

Homework 7

Harvard University

Fall 2018

Instructors: Rahul Dave

Due Date: Saturday, October 27th, 2018 at 11:59pm

Instructions:

- Upload your final answers in the form of a Jupyter notebook containing all work to Canvas.
- Structure your notebook and your work to maximize readability.

Collaborators

Michelle (Chia Chi) Ho, Jiejun Lu, Jiawen Tong

```
In [1]: 1 import numpy as np
        2 import scipy.stats
        3 import scipy.special
        4
        5 import matplotlib
        6 import matplotlib.pyplot as plt
        7 import matplotlib.mlab as mlab
        8 from matplotlib import cm
        9 import pandas as pd
       10 %matplotlib inline
       11
       12 import torch
       13 import torchvision.datasets as datasets
       14 import torchvision.transforms as transforms
       15 from torch.autograd import Variable
       16 import torch.nn as nn
       17 from torch.utils.data.sampler import SubsetRandomSampler
```

Question 1: Mon pays c'est l'MNIST. Mon cœur est brisé de Logistic Regression.

The [MNIST dataset](https://en.wikipedia.org/wiki/MNIST_database) (https://en.wikipedia.org/wiki/MNIST_database) is one of the classic datasets in Machine Learning and is often one of the first datasets against which new classification algorithms test themselves. It consists of 70,000 images of handwritten digits, each of which is 28x28 pixels. You will be using PyTorch to build a handwritten digit classifier that you will train, validate, and test with MNIST.

Your classifier MUST implement a multinomial logistic regression model (using softmax). It will take as input an array of pixel values in an image and output the images most likely digit label (i.e. 0-9). You should think of the pixel values as features of the input vector.

Using the softmax formulation, your PyTorch model should compute the cost function using an L2 regularization approach (see `optim.SGD` in PyTorch or write your own cost function) and minimize the resulting cost function using mini-batch stochastic gradient descent. We provided extensive template code in lab.

Construct and train your classifier using a batch size of 256 examples, a learning rate $\eta=0.1$, and a regularization factor $\lambda=0.01$.

1.1. Plot 10 sample images from the MNIST dataset (to develop intuition for the feature space).

1.2. Currently the MNIST dataset in Torchvision allows a Train/Test split. Use PyTorch dataloader functionality to create a Train/Validate/Test split of 50K/10K/10K samples.

Hint: Lab described a way to do it keeping within the MNIST `DataLoader` workflow: the key is to pass a `SubsetRandomSampler` to `DataLoader`

1.3. Construct a softmax formulation in PyTorch of multinomial logistic regression with Cross Entropy Loss.

1.4. Train your model using SGD to minimize the cost function. Use as many epochs as you need to achieve convergence.

1.5. Plot the cross-entropy loss on the training set as a function of iteration.

1.6. Using classification accuracy, evaluate how well your model is performing on the validation set at the end of each epoch. Plot this validation accuracy as the model trains.

1.6. Duplicate this plot for some other values of the regularization parameter λ . When should you stop the training for each of the different values of λ ? Give an approximate answer supported by using the plots.

1.7. Select what you consider the best regularization parameter and predict the labels of the test set. Compare your predictions with the given labels. What classification accuracy do you obtain on the training and test sets?

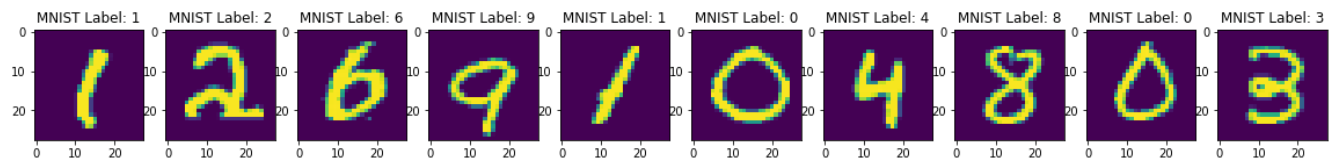
1.8. What classes are most likely to be misclassified? Plot some misclassified training and test set images.

Gratuitous Titular Reference: The recently departed French rockstar Johnny Hallyday just posthumously released what looks to be his biggest album ever "Mon pays c'est l'amour". The album sold 300,000 copies on its first day of release.

```
In [2]: 1 # Download/load MNIST data
2 train_dataset = datasets.MNIST(root='./mnist_data', train=True,
3                               transform=transforms.Compose([transforms.ToTensor(),
4                                                             transforms.Normalize((0.1307,), (0.3081,))]),
5                               download=True)
6 test_dataset = datasets.MNIST(root='./mnist_data', train=False,
7                               transform=transforms.Compose([transforms.ToTensor(),
8                                                             transforms.Normalize((0.1307,), (0.3081,))]),
9                               download=True)
```

```
In [3]: 1 # 1.1
2 n_samples = 10
3 sample_indices = np.random.choice(train_dataset.train_data.size(0), n_samples)
4
5 sample_images = train_dataset.train_data[sample_indices,:].numpy()
6 sample_labels = [train_dataset.train_labels[x] for x in sample_indices]
7
8 # plot sample images
9 fig, ax = plt.subplots(1, n_samples, figsize=(20, 4))
10 plt.suptitle("MNIST Sample Images", fontsize=20, weight='heavy')
11 for i in range(n_samples):
12     ax[i].imshow(sample_images[i])
13     ax[i].set_title("MNIST Label: {}".format(sample_labels[i], weight='bold'))
14
15 plt.show()
```

