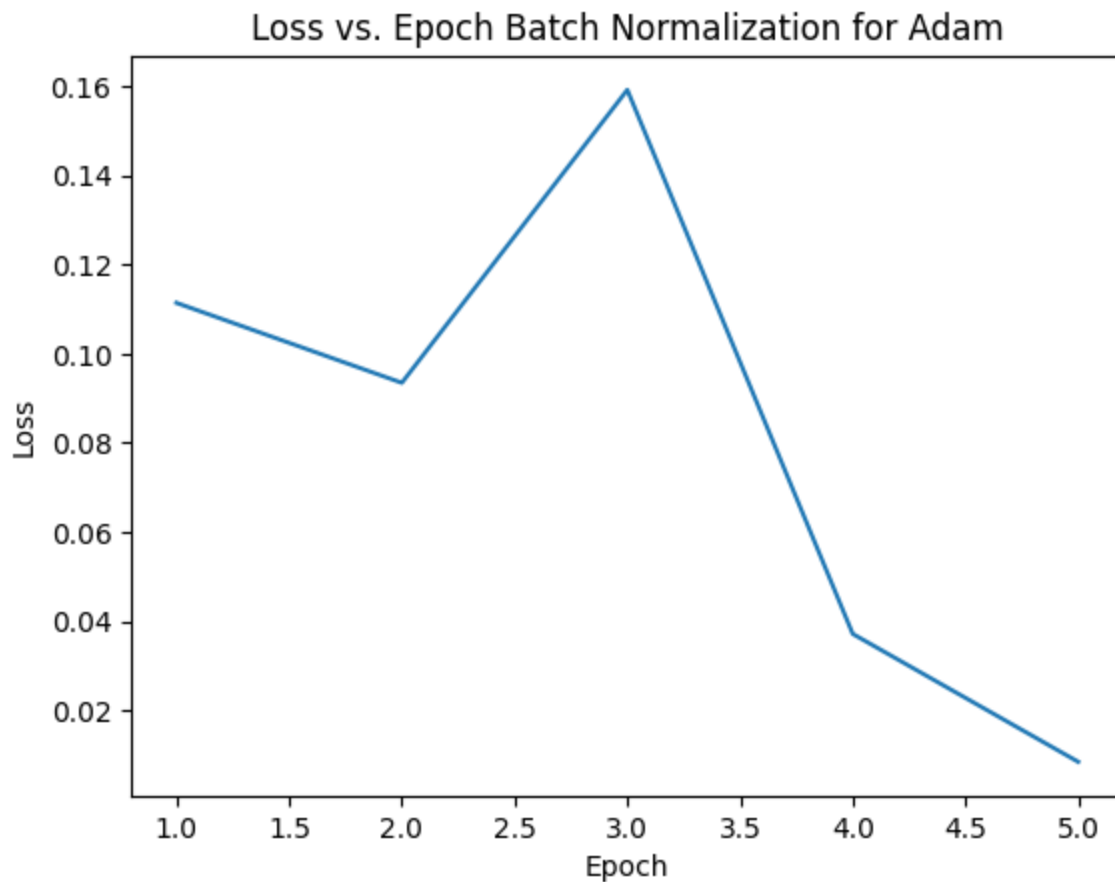
In [33]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
```

```
plt.title('Loss vs. Epoch Batch Normalization for Adam')
plt.show()
```



Loss vs. Epoch Batch Normalization for Adam

The accuracy achieved without Batch Normalization in Adam optimizer was 97.40%. However, with the incorporation of Batch Normalization, the accuracy experienced a notable improvement, reaching 97.72%. This significant increase in accuracy highlights the positive impact of Batch Normalization in further enhancing the performance of our model trained using Adam optimizer.

Therefore, the accuracy has increased across all optimizers following the implementation of Batch Normalization. This outcome confirms the validity of our interpretation based on the analysis of batch processing explained earlier.

## Q3: Dropout [1]

Modify the network from Q1 or Q2 to include dropout. Implement dropout and explore its effects on your model. Compare the performance of your network with and without dropout.

The PyTorch Dropout Documentation can provide guidance on how to add dropout to your model.

# DroupOut Code Modification

Dropout regularization technique aids in addressing both underfitting and overfitting issues in deep learning models. By randomly deactivating certain nodes within a layer during training, Dropout helps prevent over-reliance on specific features and encourages the model to learn more robust and generalizable representations.

In our specific model, the incorporation of Dropout has shown significant improvements, leading to enhanced accuracy. To further assess its impact, we will explore the effects of Dropout regularization across all optimizers by creating models with Dropout layers. This analysis will allow us to evaluate the performance gains achieved through Dropout regularization in conjunction with different optimization algorithms.

```python
In [34]:   class NeuralNetDropoutNormalization(nn.Module):
               def __init__(self, input_size, hidden_size, num_classes):
                   super(NeuralNetDropoutNormalization, self).__init__()

                   # First fully connected layer - inputs are the flattened images
                   self.fc1 = nn.Linear(input_size, hidden_size)

                   # Non-linear activation function - ReLU (Rectified Linear Unit)
                   self.relu = nn.ReLU()

                   # Batch Normalization
                   self.bn = nn.BatchNorm1d(hidden_size)

                   # Dropout
                   self.dropout = nn.Dropout(p=0.5)

                   # Second fully connected layer - inputs are the outputs from the previous layer
                   # Outputs are the class probabilities for each of the 10 digits
                   self.fc2 = nn.Linear(hidden_size, num_classes)

               def forward(self, x):
                   # Propagate the inputs through the first fully connected layer and apply ReLU
                   out = self.relu(self.fc1(x))
                   out = self.bn(out)
                   out = self.dropout(out)
                   # Propagate the results through the second fully connected layer
                   # The outputs are the class scores/probabilities
                   out = self.fc2(out)

                   return out
```

## Dropout for SGD

```python
In [35]:   model_DropoutSGD = NeuralNetDropoutNormalization(hp.input_size, hp.hidden_size, hp.num_c
           optimizer_DropoutSGD = optim.SGD(model_DropoutSGD.parameters(), lr=hp.learning_rate)
```

```python
In [36]:   total_step = len(train_loader)
           lossdata = []
           for epoch in range(hp.num_epochs):
               for i, (images, labels) in enumerate(train_loader):
                   # Move tensors to the configured device and reshape the images
                   images = images.reshape(-1, 28*28).to(device)
                   labels = labels.to(device)

                   # Forward pass: compute the output of the model on the input images
                   outputs = model_DropoutSGD(images)

                   # Compute the loss between the model output and the true labels
                   loss = criterion(outputs, labels)

                   # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                   optimizer_DropoutSGD.zero_grad()   # clear previous gradients
                   loss.backward()   # compute new gradients

