# ECE521 Winter 2017: Assignment 1

FuYuan Tee, (999295837) [*]     Chee Loong Soon, (999295793) [†]

February 8th, 2017

## Contents

---

[*]Equal Contribution (50%), fuyuan.tee@mail.utoronto.ca

[†]Equal Contribution (50%), cheeloong.soon@mail.utoronto.ca

1

# 1 k-Nearest Neighbour

## 1.1 Geometry of k-NN

### 1.1.1 Describe 1D Dataset

An example of a 1-D dataset with two classes in which k-NN produces an accuracy that is periodic with k is illustrated in Figure 1. The data point comes from the training set itself. The data points are equally distant from each adjacent data point.



Figure 1: 1-D Dataset Illustration with classification accuracy that is periodic to $k$.

For such a dataset, the classification accuracy of the data point can be summarised in the Table 1 below. The accuracy follows a periodic function of $50\% \cdot \sin\left(\frac{\pi}{2}k\right) + 50\%$.

Table 1: $k$ and prediction accuracy for two periods

| $k$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| Prediction Accuracy (%) | 100 | 50 | 0 | 50 | 100 | 50 | 0 | 50 |

### 1.1.2 Curse of Dimensionality

Proving equation 1,

$$\mathbf{var}(\frac{||x^{(i)} - x^{(j)}||_2^2}{\mathbf{E}[||x^{(i)} - x^{(j)}||_2^2]}) = \frac{N+2}{N} - 1 \tag{1}$$

We can utilise equation 2 from Probability Theory,

$$\mathbf{var}[x] = \mathbf{E}[x^2] - \mathbf{E}[x]^2 \tag{2}$$

Below are the given equations,

$$x \in \mathbb{R}^n \tag{3}$$

$$\Pr(X) \sim \prod_{n=1}^{N} \mathcal{N}(x_n, \ 0 \ \sigma^2) \tag{4}$$

where $n$ represents the $n^{th}$ dimension.

$N$ represents the number of training data.

$$d_n = x_n^i - x_n^j \tag{5}$$

where $i$ represents the $i^{th}$ training data.

$j$ represents the $j^{th}$ training data.

$$\Pr(d_n) \sim \mathcal{N}(d_n; 0, 2\sigma^2) \tag{6}$$

$$\mathbf{E}[d_n^2 d_m^2] = \mathbf{E}[d_n^2]\mathbf{E}[d_m^2] \tag{7}$$

$$\mathbf{E}[d_n^4] = 3(\sqrt{2}\sigma)^4 = 12\sigma^4 \tag{8}$$

From equations 4 and 2, it is implied that

$$\mathbf{E}[x_n] = 0 \tag{9}$$

$$\mathbf{var}[x_n] = \sigma^2 = \mathbf{E}[x_n^2] \tag{10}$$

From equations 6 and 2, it is implied

$$\mathbf{E}[d_n] = 0 \tag{11}$$

4

$$\textbf{var}[d_n] = 2\sigma^2 = \textbf{E}[d_n^2] \tag{12}$$

From equations 7 and 12,

$$\textbf{E}[d_n^2 d_m^2] = \textbf{E}[d_n^2]\textbf{E}[d_m^2] = (2\sigma^2)(2\sigma^2) = 4\sigma^4 \tag{13}$$

From equation 1 and 5,

$$||x^{(i)} - x^{(j)}||_2^2 = \sum_{n=1}^{N}(x_n^{(i)} - x_n^{(j)})^2 = \sum_{n=1}^{N} d_n^2 \tag{14}$$

Substituting equations 1, 14 into 2,

$$\textbf{var}\left(\frac{\sum_{n=1}^{N} d_n^2}{\textbf{E}[\sum_{n=1}^{N} d_n^2]}\right) = \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\textbf{E}[\sum_{n=1}^{N} d_n^2]}\right)^2\right] - \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\textbf{E}[\sum_{n=1}^{N} d_n^2]}\right)\right]^2 \tag{15}$$