MNIST Sample Images



```
In [4]: 1 # Regression Parent Class -- copied from lab
2 class Regression(object):
3     def __init__(self):
4         self.params = dict()
5
6     def get_params(self, k):
7         return self.params.get(k, None)
8
9     def set_params(self, **kwargs):
10        for k,v in kwargs.items():
11            self.params[k] = v
12
13    def fit(self, X, y):
14        raise NotImplementedError()
15
16    def predict(self, X):
17        raise NotImplementedError()
18
19    def score(self, X, y):
20        raise NotImplementedError()
21
22
23 class LRPyTorch(nn.Module): # PyTorch implementation of Logistic Regression -- copied from lab
24     def __init__(self):
25         super().__init__()
26         self.l1 = nn.Linear(784, 10)
27
28     def forward(self, x):
29         x = self.l1(x)
30         return x
31
```

```

In [5]: 1 class MNIST_Classifier(Regression): # MNIST Classifier extends Regression
2     def __init__(self, torch_model, learning_rate, batch_size, regularization, epochs):
3         super().__init__()
4
5         # define model, loss, optimizer
6         model = torch_model
7         criterion = nn.CrossEntropyLoss() # Softmax = Cross Entropy
8         optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate, weight_decay=regularization)
9
10        # Initialize the MNIST classifier's params
11        self.set_params(optimizer=optimizer,
12                        learning_rate=learning_rate,
13                        batch_size=batch_size,
14                        model=model,
15                        criterion=criterion,
16                        epochs=epochs)
17
18    def predict(self, loader, dataset_type, save_misclassified=False):
19        # dataset_type = 'Test' for test set; otherwise, load train_label or train_data
20        dataset_labels = loader.dataset.test_labels if dataset_type == 'Test' else loader.dataset.train_labels
21        dataset_images = loader.dataset.test_data if dataset_type == 'Test' else loader.dataset.train_data
22        dataset_labels = dataset_labels.numpy()
23
24        predictions = []
25        correct_count = 0
26        model = self.get_params('model')
27
28        y_true = [] # labels enumerated
29        X = [] # images enumerated
30        for inputs, labels in loader:
31            # forward pass
32            inputs = Variable(inputs)
33            inputs = inputs.view(-1, 28*28)
34            outputs = model(inputs)
35
36            # count correct
37            pred = outputs.data.max(1)[1] # predicted class label, one-hot format
38            correct_count += (pred == labels).sum()
39            predictions += list(pred)
40
41            # concatenate labels & images
42            y_true += list(labels)
43            X += list(inputs.numpy())
44
45        predictions = np.array(predictions)
46        y_true = np.array(y_true)
47        X = np.array(X).reshape((len(X), 28, 28))
48
49        # save 10 misclassified sample images
50        if save_misclassified:
51            mis_X = X[predictions != y_true]
52            mis_preds = predictions[predictions != y_true]
53            mis_y_true = y_true[predictions != y_true]
54            self.save_misclassified(mis_X, mis_preds, mis_y_true)
55
56        return predictions, y_true, correct_count
57
58    def save_misclassified(self, mis_X, mis_preds, mis_y_true):
59        # save all misclassified labels & preds
60        self.set_params(all_mis_y_true = mis_y_true)
61        self.set_params(all_mis_preds = mis_preds)
62
63        # sample 10 misclassified images X
64        sample_indices = np.random.choice(np.arange(mis_X.shape[0]), 10)
65
66        # save to params the sampled misclassified X, pred/true label
67        self.set_params(misclassified_images = mis_X[sample_indices, :, :])
68        self.set_params(misclassified_labels = mis_preds[sample_indices])
69        self.set_params(misclassified_true_labels = mis_y_true[sample_indices])
70
71    def viz_misclassified(self):
72        # get the images and labels
73        sample_images = self.get_params("misclassified_images")
74        sample_labels = self.get_params("misclassified_labels")
75        true_labels = self.get_params("misclassified_true_labels")
76
77        fig, ax = plt.subplots(1, 10, figsize=(20, 4))
78        plt.suptitle("Sample Misclassified Images", fontsize=20, weight='heavy')
79        for i in range(10):
80            ax[i].imshow(sample_images[i])
81            ax[i].set_title("True Label: {}\n Classified: {}".format(true_labels[i], sample_labels[i]))
82        plt.show()
83
84    def score(self, loader, dataset_type, data_size, save_misclassified=False, verbose=False):
85        _, _, correct_count = self.predict(loader, dataset_type, save_misclassified)
86        if verbose:
87            print('On {} set: Accuracy: {}/{ } ({:.1f}%)\n'.format(dataset_type, correct_count, data_size,
88                                                                100.0*float(correct_count)/data_size))
89        return (float(correct_count)/data_size)
90
91    def fit(self, train_loader, train_size, validation_loader, validation_size):
92        self.set_params(n_iter=int(np.ceil(train_size/self.get_params('batch_size'))))
93
94        # reclaim parameters
95        optimizer = self.get_params('optimizer')

```

```

96     model = self.get_params('model')
97     epochs = self.get_params('epochs')
98     criterion = self.get_params('criterion')
99
100     losses = []
101     val_scores = [] # validation score at the end of each epoch
102     for epoch in range(epochs):
103         for batch_index, (inputs, labels) in enumerate(train_loader): # loop thru all batches per epoch
104             # forward pass
105             inputs, labels = Variable(inputs), Variable(labels)
106             inputs = inputs.view(-1, 28*28)
107             optimizer.zero_grad()
108             outputs = model(inputs)
109
110             # record loss & backward pass
111             loss = criterion(outputs, labels)
112             losses.append(loss.data[0].numpy().reshape([1])[0])
113             loss.backward()
114             optimizer.step()
115
116         val_score = self.score(validation_loader, 'Validation', validation_size)
117         val_scores.append(val_score)
118
119         print('{} epoch fitted'.format(epoch+1), end='\r')
120
121     self.set_params(training_losses=losses) # set params: 'training_losses'
122     self.set_params(validation_scores=val_scores) # validation score by the end of each epoch
123
124     return self
125
126     def viz_training_loss(self):
127         losses = self.get_params('training_losses')
128         n_iter = self.get_params('n_iter')
129         epochs = self.get_params('epochs')
130         fig, axes = plt.subplots(nrows=1, ncols=epochs, figsize=(20, 5), sharex=True, sharey=True)
131         for i in range(epochs):
132             if i == epochs:
133                 axes[i].plot(range(len(losses[i*n_iter:])), losses[i*n_iter:])
134             else:
135                 axes[i].plot(range(n_iter), losses[i*n_iter:(i+1)*n_iter])
136                 axes[i].set_title('ep{}'.format(i))
137                 if i % 2 == 1:
138                     axes[i].axvspan(-10, n_iter, facecolor='gray', alpha=0.2)
139         plt.subplots_adjust(wspace=0)
140         plt.show()

```

```

In [6]: 1 # 1.2
2 # split out validation set
3 indices = list(range(len(train_dataset)))
4 val_split = 10000
5 batch_size = 256
6 lr = LRPyTorch()
7
8 # Random, non-contiguous split
9 validation_idx = np.random.choice(indices, size=val_split, replace=False)
10 train_idx = np.array(list(set(indices) - set(validation_idx)))
11 train_size, validation_size, test_size = train_idx.shape[0], validation_idx.shape[0], len(test_dataset)
12 print('Train size = {}, Validation size = {}, Test size = {}'.format(train_size,
13                             validation_size,
14                             test_size))
15
16 # train & validation samplers, test set does not require shuffling
17 train_sampler = SubsetRandomSampler(train_idx)
18 validation_sampler = SubsetRandomSampler(validation_idx)
19
20 # train, validation & test loader
21 train_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=batch_size, sampler=train_sampler)
22 validation_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=validation_size, sampler=validation_sampler)
23 test_loader = torch.utils.data.DataLoader(dataset=test_dataset, batch_size=test_size, shuffle=False)

```

Train size = 50000, Validation size = 10000, Test size = 10000

```

In [7]: 1 # 1.3, 1.4
2 # fit the classifier
3 mnist_clf = MNIST_Classifier(lr, learning_rate=0.1, batch_size=batch_size, regularization=0.01, epochs=30)
4 mnist_clf = mnist_clf.fit(train_loader, train_size, validation_loader, validation_size)

/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:112: UserWarning: invalid index of a 0-dim tensor. This will be an error in PyTorch 0.5. Use tensor.item() to convert a 0-dim tensor to a Python number

30 epoch fitted

