

# Assignment1

November 16, 2021

## 1 Assignment 1. Music Century Classification

**Assignment Responsible:** Natalie Lang.

In this assignment, we will build models to predict which **century** a piece of music was released. We will be using the "YearPredictionMSD Data Set" based on the Million Song Dataset. The data is available to download from the UCI Machine Learning Repository. Here are some links about the data:

- <https://archive.ics.uci.edu/ml/datasets/yearpredictionmsd>
- <http://millionsongdataset.com/pages/tasks-demos/#yearrecognition>

Note that you are not allowed to import additional packages (**especially not PyTorch**). One of the objectives is to understand how the training procedure actually operates, before working with PyTorch's autograd engine which does it all for us.

### 1.1 Question 1. Data (21%)

Start by setting up a Google Colab notebook in which to do your work. Since you are working with a partner, you might find this link helpful:

- <https://colab.research.google.com/github/googlecolab/colabtools/blob/master/notebooks/colab-github-demo.ipynb>

The recommended way to work together is pair coding, where you and your partner are sitting together and writing code together.

To process and read the data, we use the popular pandas package for data analysis.

```
[3]: import pandas
import numpy as np
import matplotlib.pyplot as plt
np.random.seed(1) # To reproduce our experiment
```

Now that your notebook is set up, we can load the data into the notebook. The code below provides two ways of loading the data: directly from the internet, or through mounting Google Drive. The first method is easier but slower, and the second method is a bit involved at first, but can save you time later on. You will need to mount Google Drive for later assignments, so we recommend figuring how to do that now.

Here are some resources to help you get started:

- <http://colab.research.google.com/notebooks/io.ipynb>

```
[4]: load_from_drive = False

if not load_from_drive:
    csv_path = "http://archive.ics.uci.edu/ml/machine-learning-databases/00203/
    ↳YearPredictionMSD.txt.zip"
else:
    from google.colab import drive
    drive.mount('/content/gdrive')
    csv_path = '/content/gdrive/My Drive/YearPredictionMSD.txt.zip' # TODO -
    ↳UPDATE ME WITH THE TRUE PATH!

t_label = ["year"]
x_labels = ["var%d" % i for i in range(1, 91)]
df = pandas.read_csv(csv_path, names=t_label + x_labels)
```

Now that the data is loaded to your Colab notebook, you should be able to display the Pandas DataFrame df as a table:

```
[5]: df
```

```
[5]:
```

	year	var1	var2	...	var88	var89	var90
0	2001	49.94357	21.47114	...	-1.82223	-27.46348	2.26327
1	2001	48.73215	18.42930	...	12.04941	58.43453	26.92061
2	2001	50.95714	31.85602	...	-0.05859	39.67068	-0.66345
3	2001	48.24750	-1.89837	...	9.90558	199.62971	18.85382
4	2001	50.97020	42.20998	...	7.88713	55.66926	28.74903
...	...	...	...	...	...	...	...
515340	2006	51.28467	45.88068	...	3.42901	-41.14721	-15.46052
515341	2006	49.87870	37.93125	...	12.96552	92.11633	10.88815
515342	2006	45.12852	12.65758	...	-6.07171	53.96319	-8.09364
515343	2006	44.16614	32.38368	...	20.32240	14.83107	39.74909
515344	2005	51.85726	59.11655	...	-5.51512	32.35602	12.17352

[515345 rows x 91 columns]

To set up our data for classification, we'll use the "year" field to represent whether a song was released in the 20-th century. In our case df["year"] will be 1 if the year was released after 2000, and 0 otherwise.

```
[6]: df["year"] = df["year"].map(lambda x: int(x > 2000))
```

```
[7]: df.head(20)
```

```
[7]:
```

	year	var1	var2	...	var88	var89	var90
0	1	49.94357	21.47114	...	-1.82223	-27.46348	2.26327
1	1	48.73215	18.42930	...	12.04941	58.43453	26.92061
2	1	50.95714	31.85602	...	-0.05859	39.67068	-0.66345
3	1	48.24750	-1.89837	...	9.90558	199.62971	18.85382
4	1	50.97020	42.20998	...	7.88713	55.66926	28.74903
5	1	50.54767	0.31568	...	5.00283	-11.02257	0.02263
6	1	50.57546	33.17843	...	4.50056	-4.62739	1.40192

7	1	48.26892	8.97526	...	-0.30633	3.98364	-3.72556
8	1	49.75468	33.99581	...	5.48708	-9.13495	6.08680
9	1	45.17809	46.34234	...	11.49326	-89.21804	-15.09719
10	1	39.13076	-23.01763	...	-2.59543	109.19723	23.36143
11	1	37.66498	-34.05910	...	-5.19009	8.83617	-7.16056
12	1	26.51957	-148.15762	...	23.00230	-164.02536	51.54138
13	1	37.68491	-26.84185	...	14.11648	-1030.99180	99.28967
14	0	39.11695	-8.29767	...	-2.14942	-211.48202	-12.81569
15	1	35.05129	-67.97714	...	20.73063	-562.07671	43.44696
16	1	33.63129	-96.14912	...	3.44539	259.10825	10.28525
17	0	41.38639	-20.78665	...	4.44627	58.16913	-0.02409
18	0	37.45034	11.42615	...	25.73235	157.22967	38.70617
19	0	39.71092	-4.92800	...	2.34002	-31.57015	1.58400

[20 rows x 91 columns]

### 1.1.1 Part (a) -- 7%

The data set description text asks us to respect the below train/test split to avoid the "producer effect". That is, we want to make sure that no song from a single artist ends up in both the training and test set.

