# APMTH 207: Advanced Scientific Computing:

## Stochastic Methods for Data Analysis, Inference and Optimization

## Homework #3

**Harvard University**
**Spring 2018**
**Instructors: Rahul Dave**
**Due Date:** Friday, Febrary 16th, 2018 at 10:00am

**Instructions:**

- Upload your final answers as an iPython notebook containing all of your work to Canvas.
- Structure your notebook and your work to maximize readability.

```
In [1]:
1  import numpy as np
2  import pandas as pd
3  from copy import deepcopy
4  import time
5
6  import matplotlib
7  import matplotlib.pyplot as plt
8  from mpl_toolkits.mplot3d import axes3d, Axes3D
9
10 import seaborn as sns
11 sns.set_style("whitegrid", {'axes.grid' : False})
12 sns.set_context('talk')
13 %matplotlib inline
```

## Problem 1: Optimization via Descent

Suppose you are building a pricing model for laying down telecom cables over a geographical region. Your model takes as input a pair of coordinates, $(x, y)$, and contains two parameters, $\lambda_1, \lambda_2$. Given a coordinate, $(x, y)$, and model parameters, the loss in revenue corresponding to the price model at location $(x, y)$ is described by

$$L(x, y, \lambda_1, \lambda_2) = 0.000045\lambda_2^2 y - 0.000098\lambda_1^2 x + 0.003926\lambda_1 x \exp\left\{\left(y^2 - x^2\right)\left(\lambda_1^2 + \lambda_2^2\right)\right\}$$

Read the data contained in `HW3_data.csv`. This is a set of coordinates configured on the curve $y^2 - x^2 = -0.1$. Given the data, find parameters $\lambda_1, \lambda_2$ that minimize the net loss over the entire dataset.

## Part A

- Visually verify that for $\lambda_1 = 2.05384, \lambda_2 = 0$, the loss function $L$ is minimized for the given data.
- Implement gradient descent for minimizing $L$ for the given data, using the learning rate of 0.001.
- Implement stochastic gradient descent for minimizing $L$ for the given data, using the learning rate of 0.001.

## Part B

- Compare the average time it takes to update the parameter estimation in each iteration of the two implementations. Which method is faster? Briefly explain why this result should be expected.
- Compare the number of iterations it takes for each algorithm to obtain an estimate accurate to `1e-3` (you may wish to set a cap for maximum number of iterations). Which method converges to the optimal point in fewer iterations? Briefly explain why this result should be expected.

## Part C

Compare the performance of stochastic gradient descent for the following learning rates: 1, 0.1, 0.001, 0.0001. Based on your observations, briefly describe the effect of the choice of learning rate on the performance of the algorithm.

## Answer to Problem 1 Part A

```
In [2]:   1  # Load data
          2
          3  data = np.genfromtxt('HW3_data.csv', delimiter=',')
          4  x = data[0, ]
          5  y = data[1, ]
          6  print('Number of data points: {}.'.format(len(x)))
```

Number of data points: 16000.

```
In [3]:    1  # In GD, we use the gradient of total loss at each iteration
           2  # In SGD, we multiply the gradient by total sample size at each iteration
           3
           4  # For the convenience of comparison, we compute the average loss
           5  # when displaying / evaluating results
           6
           7  def L(x, y, lam):
           8
           9      # Average loss
          10
          11      return np.mean(0.000045 * lam[1]**2 * y - 0.000098 * lam[0]**2 * x \
          12                   + 0.003926 * lam[0] * x * np.exp((y**2 - x**2) * (lam[0]**2 + lam[1]**2)))
          13
          14  def dL(x, y, lam):
          15
          16      # Gradient of total loss
          17
          18      z = y*y - x*x
          19      z1 = x*np.exp((lam[0]**2+lam[1]**2)*z)
          20      a = np.sum(-0.000196*lam[0]*x + (0.003926+0.007852*lam[0]**2*z)*z1)
          21      b = np.sum(0.00009*lam[1]*y + 0.007852*lam[0]*lam[1]*z*z1)
          22      return np.array([a, b])
```

**Visually verify that for $\lambda_1 = 2.05384, \lambda_2 = 0$, the loss function $L$ is minimized for the given data.**
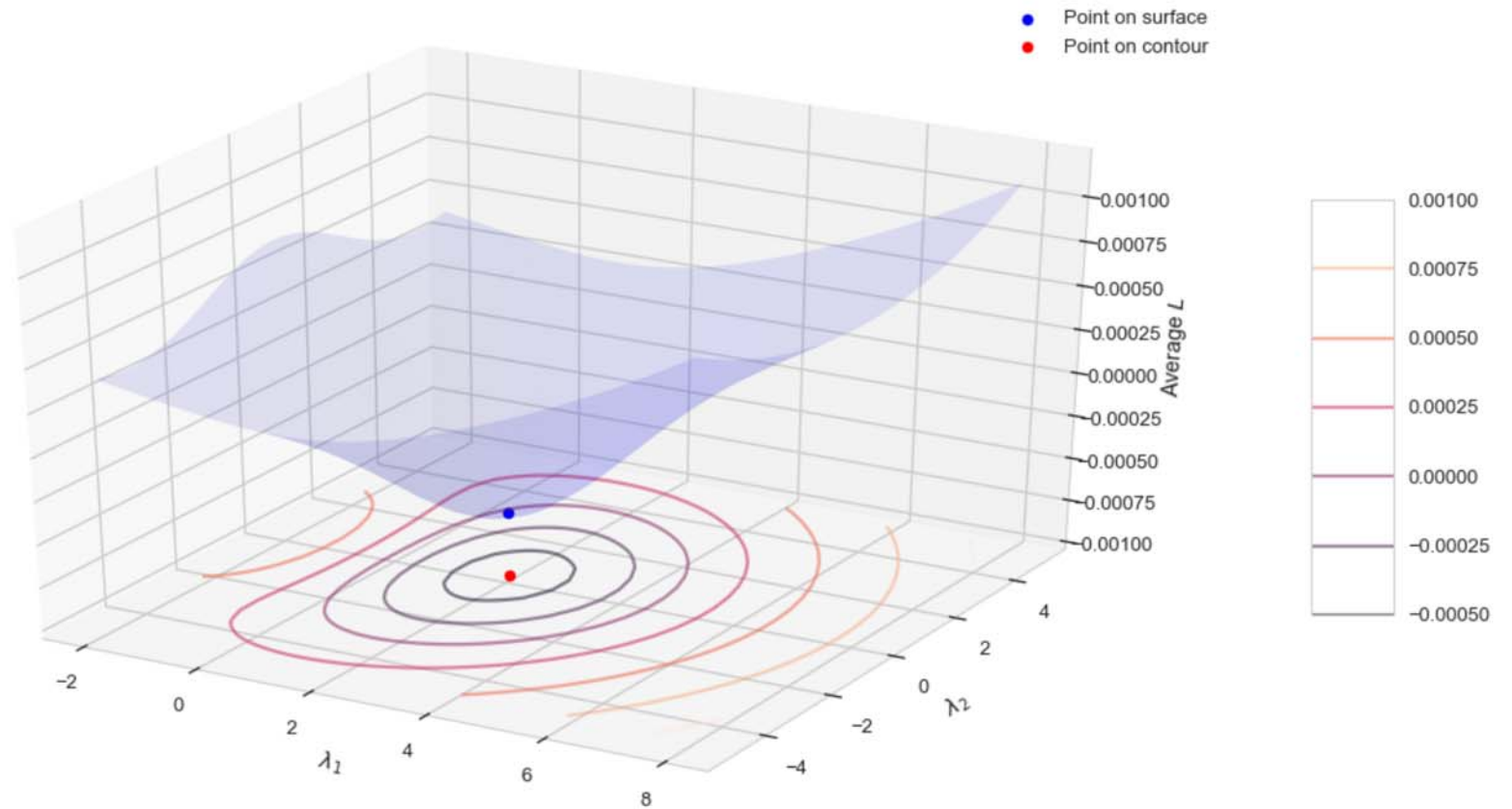
We plot the point on 3D surface as well as the contour plot.