```

```

In [8]: 1 # print scores
2 _ = mnist_clf.score(train_loader, 'Train', train_size, save_misclassified=False, verbose=True)
3 _ = mnist_clf.score(validation_loader, 'Validation', validation_size, save_misclassified=False, verbose=True)
4 _ = mnist_clf.score(test_loader, 'Test', test_size, save_misclassified=False, verbose=True)

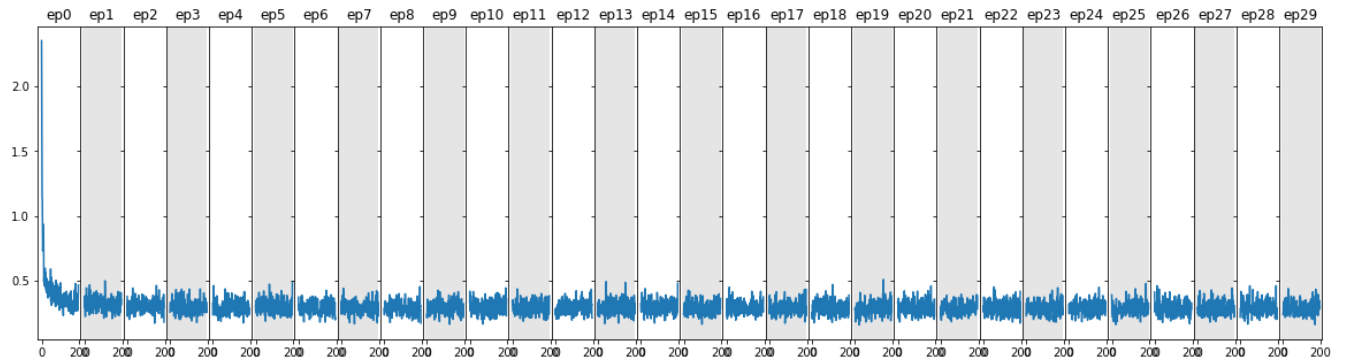
```

On Train set: Accuracy: 45704/50000 (91.4%)

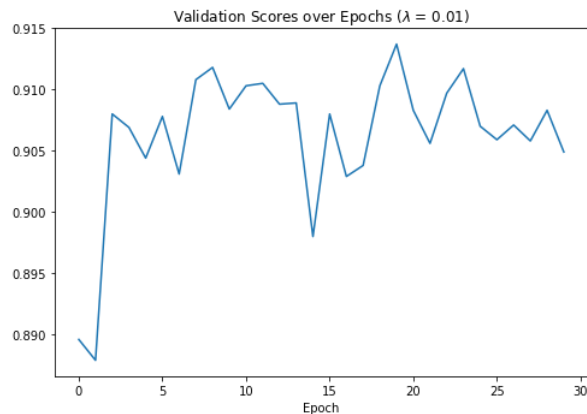
On Validation set: Accuracy: 9049/10000 (90.5%)

On Test set: Accuracy: 9122/10000 (91.2%)

```
In [9]: 1 # 1.5
2 # plot training loss
3 mnist_clf.viz_training_loss()
```



```
In [10]: 1 # 1.6
2 # plot validation score over epochs - lambda = 0.01
3 validation_scores = mnist_clf.get_params('validation_scores')
4 fig, ax = plt.subplots(1, 1, figsize=(7, 5))
5 pd.Series(validation_scores).plot(title=r'Validation Scores over Epochs ($\lambda$ = 0.01)', ax=ax)
6 ax.set_xlabel('Epoch')
7 plt.tight_layout()
```



```
In [18]: 1 # 1.6
2 # Different regularization lambdas
3 reg_params = np.array([1, 0.1, 0.01, 0.001, 0.0001, 0.00001])
4 val_scores_arr = []
5 final_val_score_arr = []
6
7 for reg in reg_params:
8     print('lambda = {}'.format(reg))
9     # fit a classifier for this lambda
10    _clf = MNIST_Classifier(lr, learning_rate=0.1, batch_size=batch_size, regularization=reg, epochs=30)
11    _clf.fit(train_loader, train_size, validation_loader, validation_size)
12
13    # record the validation scores over epochs of this classifier
14    val_scores_arr.append(_clf.get_params('validation_scores'))
15
16    # record the final validation score of this classifier
17    final_val_score_arr.append(_clf.score(validation_loader, 'Validation', validation_size))
18
```

lambda = 1.0

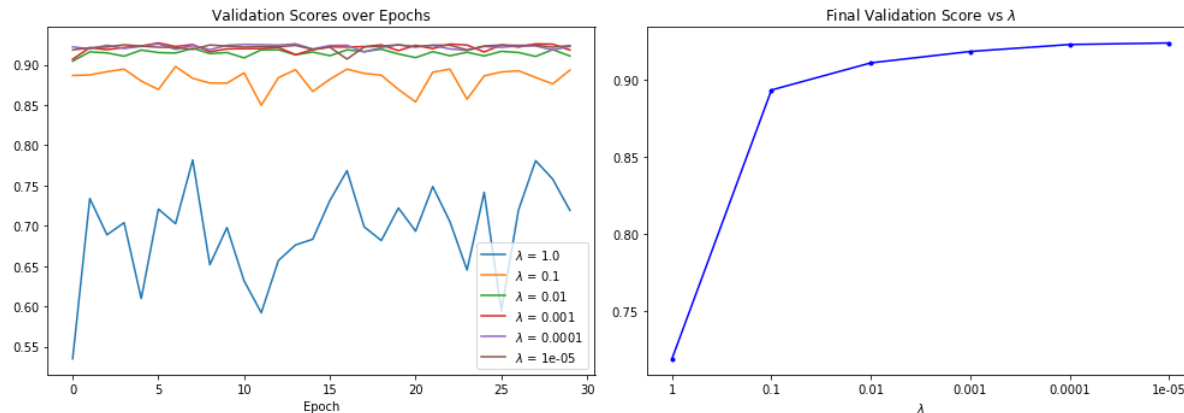
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:108: UserWarning: invalid index of a 0-dim tensor. This will be an error in PyTorch 0.5. Use tensor.item() to convert a 0-dim tensor to a Python number

lambda = 0.1ted
lambda = 0.01ed
lambda = 0.001d
lambda = 0.0001
lambda = 1e-05d
30 epoch fitted

```

In [24]: 1 # 1.7
2 # find the best regularization parameter using the validation set
3
4 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
5 for i, val_scores in enumerate(val_scores_arr):
6     # plot validation scores over epochs for this classifier
7     axes[0].plot(_val_scores, label=r'$\lambda$ = {}'.format(reg_params[i]))
8
9 axes[0].set_xlabel('Epoch')
10 axes[0].set_title('Validation Scores over Epochs')
11 axes[0].legend()
12
13 axes[1].plot(final_val_score_arr, '-.', c='b')
14 axes[1].set_xlabel(r'$\lambda$')
15 axes[1].set_title(r'Final Validation Score vs $\lambda$')
16 axes[1].set_xticklabels([-1, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001])
17
18 plt.tight_layout()

```



Answer 1.7

When to stop the training for each of the different values of λ ?

- Based on the plots of **validation scores over epochs** with different regularization parameter λ 's, all classifiers with $\lambda \leq 0.1$ converged relatively quickly (within the first 5 epochs). With $\lambda = 1$, the classifier had poorer validation performance with larger fluctuation.
- Based on the plot of **final validation accuracies vs different λ 's**, the smaller the regularization, the better validation performance we can get. But the performance improvements after $\lambda = 0.001$ became less obvious. In this case, we can see overfitting is not a major problem of logistic regression modelling. Validation/test set has almost the same level of accuracy as training set.

```

In [7]: 1 # fit a clf with the best selected lambda
2 best_lr_clf = MNIST_Classifier(lr, learning_rate=0.1, batch_size=batch_size, regularization=0.001, epochs=30)
3 best_lr_clf.fit(train_loader, train_size, validation_loader, validation_size)
4
5 # predict on the training & test set
6 train_acc = best_lr_clf.score(train_loader, 'Train', train_size, save_misclassified=False, verbose=True)
7 test_acc = best_lr_clf.score(test_loader, 'Test', test_size, save_misclassified=False, verbose=True)

```

/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:112: UserWarning: invalid index of a 0-dim tensor. This will be an error in PyTorch 0.5. Use tensor.item() to convert a 0-dim tensor to a Python number