Explain why it would be problematic to have some songs from an artist in the training set, and other songs from the same artist in the test set. (Hint: Remember that we want our test accuracy to predict how well the model will perform in practice on a song it hasn't learned about.)

```
[8]: df_train = df[:463715]
df_test = df[463715:]

# convert to numpy
train_xs = df_train[x_labels].to_numpy()
train_ts = df_train[t_label].to_numpy()
test_xs = df_test[x_labels].to_numpy()
test_ts = df_test[t_label].to_numpy()

# Write your explanation here
'''
To evaluate our model properly, we want to know how well the model will perform
→ over a new songs that it hasn't encountered before.
Using songs from the same artist for both the training set and the test set can
→ lead to a situation in which the model learns to extract features based on
→ the singer,
e.g., the singer's voice, and will decide based on those features, instead
→ based of the century.
In this scenario, the model might perform well on samples of songs in the test
→ where their corresponding singers also appeared in the train set.
But this doesn't guarantee that the model will generalize for songs in the test
→ set that their singers didn't appear in the train set.
```

```
'''
```

[8]: "\nTo evaluate our model properly, we want to know how well the model will perform over a new songs that it hasn't encountered before.\nUsing songs from the same artist for both the training set and the test set can lead to a situation in which the model learns to extract features based on the singer,\ne.g., the singer's voice, and will decide based on those features, instead based of the century.\n\nIn this scenario, the model might perform well on samples of songs in the test where their corresponding singers also appeared in the train set.\nBut this doesn't guarantee that the model will generalize for songs in the test set that their singers didn't appear in the train set.\n"

### 1.1.2 Part (b) -- 7%

It can be beneficial to **normalize** the columns, so that each column (feature) has the *same* mean and standard deviation.

```
[9]: feature_means = df_train.mean()[1:].to_numpy() # the [1:] removes the mean of
      feature_stds  = df_train.std()[1:].to_numpy()
      # the "year" field

train_norm_xs = (train_xs - feature_means) / feature_stds
test_norm_xs  = (test_xs - feature_means) / feature_stds
```

Notice how in our code, we normalized the test set using the *training data means and standard deviations*. This is *not* a bug.

Explain why it would be improper to compute and use test set means and standard deviations. (Hint: Remember what we want to use the test accuracy to measure.)

```
[10]: # Write your explanation here
      """
      We normalize the data with respect to the training set to make the input data
      used for training "more identical distributed" during training.
      We are using test data to evaluate our model performance over new unseen data.
      If we are using the average and standard deviation of this test set to
      normalize it, we enforce dependency between different samples in the test
      set.
      This harms the idea of model evaluation over new unseen data, e.g., if we
      evaluate our model for two different test sets, both contain some identical
      samples.
      Following this approach will lead to different outcomes for the same input
      sample (because the average and the standard deviation will be different),
      which demonstrates the problem with this approach.
      In addition, in real-time, there is usually a single input and not a set.
      """
```

```
[10]: '\nWe normalize the data with respect to the training set to make the input data
used for training "more identical distributed" during training.\nWe are using
test data to evaluate our model performance over new unseen data.\nIf we are
using the average and standard deviation of this test set to normalize it, we
```

enforce dependency between different samples in the test set. \nThis harms the idea of model evaluation over new unseen data, e.g., if we evaluate our model for two different test sets, both contain some identical samples.\nFollowing this approach will lead to different outcomes for the same input sample (because the average and the standard deviation will be different), which demonstrates the problem with this approach.\nIn addition, in real-time, there is usually a single input and not a set.\n'

### 1.1.3 Part (c) -- 7%

Finally, we'll move some of the data in our training set into a validation set.

Explain why we should limit how many times we use the test set, and that we should use the validation set during the model building process.

```
[11]: # shuffle the training set
reindex = np.random.permutation(len(train_xs))
train_xs = train_xs[reindex]
train_norm_xs = train_norm_xs[reindex]
train_ts = train_ts[reindex]

# use the first 50000 elements of `train_xs` as the validation set
train_xs, val_xs = train_xs[50000:], train_xs[:50000]
train_norm_xs, val_norm_xs = train_norm_xs[50000:], train_norm_xs[:50000]
train_ts, val_ts = train_ts[50000:], train_ts[:50000]

# Write your explanation here
'''
If we will not limit the times that we are using the test set, we may encounter
    ↳overfitting phenomena over the test data,
which in turn will lead to good performance over the test set, this phenomena
    ↳may occur because we are tuning our model to achieve
good performance over a specific set of test samples.
This is a problem because the main reason that we are using a test set is to
    ↳evaluate our model performance over an unseen data,
i.e., we test how well the training procedure will be able to generalize for
    ↳unseen samples.
We can avoid this issue by examining our model performance (during the building
    ↳procedure) via
another set which we call "validation" set. The validation dataset is different
    ↳from the test dataset,
both datasets are held back from the training of the model, but the validation
    ↳set used to give an unbiased performance estimation
of the final tuned model when comparing between different final models."
'''
```

```
[11]: '\nIf we will not limit the times that we are using the test set, we may
encounter overfitting phenomena over the test data, \nwhich in turn will lead to
```

good performance over the test set, this phenomena may occur because we are tuning our model to achieve good performance over a specific set of test samples. This is a problem because the main reason that we are using a test set is to evaluate our model performance over an unseen data, i.e., we test how well the training procedure will be able to generalize for unseen samples. We can avoid this issue by examining our model performance (during the building procedure) via another set which we call "validation" set. The validation dataset is different from the test dataset, both datasets are held back from the training of the model, but the validation set used to give an unbiased performance estimation of the final tuned model when comparing between different final models."