Looking into the first term of the Right Hand Side (RHS) of equation 15,

$$\textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\textbf{E}[\sum_{n=1}^{N} d_n^2]}\right)^2\right] = \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\sum_{n=1}^{N} \textbf{E}[\,d_n^2]}\right)^2\right] = \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\sum_{n=1}^{N} 2\sigma^2}\right)^2\right]$$

$$= \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{2N\sigma^2}\right)^2\right] = \textbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{2N\sigma^2}\right)^2\right] = \textbf{E}\left[\frac{(\sum_{n=1}^{N} d_n^2)^2}{(2N\sigma^2)^2}\right]$$

$$= \textbf{E}\left[\frac{(\sum_{n=1}^{N} d_n^2)^2}{4N^2\sigma^4}\right] = \textbf{E}\left[\frac{\sum_{n=1}^{N}\sum_{m=1}^{N} d_n^2 d_m^2}{4N^2\sigma^4}\right] = \textbf{E}\left[\frac{\sum_{n=1}^{N}\sum_{m=1}^{N} d_n^2 d_m^2}{4N^2\sigma^4}\right]$$

$$= \frac{\textbf{E}[\sum_{n=1}^{N}\sum_{m=1}^{N} d_n^2 d_m^2]}{4N^2\sigma^4}$$

$$= \frac{\textbf{E}[\sum_{n=1}^{N} d_n^4 + 2\sum_{n=1}^{N-1}\sum_{m=n+1}^{N} d_n^2 d_m^2]}{4N^2\sigma^4} \tag{16}$$

$$= \frac{\textbf{E}[\sum_{n=1}^{N} d_n^4] + \textbf{E}[2\sum_{n=1}^{N-1}\sum_{m=n+1}^{N} d_n^2 d_m^2]}{4N^2\sigma^4}$$

$$= \frac{\sum_{n=1}^{N} \textbf{E}[d_n^4] + \textbf{E}[2\sum_{n=1}^{N-1}\sum_{m=n+1}^{N} d_n^2 d_m^2]}{4N^2\sigma^4}$$

$$= \frac{\sum_{n=1}^{N} 12\sigma^4 + 2\sum_{n=1}^{N-1}\sum_{m=n+1}^{N} \textbf{E}[d_n^2 d_m^2]}{4N^2\sigma^4}$$

$$= \frac{12N\sigma^4 + 2\sum_{n=1}^{N-1}\sum_{m=n+1}^{N} 4\sigma^4}{4N^2\sigma^4}$$

$$= \frac{12N\sigma^4 + 2\binom{N}{2}4\sigma^4}{4N^2\sigma^4} = \frac{12N\sigma^4 + 2(\frac{N(N-1)}{2})4\sigma^4}{4N^2\sigma^4} = \frac{4N^2\sigma^4 + 8N\sigma^4}{4N^2\sigma^4} = \frac{N+2}{N}$$

Looking into the second term of the RHS of equation 15,

$$\mathbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\mathbf{E}[\sum_{n=1}^{\infty} d_n^2]}\right)\right]^2 = \mathbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\sum_{n=1}^{N} \mathbf{E}[d_n^2]}\right)\right]^2 = \mathbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{\sum_{n=1}^{N} 2\sigma^2}\right)\right]^2$$

$$= \mathbf{E}\left[\left(\frac{\sum_{n=1}^{N} d_n^2}{2N\sigma^2}\right)\right]^2 = \left(\frac{\mathbf{E}[\sum_{n=1}^{N} d_n^2]}{2N\sigma^2}\right)^2 = \left(\frac{\sum_{n=1}^{N} \mathbf{E}[d_n^2]}{2N\sigma^2}\right)^2 \qquad (17)$$

$$= \left(\frac{N2\sigma^2}{2N\sigma^2}\right)^2 = 1^2 = 1$$

Therefore, combining both terms from the RHS of equation 15 as calculated in equations 16 and 17 results in

$$\frac{N+2}{N} - 1 \qquad (18)$$

which proves equation 1.

To show that equation 18 vanishes as $N \to \infty$,

$$\lim_{N\to\infty}\left(\frac{N+2}{N} - 1\right) = \lim_{N\to\infty}\left(\frac{N+2-N}{N}\right) = \lim_{N\to\infty}\left(\frac{2}{N}\right) = 0 \qquad (19)$$

This proves that the variance vanishes which means that a test data point will be equally close to all training examples in a high dimensional space.

## 1.2 Euclidean Distance Function

### 1.2.1 Inner Product

All input vectors have same magnitude in the training set. $||x^{(1)}||_2^2 = ... = ||x^{(M)}||_2^2$

To show that in order to find nearest neighbor of a test point $x^*$ among the training set, it is sufficient to just compare and rank the negative inner product between the training and the test data, $x^{(M)^T}x^*$.

In a 2-Dimensional case, if all input vectors have the same magnitude, that defines a circle in the training set. Taking the inner product between two vectors, $x^{(M)^T}x^*$ would simply be calculating the angle between the two vectors. This is similar to performing

a cosine similarity calculation. Therefore, ranking based on negative inner product would be looking for the smallest angle between any 2 input vectors.