On Train set: Accuracy: 46437/50000 (92.9%)

On Test set: Accuracy: 9231/10000 (92.3%)

Answer 1.8 - LR misclassification on train set

```

In [8]: 1 # 1.8 - Training set
2 train_preds, train_labels, train_correct_count = best_lr_clf.predict(train_loader, 'Train',
3                                                                    save_misclassified=True)
4
5 # misclassification counts
6 df_mis_train = pd.DataFrame(pd.Series(train_labels[train_labels != train_preds], name='digit'))
7 df_mis_train['count'] = np.ones((len(df_mis_train),)).astype(int)
8 df_mis_train = df_mis_train.groupby(['
9     'digit'])
10 )['count'].agg([sum]).rename(columns={'sum': 'mis_count'}).reset_index()
11
12 # digit counts
13 df_digit_train = pd.DataFrame(pd.Series(train_labels, name='digit'))
14 df_digit_train['count'] = np.ones((len(df_digit_train),)).astype(int)
15 df_digit_train = df_digit_train.groupby(['
16     'digit'])
17 )['count'].agg([sum]).rename(columns={'sum': 'count'}).reset_index()
18
19 df_train_res = df_digit_train.merge(df_mis_train, on='digit', how='inner').set_index('digit')
20 df_train_res['mis_percent'] = df_train_res['mis_count'] / df_train_res['count'] * 100
21 df_train_res.sort_values(['mis_percent'], ascending=False)

```

```

Out[8]:
   count  mis_count  mis_percent
digit
5    4509         604    13.395431
3    5164         658    12.742060
2    4955         442     8.920283
9    4959         417     8.408953
7    5163         339     6.565950
8    4893         289     5.906397
4    4911         281     5.721849
6    4906         211     4.300856
1    5615         192     3.419412
0    4925         130     2.639594

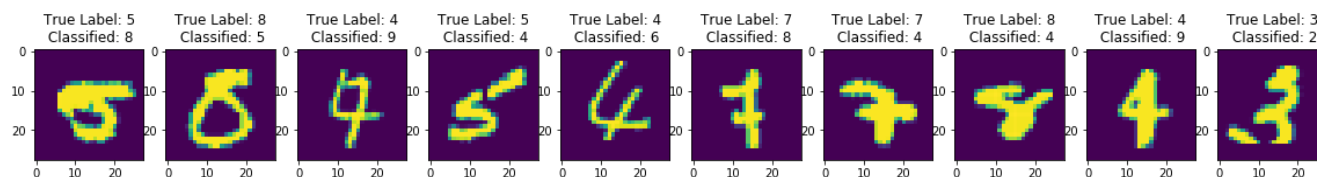
```

```

In [9]: 1 best_lr_clf.viz_misclassified()

```

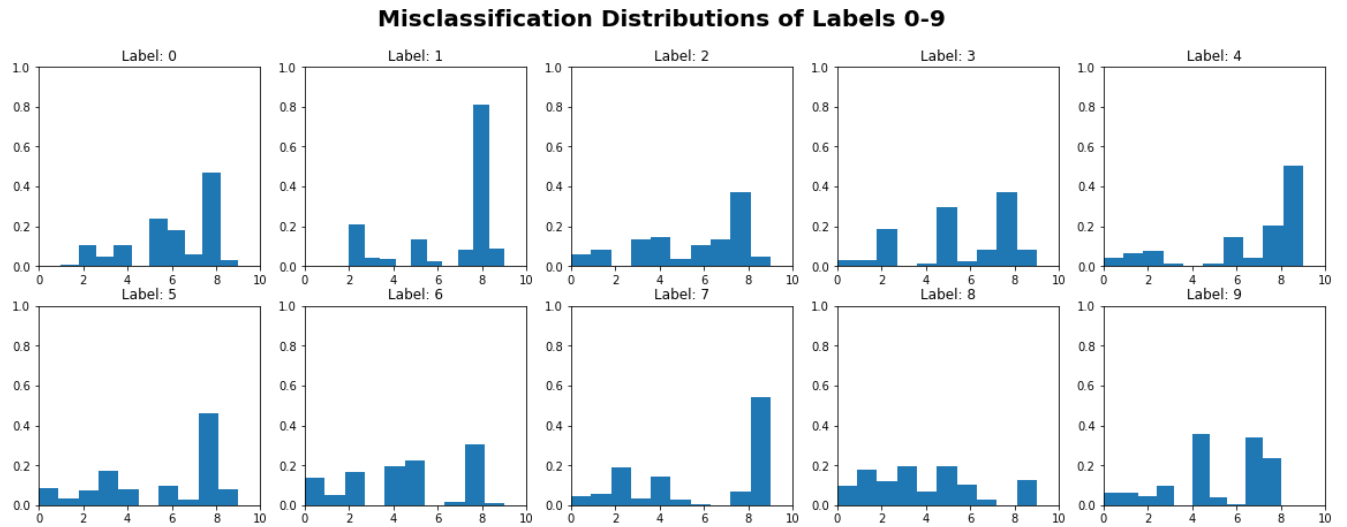
Sample Misclassified Images



```

In [10]: 1 mis_dist = [[] for i in range(10)]
2         for _i, digit in enumerate(best_lr_clf.get_params('all_mis_y_true')):
3             mis_dist[int(digit)].append(best_lr_clf.get_params('all_mis_preds')[_i])
4
5         fig, (ax1, ax2) = plt.subplots(2, 5, figsize=(20, 7))
6         plt.suptitle('Misclassification Distributions of Labels 0-9', fontsize=20, weight='heavy')
7
8         for i in range(5):
9             ax1[i].hist(mis_dist[i], normed = True)
10            ax1[i].set_title('Label: {}'.format(i))
11            ax1[i].set_xlim(0, 10)
12            ax1[i].set_ylim(0, 1)
13            ax2[i].hist(mis_dist[i + 5], normed = True)
14            ax2[i].set_title('Label: {}'.format(i+5))
15            ax2[i].set_xlim(0, 10)
16            ax2[i].set_ylim(0, 1)

```



Answer 1.8 - LR misclassification on test set


```

In [11]: 1 # 1.8 - Test set
2 test_preds, test_labels, test_correct_count = best_lr_clf.predict(test_loader, 'Test', save_misclassified=True)
3
4 # misclassification counts
5 df_mis_test = pd.DataFrame(pd.Series(test_labels[test_labels != test_preds], name='digit'))
6 df_mis_test['count'] = np.ones((len(df_mis_test),)).astype(int)
7 df_mis_test = df_mis_test.groupby([
8     'digit'
9 ])['count'].agg([sum]).rename(columns={'sum': 'mis_count'}).reset_index()
10
11 # digit counts
12 df_digit_test = pd.DataFrame(pd.Series(test_labels, name='digit'))
13 df_digit_test['count'] = np.ones((len(df_digit_test),)).astype(int)
14 df_digit_test = df_digit_test.groupby([
15     'digit'
16 ])['count'].agg([sum]).rename(columns={'sum': 'count'}).reset_index()
17
18 df_test_res = df_digit_test.merge(df_mis_test, on='digit', how='inner').set_index('digit')
19 df_test_res['mis_percent'] = df_test_res['mis_count'] / df_test_res['count'] * 100
20 df_test_res.sort_values(['mis_percent'], ascending=False)