## 1.2 Part 2. Classification (79%)

We will first build a *classification* model to perform decade classification. These helper functions are written for you. All other code that you write in this section should be vectorized whenever possible (i.e., avoid unnecessary loops).

```
[12]: def sigmoid(z):
    return 1 / (1 + np.exp(-z))

def cross_entropy(t, y):
    epsilon=1e-15 #to avoid log(0)
    return -t * np.log(y+epsilon) - (1 - t) * np.log(1 - y+epsilon)

def cost(y, t):
    return np.mean(cross_entropy(t, y))

def get_accuracy(y, t):
    acc = 0
    N = 0
    for i in range(len(y)):
        N += 1
        if (y[i] >= 0.5 and t[i] == 1) or (y[i] < 0.5 and t[i] == 0):
            acc += 1
    return acc / N
```

### 1.2.1 Part (a) -- 7%

Write a function `pred` that computes the prediction  $y$  based on logistic regression, i.e., a single layer with weights  $w$  and bias  $b$ . The output is given by:

$$y = \sigma(\mathbf{w}^T \mathbf{x} + b), \quad (1)$$

where the value of  $y$  is an estimate of the probability that the song is released in the current century, namely  $\text{year} = 1$ .

```
[13]: def pred(w, b, X):
    """
    Returns the prediction `y` of the target based on the weights `w` and scalar
    bias `b`.

    Preconditions: np.shape(w) == (90,)
                   type(b) == float
                   np.shape(X) = (N, 90) for some N

    >>> pred(np.zeros(90), 1, np.ones([2, 90]))
    array([0.73105858, 0.73105858]) # It's okay if your output differs in the
    →last
    decimals
    """
    return sigmoid(np.dot(X,w)+b)
```

### 1.2.2 Part (b) -- 7%

Write a function `derivative_cost` that computes and returns the gradients  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$  and  $\frac{\partial \mathcal{L}}{\partial b}$ . Here,  $\mathbf{X}$  is the input,  $\mathbf{y}$  is the prediction, and  $\mathbf{t}$  is the true label.

```
[14]: def derivative_cost(X, y, t):
    """
    Returns a tuple containing the gradients dLdw and dLdb.

    Precondition: np.shape(X) == (N, 90) for some N
                  np.shape(y) == (N,)
                  np.shape(t) == (N,)

    Postcondition: np.shape(dLdw) = (90,)
                   type(dLdb) = float
    """
    N = np.shape(y)[0]
    dLdb = np.mean(y - t)
    dLdw = 1/N * np.dot(y - t, X)
    return (dLdw, dLdb)
```

## 2 Explanation on Gradients

First we define our loss function as followed:

$$L(w, b) = -\frac{1}{N} \sum_{n=1}^N t_n \log(y_n) + (1 - t_n) \log(1 - y_n).$$

$$\text{Let } z_n = \mathbf{w}^T \mathbf{x}_n + b$$

where we denote bold letter as a vector.

By the chain-rule:

$$\frac{\partial L}{\partial w} = \sum_{n=1}^N \frac{\partial L}{\partial y_n} \cdot \frac{\partial y_n}{\partial z_n} \cdot \frac{\partial z_n}{\partial w}$$

$$\frac{\partial L}{\partial b} = \sum_{n=1}^N \frac{\partial L}{\partial y_n} \cdot \frac{\partial y_n}{\partial z_n} \cdot \frac{\partial z_n}{\partial b}$$

where:

$$\frac{\partial L}{\partial y_n} = -\frac{1}{N} \left( \frac{t_n}{y_n} - \frac{1-t_n}{1-y_n} \right) = \frac{y_n - t_n}{N \cdot y_n \cdot (1-y_n)},$$

$$\frac{\partial y_n}{\partial z_n} = \sigma'(z_n) = \sigma(z_n)(1 - \sigma(z_n)) = y_n(1 - y_n)$$

$$\frac{\partial z_n}{\partial b} = 1,$$

$$\frac{\partial z_n}{\partial w} = \mathbf{x}_n$$

These boils down to:

$$\frac{\partial L}{\partial b} = \sum_{n=1}^N \frac{y_n - t_n}{N \cdot y_n \cdot (1-y_n)} \cdot y_n(1 - y_n) \cdot 1 = \frac{1}{N} \sum_{n=1}^N y_n - t_n$$

$$\frac{\partial L}{\partial w} = \sum_{n=1}^N \frac{y_n - t_n}{N \cdot y_n \cdot (1-y_n)} \cdot y_n(1 - y_n) \cdot \mathbf{x}_n = \frac{1}{N} \sum_{n=1}^N (y_n - t_n) \mathbf{x}_n$$

### 2.0.1 Part (c) -- 7%

We can check that our derivative is implemented correctly using the finite difference rule. In 1D, the finite difference rule tells us that for small  $h$ , we should have

$$\frac{f(x+h) - f(x)}{h} \approx f'(x)$$

Show that  $\frac{\partial \mathcal{L}}{\partial b}$  is implemented correctly by comparing the result from `derivative_cost` with the empirical cost derivative computed using the above numerical approximation.

```
[15]: # Your code goes here
def compute_analytical_derivate(w,b,t,X,h=1e-05):
    '''
    First we calculate the derivative with respect to the bias term
    '''
    # Calculate the predications
    y = pred(w, b, X)
    # Calculate the loss of those predications
    L = cost(y,t)
    y_new = pred(w,b+h,X)
    L_new = cost(y_new,t)

    # Calculate the analytic derivative
    derivate_b = (L_new-L)/h

    '''
    Now we calculate the derivative with respect to the weights
    '''
    L_new = np.zeros(len(w))
    derivate_w = np.zeros(len(w))
    for weight_idx in range(len(w)):
        # Compute cost of w[weight_idx] + h
        w_new = np.copy(w)
        w_new[weight_idx] = w_new[weight_idx] + h
        # Calculate the predications
        y_new = pred(w_new,b,X)
        # Calculate the loss of those predications
        L_new[weight_idx]=cost(y_new,t)
```



```

    deriviate_w[weight_idx] = (L_new[weight_idx]-L)/h

    return deriviate_w, deriviate_b

```

Now we are going to check our derivative using an example

```

[16]: w=np.zeros(90)
      h=1e-05
      b=1
      t=np.array([1,1])
      X=np.ones([2, 90])
      y = pred(w, b, X)
      dw,db = derivative_cost(X,y,t)
      dw_analytic, db_analytic = compute_analytical_deriviate(w,b,t,X,h)
      print("The bias analytical derivative results is -", db_analytic)
      print("The bias algorithm derivative results is - ", db)
      print(f"The error is {np.abs(db-db_analytic)}")