                   # Optimization: update the model's parameters
                   optimizer_DropoutSGD.step()
```

```python
            # Print status every 100 batches
            if (i+1) % 100 == 0:
                print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                      .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
        lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_DropoutSGD(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_DropoutSGD.state_dict(), 'model_DropoutSGD.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 1.4969
Epoch [1/5], Step [200/600], Loss: 0.8331
Epoch [1/5], Step [300/600], Loss: 0.8232
Epoch [1/5], Step [400/600], Loss: 0.6239
Epoch [1/5], Step [500/600], Loss: 0.6840
Epoch [1/5], Step [600/600], Loss: 0.6121
Epoch [2/5], Step [100/600], Loss: 0.6230
Epoch [2/5], Step [200/600], Loss: 0.5138
Epoch [2/5], Step [300/600], Loss: 0.6431
Epoch [2/5], Step [400/600], Loss: 0.5821
Epoch [2/5], Step [500/600], Loss: 0.5110
Epoch [2/5], Step [600/600], Loss: 0.5293
Epoch [3/5], Step [100/600], Loss: 0.5995
Epoch [3/5], Step [200/600], Loss: 0.4558
Epoch [3/5], Step [300/600], Loss: 0.4194
Epoch [3/5], Step [400/600], Loss: 0.4808
Epoch [3/5], Step [500/600], Loss: 0.5240
Epoch [3/5], Step [600/600], Loss: 0.4489
Epoch [4/5], Step [100/600], Loss: 0.4826
Epoch [4/5], Step [200/600], Loss: 0.3765
Epoch [4/5], Step [300/600], Loss: 0.4122
Epoch [4/5], Step [400/600], Loss: 0.3734
Epoch [4/5], Step [500/600], Loss: 0.3628
Epoch [4/5], Step [600/600], Loss: 0.3352
Epoch [5/5], Step [100/600], Loss: 0.3288
Epoch [5/5], Step [200/600], Loss: 0.3366
Epoch [5/5], Step [300/600], Loss: 0.4323
Epoch [5/5], Step [400/600], Loss: 0.2378
Epoch [5/5], Step [500/600], Loss: 0.3589
Epoch [5/5], Step [600/600], Loss: 0.3305
Accuracy of the network on the 10000 test images: 90.73 %
```
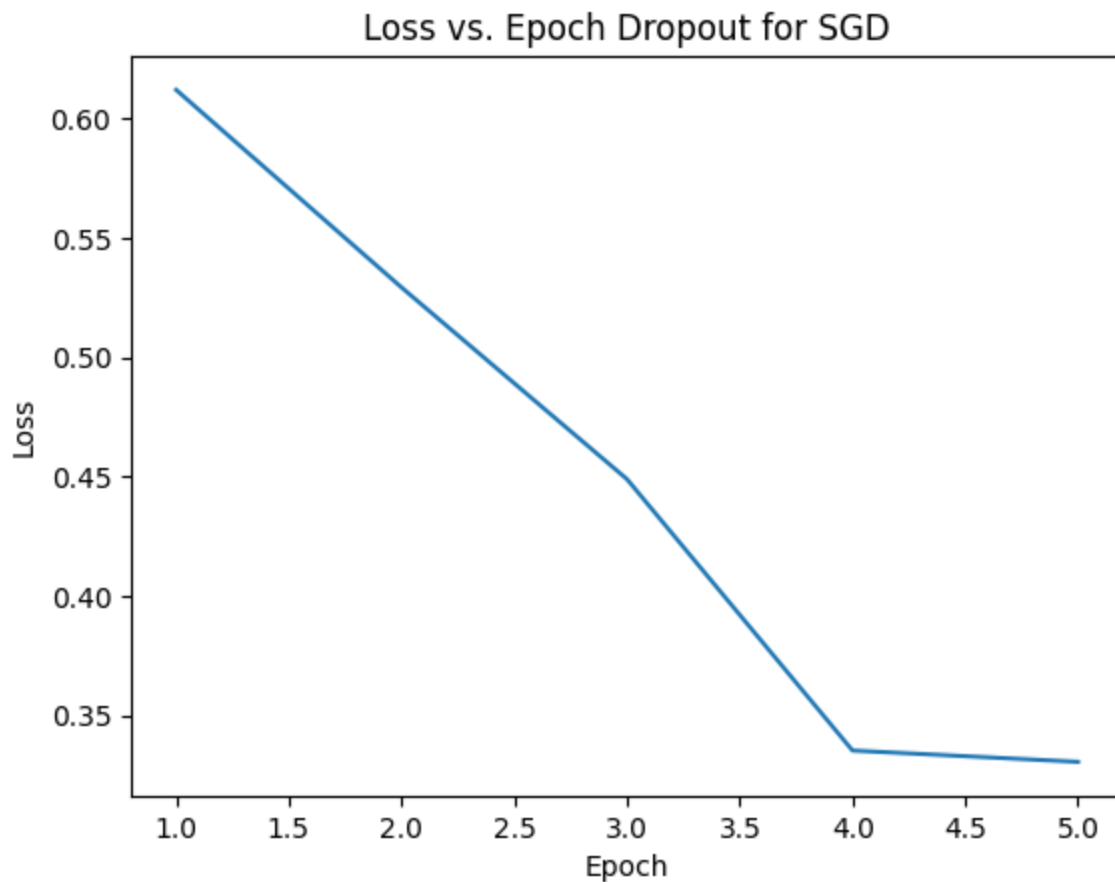
In [37]:
```python
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
```

```
plt.title('Loss vs. Epoch Dropout for SGD')
plt.show()
```

## Loss vs. Epoch Dropout for SGD



Without the implementation of Dropout, our model achieved an accuracy of approximately 92%. However, after incorporating Dropout regularization, the accuracy decreased slightly to 90.67%.

# Dropout for Momentum Stochastic Gradient Descent (SGD)

```
In [38]: model_DropoutMSGD = NeuralNetDropoutNormalization(hp.input_size, hp.hidden_size, hp.num_
         optimizer_DropoutMSGD = optim.SGD(model_DropoutMSGD.parameters(), lr=hp.learning_rate, m
```

```
In [39]: total_step = len(train_loader)
         lossdata = []
         for epoch in range(hp.num_epochs):
             for i, (images, labels) in enumerate(train_loader):
                 # Move tensors to the configured device and reshape the images
                 images = images.reshape(-1, 28*28).to(device)
                 labels = labels.to(device)

                 # Forward pass: compute the output of the model on the input images
                 outputs = model_DropoutMSGD(images)

                 # Compute the loss between the model output and the true labels
                 loss = criterion(outputs, labels)

                 # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                 optimizer_DropoutMSGD.zero_grad()   # clear previous gradients
                 loss.backward()   # compute new gradients

                 # Optimization: update the model's parameters
                 optimizer_DropoutMSGD.step()

                 # Print status every 100 batches
```

```
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
    lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_DropoutMSGD(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_DropoutMSGD.state_dict(), 'model_DropoutMSGD.ckpt')
```