```

```

Out[11]:
   count  mis_count  mis_percent
digit
5      892         127    14.237668
3     1010         118    11.683168
2     1032         112    10.852713
9     1009          89     8.820614
7     1028          86     8.365759
8      974          71     7.289528
4      982          65     6.619145
6      958          48     5.010438
1     1135          35     3.083700
0      980          18     1.836735

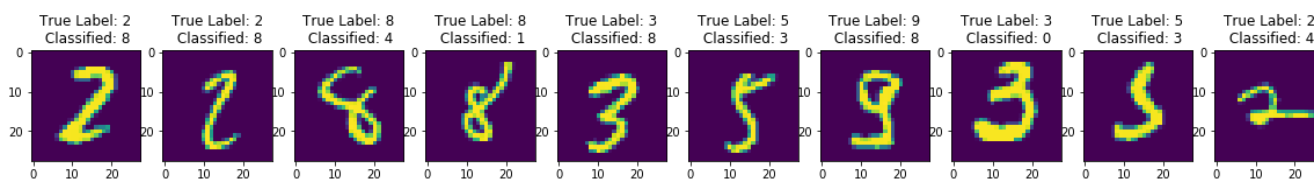
```

```

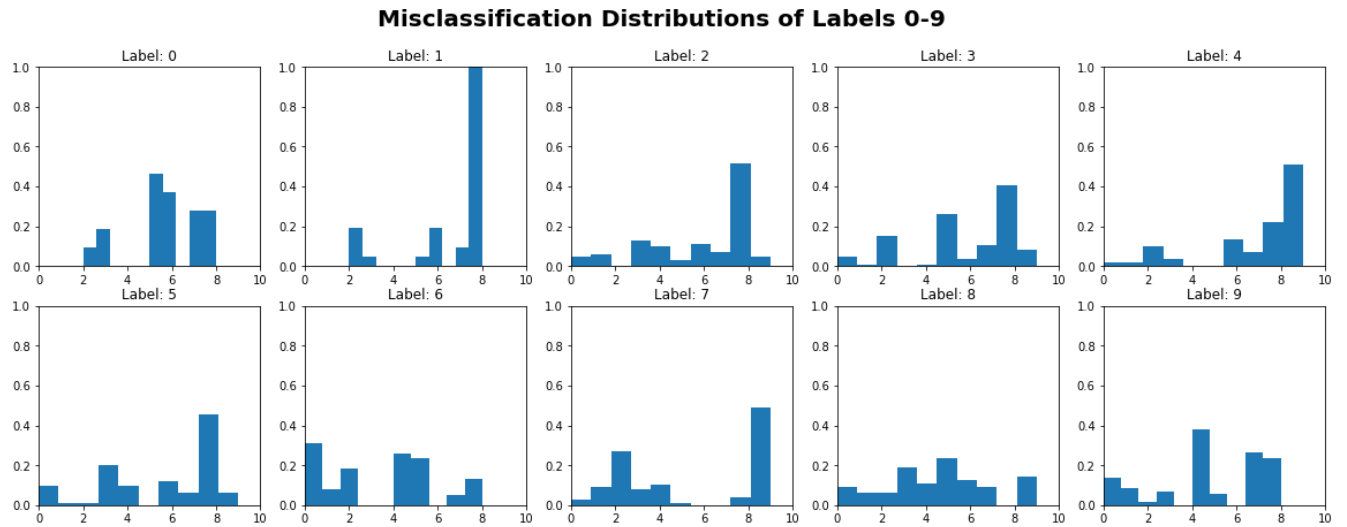
In [12]: 1 best_lr_clf.viz_misclassified()

```

Sample Misclassified Images



```
In [13]: 1 mis_dist = [[] for i in range(10)]
2         for _i, digit in enumerate(best_lr_clf.get_params('all_mis_y_true')):
3             mis_dist[int(digit)].append(best_lr_clf.get_params('all_mis_preds')[_i])
4
5         fig, (ax1, ax2) = plt.subplots(2, 5, figsize=(20, 7))
6         plt.suptitle('Misclassification Distributions of Labels 0-9', fontsize=20, weight='heavy')
7
8         for i in range(5):
9             ax1[i].hist(mis_dist[i], normed = True)
10            ax1[i].set_title('Label: {}'.format(i))
11            ax1[i].set_xlim(0, 10)
12            ax1[i].set_ylim(0, 1)
13            ax2[i].hist(mis_dist[i + 5], normed = True)
14            ax2[i].set_title('Label: {}'.format(i+5))
15            ax2[i].set_xlim(0, 10)
16            ax2[i].set_ylim(0, 1)
```

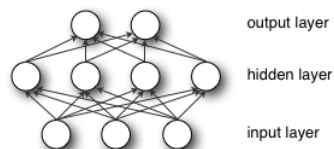


Answer 1.8

The top 3 classes that are most likely to be misclassified are digits 5, 3, and 2. On the test set, more than 10% of images in these class labels are misclassified.

Question 2: MNIST MLP! Find out what that means to me. MNIST MLP! Take care, TCB!

The multilayer perceptron can be understood as a logistic regression classifier in which the input is first transformed using a learnt non-linear transformation. The non-linear transformation is often chosen to be either the logistic function or the tanh function or the RELU function, and its purpose is to project the data into a space where it becomes linearly separable. The output of this so-called hidden layer is then passed to the logistic regression graph that we have constructed in the first problem.



We'll construct a model with 1 **hidden layer**. That is, you will have an input layer, then a hidden layer with the nonlinearity, and finally an output layer with cross-entropy loss (or equivalently log-softmax activation with a negative log likelihood loss).

2.1. Using a similar architecture as in Question 1 and the same training, validation and test sets, build a PyTorch model for the multilayer perceptron. Use the \tanh function as the non-linear activation function.

2.2. The initialization of the weights matrix for the hidden layer must assure that the units (neurons) of the perceptron operate in a regime where information gets propagated. For the \tanh function, you may find it advisable to initialize with the interval $\left[-\sqrt{\frac{6}{\text{fan}_{in} + \text{fan}_{out}}}, \sqrt{\frac{6}{\text{fan}_{in} + \text{fan}_{out}}}\right]$, where fan_{in} is the number of units in the $(i - 1)$ -th layer, and fan_{out} is the number of units in the i -th layer. This is known as **Xavier Initialization**. Use Xavier Initialization to initialize your MLP. Feel free to use PyTorch's in-built Xavier Initialization methods.

2.3. Using $\lambda = 0.01$ to compare with Question 1, experiment with the learning rate (try 0.1 and 0.01 for example), batch size (use 64, 128 and 256) and the number of units in your hidden layer (use between 25 and 200 units). For what combination of these parameters do you obtain the highest validation accuracy? You may want to start with 20 epochs for running time and experiment a bit to make sure that your models reach convergence.

2.4. For your best combination plot the cross-entropy loss on the training set as a function of iteration.

2.5. For your best combination use classification accuracy to evaluate how well your model is performing on the validation set at the end of each epoch. Plot this validation accuracy as the model trains.

2.6. Select what you consider the best set of parameters and predict the labels of the test set. Compare your predictions with the given labels. What classification accuracy do you obtain on the training and test sets?

2.7. How does your test accuracy compare to that of the logistic regression classifier in Question 1? Compare best parameters for both models.