```

```

The bias analytical derivative results is - -0.26894043830827385
The bias algorithm derivative results is - -0.2689414213699951
The error is 9.830617212491788e-07

```

## 2.0.2 Part (d) -- 7%

Show that  $\frac{\partial \mathcal{L}}{\partial w}$  is implement correctly.

```

[17]: print("The weights analytical derivative results is -", dw_analytic)
      print("The weights algorithm derivative results is - ", dw)
      print(f"The maximal error is {np.max(np.abs(dw-dw_analytic))}")

```

```

The weights analytical derivative results is - [-0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044
-0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044 -0.26894044]
The weights algorithm derivative results is - [-0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142]

```

```

-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142
-0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142 -0.26894142

```

The maximal error is 9.830617212491788e-07

Now we will check over couple of real examples from the training set

```

[18]: h = 1e-9;
      # First we initialize the weights by random numbers
      w=np.random.rand(90)
      b=np.random.rand(1)

      # Now we extract the actual samples from the dataset
      t = (train_ts[:1]).squeeze()
      X = train_norm_xs[:1]
      y = pred(w, b, X)

      dw,db = derivative_cost(X,y,t)
      dw_analytic, db_analytic = compute_analytical_derivate(w,b,t,X,h)
      print("The analytical derivative bias results is -", db_analytic)
      print("The algorithm derivative bias results is - ", db)
      print(f"Bias Error is {np.abs(db-db_analytic)}")

      print("The analytical derivative weights results is -", dw_analytic)
      print("The algorithm derivative weights results is - ", dw)
      print(f"Maximal weights Error is {np.max(np.abs(dw-dw_analytic))}")

```

The analytical derivative bias results is - -0.9672760370449395

The algorithm derivative bias results is - -0.9672759967800327

Bias Error is 4.0264906742137896e-08

The analytical derivative weights results is - [ 0.97880459 1.33712197  
1.5150321 1.17287824 -0.1407332 -0.28086733

```

0.99645137  1.66069425  0.34143222  2.53612864 -0.60000049  0.2660081
-0.16358959  0.27339908 -0.18315438 -0.02620393  0.57133454  0.29290792
-0.22191848 -0.26613023  0.32624214 -0.12356471  0.2247269 -0.1315108
0.79538109  0.54556981  0.37574432  0.19421131  0.50898086  0.31249137
-0.95565378 -0.54595706  0.15626922  0.11944934 -0.86183682 -0.96469277
-0.30012703 -0.25347191 -0.2640399 -0.20129987  0.12283019  0.02585665
-0.75847728  0.77280982  0.88499963 -0.2634799  0.25427704  0.86388496

```

```

-0.74310602 -0.05215162 0.02074074 -0.18614754 0.64476735 0.68632078
-0.79160367 -0.5105103 -0.51600768 0.62609251 -0.36517989 1.05336495
0.91098018 0.16509771 0.39282799 0.36208059 -1.05262732 0.20108226
0.01534106 -0.28936853 0.35263614 -0.85813934 0.51580074 -0.01067191
-0.41914161 0.52804738 -0.2147611 -0.83647089 0.07340262 -1.70637771
-0.25002356 0.29500091 0.20058133 -0.08705836 0.93793728 0.12571144
-1.7977615 -0.31657921 0.64228134 0.73033668 -0.51113247 0.58821525]
The algorithm derivative weights results is - [ 0.97880513 1.33712198
1.51503252 1.17287858 -0.14073285 -0.28086701
0.99645136 1.66069465 0.34143218 2.53612882 -0.60000029 0.26600803
-0.16358953 0.27339904 -0.18315388 -0.02620335 0.57133477 0.29290862
-0.22191785 -0.26613003 0.3262425 -0.12356409 0.22472752 -0.13151021
0.79538094 0.54556964 0.37574461 0.19421122 0.5089815 0.31249165
-0.95565349 -0.54595737 0.15626943 0.11944995 -0.86183685 -0.96469248
-0.30012698 -0.25347199 -0.26404002 -0.20129932 0.12282997 0.02585716
-0.75847651 0.77280993 0.88499994 -0.26347972 0.25427681 0.86388554
-0.74310507 -0.05215134 0.02074135 -0.18614762 0.64476718 0.68632109
-0.79160365 -0.5105102 -0.51600786 0.62609255 -0.36517956 1.05336573
0.91098042 0.16509771 0.39282848 0.3620809 -1.05262721 0.20108315
0.01534155 -0.28936764 0.35263587 -0.8581389 0.51580074 -0.01067161
-0.41914165 0.52804727 -0.21476141 -0.83647112 0.07340269 -1.70637775
-0.25002302 0.29500102 0.20058141 -0.08705784 0.93793721 0.12571158
-1.79776144 -0.31657967 0.64228073 0.73033717 -0.51113169 0.58821569]
Maximal weights Error is 9.504835127849276e-07

```

### 2.0.3 Part (e) -- 7%

Now that you have a gradient function that works, we can actually run gradient descent. Complete the following code that will run stochastic gradient descent training:

```

[19]: def run_gradient_descent(w0, b0, mu=0.1, batch_size=100, max_iters=100):
    """Return the values of (w, b) after running gradient descent for max_iters.
    We use:
        - train_norm_xs and train_ts as the training set
        - val_norm_xs and val_ts as the test set
        - mu as the learning rate
        - (w0, b0) as the initial values of (w, b)

    Precondition: np.shape(w0) == (90,)
                  type(b0) == float

    Postcondition: np.shape(w) == (90,)
                   type(b) == float

    """
    w = w0
    b = b0
    iter = 0
    val_loss = []