In [4]:
```python
def plot_3d(ms=np.linspace(-2, 8, 100), bs = np.linspace(-5, 5, 20), z_offset=-1e-3, lam=None):

    # reference:
    # https://am207.github.io/2018spring/wiki/gradientdescent.html

    M, B = np.meshgrid(ms, bs)
    zs = np.array([L(x, y, l) for l in zip(np.ravel(M), np.ravel(B))])
    Z = zs.reshape(M.shape)
    fig = plt.figure(figsize=(20, 10))
    ax = fig.gca(projection='3d')

    ax.plot_surface(M, B, Z, rstride=1, cstride=1, color='b', alpha=0.1)
    c = ax.contour(M, B, Z, alpha=0.5, offset=z_offset, stride=30)

    ax.set_zlim(z_offset, np.max(Z) * 1.1)
    ax.set_xlabel('$\lambda_1$', labelpad=15)
    ax.set_ylabel('$\lambda_2$', labelpad=15)
    ax.set_zlabel('Average $L$', labelpad=15)

    fig.colorbar(c, shrink=0.5, aspect=5)

    if lam is not None:
        ax.scatter([lam[0]], [lam[1]], [L(x, y, lam)], c='b', s=50, label='Point on surface')
        ax.scatter([lam[0]], [lam[1]], [z_offset], c='r', s=50, label='Point on contour')
        plt.legend()

lam_best = np.array([2.05384, 0])
print('Average L(x, y, lambda_1={}, lambda_2={}) = {}'.format(lam_best[0], lam_best[1], L(x, y, lam_best)))
print('Gradient is {}'.format(dL(x, y, lam_best)))
plot_3d(lam=lam_best)
```

```
Average L(x, y, lambda_1=2.05384, lambda_2=0.0) = -0.00062088815015901337
Gradient is [ -4.53213012e-05   0.00000000e+00]
```

From the visual check above, we know for $\lambda_1 = 2.05384, \lambda_2 = 0$, the loss function $L$ is minimized for the given data.

**Implement gradient descent for minimizing $L$ for the given data, using the learning rate of 0.001.**

```
In [5]:  1  class GD:
         2      def __init__(self, x, y, lam_init, step=0.001, max_iter=10000, tol=0.001):
         3          self.x = deepcopy(x)
         4          self.y = deepcopy(y)
         5          self.m = x.size
         6          self.lam_init = lam_init
         7          self.step = step
         8          self.max_iter = max_iter
         9          self.tol = tol
        10          self.costs = []
        11          self.time_ = []
        12          self.total_time = 0
        13          self.history = []
        14          self.iter_ = 0
        15
        16      def run_gd(self):
        17
        18          # Run max_iter iterations
        19
        20          total_start = time.time()
        21          self.history.append(self.lam_init)
        22          self.costs.append(L(self.x, self.y, self.lam_init))
        23          for _ in range(self.max_iter):
        24              start = time.time()
        25              self.iter_ += 1
        26              self.history.append(self.history[-1] - self.step * dL(self.x, self.y, self.history[-1]))
        27              self.costs.append(L(self.x, self.y, self.history[-1]))
        28              self.time_.append(time.time() - start)
        29          self.total_time = time.time() - total_start
        30          return self
        31
        32      def run_gd_test(self, actual=np.array([2.05384, 0])):
        33
        34          # Run until approaching actual within tol or reaching max_iter
        35
        36          total_start = time.time()
        37          self.history.append(self.lam_init)
        38          self.costs.append(L(self.x, self.y, self.lam_init))
        39          for _ in range(self.max_iter):
        40              start = time.time()
        41              self.iter_ += 1
        42              self.history.append(self.history[-1] - self.step * dL(self.x, self.y, self.history[-1]))
```

```
43              self.costs.append(L(self.x, self.y, self.history[-1]))
44              if np.linalg.norm(self.history[-1] - actual) <= self.tol:
45                  self.time_.append(time.time() - start)
46                  break
47              self.time_.append(time.time() - start)
48          self.total_time = time.time() - total_start
49          return self
```

In [6]:
```
1 gd = GD(x, y, np.array([1, 1]), tol=1e-3).run_gd_test()
```

In [7]:
```
1 print('Gradient descent obtains an estimate accurate to 1e-3 using {} iterations.'.format(gd.iter_))
```

Gradient descent obtains an estimate accurate to 1e-3 using 2753 iterations.

**Implement stochastic gradient descent for minimizing $L$ for the given data, using the learning rate of 0.001.**

```
In [8]:   1  class SGD:
          2      def __init__(self, x, y, lam_init, step=0.001, max_epoch=5, tol=0.001):
          3          self.x = deepcopy(x)
          4          self.y = deepcopy(y)
          5          self.m = x.size
          6          self.lam_init = lam_init
          7          self.step = step
          8          self.max_epoch = max_epoch
          9          self.tol = tol
         10          self.costs = []
         11          self.total_cost = 0
         12          self.time_ = []
         13          self.total_time = 0
         14          self.history = []
         15          self.iter_ = 0
         16
         17      def run_sgd(self):
         18
         19          # Run until reaching max_epoch
         20
         21          total_start = time.time()
         22          self.costs.append(L(self.x[0], self.y[0], self.lam_init))
         23          self.history.append(self.lam_init)
         24          for _ in range(self.max_epoch):
         25              for i in range(self.m):
         26                  start = time.time()
         27                  self.iter_ += 1
         28                  self.history.append(self.history[-1]\
         29                                  - self.step * self.m* dL(self.x[i], self.y[i], self.history[-1]))
         30                  self.total_cost += L(self.x[i], self.y[i], self.history[-1])
         31                  self.costs.append(self.total_cost / self.iter_)
         32                  self.time_.append(time.time() - start)
         33              neworder = np.random.permutation(self.m)
         34              self.x = self.x[neworder]
         35              self.y = self.y[neworder]
         36          self.total_time = time.time() - total_start
         37          return self
         38
         39      def run_sgd_test(self, actual=np.array([2.05384, 0])):
         40
         41          # Run until approaching actual within tol or reaching max_epoch
         42
```

```
43              total_start = time.time()
44              self.costs.append(L(self.x[0], self.y[0], self.lam_init))
45              self.history.append(self.lam_init)
46              done = False
47              for _ in range(self.max_epoch):
48                  for i in range(self.m):
49                      start = time.time()
50                      self.iter_ += 1
51                      self.history.append(self.history[-1]\
52                                     - self.step * self.m * dL(self.x[i], self.y[i], self.history[-1]))
53                      self.total_cost += L(self.x[i], self.y[i], self.history[-1])
54                      self.costs.append(self.total_cost / self.iter_)
55                      if np.linalg.norm(self.history[-1] - actual) <= self.tol:
56                          done = True
57                          self.time_.append(time.time() - start)
58                          break
59                      self.time_.append(time.time() - start)
60                  if done:
61                      break
62                  neworder = np.random.permutation(self.m)
63                  self.x = self.x[neworder]
64                  self.y = self.y[neworder]
65              self.total_time = time.time() - total_start
66              return self
```

In [9]:
```
1 sgd = SGD(x, y, np.array([1, 1]), step=0.001, tol=1e-3, max_epoch=5).run_sgd_test()
```

In [10]:
```
1 print('Stochastic gradient descent obtains an estimate accurate to 1e-3 using {} iterations.'.format(sgd.ite
```

Stochastic gradient descent obtains an estimate accurate to 1e-3 using 8054 iterations.

## Answer to Problem 1 Part B

**Compare the average time it takes to update the parameter estimation in each iteration of the two implementations. Which method is faster? Briefly explain why this result should be expected.**

In [11]:
```
1 print('Gradient descent (GD): total run time: {:7f} s; average time in each iteration {:7f} s.'\
2       .format(gd.total_time, gd.total_time / gd.iter_))
3 print('Stochastic gradient descent (SGD): total run time: {:7f} s; average time in each iteration {:7f} s.'\
4       .format(sgd.total_time, sgd.total_time / sgd.iter_))
```

Gradient descent (GD): total run time: 2.230600 s; average time in each iteration 0.000810 s.
Stochastic gradient descent (SGD): total run time: 0.516364 s; average time in each iteration 0.000064 s.