Now, extending this intuition to a M-Dimensional case. It is a M-dimensional hypersphere where all the vectors extend to the surface of the hypersphere as they have the same magnitude. Taking the negative inner product, $x^{(M)^T} x^*$ and ranking them would be finding the minimum angle between any of the 2 points on the surface, which is its nearest neighbor.

This can be illustrated below in equation 20, where $R$ is the radius of the hypersphere and all vectors have a magnitude equal to this radius.

$$
\begin{aligned}
||x^{(*)}||_2^2 &= \sum_{n=1}^{N} (x_n^{(*)})^2 = R^2 \\
||x^{(1)}||_2^2 &= \sum_{n=1}^{N} (x_n^{(1)})^2 = R^2 \\
&\vdots \\
||x^{(M)}||_2^2 &= \sum_{n=1}^{N} (x_n^{(M)})^2 = R^2
\end{aligned}
\tag{20}
$$

### 1.2.2 Pairwise Distances

Code snippets are included below:

```python
{python}
# 1.2 Euclidean Distance Function
# 1.2.2 Pairwise Distances
# Write a vectorized Tensorflow Python function that
#   implements
# the pairwise squared Euclidean distance function
# for two input matrices.
# No Loops and makes use of Tensorflow broadcasting.
def PairwiseDistances(X, Z):
    """
```

```python
    input:
        X is a matrix of size (B x N)
        Z is a matrix of size (C x N)
    output:
        D = matrix of size (B x C) containing
        the pairwise Euclidean distances
    """
    B = X.get_shape().as_list()[0]
    N = X.get_shape().as_list()[1]
    C = Z.get_shape().as_list()[0]
    # Ensure the N dimensions are consistent
    assert  N == Z.get_shape().as_list()[1]
    # Reshape to make use of broadcasting in Python
    X = tf.reshape(X, [B, 1, N])
    Z = tf.reshape(Z, [1, C, N])
    # The code below automatically does broadcasting.
    # Calculates the
    # pairwise squared Euclidean distance function
    D = tf.reduce_sum(tf.square(tf.sub(X, Z)), 2)
    return D
```

## 1.3  Making Predictions

### 1.3.1  Choosing Nearest Neighbour

Code snippets are included below:

```python
{python}
# 1.3 Making Predictions
# 1.3.1 Choosing nearest neighbours
```

```python
# Write a vectorized Tensorflow Python function that
# takes a pairwise distance matrix
# and returns the responsibilities of the
# training examples to a new test data point.
# It should not contain loops.
# Use tf.nn.top_k
def ChooseNearestNeighbours(D, K):
    """
    input:
        D is a matrix of size (B x C)
        K is the top K responsibilities for each test input
    output:
        topK are the value of the squared distances
        for the top K values
        indices are the index of the location of
        these top K squared distances.
    """
    # Take topK of negative distances since
    # closer data ranks higher.
    topK, indices = tf.nn.top_k(tf.neg(D), K)
    return topK, indices
```

### 1.3.2 Prediction

Code snippets are included below:

```python
{python}
# 1.3.2 Prediction
# Compute the k-NN prediction with K = {1, 3, 5, 50}
# For each value of K, compute and report:
```

```python
 5      # training MSE loss
 6      # validation MSE loss
 7      # test MSE loss
 8  # Choose best k using validation error = 50
 9  def PredictKnn(trainData , testData, trainTarget,  testTarget,
    ↪  K):
10      """
11      input:
12          trainData: Data for training KNN
13          testData: Data used in testing
14          trainTarget: Targets used to create prediction.
15          testTarget: Targets used to calculate loss.
16      output:
17          loss: The mean squared loss of the prediction.
18      """
19      D = PairwiseDistances(testData, trainData)
20      topK, indices = ChooseNearestNeighbours(D, K)
21      # Select the proper outputs to be averaged
22      # from the target values and average them
23      trainTargetSelectedAveraged = tf.reduce_mean( \
24              tf.gather(trainTarget, indices), 1)
25      # Calculate the loss from the actual values
26      loss = tf.reduce_mean(tf.square(tf.sub( \
27              trainTargetSelectedAveraged, testTarget)))
28      return loss
29
30  # Plot the prediction function for x = [0, 11]
31  def PredictedValues(x, trainData, trainTarget, K):
32      """
33      Plot the predicted values
```

```
34      input:
35          x = test target to plot and predict
36      """
37      D = PairwiseDistances(x, trainData)
38      topK, indices = ChooseNearestNeighbours(D, K)
39      predictedValues = tf.reduce_mean( \
40              tf.gather(trainTarget, indices), 1)
41      return predictedValues
```