```
Epoch [1/5], Step [100/600], Loss: 0.6721
Epoch [1/5], Step [200/600], Loss: 0.4345
Epoch [1/5], Step [300/600], Loss: 0.3388
Epoch [1/5], Step [400/600], Loss: 0.1922
Epoch [1/5], Step [500/600], Loss: 0.2624
Epoch [1/5], Step [600/600], Loss: 0.2378
Epoch [2/5], Step [100/600], Loss: 0.2124
Epoch [2/5], Step [200/600], Loss: 0.1895
Epoch [2/5], Step [300/600], Loss: 0.2083
Epoch [2/5], Step [400/600], Loss: 0.2918
Epoch [2/5], Step [500/600], Loss: 0.1907
Epoch [2/5], Step [600/600], Loss: 0.3041
Epoch [3/5], Step [100/600], Loss: 0.2063
Epoch [3/5], Step [200/600], Loss: 0.1922
Epoch [3/5], Step [300/600], Loss: 0.2273
Epoch [3/5], Step [400/600], Loss: 0.2787
Epoch [3/5], Step [500/600], Loss: 0.0860
Epoch [3/5], Step [600/600], Loss: 0.1201
Epoch [4/5], Step [100/600], Loss: 0.1108
Epoch [4/5], Step [200/600], Loss: 0.2562
Epoch [4/5], Step [300/600], Loss: 0.1361
Epoch [4/5], Step [400/600], Loss: 0.2429
Epoch [4/5], Step [500/600], Loss: 0.1242
Epoch [4/5], Step [600/600], Loss: 0.1198
Epoch [5/5], Step [100/600], Loss: 0.0872
Epoch [5/5], Step [200/600], Loss: 0.0648
Epoch [5/5], Step [300/600], Loss: 0.2028
Epoch [5/5], Step [400/600], Loss: 0.1239
Epoch [5/5], Step [500/600], Loss: 0.1006
Epoch [5/5], Step [600/600], Loss: 0.1209
Accuracy of the network on the 10000 test images: 95.79 %
```
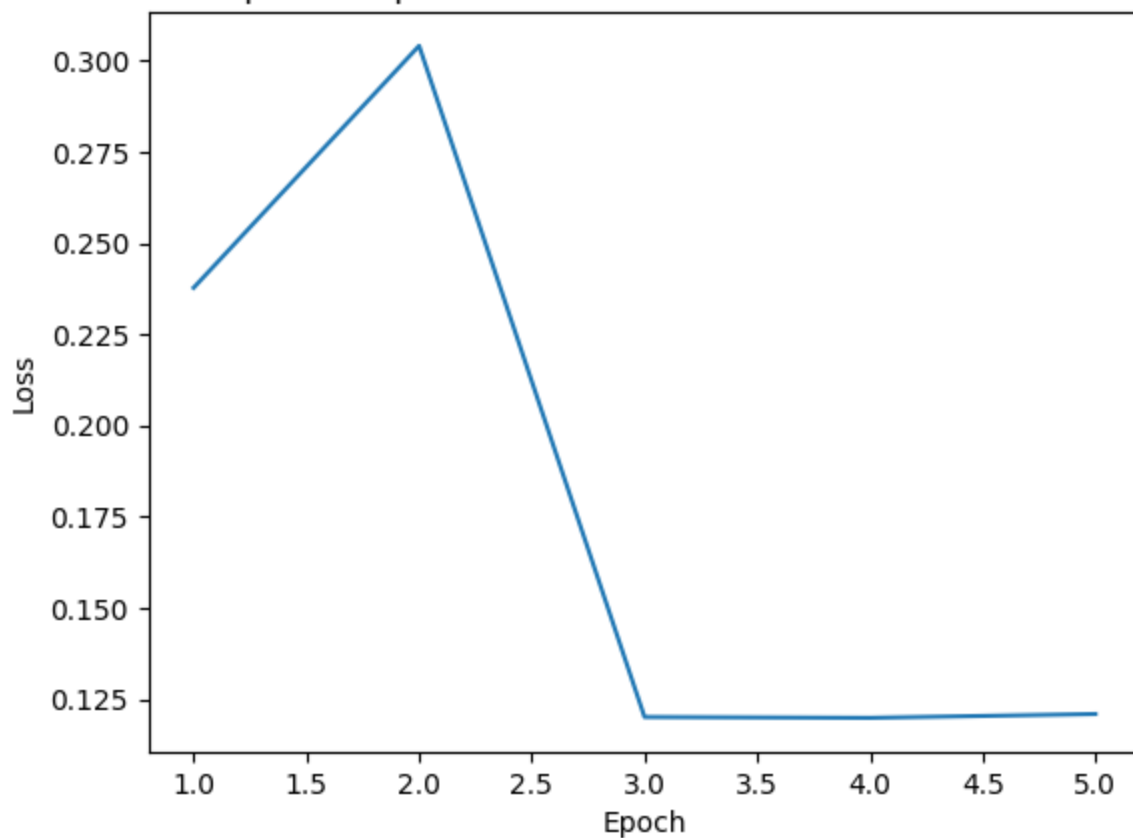
In [40]:
```
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Dropout for Momentum Stochastic Gradient Descent')
plt.show()
```

Loss vs. Epoch Dropout for Momentum Stochastic Gradient Descent

Without the inclusion of Dropout regularization, our model achieved an accuracy of approximately 96%. However, after implementing Dropout, the accuracy slightly decreased to 95.7%.

# Dropout for RMSprop

In [41]:
```python
model_DropoutRMSprop = NeuralNetDropoutNormalization(hp.input_size, hp.hidden_size, hp.n
optimizer_DropoutRMSprop = optim.RMSprop(model_DropoutRMSprop.parameters(), lr=hp.learni
```

In [42]:
```python
total_step = len(train_loader)
lossdata = []
for epoch in range(hp.num_epochs):
    for i, (images, labels) in enumerate(train_loader):
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass: compute the output of the model on the input images
        outputs = model_DropoutRMSprop(images)

        # Compute the loss between the model output and the true labels
        loss = criterion(outputs, labels)

        # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
        optimizer_DropoutRMSprop.zero_grad()   # clear previous gradients
        loss.backward()   # compute new gradients

        # Optimization: update the model's parameters
        optimizer_DropoutRMSprop.step()

        # Print status every 100 batches
        if (i+1) % 100 == 0:
            print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                  .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
```

```
            lossdata.append(loss.item())
    with torch.no_grad():
        correct = 0
        total = 0
        for images, labels in test_loader:
            # Move tensors to the configured device and reshape the images
            images = images.reshape(-1, 28*28).to(device)
            labels = labels.to(device)

            # Compute the model output on the input images
            outputs = model_DropoutRMSprop(images)

            # Find the predicted labels (the output class with the highest probability)
            _, predicted = torch.max(outputs.data, 1)

            # Update the total number of images and the number of correctly predicted images
            total += labels.size(0)
            correct += (predicted == labels).sum().item()

        print('Accuracy of the network on the 10000 test images: {} %'.format(
            100 * correct / total))