SGD is faster than GD. While GD needs to calculate the gradient of all samples in each iteration, SGD only needs to calculate the gradient of 1 sample in each iteration. As a result, the average run time in each iteration is much shorter for SGD than that for GD. Although SGD takes more iterations to "find" the optimal point, the total run time of SGD is shorter due to significantly shorter run time in each iteration.

**Compare the number of iterations it takes for each algorithm to obtain an estimate accurate to 1e-3 (you may wish to set a cap for maximum number of iterations). Which method converges to the optimal point in fewer iterations? Briefly explain why this result should be expected.**

In [12]:
```
1 print('Gradient descent obtains an estimate accurate to 1e-3 using {} iterations.'.format(gd.iter_))
2 print('Stochastic gradient descent obtains an estimate accurate to 1e-3 using {} iterations.'.format(sgd.ite
```

Gradient descent obtains an estimate accurate to 1e-3 using 2753 iterations.
Stochastic gradient descent obtains an estimate accurate to 1e-3 using 8054 iterations.

It takes more iterations for SGD. We can investigate the behaviour of GD and SGD through some visualizations.
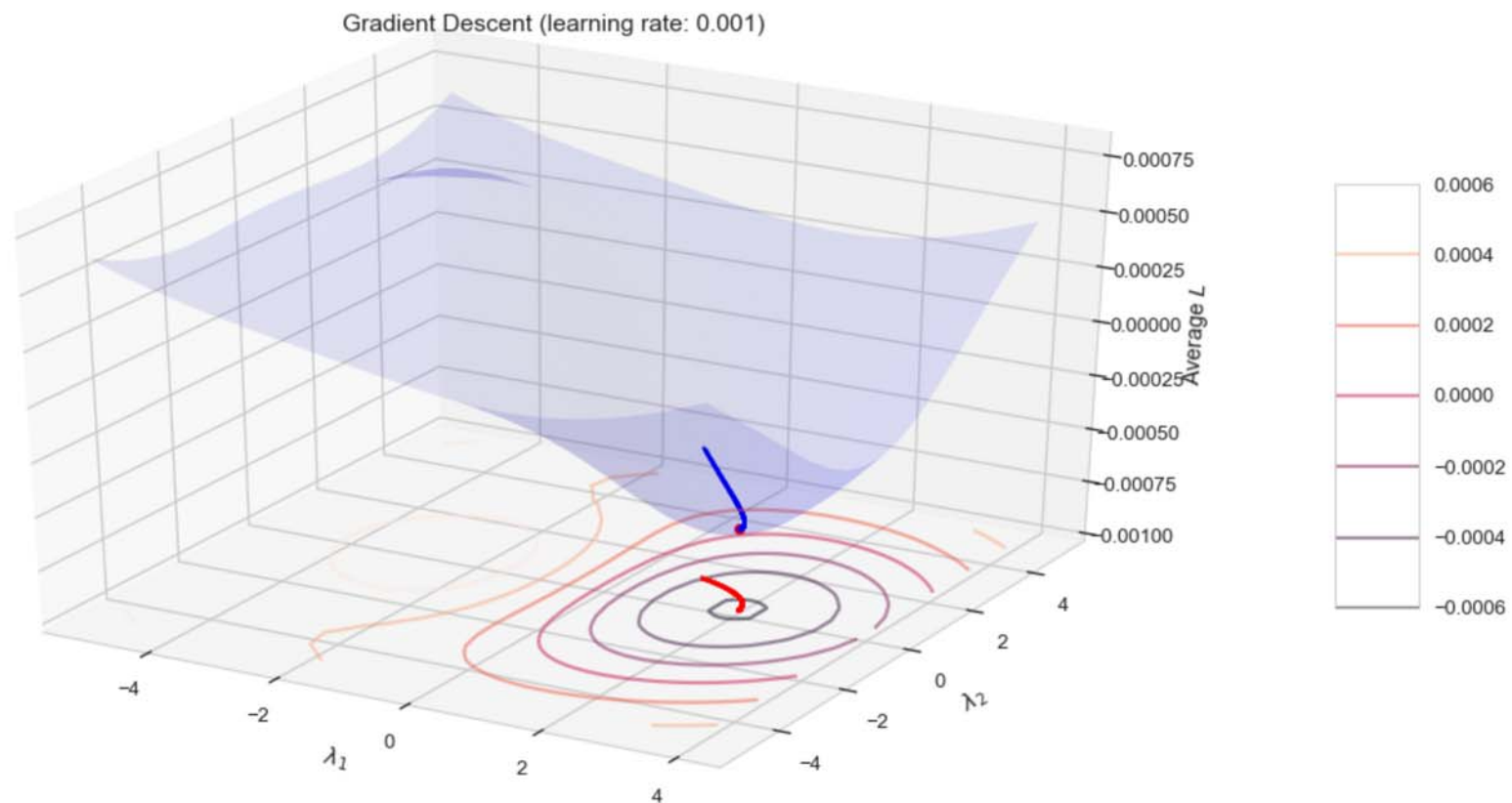
```
In [13]:    1  def plot_3d_hist(history, costs, \
            2                  ms=np.linspace(-5, 4, 100), bs = np.linspace(-5, 5, 20), z_offset=-1e-3):
            3
            4      # reference:
            5      # https://am207.github.io/2018spring/wiki/gradientdescent.html
            6
            7      M, B = np.meshgrid(ms, bs)
            8      zs = np.array([L(x, y, l) for l in zip(np.ravel(M), np.ravel(B))])
            9      Z = zs.reshape(M.shape)
           10      fig = plt.figure(figsize=(20, 10))
           11      ax = fig.gca(projection='3d')
           12
           13      ax.plot_surface(M, B, Z, rstride=1, cstride=1, color='b', alpha=0.1)
           14      c = ax.contour(M, B, Z, alpha=0.5, offset=z_offset, stride=30)
           15
           16      ax.set_zlim(z_offset, np.max(Z) * 1.1)
           17      ax.set_xlabel('$\lambda_1$', labelpad=15)
           18      ax.set_ylabel('$\lambda_2$', labelpad=15)
           19      ax.set_zlabel('Average $L$', labelpad=15)
           20
           21      fig.colorbar(c, shrink=0.5, aspect=5)
           22
           23      ax.plot([history[-1][0]], [history[-1][1]], [costs[-1]], \
           24              markerfacecolor='r', markeredgecolor='r', marker='o', markersize=7)
           25      ax.plot([t[0] for t in history], [t[1] for t in history], costs, alpha=0.5, \
           26              markerfacecolor='b', markeredgecolor='b', marker='.', markersize=5)
           27      ax.plot([t[0] for t in history], [t[1] for t in history], z_offset, alpha=0.5, \
           28              markerfacecolor='r', markeredgecolor='r', marker='.', markersize=5)
           29
           30  def plot_summary(gd, actual=np.array([2.05384, 0])):
           31      costs = np.array(gd.costs)
           32      costs = costs[~np.isnan(costs)]
           33      l = len(costs)
           34      history = np.array(gd.history)
           35      history = history[:l, :]
           36      plt.figure(figsize=(12, 10))
           37
           38
           39
           40      plt.subplot(2, 2, 1)
           41      plt.plot(range(l), costs, 'o-', markersize=5, alpha=0.5)
           42      plt.title('Learning rate: {}'.format(gd.step))
```
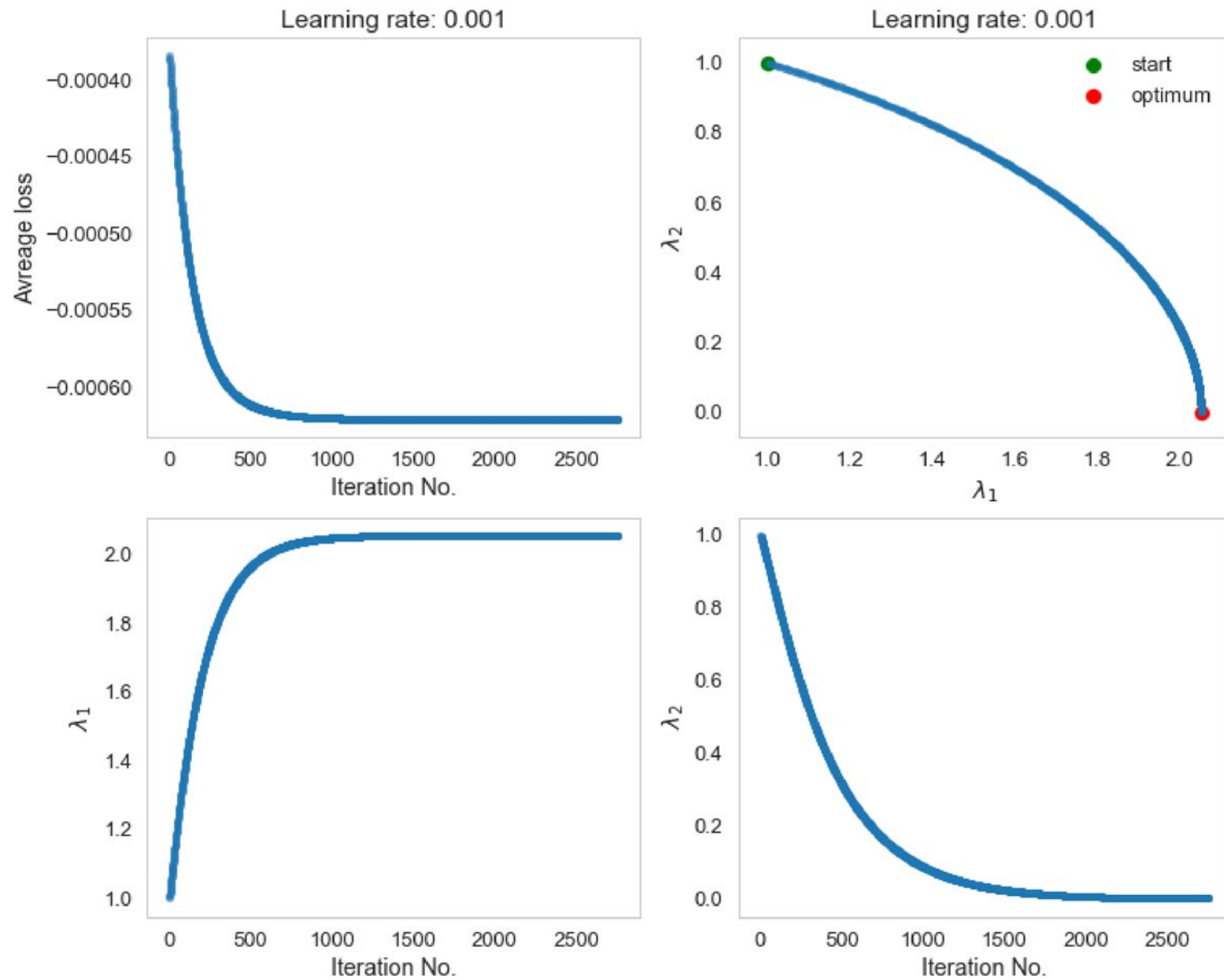
```
43      plt.xlabel('Iteration No.')
44      plt.ylabel('Avreage loss')
45
46      plt.subplot(2, 2, 2)
47      plt.plot(history[:, 0], history[:, 1], 'o-',markersize=5, alpha=0.5)
48      plt.title('Learning rate: {}'.format(gd.step))
49      plt.xlabel('$\lambda_1$')
50      plt.ylabel('$\lambda_2$')
51      plt.scatter(gd.lam_init[0], gd.lam_init[1], color='g', label='start')
52      plt.scatter(actual[0], actual[1], color='r', label='optimum')
53      plt.legend()
54
55      plt.subplot(2, 2, 3)
56      plt.plot(range(l), history[:, 0], 'o-',markersize=5, alpha=0.5)
57      plt.xlabel('Iteration No.')
58      plt.ylabel('$\lambda_1$')
59
60      plt.subplot(2, 2, 4)
61      plt.plot(range(l), history[:, 1], 'o-',markersize=5, alpha=0.5)
62      plt.xlabel('Iteration No.')
63      plt.ylabel('$\lambda_2$')
```