```

```

while iter < max_iters:
    # shuffle the training set (there is code above for how to do this)
    new_indices = np.random.permutation(len(train_norm_xs))
    train_norm_xs_temp = train_norm_xs[new_indices]
    train_ts_temp = train_ts[new_indices]

    for i in range(0, len(train_norm_xs), batch_size): # iterate over each
→minibatch
        # minibatch that we are working with:
        X = train_norm_xs_temp[i:(i + batch_size)]
        t = train_ts_temp[i:(i + batch_size), 0]

        # since len(train_norm_xs) does not divide batch_size evenly, we will
→skip over
        # the "last" minibatch
        if np.shape(X)[0] != batch_size:
            continue

        # compute the prediction
        y = pred(w, b, X)
        # update w and b
        dw,db = derivative_cost(X,y,t)
        w -= mu*dw
        b -= mu*db

        # increment the iteration count
        iter += 1
        # compute and print the *validation* loss and accuracy
        if (iter % 10 == 0):
            # Calculate the predications over the validation data
            y_val = pred(w,b,val_norm_xs)
            # Calculate the cost over the predications
            val_cost = cost(y_val,val_ts.squeeze())
            val_acc = get_accuracy(y_val, val_ts)
            print("Iter %d. [Val Acc %.0f%%, Loss %f]" % (iter, val_acc * 100,
→val_cost))
            val_loss.append(val_cost)
        if iter >= max_iters:
            break

        # Think what parameters you should return for further use
        # We choose to return in addition to w and b the validation loss in order
→to evaluate our model training procedure over unseen data.
    return w, b, val_loss

```

#### 2.0.4 Part (f) -- 7%

Call `run_gradient_descent` with the weights and biases all initialized to zero. Show that if the learning rate  $\mu$  is too small, then convergence is slow. Also, show that if  $\mu$  is too large, then the optimization algorithm does not converge. The demonstration should be made using plots showing these effects.

```
[20]: w0 = np.zeros(90)
      b0 = 0
      Iterations = 100
      mus=[0.005,0.1,5]
      for mu in mus:
          print("mu=",mu,":")
          w,b,losses=run_gradient_descent(w0,b0,mu, max_iters=Iterations)
          print("-----")
          plt.semilogy(range(0,Iterations,10),losses)

      plt.legend(mus,loc='best')
      plt.xlabel("Iteration")
      plt.ylabel("Loss")
      plt.grid()
      plt.title("Learning rate convergance properties")
```

```
mu= 0.005 :
Iter 10. [Val Acc 64%, Loss 0.689537]
Iter 20. [Val Acc 64%, Loss 0.686243]
Iter 30. [Val Acc 65%, Loss 0.683223]
Iter 40. [Val Acc 66%, Loss 0.680522]
Iter 50. [Val Acc 66%, Loss 0.677757]
Iter 60. [Val Acc 66%, Loss 0.675389]
Iter 70. [Val Acc 66%, Loss 0.673272]
Iter 80. [Val Acc 66%, Loss 0.670915]
Iter 90. [Val Acc 66%, Loss 0.668803]
Iter 100. [Val Acc 67%, Loss 0.666691]
```

```
-----
mu= 0.1 :
Iter 10. [Val Acc 67%, Loss 0.640029]
Iter 20. [Val Acc 68%, Loss 0.626906]
Iter 30. [Val Acc 69%, Loss 0.614803]
Iter 40. [Val Acc 70%, Loss 0.608335]
Iter 50. [Val Acc 70%, Loss 0.601604]
Iter 60. [Val Acc 70%, Loss 0.600578]
Iter 70. [Val Acc 71%, Loss 0.593426]
Iter 80. [Val Acc 71%, Loss 0.589205]
Iter 90. [Val Acc 70%, Loss 0.591799]
Iter 100. [Val Acc 71%, Loss 0.584431]
```

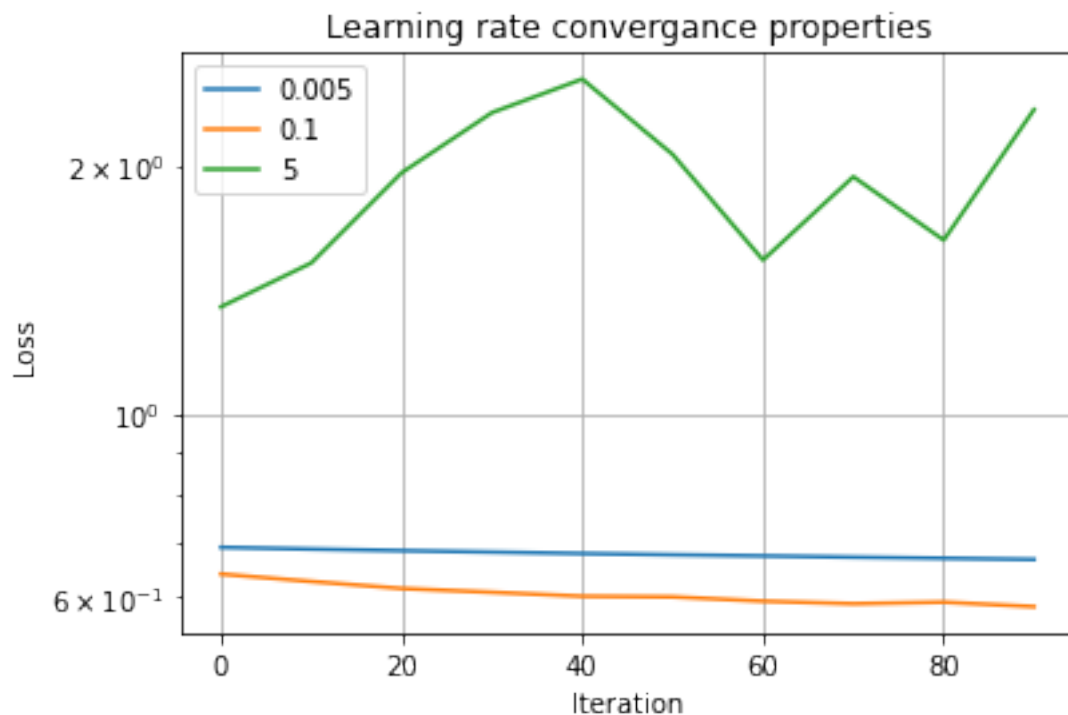
```
-----
mu= 5 :
Iter 10. [Val Acc 67%, Loss 1.348408]
```

```

Iter 20. [Val Acc 60%, Loss 1.523874]
Iter 30. [Val Acc 61%, Loss 1.959523]
Iter 40. [Val Acc 64%, Loss 2.317599]
Iter 50. [Val Acc 64%, Loss 2.544981]
Iter 60. [Val Acc 65%, Loss 2.062478]
Iter 70. [Val Acc 69%, Loss 1.536574]
Iter 80. [Val Acc 68%, Loss 1.938252]
Iter 90. [Val Acc 65%, Loss 1.624422]
Iter 100. [Val Acc 61%, Loss 2.337443]
-----