Table 2: KNN and Loss

| $k$ | Training MSE Loss | Validation MSE Loss | Test MSE Loss |
|---|---|---|---|
| 1 | 0.000 | 0.272 | 0.311 |
| 3 | 0.105 | 0.326 | 0.145 |
| 5 | 0.119 | 0.310 | 0.178 |
| 50 | 1.248 | 1.229 | 0.707 |

The best value of $k$ is based on one that gives the lowest validation MSE loss. In this case, the best $k$ is found to be $k = 1$.

By inspecting the plot, a $k$ value of 3 or 5 would be picked, instead of $k = 1$. The reason for this is simply because $k = 1$ is overfitting the training data, and would have modeled noises in the data. This is confirmed when comparing between the Test MSE losses for different values of $k$.

Figure 2: k-NN regression on data1D for various values of k

## 1.4 Soft kNN and Gaussian Processes

As shown in Table 3 , the Soft Decision performs better than the Gaussian Process Regression Model as it has a lower Test Mean Squared Error Loss. The algorithm was run on the test set.

Table 3: Loss on Test Set

| *Algorithm* | Test MSE Loss |
|---|---|
| Soft Decision | 0.159 |
| Gaussian Process Regression | 0.380 |



Figure 3: Soft Decision kNN on Test Set

Figure 4: Gaussian Process Regression on Test Set

```python
{python}
def SortData(inputVal, outputVal):
    """
    This sorts a given test set by the dataValue before
    plotting it.
    """
    # Sort across the data values on both pairs of sets
    # Sort in numpy itself to not lose precision.
    p = np.argsort(inputVal, axis=0)
    inputVal = np.array(inputVal)[p]
    outputVal = np.array(outputVal)[p]
```

```python
11    # Get rid of extra dimensions from np.argsort
12    inputVal = inputVal[:, :,0]
13    outputVal = outputVal[:, :,0]
14    return inputVal, outputVal
```