In [14]:
```
1 plot_3d_hist(gd.history, gd.costs)
2 plt.title('Gradient Descent (learning rate: {})'.format(gd.step));
```

Gradient Descent (learning rate: 0.001)

```
In [15]:   1 print('Gradient Descent')
           2 plot_summary(gd)
```

Gradient Descent

In [16]:
```
1 plot_3d_hist(sgd.history, [L(x, y, h) for h in sgd.history])
2 plt.title('Stochastic Gradient Descent (learning rate: {})'.format(sgd.step));
```



Stochastic Gradient Descent (learning rate: 0.001)

In [17]:
```
1  print('Stochastic Gradient Descent')
2  plot_summary(sgd)
```

Stochastic Gradient Descent

The path GD takes seems much smoother than that SGD takes. While batch gradient ensures GD approaching local optimum in each iteration, SGD takes some detour in "finding" the optimum, which results in more iterations than GD.

In this case, it seems SGD "bumped into" the optimum we were waiting for rather than finding the optimum. If we don't know the optimum beforehand, SGD might not converge there. SGD usually gets close to the optimum faster than batch method, but never fully converge to the optimum. Also, in this case although SGD takes more iterations than GD, it takes much less computation.

## Answer to Problem 1 Part C

In [18]:
```python
steps = [1, 0.1, 0.001, 0.0001, 0.00001]
sgds = []
for s in steps:
    print('Learning rate: {}.'.format(s))
    sgds.append(SGD(x, y, np.array([1, 1]), step=s, tol=1e-3, max_epoch=5).run_sgd_test())
    print('Number of iterations: {}.'.format(sgds[-1].iter_))
    print('Final lambda: {}'.format(sgds[-1].history[-1]))
    print('L2 distance to the optimum: {}.'.format(np.linalg.norm(sgds[-1].history[-1] - lam_best)))
    print('Average loss along the path: {}.'.format(sgds[-1].costs[-1]))
    print('Average loss on the dataset: {}.'.format(L(x, y, sgds[-1].history[-1])))
    print('----------------')
    print()
```

Learning rate: 1.

C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:3: RuntimeWarning: overflow encountered in dou
ble_scalars
  This is separate from the ipykernel package so we can avoid doing imports until
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:7: RuntimeWarning: overflow encountered in dou
ble_scalars
  import sys
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:8: RuntimeWarning: overflow encountered in dou
ble_scalars

C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:8: RuntimeWarning: invalid value encountered i
n double_scalars