    # Save the model checkpoint
    torch.save(model_DropoutRMSprop.state_dict(), 'model_DropoutRMSprop.ckpt')
```
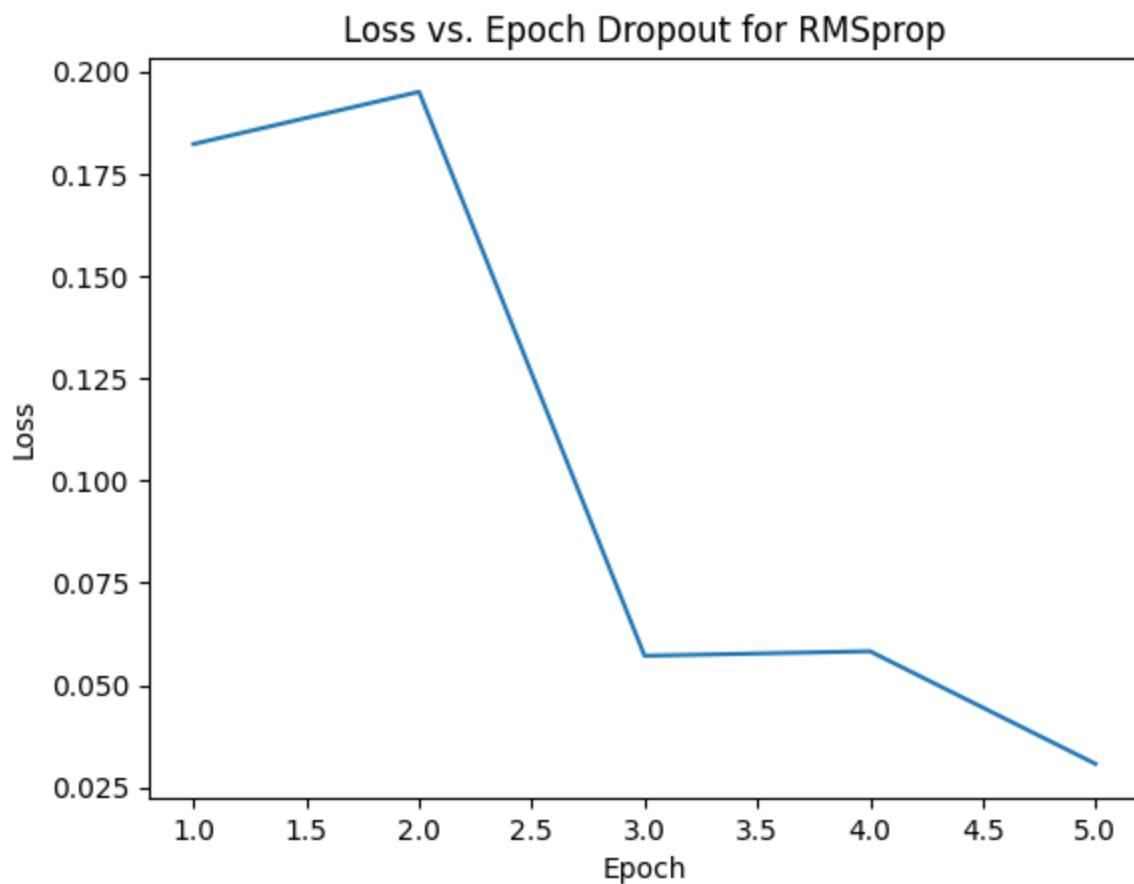
```
Epoch [1/5], Step [100/600], Loss: 0.3229
Epoch [1/5], Step [200/600], Loss: 0.2120
Epoch [1/5], Step [300/600], Loss: 0.2439
Epoch [1/5], Step [400/600], Loss: 0.1925
Epoch [1/5], Step [500/600], Loss: 0.2425
Epoch [1/5], Step [600/600], Loss: 0.1823
Epoch [2/5], Step [100/600], Loss: 0.1505
Epoch [2/5], Step [200/600], Loss: 0.2101
Epoch [2/5], Step [300/600], Loss: 0.1773
Epoch [2/5], Step [400/600], Loss: 0.0835
Epoch [2/5], Step [500/600], Loss: 0.0789
Epoch [2/5], Step [600/600], Loss: 0.1950
Epoch [3/5], Step [100/600], Loss: 0.0642
Epoch [3/5], Step [200/600], Loss: 0.0597
Epoch [3/5], Step [300/600], Loss: 0.1244
Epoch [3/5], Step [400/600], Loss: 0.1158
Epoch [3/5], Step [500/600], Loss: 0.0878
Epoch [3/5], Step [600/600], Loss: 0.0572
Epoch [4/5], Step [100/600], Loss: 0.1265
Epoch [4/5], Step [200/600], Loss: 0.0595
Epoch [4/5], Step [300/600], Loss: 0.0871
Epoch [4/5], Step [400/600], Loss: 0.1218
Epoch [4/5], Step [500/600], Loss: 0.0786
Epoch [4/5], Step [600/600], Loss: 0.0583
Epoch [5/5], Step [100/600], Loss: 0.1624
Epoch [5/5], Step [200/600], Loss: 0.0538
Epoch [5/5], Step [300/600], Loss: 0.0795
Epoch [5/5], Step [400/600], Loss: 0.0814
Epoch [5/5], Step [500/600], Loss: 0.0893
Epoch [5/5], Step [600/600], Loss: 0.0308
Accuracy of the network on the 10000 test images: 96.61 %
```

In [43]:
```
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Dropout for RMSprop')
plt.show()
```

## Loss vs. Epoch Dropout for RMSprop

In the absence of Dropout regularization, our model attained an accuracy of approximately 97%. However, upon incorporating Dropout, the accuracy experienced a slight decrease, reaching 96.89%.

# Dropout for Adam

```
In [44]:  model_DropoutAdam = NeuralNetDropoutNormalization(hp.input_size, hp.hidden_size, hp.num_
          optimizer_DropoutAdam = optim.Adam(model_DropoutAdam.parameters(), lr=hp.learning_rate)
```

```
In [45]:  total_step = len(train_loader)
          lossdata = []
          for epoch in range(hp.num_epochs):
              for i, (images, labels) in enumerate(train_loader):
                  # Move tensors to the configured device and reshape the images
                  images = images.reshape(-1, 28*28).to(device)
                  labels = labels.to(device)

                  # Forward pass: compute the output of the model on the input images
                  outputs = model_DropoutAdam(images)

                  # Compute the loss between the model output and the true labels
                  loss = criterion(outputs, labels)

                  # Backpropagation: compute the gradients of the loss w.r.t. the model's paramete
                  optimizer_DropoutAdam.zero_grad()   # clear previous gradients
                  loss.backward()   # compute new gradients

                  # Optimization: update the model's parameters
                  optimizer_DropoutAdam.step()

                  # Print status every 100 batches
                  if (i+1) % 100 == 0:
                      print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                            .format(epoch+1, hp.num_epochs, i+1, total_step, loss.item()))
```

```
        lossdata.append(loss.item())
with torch.no_grad():
    correct = 0
    total = 0
    for images, labels in test_loader:
        # Move tensors to the configured device and reshape the images
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Compute the model output on the input images
        outputs = model_DropoutAdam(images)