### 1.4.1  Soft Decisions and Gaussian Process Regression

```python
1  {python}
2  # Predict values using soft decision
3  # 1.4.1.1 Soft Knn Decision
4
5  def PredictedValuesSoftDecision(x, trainData, trainTarget):
6      # use hyper parameter of 100 as given by Jimmy.
7      hyperParam = 100
8      # Compute pairwise differences.
9      D1 = PairwiseDistances(x, trainData)
10     K1 =  tf.exp(-hyperParam*D1)
11     # Get the sum term used for normalization
12     sum1 = tf.reduce_sum(tf.transpose(K1), axis=0)
13     # Reshape to enable broadcasting during division.
14     N = sum1.get_shape().as_list()[0]
15     sum1 = tf.reshape(sum1, [N,1])
16     # Normalize the data using broadcast
17     # to calculate the final responsibility values
18     rStar = tf.div(K1, sum1)
19     # Calculate the predicted value
20     # using the new responsibilities, rStar
21     predictedValues = tf.matmul(rStar,trainTarget)
22     return predictedValues
```
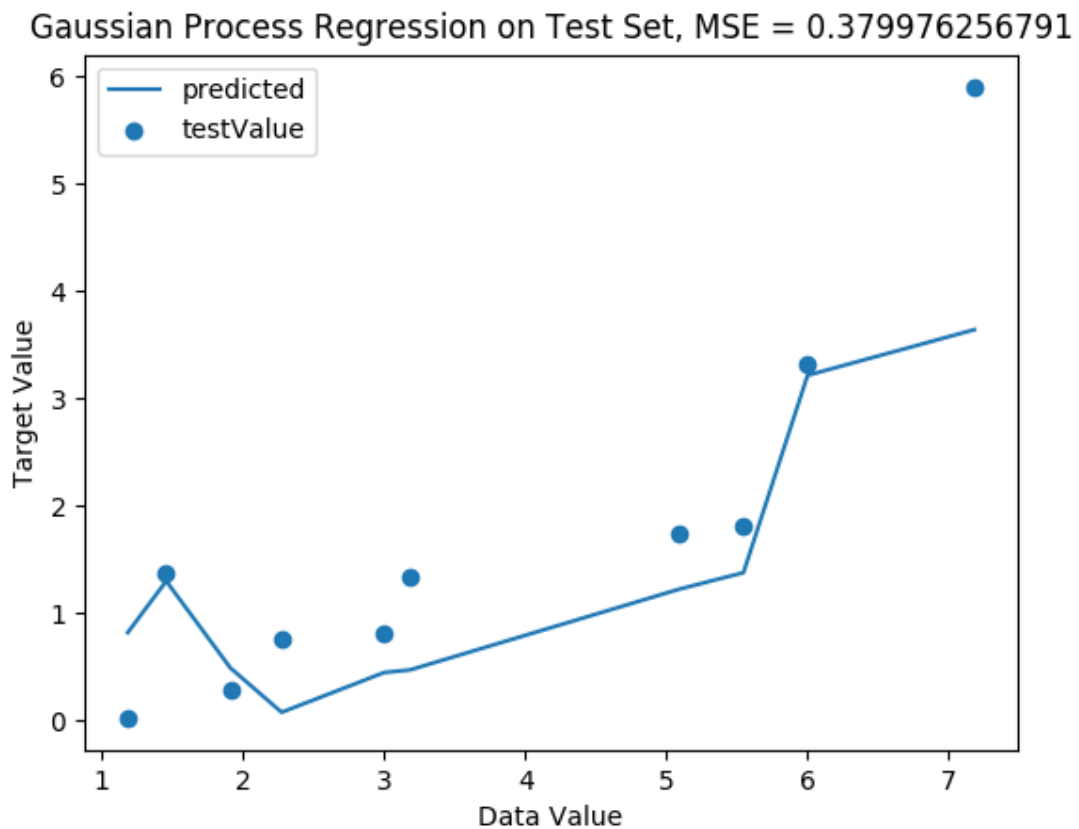
```python
{python}
# Predict values using Gaussian
# 1.4.1.1 Gaussian Processes
def PredictedValuesGaussianProcesses(x, trainData,
↪ trainTarget):
    # use hyper parameter of 100 as given by Jimmy.
    hyperParam = 100
    D1 = PairwiseDistances(x, trainData)
    K1 =  tf.exp(-hyperParam*D1)
    D2 = PairwiseDistances(trainData, trainData)
    K2 =  tf.matrix_inverse(tf.exp(-hyperParam*D2))
    # Calculate the responsibilites, rStar
    # by normalizing using the inverse ofK2.
    rStar = tf.matmul(K1, K2)
    # Calculate the predicted value
    # using the new responsibilities, rStar
    predictedValues = tf.matmul(rStar,trainTarget)
    return predictedValue
```

### 1.4.2 Conditional Distribution of a Gaussian

Given an $M + 1$ Gaussian random vector:

$$\boldsymbol{y} = \begin{bmatrix} y^* \\ y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \Sigma\right) \tag{21}$$

Splitting the vector into two parts as follows and describing the resulting distribution using stacked block notation:

$$\boldsymbol{y} = \begin{bmatrix} y^* \\ \boldsymbol{y}_{train} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} 0 \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{y^*y^*} & \Sigma_{y^*\boldsymbol{y}_{train}} \\ \Sigma_{\boldsymbol{y}_{train}y^*} & \Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} \end{bmatrix} \right) \tag{22}$$

$$\sim \mathcal{N}\left( \begin{bmatrix} 0 \\ \mathbf{0} \end{bmatrix}, \Sigma = \Lambda^{-1} = \begin{bmatrix} \Lambda_{y^*y^*} & \Lambda_{y^*\boldsymbol{y}_{train}} \\ \Lambda_{\boldsymbol{y}_{train}y^*} & \Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} \end{bmatrix}^{-1} \right) \tag{23}$$

Given that $\boldsymbol{y}_{train}$ is observed, find $P(y^*|\boldsymbol{y}_{train}) \sim \mathcal{N}\left(y^*; \mu^*, \Sigma^*\right)$.