Number of iterations: 80000.
Final lambda: [ nan  nan]
L2 distance to the optimum: nan.
Average loss along the path: nan.
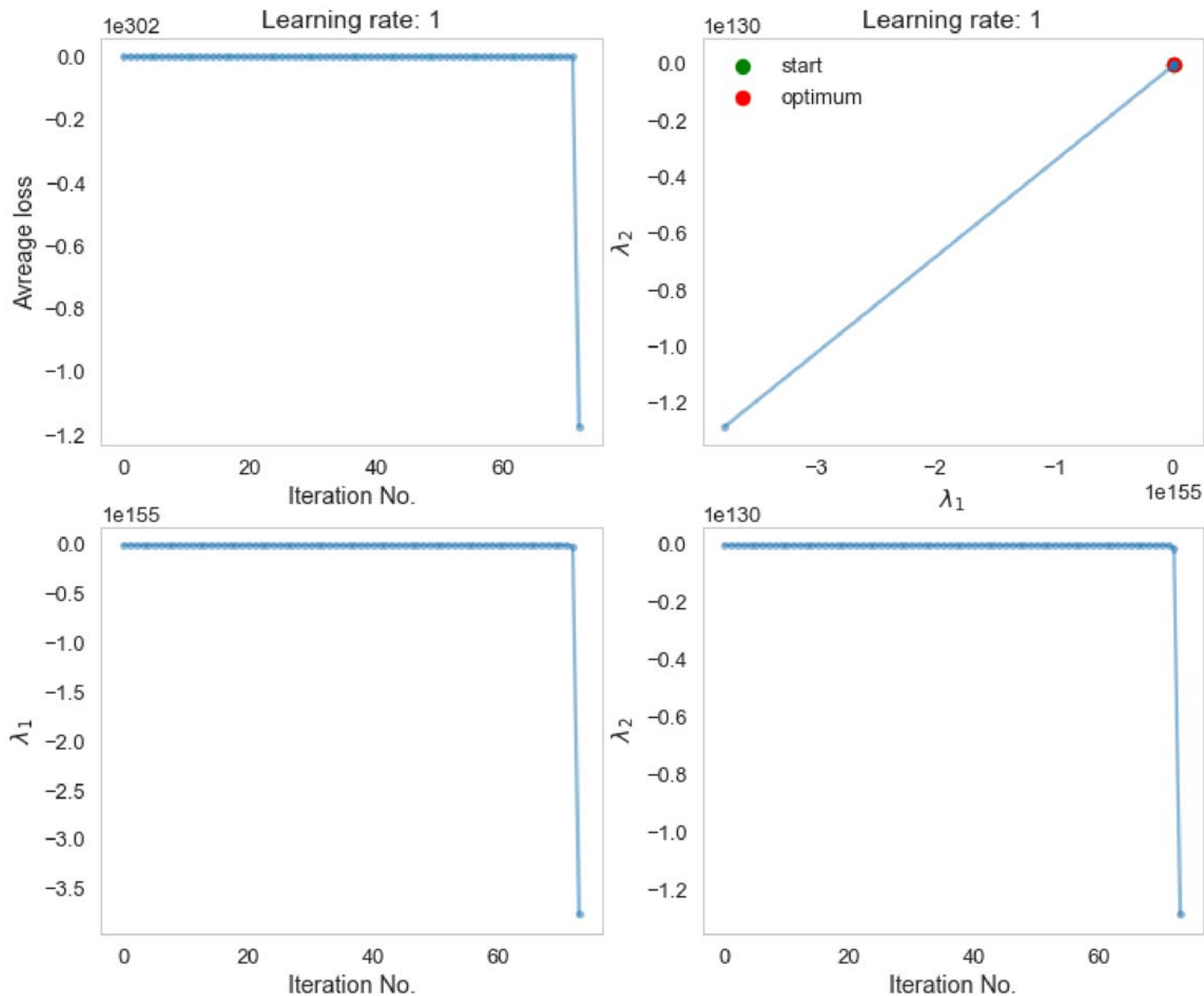Average loss on the dataset: nan.
----------------

Learning rate: 0.1.
Number of iterations: 80000.
Final lambda: [ nan  nan]
L2 distance to the optimum: nan.
Average loss along the path: nan.
Average loss on the dataset: nan.
----------------

Learning rate: 0.001.
Number of iterations: 8054.
Final lambda: [  2.05364425e+000   6.32404027e-322]
L2 distance to the optimum: 0.0001957475651179763.
Average loss along the path: -0.2918589239125654.
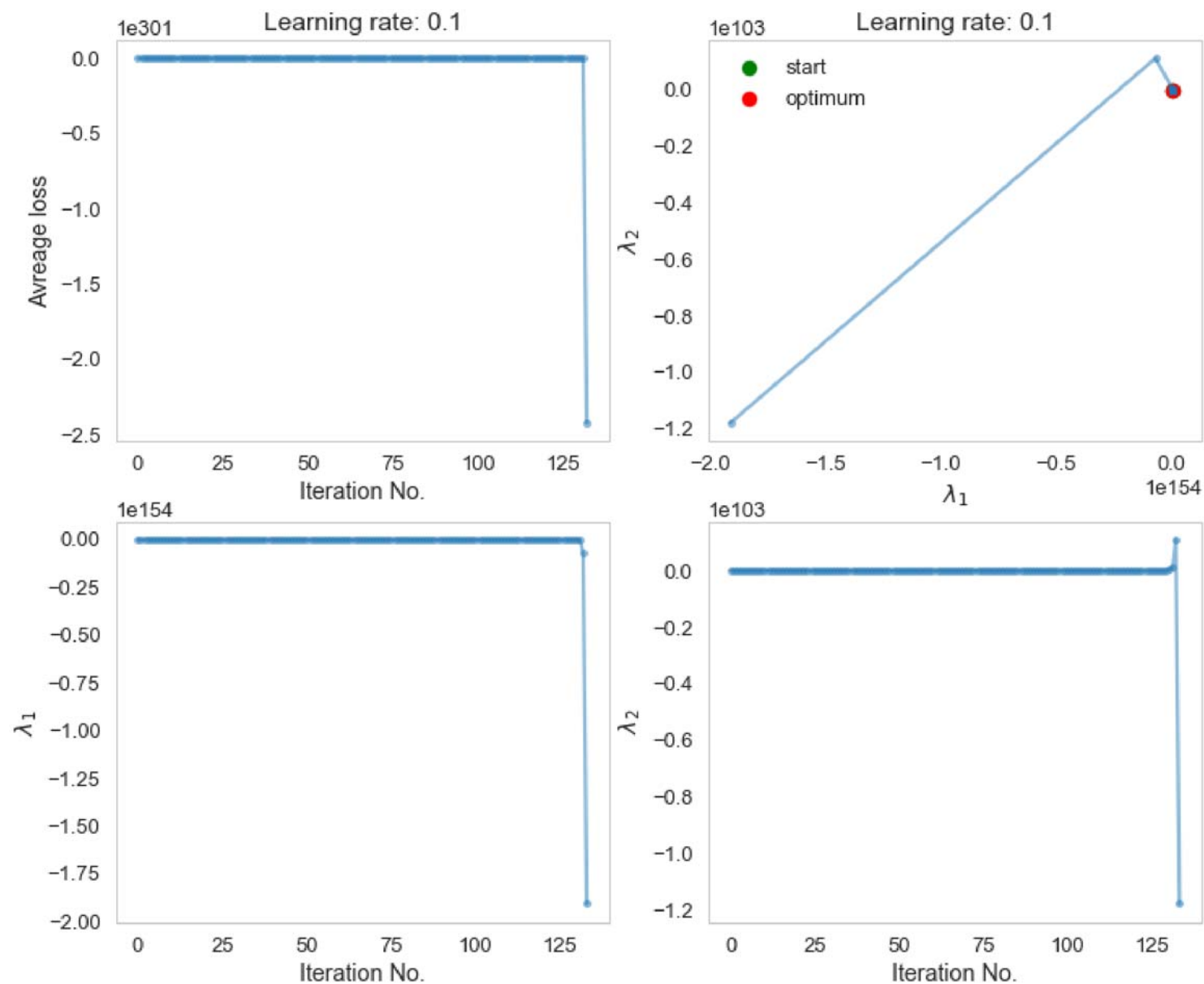Average loss on the dataset: -0.0006208814949922327.
----------------

Learning rate: 0.0001.
Number of iterations: 8270.
Final lambda: [  2.05287932e+000   7.60692890e-306]
L2 distance to the optimum: 0.00096067954052792.
Average loss along the path: -0.27727466193866285.
Average loss on the dataset: -0.0006208813532968805.
----------------

Learning rate: 1e-05.
Number of iterations: 80000.
Final lambda: [ -5.36324926e+00   4.35840383e-17]
L2 distance to the optimum: 7.417089260285739.
Average loss along the path: -0.005923523951164853.
Average loss on the dataset: 0.0005100955437239491.
----------------