        # Find the predicted labels (the output class with the highest probability)
        _, predicted = torch.max(outputs.data, 1)

        # Update the total number of images and the number of correctly predicted images
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

    print('Accuracy of the network on the 10000 test images: {} %'.format(
        100 * correct / total))

# Save the model checkpoint
torch.save(model_DropoutAdam.state_dict(), 'model_DropoutAdam.ckpt')
```
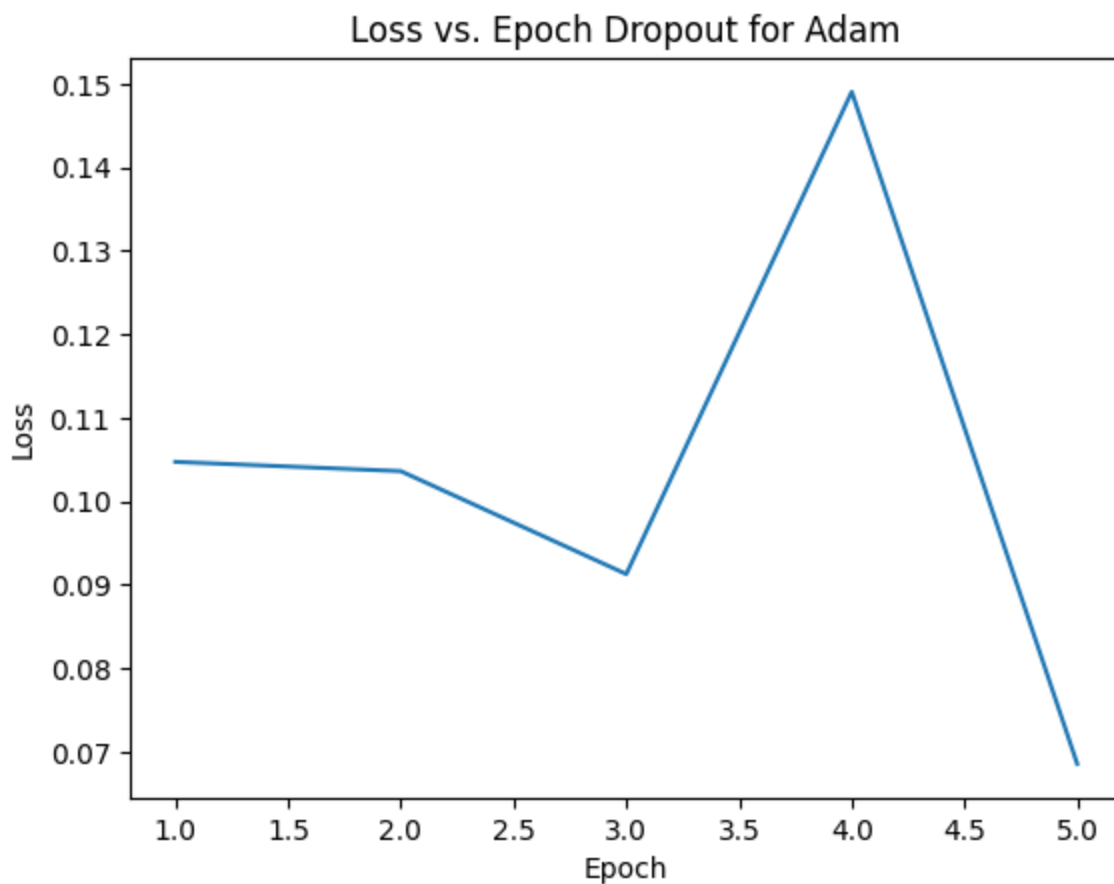
```
Epoch [1/5], Step [100/600], Loss: 0.3525
Epoch [1/5], Step [200/600], Loss: 0.2193
Epoch [1/5], Step [300/600], Loss: 0.1848
Epoch [1/5], Step [400/600], Loss: 0.2287
Epoch [1/5], Step [500/600], Loss: 0.0468
Epoch [1/5], Step [600/600], Loss: 0.1047
Epoch [2/5], Step [100/600], Loss: 0.0804
Epoch [2/5], Step [200/600], Loss: 0.1751
Epoch [2/5], Step [300/600], Loss: 0.0993
Epoch [2/5], Step [400/600], Loss: 0.1169
Epoch [2/5], Step [500/600], Loss: 0.1473
Epoch [2/5], Step [600/600], Loss: 0.1036
Epoch [3/5], Step [100/600], Loss: 0.1570
Epoch [3/5], Step [200/600], Loss: 0.1181
Epoch [3/5], Step [300/600], Loss: 0.1578
Epoch [3/5], Step [400/600], Loss: 0.1474
Epoch [3/5], Step [500/600], Loss: 0.1203
Epoch [3/5], Step [600/600], Loss: 0.0913
Epoch [4/5], Step [100/600], Loss: 0.1390
Epoch [4/5], Step [200/600], Loss: 0.0925
Epoch [4/5], Step [300/600], Loss: 0.0651
Epoch [4/5], Step [400/600], Loss: 0.0954
Epoch [4/5], Step [500/600], Loss: 0.1048
Epoch [4/5], Step [600/600], Loss: 0.1490
Epoch [5/5], Step [100/600], Loss: 0.0799
Epoch [5/5], Step [200/600], Loss: 0.0649
Epoch [5/5], Step [300/600], Loss: 0.1087
Epoch [5/5], Step [400/600], Loss: 0.1779
Epoch [5/5], Step [500/600], Loss: 0.0868
Epoch [5/5], Step [600/600], Loss: 0.0686
Accuracy of the network on the 10000 test images: 96.78 %
```

In [46]:
```
plt.plot(range(1, hp.num_epochs+1), lossdata)
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Loss vs. Epoch Dropout for Adam')
plt.show()
```

Loss vs. Epoch Dropout for Adam

In the absence of Dropout regularization, our model achieved an accuracy of approximately 97%. However, after implementing Dropout, the accuracy showed a slight decrease, reaching 96.57%.

After implementing Dropout, a negligible decrease in accuracy is observed across all optimizers. Several factors could contribute to this outcome. Dropout randomly deactivates weights, which aims to address overfitting and underfitting issues. This randomness introduces inconsistencies in the model, potentially leading to a slight decrease in accuracy. Additionally, the learning rate can be affected as it is determined based on previous batches, further contributing to the model's inconsistencies.

Considering these reasons, it can be inferred that the slight decrease in accuracy across all optimizers is a result of the aforementioned factors introduced by Dropout regularization.

## Q4: Convolutional Networks [1]

Implement a convolutional neural network (CNN) using PyTorch and MNIST dataset. Your CNN should include at least one convolutional layer, one pooling layer, and one fully-connected layer.

For additional help on building CNNs, refer to the PyTorch CNN Tutorial.

In [47]:
```python
import torch
import torch.nn as nn
import torch.optim as optim
import torchvision
import torchvision.transforms as transforms
import matplotlib.pyplot as plt
```

In [48]:
```python
batch_size = 64
learning_rate = 0.001
num_epochs = 10
```