By completing the squares on the quadratic term of the exponent for the multivariate Gaussian equation,

$$-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\boldsymbol{y} - \boldsymbol{\mu})$$

$$= -\frac{1}{2}\boldsymbol{y}^T \Sigma^{-1} \boldsymbol{y}^T$$

$$= -\frac{1}{2}\begin{bmatrix} y^{*T} & \boldsymbol{y}_{train}^T \end{bmatrix} \begin{bmatrix} \Lambda_{y^*y^*} & \Lambda_{y^*\boldsymbol{y}_{train}} \\ \Lambda_{\boldsymbol{y}_{train}y^*} & \Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} \end{bmatrix} \begin{bmatrix} y^* \\ \boldsymbol{y}_{train} \end{bmatrix} \tag{24}$$

$$= -\frac{1}{2}\begin{bmatrix} y^{*T}\Lambda_{y^*y^*} + \boldsymbol{y}_{train}\Lambda_{\boldsymbol{y}_{train}y^*} & y^{*T}\Lambda_{y^*\boldsymbol{y}_{train}} + \boldsymbol{y}_{train}\Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} \end{bmatrix} \begin{bmatrix} y^* \\ \boldsymbol{y}_{train} \end{bmatrix}$$

$$= -\frac{1}{2}\left( y^{*T}\Lambda_{y^*y^*}y^* + \boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}y^*}y^* + y^{*T}\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train} + \boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}\boldsymbol{y}_{train} \right)$$

Since $\boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}y^*}y^*$ is a scalar and $\Lambda_{\boldsymbol{y}_{train}y^*}^T = \Lambda_{y^*\boldsymbol{y}_{train}}$,

$$\left( \boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}y^*}y^* \right)^T = y^{*T}\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train} \tag{25}$$

Hence, equation 24 simplifies to:

$$-\frac{1}{2}\left( y^{*T}\Lambda_{y^*y^*}y^* + 2y^{*T}\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train} + \boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}\boldsymbol{y}_{train} \right)$$

$$= -\frac{1}{2}y^{*T}\Lambda_{y^*y^*}y^* - y^{*T}\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train} - \frac{1}{2}\boldsymbol{y}_{train}^T\Lambda_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}\boldsymbol{y}_{train} \tag{26}$$

Completing the square for a multivariate Gaussian quadratic term with $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ yields:

$$-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\boldsymbol{x} - \boldsymbol{\mu})$$

$$= -\frac{1}{2}\boldsymbol{x}^T\Sigma^{-1}\boldsymbol{x} + \boldsymbol{x}^T\Sigma^{-1}\boldsymbol{\mu} - \frac{1}{2}\boldsymbol{\mu}^T\Sigma^{-1}\boldsymbol{\mu} \tag{27}$$

Given that $\boldsymbol{y}_{train}$ and $\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}$ are known, the conditional Gaussian Distribution $P(y^*|\boldsymbol{y}_{train})$ can be inferred by performing a term-by-term comparison between the terms of equation 26 and those of equation 27.

Comparing terms that are of second order with respect to $y^*$,

$$\Sigma^* = \Lambda_{y^*\boldsymbol{y}_{train}}^{-1} \tag{28}$$

Comparing terms that are of first order with respect to $y^*$,

$$\Sigma^{*-1}\mu^* = \Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train}$$

$$\mu* = \Sigma^*\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train} \tag{29}$$

$$=\Lambda_{y^*y^*}^{-1}\Lambda_{y^*\boldsymbol{y}_{train}}\boldsymbol{y}_{train}$$

Using the matrix-inverse identity provided in Tutorial 3 (pg. 41), the terms $\Lambda_{y^*y^*}$ and $\Lambda_{y^*\boldsymbol{y}_{train}}$ can be expressed using terms in $\Sigma = \begin{bmatrix} \Sigma_{y^*y^*} & \Sigma_{y^*\boldsymbol{y}_{train}} \\ \Sigma_{\boldsymbol{y}_{train}y^*} & \Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} \end{bmatrix}$,

$$\Lambda_{y^*y^*} = \left(\Sigma_{y^*y^*} - \Sigma_{y^*\boldsymbol{y}_{train}}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1}\Sigma_{\boldsymbol{y}_{train}y^*}\right)^{-1}$$
$$\Lambda_{y^*\boldsymbol{y}_{train}} = - \left(\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}} - \Sigma_{\boldsymbol{y}_{train}y^*}\Sigma_{y^*y^*}^{-1}\Sigma_{y^*\boldsymbol{y}_{train}}\right)^{-1}\Sigma_{y^*\boldsymbol{y}_{train}}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1} \tag{30}$$