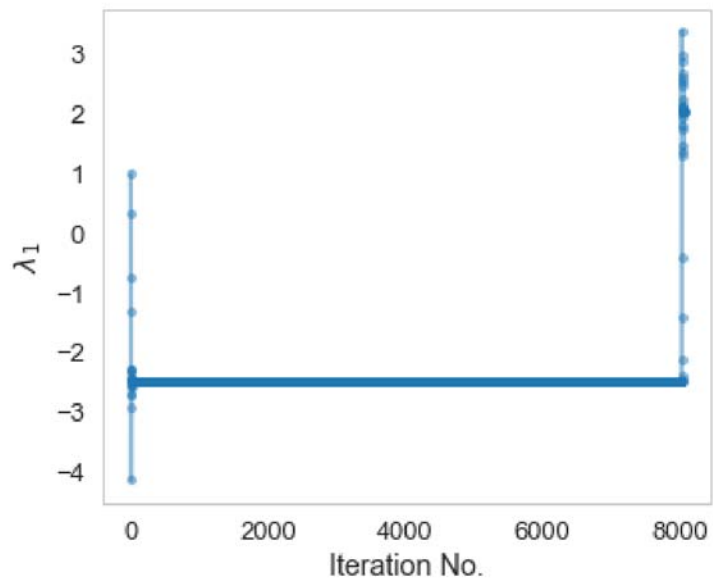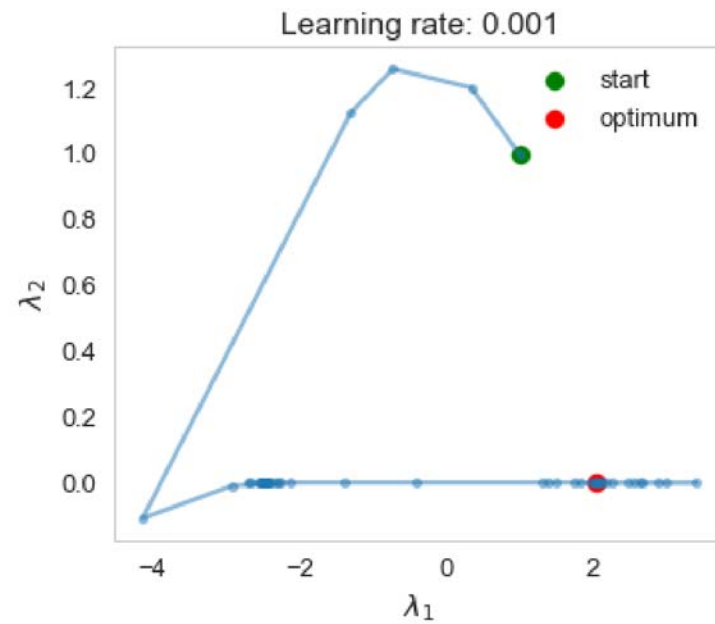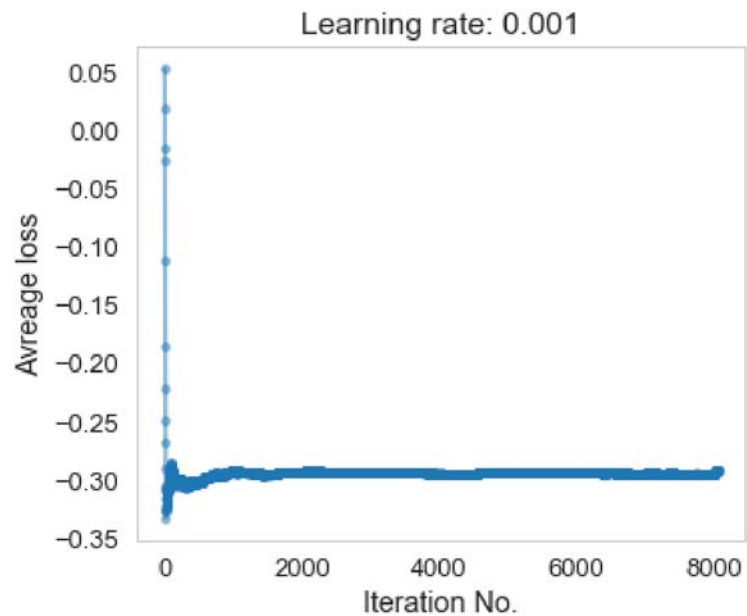We can visualize the path in each case as follows.

```
In [19]:  1  plot_summary(sgds[0])
          2  plot_summary(sgds[1])
```
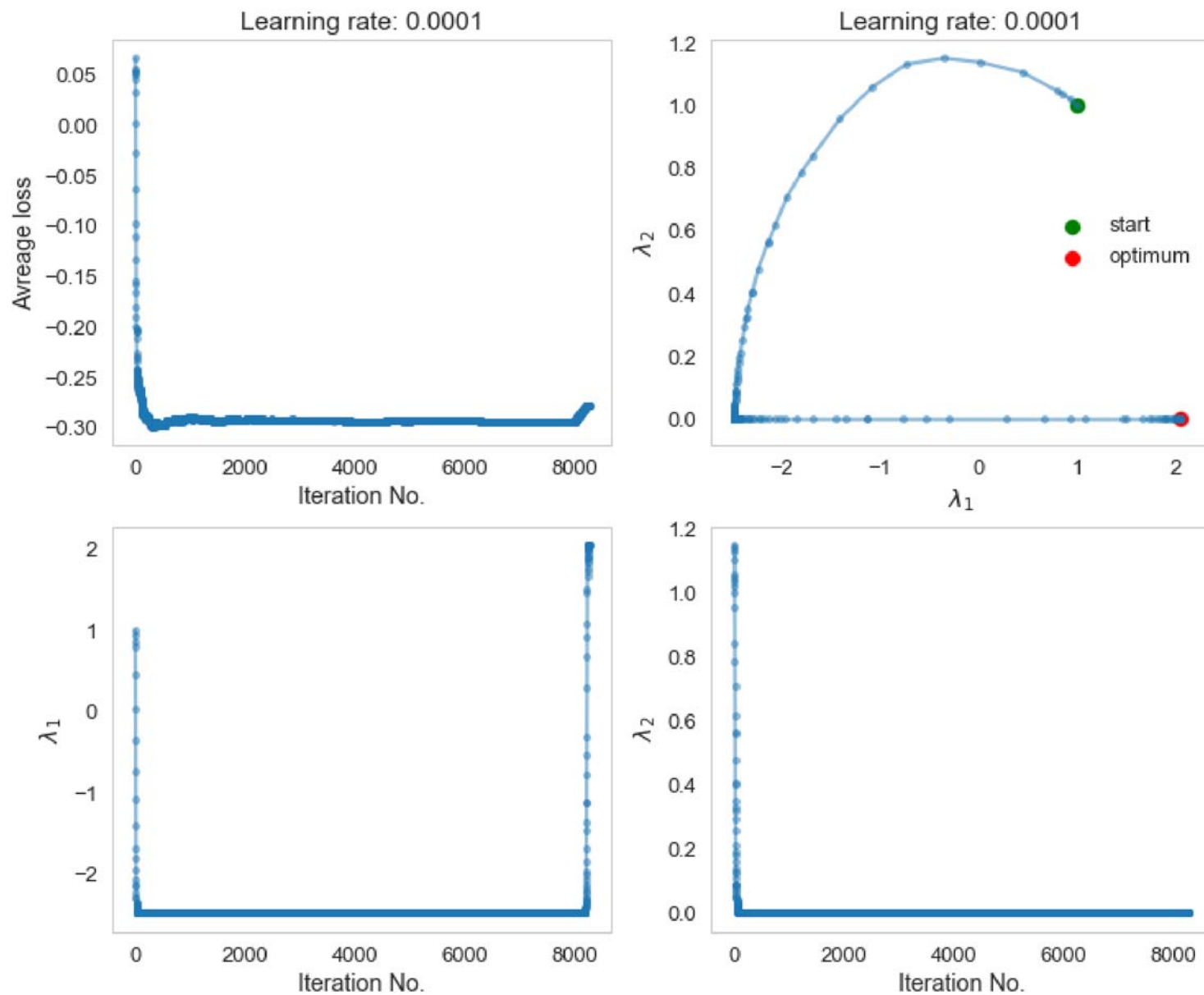
As we can see, when we set the learning rate at 1 or 0.1, the update is too large at each step; the loss function and the gradient blow up, and we finally encounter overflow issue.
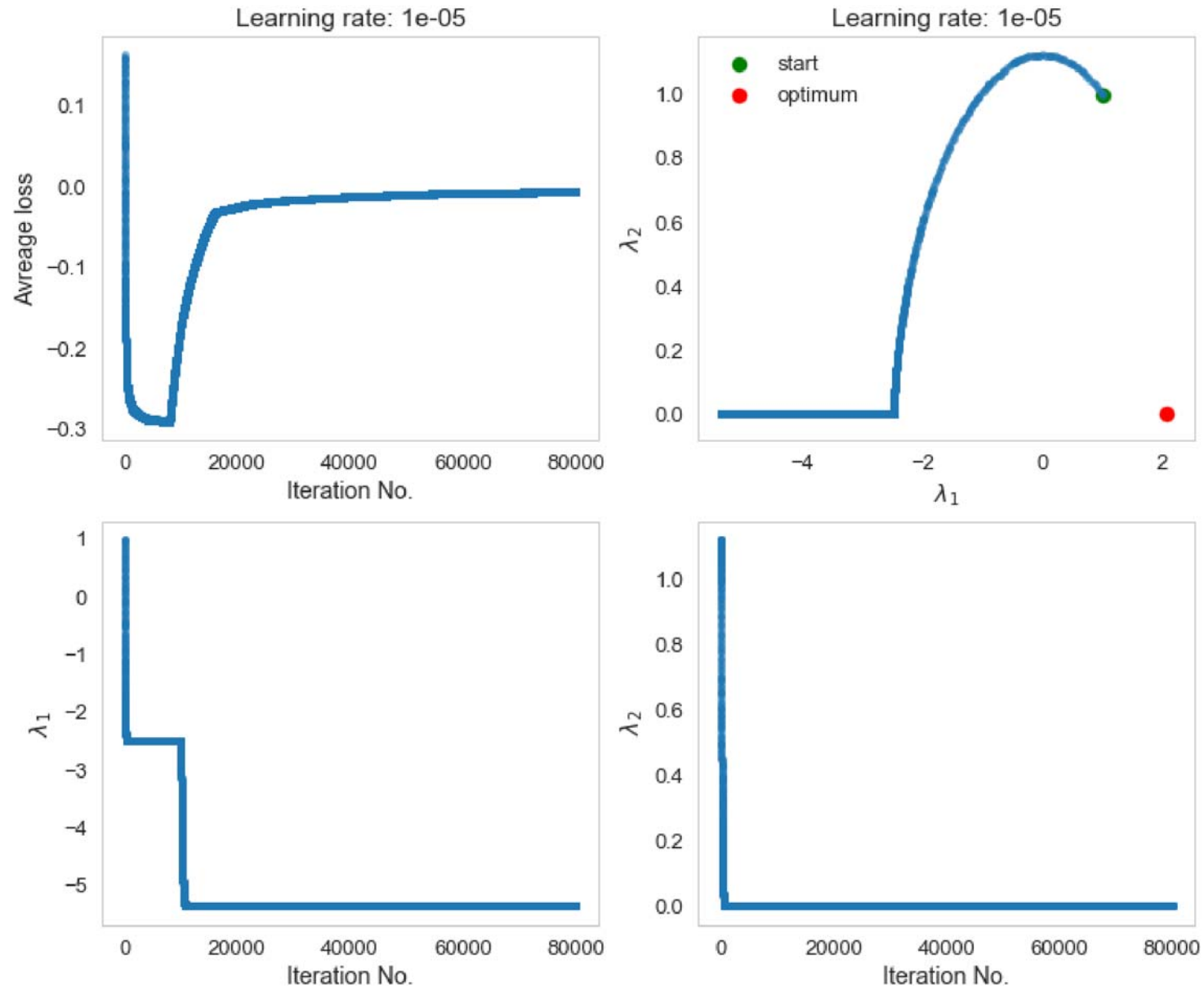
In [20]:
```
1  plot_summary(sgds[2])
2  plot_summary(sgds[3])
```