2.8. What classes are most likely to be misclassified? Plot some misclassified training and test set images.

Gratuitous Titular Reference: Respect, originally performed by Otis Redding, became a huge hit and an anthem for the recently departed "Queen of Soul" Aretha Franklin. Respect is often credited with popularizing the word usages "probers" (a synonym for respect) and "sock it to me".

```
In [14]: 1 # 2.1, 2.2
2 class MLPPyTorch(nn.Module):
3     def __init__(self, N_hidden):
4         super().__init__()
5         self.l1 = nn.Linear(784, N_hidden)
6
7         # Xavier Initialization
8         torch.nn.init.xavier_uniform_(self.l1.weight)
9         torch.nn.init.constant_(self.l1.bias, 0.)
10
11         self.l2 = nn.Linear(N_hidden, 10)
12
13     def forward(self, x):
14         x = self.l1(x)
15         x = torch.tanh(x)
16         x = self.l2(x)
17         return x
```

```

In [18]: 1 # 2.3 grid search best hyper parameters
2 learning_rates = [0.1, 0.01]
3 batch_sizes = [64, 128, 256]
4 N_hiddens = [25, 50, 100, 150, 200]
5 val_scores_res = []
6
7 for eta in learning_rates:
8     for bs in batch_sizes:
9         for nh in N_hiddens:
10             print('eta = {}, bs = {}, nh = {}'.format(eta, bs, nh))
11             mlp = MLPPyTorch(nh)
12             mlp_clf = MNIST_Classifier(mlp, learning_rate=eta, batch_size=bs, regularization=0.01, epochs=30)
13             batch_train_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=bs, sampler=train_sampler)
14             mlp_clf.fit(batch_train_loader, train_size, validation_loader, validation_size)
15             val_scores_res.append(mlp_clf.get_params('validation_scores'))
16

```

eta = 0.1, bs = 64, nh = 25

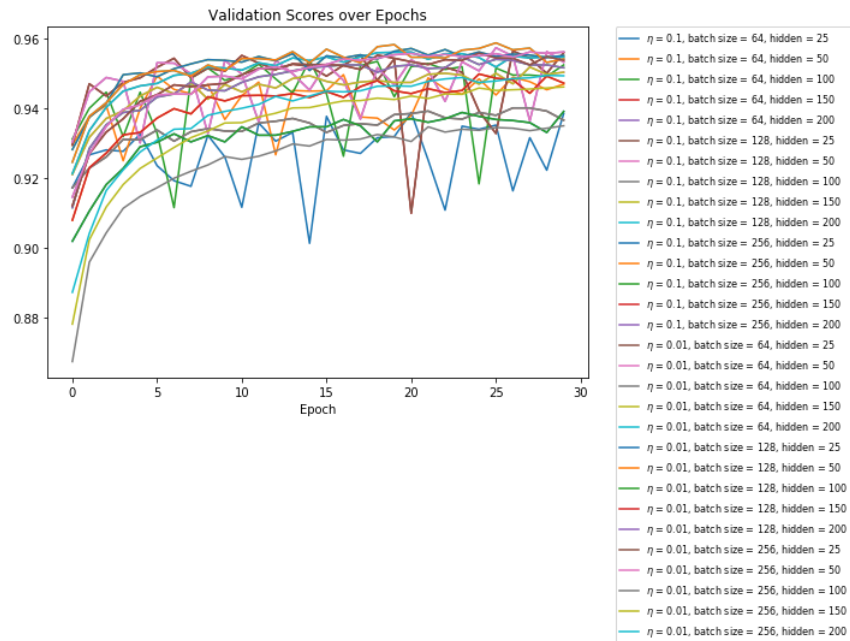
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:112: UserWarning: invalid index of a 0-dim tensor. This will be an error in PyTorch 0.5. Use tensor.item() to convert a 0-dim tensor to a Python number

eta = 0.1, bs = 64, nh = 50
eta = 0.1, bs = 64, nh = 100
eta = 0.1, bs = 64, nh = 150
eta = 0.1, bs = 64, nh = 200
eta = 0.1, bs = 128, nh = 25
eta = 0.1, bs = 128, nh = 50
eta = 0.1, bs = 128, nh = 100
eta = 0.1, bs = 128, nh = 150
eta = 0.1, bs = 128, nh = 200
eta = 0.1, bs = 256, nh = 25
eta = 0.1, bs = 256, nh = 50
eta = 0.1, bs = 256, nh = 100
eta = 0.1, bs = 256, nh = 150
eta = 0.1, bs = 256, nh = 200
eta = 0.01, bs = 64, nh = 25
eta = 0.01, bs = 64, nh = 50
eta = 0.01, bs = 64, nh = 100
eta = 0.01, bs = 64, nh = 150
eta = 0.01, bs = 64, nh = 200
eta = 0.01, bs = 128, nh = 25
eta = 0.01, bs = 128, nh = 50
eta = 0.01, bs = 128, nh = 100
eta = 0.01, bs = 128, nh = 150
eta = 0.01, bs = 128, nh = 200
eta = 0.01, bs = 256, nh = 25
eta = 0.01, bs = 256, nh = 50
eta = 0.01, bs = 256, nh = 100
eta = 0.01, bs = 256, nh = 150
eta = 0.01, bs = 256, nh = 200
30 epoch fitted

```

In [19]: 1 fig, ax = plt.subplots(1, 1, figsize=(7, 5))
2         for _a, eta in enumerate(learning_rates):
3             for _b, bs in enumerate(batch_sizes):
4                 for _c, nh in enumerate(N_hidden):
5                     ax.plot(val_scores_res[_a*len(batch_sizes) + _b*len(N_hidden) + _c],
6                             label=r'$\eta$ = {}, batch size = {}, hidden = {}'.format(eta, bs, nh))
7
8         ax.set_xlabel('Epoch')
9         ax.set_title('Validation Scores over Epochs')
10        ax.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0., fontsize=8)
11
12        plt.tight_layout()

```



```

In [20]: 1 np.array(val_scores_res).shape

```

```

Out[20]: (30, 30)

```

```

In [21]: 1 # best params combination's 1D index
2         np.argmax(np.array(val_scores_res)[: , -10: ].mean(axis=1))

```

```

Out[21]: 9

```

Answer 2.3

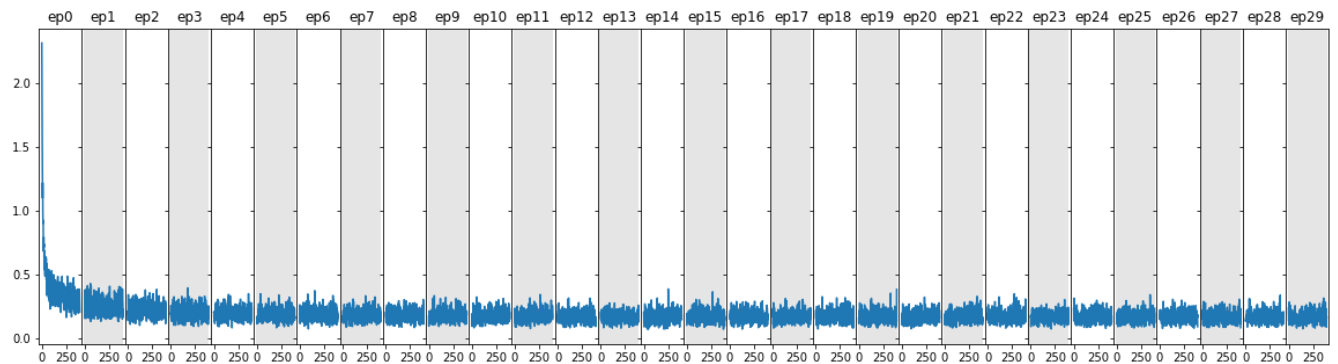
The highest validation accuracy with $\lambda = 0.01$ was achieved with the following combination of hyper parameters:

- learning_rate = 0.1
- batch_size = 128
- N_hidden = 200

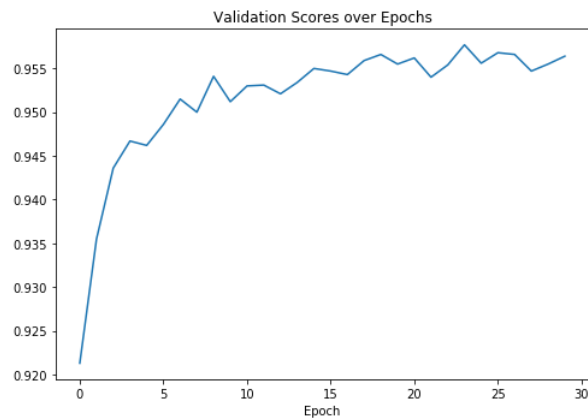
```
In [22]: 1 # 2.4
2 # best model cross-entropy loss over iterations
3 best_mlp_clf = MNIST_Classifier(MLPPyTorch(200), learning_rate=0.1, batch_size=128,
4                                   regularization=0.01, epochs=30)
5 batch_train_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=128, sampler=train_sampler)
6 best_mlp_clf.fit(batch_train_loader, train_size, validation_loader, validation_size)
7 best_mlp_clf.viz_training_loss()
```

/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:112: UserWarning: invalid index of a 0-dim tensor. This will be an error in PyTorch 0.5. Use tensor.item() to convert a 0-dim tensor to a Python number

30 epoch fitted



```
In [23]: 1 # 2.5
2 # plot final iteration validation score over epochs
3 validation_scores = best_mlp_clf.get_params('validation_scores')
4
5 fig, ax = plt.subplots(1, 1, figsize=(7, 5))
6 pd.Series(validation_scores).plot(title='Validation Scores over Epochs', ax=ax)
7 ax.set_xlabel('Epoch')
8 plt.tight_layout()
```



```
In [24]: 1 # 2.6
2 # accuracy on train & test set
3 train_acc = best_mlp_clf.score(train_loader, 'Train', train_size, save_misclassified=False, verbose=True)
4 test_acc = best_mlp_clf.score(test_loader, 'Test', test_size, save_misclassified=False, verbose=True)
```

On Train set: Accuracy: 48086/50000 (96.2%)

On Test set: Accuracy: 9589/10000 (95.9%)

Answer 2.7

With $\lambda = 0.01$, the best test accuracies of the 2 models are:

- Logistic regression: 91.2%
- Multi-layer perceptron: 95.9%

Answer 2.8 - MLP misclassification on train set:

```

In [25]: 1 # 2.8 - Training set
2 train_preds, train_labels, train_correct_count = best_mlp_clf.predict(train_loader, 'Train',
3                                     save_misclassified=True)
4
5 # misclassification counts
6 df_mis_train = pd.DataFrame(pd.Series(train_labels[train_labels != train_preds], name='digit'))
7 df_mis_train['count'] = np.ones((len(df_mis_train),)).astype(int)
8 df_mis_train = df_mis_train.groupby(['
9     'digit'])
10 )['count'].agg([sum]).rename(columns={'sum': 'mis_count'}).reset_index()
11
12 # digit counts
13 df_digit_train = pd.DataFrame(pd.Series(train_labels, name='digit'))
14 df_digit_train['count'] = np.ones((len(df_digit_train),)).astype(int)
15 df_digit_train = df_digit_train.groupby(['
16     'digit'])
17 )['count'].agg([sum]).rename(columns={'sum': 'count'}).reset_index()
18
19 df_train_res = df_digit_train.merge(df_mis_train, on='digit', how='inner').set_index('digit')
20 df_train_res['mis_percent'] = df_train_res['mis_count'] / df_train_res['count'] * 100
21 df_train_res.sort_values(['mis_percent'], ascending=False)

```

```

Out[25]:
count  mis_count  mis_percent
digit
3      5164      298      5.770720
5      4509      253      5.611000
8      4893      260      5.313713
9      4959      230      4.638032
4      4911      216      4.398290
2      4955      184      3.713421
7      5163      166      3.215185
6      4906      104      2.119853
0      4925       98      1.989848
1      5615      105      1.869991

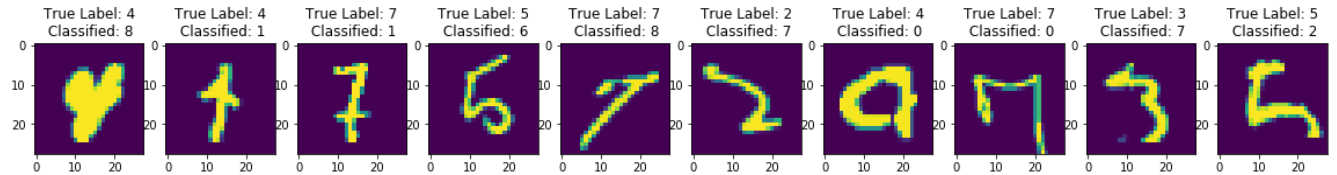
```

```

In [26]: 1 best_mlp_clf.viz_misclassified()

```

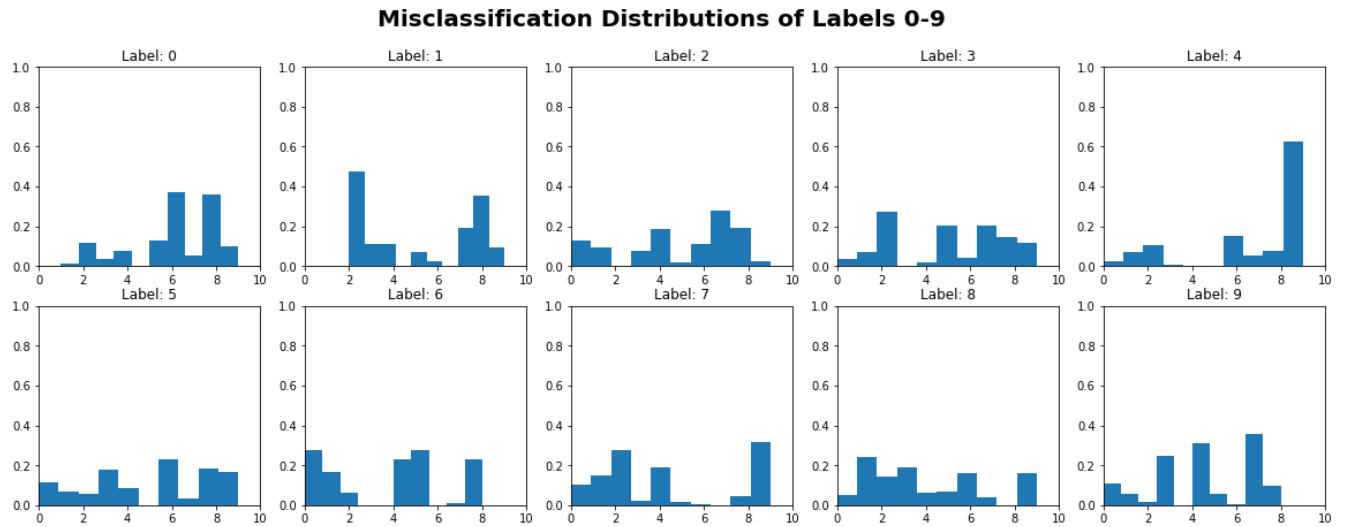
Sample Misclassified Images



```

In [27]: 1 mis_dist = [[] for i in range(10)]
2         for _i, digit in enumerate(best_mlp_clf.get_params('all_mis_y_true')):
3             mis_dist[int(digit)].append(best_mlp_clf.get_params('all_mis_preds')[_i])
4
5         fig, (ax1, ax2) = plt.subplots(2, 5, figsize=(20, 7))
6         plt.suptitle('Misclassification Distributions of Labels 0-9', fontsize=20, weight='heavy')
7
8         for i in range(5):
9             ax1[i].hist(mis_dist[i], normed = True)
10            ax1[i].set_title('Label: {}'.format(i))
11            ax1[i].set_xlim(0, 10)
12            ax1[i].set_ylim(0, 1)
13            ax2[i].hist(mis_dist[i + 5], normed = True)
14            ax2[i].set_title('Label: {}'.format(i+5))
15            ax2[i].set_xlim(0, 10)
16            ax2[i].set_ylim(0, 1)