The results above allow equations 28 and 29 to be expressed in terms of the original $\Sigma$ block terms:

$$\mu^* = - \Sigma_{y^*\boldsymbol{y}_{train}}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1}\boldsymbol{y}_{train} \tag{31}$$

$$\Sigma^* = \Sigma_{y^*y^*} - \Sigma_{\boldsymbol{y}_{train}y^*}^{T}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1}\Sigma_{\boldsymbol{y}_{train}y^*} \tag{32}$$

$$\therefore P(y^*|\boldsymbol{y}_{train}) \sim \mathcal{N}\left(-\Sigma_{y^*\boldsymbol{y}_{train}}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1}\boldsymbol{y}_{train}, \Sigma_{y^*y^*} - \Sigma_{\boldsymbol{y}_{train}y^*}^{T}\Sigma_{\boldsymbol{y}_{train}\boldsymbol{y}_{train}}^{-1}\Sigma_{\boldsymbol{y}_{train}y^*}\right) \tag{33}$$

# 2 Linear and Logistic Regression

## 2.1 Geometry of Linear Regression

### 2.1.1 Convex Function

Need to show if equation 34 is a convex function of W using Jensen Inequality given by equation 35.

$M$ is the total number of training data. $N$ is the number of dimension for each training data.

$$
\begin{aligned}
\mathcal{L} &= \mathcal{L}_{\mathcal{D}} + \mathcal{L}_{\mathcal{W}} \\
&= \sum_{m=1}^{M} \frac{1}{2M} ||W^T x^{(m)} + b - y^{(m)}||_2^2 + \frac{\lambda}{2} ||W||_2^2 \\
&= \sum_{m=1}^{M} \frac{1}{2M} \left[ \sum_{n=1}^{N} (W_n x_n^{(m)}) + b - y^{(m)} \right]^2 + \frac{\lambda}{2} W^T W \\
&= \sum_{m=1}^{M} \frac{1}{2M} \left[ (W^T x^{(m)} + b - y^{(m)})^2 \right] + \frac{\lambda}{2} W^T W
\end{aligned}
\tag{34}
$$

$$
f(\alpha W_1 + (1-\alpha)W_2) \leq \alpha f(W_1) + (1-\alpha)f(W_2) \tag{35}
$$

Since the sum of two convex function is convex, we can prove each sum term, $\mathcal{L}_{\mathcal{D}}$ and $\mathcal{L}_{\mathcal{W}}$ on equation 34 separately.

You can easily prove that the sum of two convex is still convex by summing both sides of Jensen's Inequality for both convex functions and showing that Jensen's Inequality would still hold, indicating that the sum of the convex functions is still convex. This enables us to simplify the analysis by proving that each term is convex separately

As a convex function divided by a positive value is still a convex function, we can ignore the division by $2M$ since $M > 0$ means that $\frac{1}{2M} > 0$. Similarly, $\frac{\lambda}{2} \geq 0$ and it only scales $W^T W$. You can divide both sides of Jensen's inequality by a positive constant and the Jensen's inequality would still hold, implying that the function is still

convex when divided by a positive constant. This enables us to simplify the analysis by ignoring the positive constants.

This means we just have to prove that equation 36 is convex and equation 37 is convex. By doing so, we would have proven that equation 34 is convex.

$$\mathcal{L}_{\mathcal{D}} = \left[ W_n^T x_n^{(m)} + b_n - y_n^{(m)} \right]^2 \tag{36}$$

$$\mathcal{L}_{\mathcal{W}} = W^T W \tag{37}$$

**2.1.1.1 Proof of Convexity of $\mathcal{L}_{\mathcal{D}}$ with respect to $W$**

Rearranging equation 35 for $\mathcal{L}_{\mathcal{D}}$,

$$\mathcal{L}_{\mathcal{D}}(\alpha W_1 + (1 - \alpha)W_2) - \alpha \mathcal{L}_{\mathcal{D}}(W_1) - (1 - \alpha)\mathcal{L}_{\mathcal{D}} \leq 0 \tag{38}$$

$$
\begin{aligned}
LHS \equiv \; & \mathcal{L}_{\mathcal{D}}[\alpha W_1 + (1 - \alpha)W_2] - \alpha \mathcal{L}_{\mathcal{D}}(W_1) - (1 - \alpha)\mathcal{L}_{\mathcal{D}}(W_2) \\
= \; & \{[\alpha W_1 + (1 - \alpha)W_2]^T x + (b - y)\}^T \{[\alpha W_1 + (1 - \alpha)W_2]^T x + (b - y)\} \\
& - \alpha[W_1^T x + (b - y)]^T[W_1^T x + (b - y)] \\
& \qquad - (1 - \alpha)[W_2^T x + (b - y)]^T[W_2^T x + (b - y)] \quad (39)
\end{aligned}
$$