```

[20]: Text(0.5, 1.0, 'Learning rate convergence properties')



We conclude from the figure above that for a small learning rate, the algorithm may be capable to converge but very slowly due to the small changes in the parameters space. On the other hand, we see that by using too large learning rate, the algorithm is incapable to converge due to the large changes in the parameters space. Thus, we conclude from the empirical results that choosing a learning rate somewhere in between those learning rates leads to the best performance. We explain this phenomenon by the fact that the algorithm can change its parameters efficiently towards a good optimum. Meaning that we have found a good trade-off between a small enough step size that allows convergence but not too small such that it leads to faster convergence rate to a good enough optimum.

## 2.0.5 Part (g) -- 7%

Find the optimal value of  $w$  and  $b$  using your code. Explain how you chose the learning rate  $\mu$  and the batch size. Show plots demonstrating good and bad behaviours.

```
[21]: '''
Hyperparameters search
'''

mus = [0.1, 0.5, 1]
batch_sizes = [64, 512, 1024]
Iterations = 100
loss_results = np.zeros((len(mus),len(batch_sizes)))
w_history = np.zeros((len(mus),len(batch_sizes), 90))
b_history = np.zeros((len(mus),len(batch_sizes), 1))

w0 = np.random.randn(90)
b0 = np.random.randn(1)[0]

for mu_idx,mu in enumerate(mus):
    for minibatch_idx, minibatch_size in enumerate(batch_sizes):
        print("mu=",mu, 'batch_size=', minibatch_size, ':')
        w,b,losses=run_gradient_descent(w0,b0,mu, max_iters=Iterations)
        print("-----")
        plt.plot(list(range(0,Iterations,10)), losses, label=r'$\mu$='+str(mu)+r',
        ↳$BS$='+str(minibatch_size))
        loss_results[mu_idx][minibatch_idx] = losses[-1]
        w_history[mu_idx][minibatch_idx] = w
        b_history[mu_idx][minibatch_idx] = b
plt.legend()
plt.grid()
plt.xlabel('Iteration')
plt.ylabel('Label')
plt.title('Loss VS Iteration for various learning rates and batch sizes')
plt.ylim([0.5,2.5])
best_mu_idx, best_batch_size_idx = np.unravel_index(np.argmin(loss_results,
↳axis=None), loss_results.shape)
print("optimal mu=",mus[best_mu_idx], 'optimal batch_size=',
↳batch_sizes[best_batch_size_idx])
'''

Save the optimal weight and bias
'''

opt_w = w_history[best_mu_idx][best_batch_size_idx]
opt_b = b_history[best_mu_idx][best_batch_size_idx]
```

```
mu= 0.1 batch_size= 64 :
Iter 10. [Val Acc 48%, Loss 2.504536]
Iter 20. [Val Acc 48%, Loss 2.269788]
Iter 30. [Val Acc 49%, Loss 2.106893]
```

```
Iter 40. [Val Acc 50%, Loss 1.962389]
Iter 50. [Val Acc 51%, Loss 1.824749]
Iter 60. [Val Acc 52%, Loss 1.711386]
Iter 70. [Val Acc 53%, Loss 1.615538]
Iter 80. [Val Acc 54%, Loss 1.521306]
Iter 90. [Val Acc 55%, Loss 1.436870]
Iter 100. [Val Acc 55%, Loss 1.372454]
```

-----  
mu= 0.1 batch\_size= 512 :

```
Iter 10. [Val Acc 57%, Loss 1.297985]
Iter 20. [Val Acc 57%, Loss 1.237497]
Iter 30. [Val Acc 58%, Loss 1.182943]
Iter 40. [Val Acc 58%, Loss 1.140016]
Iter 50. [Val Acc 59%, Loss 1.088942]
Iter 60. [Val Acc 59%, Loss 1.044982]
Iter 70. [Val Acc 60%, Loss 1.009848]
Iter 80. [Val Acc 60%, Loss 0.978319]
Iter 90. [Val Acc 61%, Loss 0.948770]
Iter 100. [Val Acc 62%, Loss 0.920177]
```

-----  
mu= 0.1 batch\_size= 1024 :

```
Iter 10. [Val Acc 62%, Loss 0.901723]
Iter 20. [Val Acc 63%, Loss 0.875667]
Iter 30. [Val Acc 63%, Loss 0.856669]
Iter 40. [Val Acc 64%, Loss 0.838742]
Iter 50. [Val Acc 64%, Loss 0.821935]
Iter 60. [Val Acc 65%, Loss 0.804251]
Iter 70. [Val Acc 65%, Loss 0.791927]
Iter 80. [Val Acc 65%, Loss 0.773395]
Iter 90. [Val Acc 66%, Loss 0.763514]
Iter 100. [Val Acc 66%, Loss 0.753636]
```

-----  
mu= 0.5 batch\_size= 64 :

```
Iter 10. [Val Acc 66%, Loss 0.727669]
Iter 20. [Val Acc 67%, Loss 0.709271]
Iter 30. [Val Acc 68%, Loss 0.693618]
Iter 40. [Val Acc 69%, Loss 0.660371]
Iter 50. [Val Acc 69%, Loss 0.653615]
Iter 60. [Val Acc 69%, Loss 0.643698]
Iter 70. [Val Acc 70%, Loss 0.639843]
Iter 80. [Val Acc 71%, Loss 0.617782]
Iter 90. [Val Acc 70%, Loss 0.634432]
Iter 100. [Val Acc 71%, Loss 0.609879]
```