When we set the learning rate at 0.001 or 0.0001, we can reach the optimum. The number of iterations corresponding to 0.0001 is slightly larger than that of 0.001.

In [21]: `1 plot_summary(sgds[4])`



When we set the learning rate at 0.00001, SGD didn't "find" the optimum when it reaches `max_epoch`; it seems it "converges" to another answer.

In summary, when the learning rate is too large, SGD would take large random steps in each iteration and is likely to blow up and runs into overflow issue.

When the learning rate is neither too large nor too small, SGD could "reach" the optimum although it might not be able to "realize" it without evaluating the batch gradient or loss function. Smaller learning rate would usually take more iterations.

When the learning rate is too small, it takes longer for SGD to converge and SGD might not be able to "reach" the optimum within a fixed amount of iterations.

# Problem 2. SGD for Multinomial Logistic Regression on MNIST

The _MNIST dataset (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 and test with MNIST.

**The MNIST dataset (including a train/test split which you must use) is part of PyTorch in the torchvision module. The Lab will have details of how to load it.**

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.

1. Plot 10 sample images from the MNIST dataset (to develop intuition for the feature space).
2. Construct a softmax formulation in PyTorch of multinomial logistic regression with Cross Entropy Loss.
3. Train your model using SGD to minimize the cost function. _Use a batch size of 64, a learning rate $\eta = 0.01$, and 10 epochs._
4. Plot the cross-entropy loss on the training set as a function of iteration.
5. What are the training and test set accuracies?
6. Plot some (around 5) examples of misclassifications.

## Answer to Problem 2

```
In [2]:  1  import torch
         2  import torchvision
         3  from torchvision import datasets
         4  from torchvision import transforms
         5  from torch.autograd import Variable
         6  import torch.nn as nn
         7  import torch.nn.functional as F
         8  import torch.optim as optim
```

```
In [3]:  1  # Load data
         2
         3  transform = transforms.Compose([transforms.ToTensor(), \
         4                                    transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])
         5
         6  trainset = datasets.MNIST(root='', train=True, download=True, transform=transform)
         7  testset = datasets.MNIST(root='', train=False, transform=transform)
         8
```

**1. Plot 10 sample images from the MNIST dataset**

```
In [4]:    1  # Reference:
           2  # http://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html#sphx-glr-beginner-blitz-cifar10-tutorial
           3
           4  trainloader = torch.utils.data.DataLoader(trainset, batch_size=10, shuffle=False)
           5
           6  def imshow(img):
           7      img = img / 2 + 0.5      # unnormalize
           8      npimg = img.numpy()
           9      plt.imshow(np.transpose(npimg, (1, 2, 0)))
          10
          11  dataiter = iter(trainloader)
          12  images, labels = dataiter.next()
          13
          14  # show images
          15  imshow(torchvision.utils.make_grid(images))
          16  # print labels
          17  print(' '.join([str(l) for l in labels]))
```