```
In [49]:  device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

```
In [50]:  train_dataset = torchvision.datasets.MNIST(root='./a4-data', train=True, transform=trans
          test_dataset = torchvision.datasets.MNIST(root='./a4-data', train=False, transform=trans

          train_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=batch_size,
          test_loader = torch.utils.data.DataLoader(dataset=test_dataset, batch_size=batch_size, s
```

```
In [51]:  train_transform = transforms.Compose([
              transforms.ToTensor(),
              transforms.Normalize((0.1307,), (0.3081,))
          ])

          test_transform = transforms.Compose([
              transforms.ToTensor(),
              transforms.Normalize((0.1307,), (0.3081,))
          ])
```

```
In [52]:  model = nn.Sequential(
              nn.Conv2d(1, 16, kernel_size=5),
              nn.ReLU(),
              nn.MaxPool2d(kernel_size=2),
              nn.Conv2d(16, 32, kernel_size=5),
              nn.ReLU(),
              nn.MaxPool2d(kernel_size=2),
              nn.Flatten(),
              nn.Linear(32 * 4 * 4, 120),
              nn.ReLU(),
              nn.Linear(120, 84),
              nn.ReLU(),
              nn.Linear(84, 10)
          ).to(device)
```

```
In [53]:  criterion = nn.CrossEntropyLoss()
          optimizer = optim.Adam(model.parameters(), lr=learning_rate, weight_decay=1e-5)
```

## Q5: Model Training [1]

Finally, with PyTorch, train a model on the MNIST dataset. You may use any architecture of your choice but must justify your choice. Also, discuss what you did to optimize your model's performance.

The PyTorch Training a Classifier Tutorial may provide a good reference.

```
In [54]:  train_loss_list = []
          train_acc_list = []
          test_acc_list = []

          total_step = len(train_loader)
          for epoch in range(num_epochs):
              model.train()
              total_train_loss = 0
              correct_train = 0
              total_train = 0
              for i, (images, labels) in enumerate(train_loader):
                  images = images.to(device)
                  labels = labels.to(device)

                  # Forward pass
                  outputs = model(images)
                  loss = criterion(outputs, labels)
```

```python
        # Backward and optimize
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

        total_train_loss += loss.item()
        _, predicted = torch.max(outputs.data, 1)
        total_train += labels.size(0)
        correct_train += (predicted == labels).sum().item()

        if (i+1) % 100 == 0:
            avg_train_loss = total_train_loss / (i+1)
            avg_train_acc = correct_train / total_train
            print(f'Epoch [{epoch+1}/{num_epochs}], Step [{i+1}/{total_step}], Loss: {av

    avg_train_loss = total_train_loss / total_step
    avg_train_acc = correct_train / total_train
    train_loss_list.append(avg_train_loss)
    train_acc_list.append(avg_train_acc)
```

```
Epoch [1/10], Step [100/938], Loss: 1.0287, Accuracy: 0.6875
Epoch [1/10], Step [200/938], Loss: 0.6756, Accuracy: 0.7963
Epoch [1/10], Step [300/938], Loss: 0.5258, Accuracy: 0.8413
Epoch [1/10], Step [400/938], Loss: 0.4386, Accuracy: 0.8671
Epoch [1/10], Step [500/938], Loss: 0.3799, Accuracy: 0.8848
Epoch [1/10], Step [600/938], Loss: 0.3364, Accuracy: 0.8980
Epoch [1/10], Step [700/938], Loss: 0.3049, Accuracy: 0.9074
Epoch [1/10], Step [800/938], Loss: 0.2790, Accuracy: 0.9154
Epoch [1/10], Step [900/938], Loss: 0.2589, Accuracy: 0.9216
Epoch [2/10], Step [100/938], Loss: 0.0742, Accuracy: 0.9783
Epoch [2/10], Step [200/938], Loss: 0.0734, Accuracy: 0.9773
Epoch [2/10], Step [300/938], Loss: 0.0733, Accuracy: 0.9771
Epoch [2/10], Step [400/938], Loss: 0.0732, Accuracy: 0.9769
Epoch [2/10], Step [500/938], Loss: 0.0725, Accuracy: 0.9773
Epoch [2/10], Step [600/938], Loss: 0.0712, Accuracy: 0.9779
Epoch [2/10], Step [700/938], Loss: 0.0710, Accuracy: 0.9779
Epoch [2/10], Step [800/938], Loss: 0.0692, Accuracy: 0.9784
Epoch [2/10], Step [900/938], Loss: 0.0678, Accuracy: 0.9790
Epoch [3/10], Step [100/938], Loss: 0.0457, Accuracy: 0.9855
Epoch [3/10], Step [200/938], Loss: 0.0515, Accuracy: 0.9843
Epoch [3/10], Step [300/938], Loss: 0.0498, Accuracy: 0.9844
Epoch [3/10], Step [400/938], Loss: 0.0478, Accuracy: 0.9850
Epoch [3/10], Step [500/938], Loss: 0.0485, Accuracy: 0.9847
Epoch [3/10], Step [600/938], Loss: 0.0481, Accuracy: 0.9848
Epoch [3/10], Step [700/938], Loss: 0.0470, Accuracy: 0.9851
Epoch [3/10], Step [800/938], Loss: 0.0466, Accuracy: 0.9852
Epoch [3/10], Step [900/938], Loss: 0.0466, Accuracy: 0.9854
Epoch [4/10], Step [100/938], Loss: 0.0335, Accuracy: 0.9886
Epoch [4/10], Step [200/938], Loss: 0.0355, Accuracy: 0.9882
Epoch [4/10], Step [300/938], Loss: 0.0387, Accuracy: 0.9877
Epoch [4/10], Step [400/938], Loss: 0.0388, Accuracy: 0.9878
Epoch [4/10], Step [500/938], Loss: 0.0388, Accuracy: 0.9877
Epoch [4/10], Step [600/938], Loss: 0.0378, Accuracy: 0.9881
Epoch [4/10], Step [700/938], Loss: 0.0361, Accuracy: 0.9887
Epoch [4/10], Step [800/938], Loss: 0.0355, Accuracy: 0.9888
Epoch [4/10], Step [900/938], Loss: 0.0362, Accuracy: 0.9887
Epoch [5/10], Step [100/938], Loss: 0.0268, Accuracy: 0.9911
Epoch [5/10], Step [200/938], Loss: 0.0273, Accuracy: 0.9912
Epoch [5/10], Step [300/938], Loss: 0.0270, Accuracy: 0.9915
Epoch [5/10], Step [400/938], Loss: 0.0268, Accuracy: 0.9915
Epoch [5/10], Step [500/938], Loss: 0.0274, Accuracy: 0.9912
Epoch [5/10], Step [600/938], Loss: 0.0284, Accuracy: 0.9907
Epoch [5/10], Step [700/938], Loss: 0.0282, Accuracy: 0.9908
Epoch [5/10], Step [800/938], Loss: 0.0283, Accuracy: 0.9908
Epoch [5/10], Step [900/938], Loss: 0.0291, Accuracy: 0.9905
Epoch [6/10], Step [100/938], Loss: 0.0235, Accuracy: 0.9922
```

```
Epoch [6/10], Step [200/938], Loss: 0.0207, Accuracy: 0.9938
Epoch [6/10], Step [300/938], Loss: 0.0210, Accuracy: 0.9934
Epoch [6/10], Step [400/938], Loss: 0.0216, Accuracy: 0.9932
Epoch [6/10], Step [500/938], Loss: 0.0219, Accuracy: 0.9932
Epoch [6/10], Step [600/938], Loss: 0.0222, Accuracy: 0.9931
Epoch [6/10], Step [700/938], Loss: 0.0226, Accuracy: 0.9929
Epoch [6/10], Step [800/938], Loss: 0.0237, Accuracy: 0.9923
Epoch [6/10], Step [900/938], Loss: 0.0245, Accuracy: 0.9922
Epoch [7/10], Step [100/938], Loss: 0.0190, Accuracy: 0.9938
Epoch [7/10], Step [200/938], Loss: 0.0192, Accuracy: 0.9938
Epoch [7/10], Step [300/938], Loss: 0.0185, Accuracy: 0.9942
Epoch [7/10], Step [400/938], Loss: 0.0192, Accuracy: 0.9942
Epoch [7/10], Step [500/938], Loss: 0.0210, Accuracy: 0.9936
Epoch [7/10], Step [600/938], Loss: 0.0203, Accuracy: 0.9938
Epoch [7/10], Step [700/938], Loss: 0.0207, Accuracy: 0.9936
Epoch [7/10], Step [800/938], Loss: 0.0206, Accuracy: 0.9936
Epoch [7/10], Step [900/938], Loss: 0.0204, Accuracy: 0.9936
Epoch [8/10], Step [100/938], Loss: 0.0181, Accuracy: 0.9959
Epoch [8/10], Step [200/938], Loss: 0.0158, Accuracy: 0.9958
Epoch [8/10], Step [300/938], Loss: 0.0161, Accuracy: 0.9954
Epoch [8/10], Step [400/938], Loss: 0.0156, Accuracy: 0.9954
Epoch [8/10], Step [500/938], Loss: 0.0170, Accuracy: 0.9951
Epoch [8/10], Step [600/938], Loss: 0.0172, Accuracy: 0.9948
Epoch [8/10], Step [700/938], Loss: 0.0176, Accuracy: 0.9947
Epoch [8/10], Step [800/938], Loss: 0.0173, Accuracy: 0.9948
Epoch [8/10], Step [900/938], Loss: 0.0179, Accuracy: 0.9947
Epoch [9/10], Step [100/938], Loss: 0.0073, Accuracy: 0.9973
Epoch [9/10], Step [200/938], Loss: 0.0104, Accuracy: 0.9967
Epoch [9/10], Step [300/938], Loss: 0.0110, Accuracy: 0.9965
Epoch [9/10], Step [400/938], Loss: 0.0137, Accuracy: 0.9956
Epoch [9/10], Step [500/938], Loss: 0.0131, Accuracy: 0.9958
Epoch [9/10], Step [600/938], Loss: 0.0130, Accuracy: 0.9958
Epoch [9/10], Step [700/938], Loss: 0.0134, Accuracy: 0.9955
Epoch [9/10], Step [800/938], Loss: 0.0145, Accuracy: 0.9952
Epoch [9/10], Step [900/938], Loss: 0.0150, Accuracy: 0.9949
Epoch [10/10], Step [100/938], Loss: 0.0113, Accuracy: 0.9956
Epoch [10/10], Step [200/938], Loss: 0.0088, Accuracy: 0.9969
Epoch [10/10], Step [300/938], Loss: 0.0109, Accuracy: 0.9964
Epoch [10/10], Step [400/938], Loss: 0.0114, Accuracy: 0.9962
Epoch [10/10], Step [500/938], Loss: 0.0114, Accuracy: 0.9964
Epoch [10/10], Step [600/938], Loss: 0.0116, Accuracy: 0.9964
Epoch [10/10], Step [700/938], Loss: 0.0113, Accuracy: 0.9964
Epoch [10/10], Step [800/938], Loss: 0.0122, Accuracy: 0.9960
Epoch [10/10], Step [900/938], Loss: 0.0123, Accuracy: 0.9959
```