-----  
mu= 0.5 batch\_size= 512 :

```
Iter 10. [Val Acc 71%, Loss 0.611268]
Iter 20. [Val Acc 71%, Loss 0.608491]
Iter 30. [Val Acc 72%, Loss 0.599843]
```



```
Iter 40. [Val Acc 72%, Loss 0.601804]
Iter 50. [Val Acc 72%, Loss 0.591962]
Iter 60. [Val Acc 72%, Loss 0.589068]
Iter 70. [Val Acc 72%, Loss 0.603013]
Iter 80. [Val Acc 72%, Loss 0.585173]
Iter 90. [Val Acc 71%, Loss 0.596688]
Iter 100. [Val Acc 71%, Loss 0.599163]
```

```
-----
mu= 0.5 batch_size= 1024 :
```

```
Iter 10. [Val Acc 70%, Loss 0.609293]
Iter 20. [Val Acc 72%, Loss 0.582921]
Iter 30. [Val Acc 71%, Loss 0.590241]
Iter 40. [Val Acc 72%, Loss 0.582173]
Iter 50. [Val Acc 72%, Loss 0.590123]
Iter 60. [Val Acc 71%, Loss 0.592309]
Iter 70. [Val Acc 71%, Loss 0.590857]
Iter 80. [Val Acc 72%, Loss 0.589254]
Iter 90. [Val Acc 72%, Loss 0.580167]
Iter 100. [Val Acc 72%, Loss 0.595504]
```

```
-----
mu= 1 batch_size= 64 :
```

```
Iter 10. [Val Acc 69%, Loss 0.643439]
Iter 20. [Val Acc 71%, Loss 0.619565]
Iter 30. [Val Acc 70%, Loss 0.629722]
Iter 40. [Val Acc 71%, Loss 0.639629]
Iter 50. [Val Acc 70%, Loss 0.634957]
Iter 60. [Val Acc 71%, Loss 0.635613]
Iter 70. [Val Acc 68%, Loss 0.695302]
Iter 80. [Val Acc 71%, Loss 0.637624]
Iter 90. [Val Acc 70%, Loss 0.624395]
Iter 100. [Val Acc 67%, Loss 0.686641]
```

```
-----
mu= 1 batch_size= 512 :
```

```
Iter 10. [Val Acc 69%, Loss 0.663816]
Iter 20. [Val Acc 70%, Loss 0.650512]
Iter 30. [Val Acc 69%, Loss 0.657354]
Iter 40. [Val Acc 72%, Loss 0.610172]
Iter 50. [Val Acc 71%, Loss 0.600712]
Iter 60. [Val Acc 72%, Loss 0.600045]
Iter 70. [Val Acc 71%, Loss 0.630042]
Iter 80. [Val Acc 72%, Loss 0.616479]
Iter 90. [Val Acc 70%, Loss 0.631277]
Iter 100. [Val Acc 70%, Loss 0.633159]
```

```
-----
mu= 1 batch_size= 1024 :
```

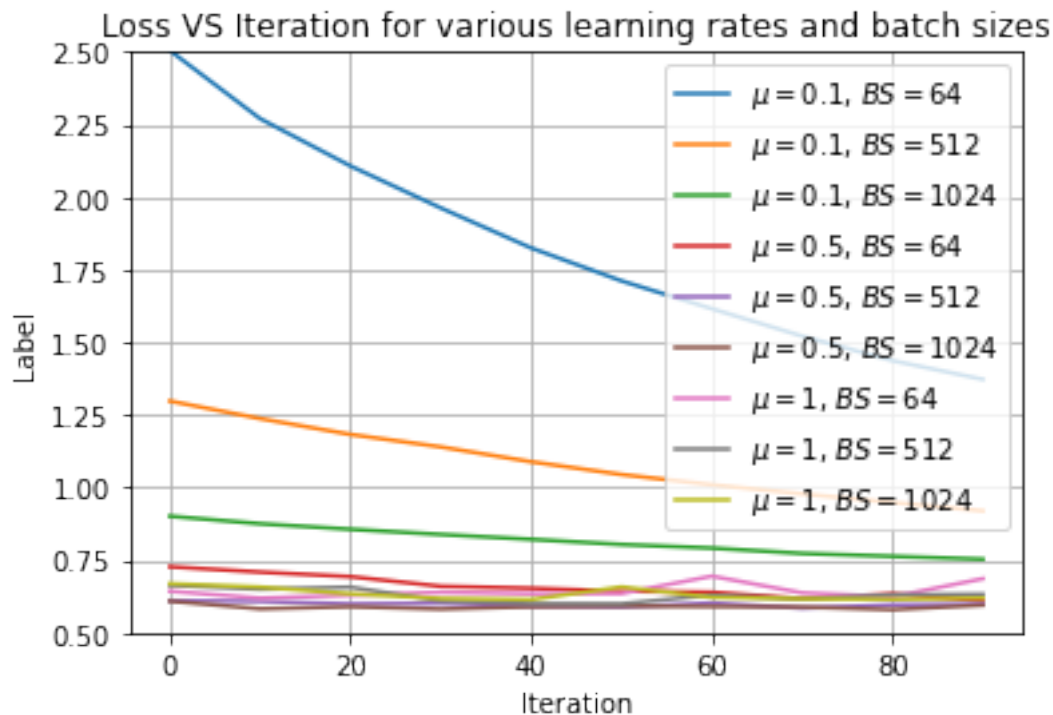
```
Iter 10. [Val Acc 68%, Loss 0.669234]
Iter 20. [Val Acc 69%, Loss 0.657841]
Iter 30. [Val Acc 70%, Loss 0.634529]
```

```

Iter 40. [Val Acc 71%, Loss 0.618299]
Iter 50. [Val Acc 71%, Loss 0.614596]
Iter 60. [Val Acc 69%, Loss 0.658339]
Iter 70. [Val Acc 70%, Loss 0.622178]
Iter 80. [Val Acc 71%, Loss 0.617426]
Iter 90. [Val Acc 70%, Loss 0.615538]
Iter 100. [Val Acc 71%, Loss 0.619879]