5 0 4 1 9 2 1 3 1 4



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

In [5]:
```python
class Model(nn.Module):
    def __init__(self):
        super(Model, self).__init__()
        self.linear = nn.Linear(28 * 28, 10)
        self.logsoftmax = nn.LogSoftmax(dim=1)

    def forward(self, x):
        x = x.view(x.shape[0], 28*28)
        return self.logsoftmax(self.linear(x))


class MLR:
    def __init__(self, lr=0.01, max_epoch=10):
        self.max_epoch = max_epoch
        self.model = Model()
        self.criterion = nn.NLLLoss()
        self.optimizer = optim.SGD(self.model.parameters(), lr=lr)
        self.loss_ = []

    def fit(self, trainloader):
        for epoch in range(self.max_epoch):
            running_loss = 0
            for i, data in enumerate(trainloader, 0):
                inputs, labels = data
                inputs, labels = Variable(inputs), Variable(labels)
                self.optimizer.zero_grad()
                outputs = self.model(inputs)
                loss = self.criterion(outputs, labels)
                loss.backward()
                self.optimizer.step()
                running_loss += loss.data[0]
                self.loss_.append(loss.data[0])
            print('Epoch {} loss: {}'.format(epoch + 1, running_loss / len(trainloader)))
        print('Finished Training.')
        return self

    def predict(self, x):
        outputs = self.model(Variable(deepcopy(x)))
        _, pred = torch.max(outputs.data, 1)
        return pred

def getData(testloader):
```

```
43        return iter(testloader).next()
```

**3. Train your model using SGD to minimize the cost function.** *Use a batch size of 64, a learning rate $\eta = 0.01$, and 10 epochs.*
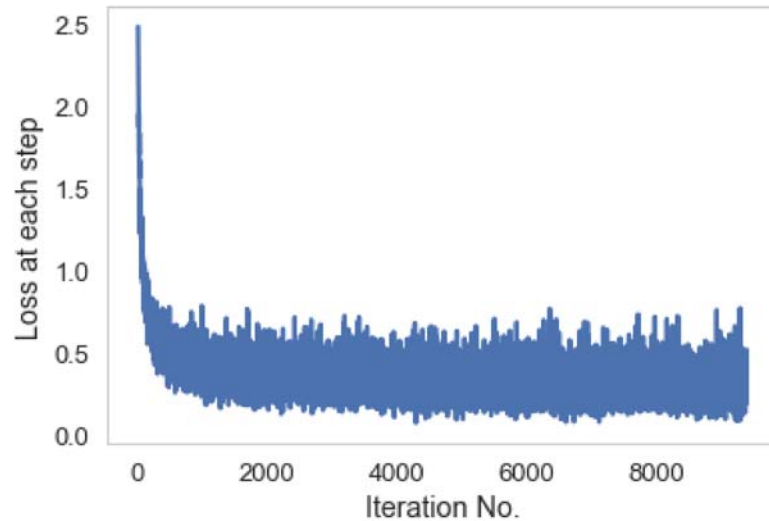
In [6]:
```
1 trainloader = torch.utils.data.DataLoader(trainset, batch_size=64, shuffle=True)
2 trainloader2 = torch.utils.data.DataLoader(trainset, batch_size=60000, shuffle=False)
3 testloader = torch.utils.data.DataLoader(testset, batch_size=10000, shuffle=False)
4
5 mlr = MLR(lr=0.01, max_epoch=10).fit(trainloader)
```

```
Epoch 1 loss: 0.6085820051112663
Epoch 2 loss: 0.3873185949729704
Epoch 3 loss: 0.35267389373484453
Epoch 4 loss: 0.3349054019842575
Epoch 5 loss: 0.3236767000067971
Epoch 6 loss: 0.3159288539370494
Epoch 7 loss: 0.30913406163295193
Epoch 8 loss: 0.30494309202440256
Epoch 9 loss: 0.30115880849741417
Epoch 10 loss: 0.2975056896220519
Finished Training.
```
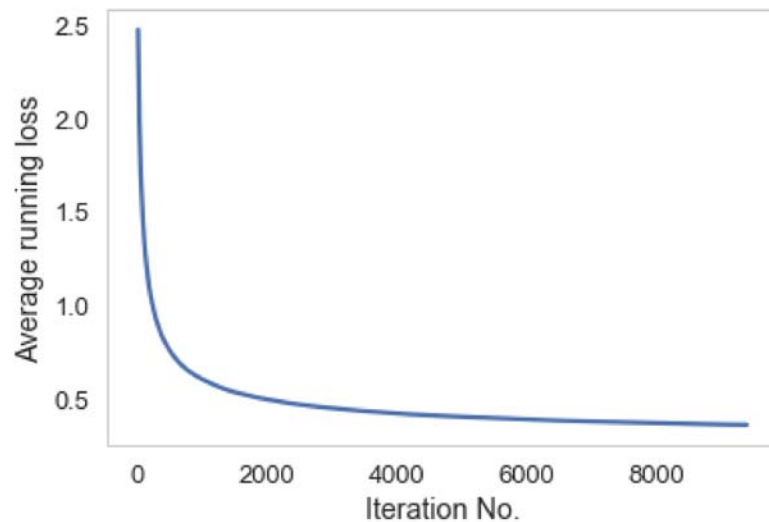
**4. Plot the cross-entropy loss on the training set as a function of iteration.**

We can plot the loss at each step, as well as average cross-entropy loss after each iteration as follows.

In [7]:
```
1  plt.plot(np.arange(1, len(mlr.loss_) + 1), mlr.loss_);
2  plt.xlabel('Iteration No.');
3  plt.ylabel('Loss at each step');
```

In [8]:
```
1  mean_loss = np.cumsum(mlr.loss_) / np.arange(1, len(mlr.loss_) + 1)
2  plt.plot(np.arange(1, len(mean_loss) + 1), mean_loss);
3  plt.xlabel('Iteration No.');
4  plt.ylabel('Average running loss');
```

### 5. What are the training and test set accuracies?

```
In [9]:    1  def accuracy_score(y_true, y_pred, verbose=False):
           2      if not verbose:
           3          return np.mean(y_true == y_pred)
           4      else:
           5          return np.array([np.mean(y_pred_test[y_test == i] == i) for i in range(10)])
```

```
In [10]:   1  x_train, y_train = getData(trainloader2)
           2  x_test, y_test = getData(testloader)
```

```
In [11]:   1  y_pred_train = mlr.predict(x_train)
           2  y_pred_test = mlr.predict(x_test)
           3
           4  accu_train = accuracy_score(y_train, y_pred_train, verbose=True)
           5  accu_test = accuracy_score(y_test, y_pred_test, verbose=True)
           6
           7  print('Training set:')
           8  print('Overall accuracy : {:.4f}'.format(accuracy_score(y_train, y_pred_train)))
           9  for i in range(10):
          10      print('Accuracy of {} : {:.4f}'.format(i, accu_train[i]))
          11
          12  print('----------------')
          13  print()
          14
          15  print('Test set:')
          16  print('Overall accuracy : {:.4f}'.format(accuracy_score(y_test, y_pred_test)))
          17  for i in range(10):
          18      print('Accuracy of {} : {:.4f}'.format(i, accu_test[i]))
```

```
Training set:
Overall accuracy : 0.9162
Accuracy of 0 : 0.9755
Accuracy of 1 : 0.9771
Accuracy of 2 : 0.8769
Accuracy of 3 : 0.9139
Accuracy of 4 : 0.9257
Accuracy of 5 : 0.8296
Accuracy of 6 : 0.9614
Accuracy of 7 : 0.9105
Accuracy of 8 : 0.9138
Accuracy of 9 : 0.8940
----------------

Test set:
Overall accuracy : 0.9191
Accuracy of 0 : 0.9755
Accuracy of 1 : 0.9771
Accuracy of 2 : 0.8769
Accuracy of 3 : 0.9139
Accuracy of 4 : 0.9257
Accuracy of 5 : 0.8296
Accuracy of 6 : 0.9614
Accuracy of 7 : 0.9105
```

```
Accuracy of 8 : 0.9138
Accuracy of 9 : 0.8940
```

### 6. Plot some (around 5) examples of misclassifications.

In [12]:
```python
1 imshow(torchvision.utils.make_grid(x_test[[np.where(y_test != y_pred_test)[0][:8]]]))
2 print('Actual labels: \n' + ' '.join([str(l) for l in y_test[y_test != y_pred_test][:8]]))
3 print()
4 print('Predicted labels: \n' + ' '.join([str(l) for l in y_pred_test[y_test != y_pred_test][:8]]))
```

```
Actual labels:
5 4 3 2 9 7 7 2

Predicted labels:
6 6 2 7 4 1 4 9
```