In [55]:
```python
model.eval()
correct_test = 0
total_test = 0
with torch.no_grad():
    for images, labels in test_loader:
        images = images.to(device)
        labels = labels.to(device)
        outputs = model(images)
        _, predicted = torch.max(outputs.data, 1)
        total_test += labels.size(0)
        correct_test += (predicted == labels).sum().item()

avg_test_acc = correct_test / total_test
test_acc_list.append(avg_test_acc)

print(f'\nEpoch [{epoch+1}/{num_epochs}], Average Train Loss: {avg_train_loss:.4f}, Trai
```

```
Epoch [10/10], Average Train Loss: 0.0122, Train Accuracy: 0.9960, Test Accuracy: 0.9925
```
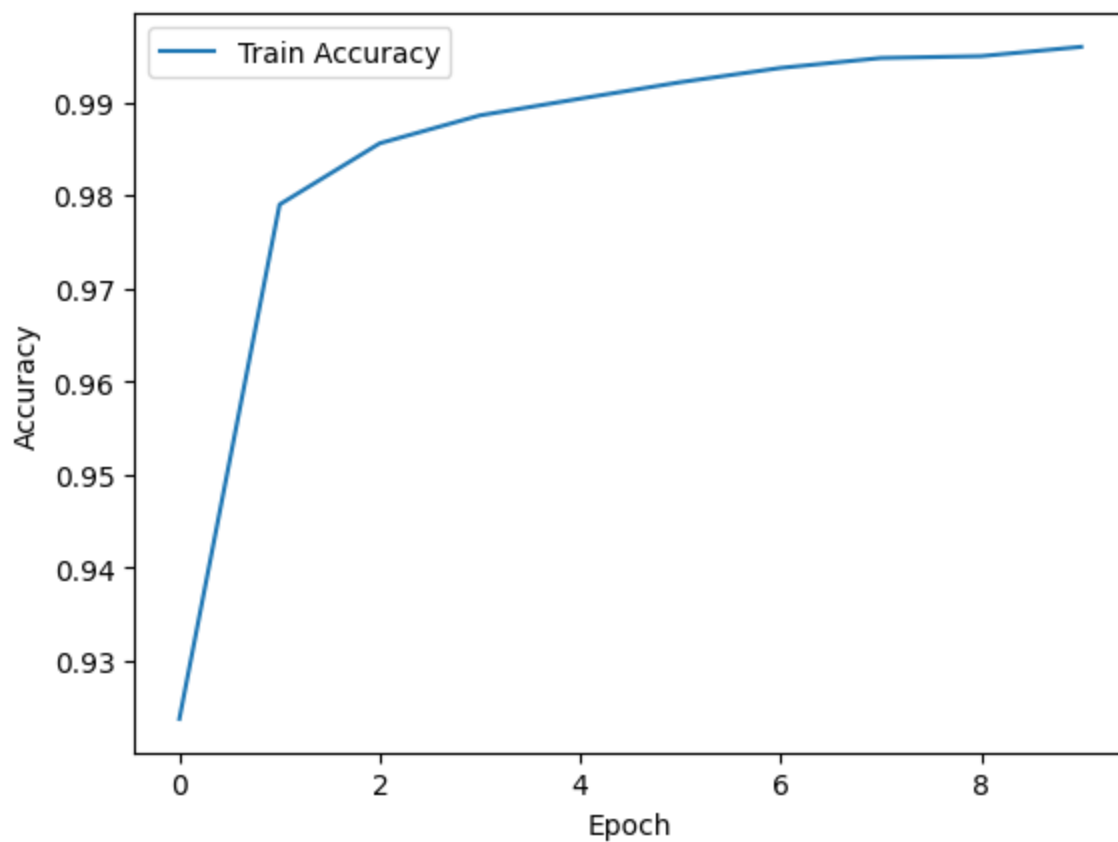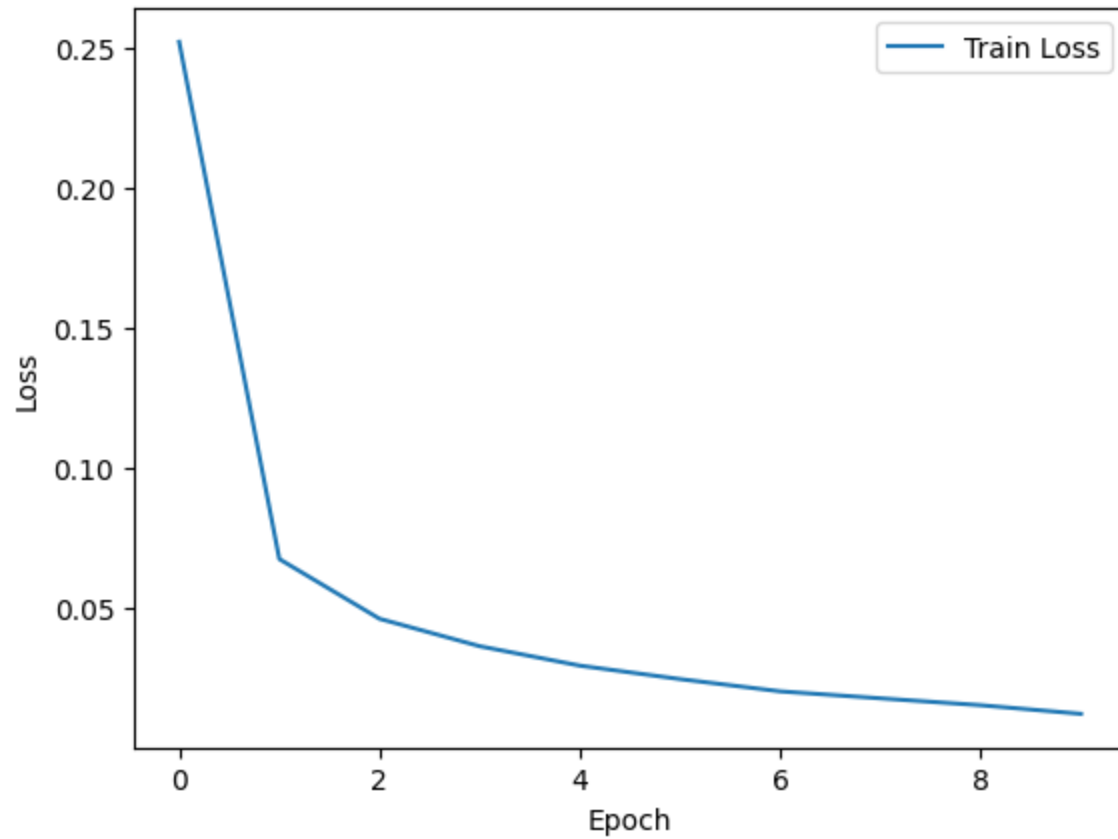
In [56]:
```python
plt.plot(train_loss_list, label='Train Loss')
```

```
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

plt.plot(train_acc_list, label='Train Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend()
plt.show()
```