```

-----  
optimal mu= 0.5 optimal batch\_size= 1024



We found the optimal weights and biases by scanning over a range of different learning rates and batch sizes for 100 iterations each and then picking the weights and biases in which the loss over the validation set was minimal. The best weights and biases were observed when  $\mu = 0.5$  and batch\_size = 1024. The best results are achieved when the batch size is maximal. The reason is that the number of iterations is the number of update steps. Because this number is constant (100), we have the same number of updates independently of the batch size. Thus, when the batch size is large, as we learned in the lectures, we get a better approximation of the true gradient, and hence the results are better. It is important to mention that we should not conclude that the best choice is always picking the largest batch size. The reason is that if we measure the validation loss per epoch (for a fixed number of epochs), the number of update steps with a large batch size will be smaller than the number of updates when the batch size is smaller which may lead to worse performance. Furthermore, a smaller batch size can assist the algorithm in avoiding it from stuck in a local minimum or a saddle point due to the noisier gradient.

### 2.0.6 Part (h) -- 15%

Using the values of  $w$  and  $b$  from part (g), compute your training accuracy, validation accuracy, and test accuracy. Are there any differences between those three values? If so, why?

```
[22]: train_ys=pred(opt_w,opt_b,train_norm_xs)
      test_ys=pred(opt_w,opt_b,test_norm_xs)
      val_ys=pred(opt_w,opt_b,val_norm_xs)

      train_acc = get_accuracy(train_ys,train_ts)
      val_acc = get_accuracy(val_ys,val_ts)
      test_acc = get_accuracy(test_ys,test_ts)

      print('train_acc = ', train_acc, ' val_acc = ', val_acc, ' test_acc = ',
            →test_acc)
```

```
train_acc = 0.7107259828625987  val_acc = 0.71514  test_acc =
0.7072825876428434
```

Both train and validation accuracies are better than the test accuracy, we explain this phenomena by the fact that we didn't see the test data during the model building process, so we expect to encounter minor degradation in the model performance over the new unseen data.

Moreover, we can see that the gap between the train and test data is not significant and therefore we conclude that the model can extract meaningful information from the training dataset. Meaning that the model is capable to perform similarly over unseen data, e.g., the model is well generalized.

Regarding the validation accuracy, we have tuned our model hyperparameters (Learning rate and Batch size) such that it achieves the best performance over the validation data, therefore, we shall expect to achieve better results in the validation data than the test data, and this is indeed the result that we've achieved.

### 2.0.7 Part (i) -- 15%

Writing a classifier like this is instructive, and helps you understand what happens when we train a model. However, in practice, we rarely write model building and training code from scratch. Instead, we typically use one of the well-tested libraries available in a package.

Use `sklearn.linear_model.LogisticRegression` to build a linear classifier, and make predictions about the test set. Start by reading the [API documentation here](#).

Compute the training, validation and test accuracy of this model.

```
[23]: import sklearn.linear_model

      model = sklearn.linear_model.LogisticRegression()
      model.fit(train_norm_xs, train_ts.squeeze())
      train_ys = model.predict(train_norm_xs)
      test_ys = model.predict(test_norm_xs)
      val_ys = model.predict(val_norm_xs)
      train_acc = get_accuracy(train_ys,train_ts)
      val_acc = get_accuracy(val_ys,val_ts)
```

```
test_acc = get_accuracy(test_ys, test_ts)

print('train_acc = ', train_acc, ' val_acc = ', val_acc, ' test_acc = ',
      test_acc)
```

```
train_acc = 0.7323979067715698 val_acc = 0.73642 test_acc =
0.7265736974627155
```

This part helps by checking if the code worked. Check if you get similar results, if not repair your code

### Covertng the Notebook to PDF

```
[!]: !apt-get install texlive texlive-xetex texlive-latex-extra pandoc
      !pip install py pandoc
```

```
Reading package lists... Done
Building dependency tree
Reading state information... Done
pandoc is already the newest version (1.19.2.4~dfsg-1build4).
pandoc set to manually installed.
texlive-latex-extra is already the newest version (2017.20180305-2).
texlive-latex-extra set to manually installed.
texlive-xetex is already the newest version (2017.20180305-1).
The following NEW packages will be installed:
  texlive
0 upgraded, 1 newly installed, 0 to remove and 37 not upgraded.
Need to get 14.4 kB of archives.
After this operation, 70.7 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive all
2017.20180305-1 [14.4 kB]
Fetched 14.4 kB in 1s (9,916 B/s)
```

```
[29]: from google.colab import drive
      drive.mount('/content/drive')
```

Mounted at /content/drive

```
[30]: !cp /content/drive/My\ Drive/Colab\ Notebooks/Assignment1.ipynb ./
```

```
[34]: !jupyter nbconvert --to PDF "Assignment1.ipynb"
```

```
[NbConvertApp] Converting notebook Assignment1.ipynb to PDF
[NbConvertApp] Support files will be in Assignment1_files/
[NbConvertApp] Making directory ./Assignment1_files
[NbConvertApp] Making directory ./Assignment1_files
[NbConvertApp] Writing 89561 bytes to ./notebook.tex
[NbConvertApp] Building PDF
```

```
[NbConvertApp] Running xelatex 3 times: [u'xelatex', u'./notebook.tex',  
'-quiet']  
[NbConvertApp] Running bibtex 1 time: [u'bibtex', u'./notebook']  
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no  
citations  
[NbConvertApp] PDF successfully created  
[NbConvertApp] Writing 129571 bytes to Assignment1.pdf
```