```



Answer 2.8 - MLP misclassification on test set


```

In [28]: 1 # 2.8 - Test set
2 test_preds, test_labels, test_correct_count = best_mlp_clf.predict(test_loader, 'Test', save_misclassified=True)
3
4 # misclassification counts
5 df_mis_test = pd.DataFrame(pd.Series(test_labels[test_labels != test_preds], name='digit'))
6 df_mis_test['count'] = np.ones((len(df_mis_test),)).astype(int)
7 df_mis_test = df_mis_test.groupby(['
8     'digit'])
9     ][ 'count'].agg([sum]).rename(columns={'sum': 'mis_count'}).reset_index()
10
11 # digit counts
12 df_digit_test = pd.DataFrame(pd.Series(test_labels, name='digit'))
13 df_digit_test['count'] = np.ones((len(df_digit_test),)).astype(int)
14 df_digit_test = df_digit_test.groupby(['
15     'digit'])
16     ][ 'count'].agg([sum]).rename(columns={'sum': 'count'}).reset_index()
17
18 df_test_res = df_digit_test.merge(df_mis_test, on='digit', how='inner').set_index('digit')
19 df_test_res['mis_percent'] = df_test_res['mis_count'] / df_test_res['count'] * 100
20 df_test_res.sort_values(['mis_percent'], ascending=False)

```

```

Out[28]:      count  mis_count  mis_percent
digit
5      892         60      6.726457
8      974         53      5.441478
9     1009         52      5.153617
4      982         50      5.091650
7     1028         49      4.766537
2     1032         46      4.457364
3     1010         43      4.257426
6      958         31      3.235908
1     1135         15      1.321586
0      980         12      1.224490

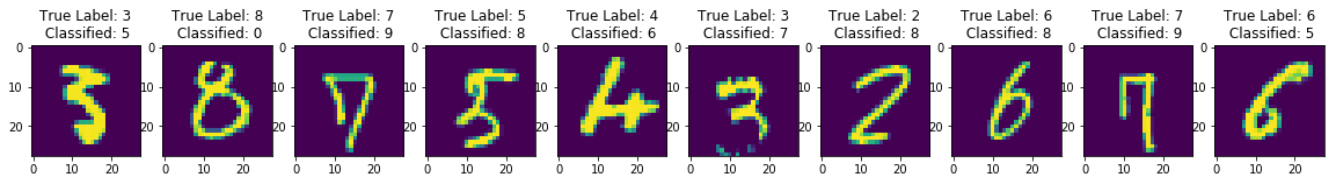
```

```

In [29]: 1 best_mlp_clf.viz_misclassified()

```

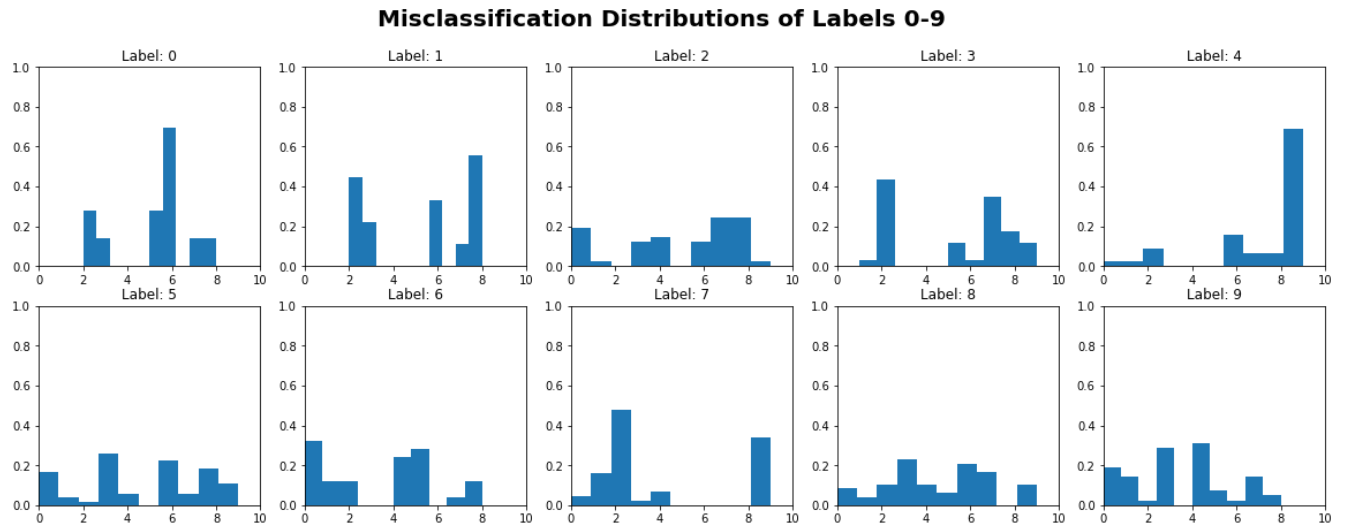
Sample Misclassified Images



```

In [30]: 1 mis_dist = [[] for i in range(10)]
2         for _i, digit in enumerate(best_mlp_clf.get_params('all_mis_y_true')):
3             mis_dist[int(digit)].append(best_mlp_clf.get_params('all_mis_preds')[_i])
4
5         fig, (ax1, ax2) = plt.subplots(2, 5, figsize=(20, 7))
6         plt.suptitle('Misclassification Distributions of Labels 0-9', fontsize=20, weight='heavy')
7
8         for i in range(5):
9             ax1[i].hist(mis_dist[i], normed = True)
10            ax1[i].set_title('Label: {}'.format(i))
11            ax1[i].set_xlim(0, 10)
12            ax1[i].set_ylim(0, 1)
13            ax2[i].hist(mis_dist[i + 5], normed = True)
14            ax2[i].set_title('Label: {}'.format(i+5))
15            ax2[i].set_xlim(0, 10)
16            ax2[i].set_ylim(0, 1)

```



Answer 2.8

With the MLP model, the top 3 classes that are most likely to be misclassified are digits 5, 8, and 9. The percentage of misclassified for these classes are 6.7%, 5.4% and 5.2%, respectively, in the test set. These percentages are significantly (almost 2x) lower than the highest misclassification percentages using a single layer Logistic Regression model, which are all above 10%.