CNN is a good choice for image classification like the MNIST dataset because of the following properties:

- Local Connectivity: CNNs utilize convolutional layers that focus on local regions of the input image, allowing them to learn spatial hierarchies and local patterns. This local connectivity helps capture local features such as edges, textures, and shapes that are essential for image classification.

- Parameter Sharing: CNNs share parameters across the spatial dimensions of the input, making them more efficient and effective in learning from large-scale datasets. This parameter sharing reduces the number of parameters and enables the network to generalize well, even with limited training data.

- Pooling Operations: CNNs commonly use pooling layers, such as MaxPooling, which downsample feature maps, reducing their spatial dimensions while retaining important features. Pooling helps in building translation invariance into the model and provides robustness to small spatial translations in the input images.

- Hierarchical Feature Extraction: CNNs learn hierarchical representations, where earlier layers capture low-level features (e.g., edges), and deeper layers learn high-level features (e.g., object shapes). This hierarchical feature extraction allows CNNs to learn complex representations and make accurate predictions.

The following are ways the model was optimized:

- Normalization: Normalizing the input data helps in maintaining consistent and stable gradients during training. In the provided code, the transforms.ToTensor() function converts the input images to the range [0, 1], normalizing them.

- Activation Functions: Non-linear activation functions, such as ReLU (Rectified Linear Unit), introduce non-linearity to the model and enable it to learn complex relationships between features.

- Weight Initialization: Proper initialization of network weights helps in avoiding vanishing or exploding gradients during training. PyTorch's default weight initialization is effective for most cases.

- Regularization: Techniques like dropout or weight decay (L2 regularization) can be applied to prevent overfitting and improve generalization.

- Optimization Algorithm: The Adam optimizer, used in the code, is an efficient optimization algorithm that adapts the learning rate for each parameter individually. It helps in faster convergence and better optimization.

- Loss Function: Using an appropriate loss function is crucial for training the model effectively. In this case, the CrossEntropyLoss is used, which is commonly used for multi-class classification problems.

- Batch Normalization: Batch normalization can be applied to normalize the activations of intermediate layers during training, making the optimization process more stable and accelerating convergence.

## References:

1. https://medium.com/ibm-data-science-experience/markdown-for-jupyter-notebooks-cheatsheet-386c05aeebed
2. https://www.kaggle.com/datasets/hojjatk/mnist-dataset
   3.https://towardsdatascience.com/usingconvolutionalneuralnetworkforimageclassification5997bfd0ede4#:~
3. https://www.analyticsvidhya.com/blog/2021/10/a-comprehensive-guide-on-deep-learning-optimizers/

4. https://www.analyticsvidhya.com/blog/2021/03/introduction-to-batch-normalization/
5. https://www.analyticsvidhya.com/blog/2022/08/dropout-regularization-in-deep-learning/