$$
\begin{aligned}
= \; & [x^T(\alpha W_1 + (1 - \alpha)W_2) + (b - y)^T]\{[\alpha W_1 + (1 - \alpha)W_2]^T x + (b - y)\} \\
& - \alpha[x^T W_1 + (b - y)^T][W_1^T x + (b - y)] \\
& \qquad - (1 - \alpha)[x^T W_2 + (b - y)^T][W_2^T x + (b - y)] \quad (40)
\end{aligned}
$$

$$
\begin{aligned}
= \; & \{x^T[\alpha W_1 + (1 - \alpha)W_2][\alpha W_1^T + (1 - \alpha)W_2^T]x \\
& + 2x^T[\alpha W_1 + (1 - \alpha)W_2](b - y) + (b - y)^T(b - y)\} \\
& - \alpha[x^T W_1 W_1^T x + 2x^T W_1(b - y) + (b - y)^T(b - y)] \\
& \qquad - (1 - \alpha)[x^T W_2 W_2^T x + 2x^T W_2(b - y) + (b - y)^T(b - y)] \quad (41)
\end{aligned}
$$

Rearranging similar terms together,

$$
\begin{aligned}
= \{ x^T & \left[ \alpha^2 W_1 W_1^T + 2\alpha(1-\alpha)W_1^T W_2 + (1-\alpha)^2 W_2^T W_2 \right] x \\
& - \alpha(x^T W_1 W_1^T x) - (1-\alpha)(x^T W_2 W_2 Tx) \} \\
& + \{ 2\alpha x^T W_1(b-y) - 2\alpha x^T W_1(b-y) \} \\
& + \{ 2(1-\alpha)x^T W_2(b-y) - 2(1-\alpha)x^T W_2(b-y) \} \\
& + \{ 1 - \alpha - (1-\alpha) \}(b-y)^T(b-y) \quad \text{(42)}
\end{aligned}
$$

$$
\begin{aligned}
&= (\alpha^2 - \alpha)x^T W_1 W_1^T x + 2\alpha(1-\alpha)x^T W_1^T W_2 x + [(1-\alpha)^2 - (1-\alpha)]x^T W_2 W_2^T x \\
&= -\alpha(1-\alpha)x^T W_1 W_1^T x + 2\alpha(1-\alpha)x^T W_1^T W_2 x - \alpha(1-\alpha)x^T W_2 W_2^T x \\
&= -\alpha(1-\alpha)[x^T W_1 W_1^T x - 2x^T W_1^T W_2 x + x^T W_2 W_2^T x] \quad \text{(43)} \\
&= -\alpha(1-\alpha)[W_1^T x - W_2^T x]^T [W_1^T x - W_2^T x] \\
&\leq 0 \equiv RHS
\end{aligned}
$$

Equation 43 is less than or equal to zero as $-\alpha(1-\alpha) \leq 0$; $\forall \alpha \in [0,1]$. Furthermore, the remaining quadratic term, $[W_1^T x - W_2^T x]^T [W_1^T x - W_2^T x] \geq 0$ since it is a square of the term $\forall \left( W_1^T x - W_2^T x \right) \in \mathbb{R}^N$.

Hence, it has been shown that $\mathcal{L}_\mathcal{D}$ is convex.

**2.1.1.2 Proof of Convexity of $\mathcal{L}_\mathcal{W}$ to $W$**

From the Left Hand Side of Equation 38,

$$
\begin{aligned}
LHS \equiv{} & \mathcal{L}_\mathcal{W}[\alpha W_1 + (1-\alpha)W_2] - \alpha\mathcal{L}_\mathcal{W}(W_1) - (1-\alpha)\mathcal{L}_\mathcal{W}(W_2) \\
={} & [\alpha W_1 + (1-\alpha)W_2]^T[\alpha W_1 + (1-\alpha)W_2] - \alpha W_1^T W_1 - (1-\alpha)W_2^T W_2 \\
={} & [\alpha W_1^T + (1-\alpha)W_2^T][\alpha W_1 + (1-\alpha)W_2] - \alpha W_1^T W_1 - (1-\alpha)W_2^T W_2 \\
={} & [\alpha^2 W_1^T W_1 + 2\alpha(1-\alpha)W_1^T W_2 + (1-\alpha)^2 W_2^T W_2] - \alpha W_1^T W_1 - (1-\alpha)W_2^T W_2 \\
={} & (\alpha^2 - \alpha)x^T W_1 W_1^T x + 2\alpha(1-\alpha)x^T W_1^T W_2 x + [(1-\alpha)^2 - (1-\alpha)]x^T W_2 W_2^T x \\
={} & -\alpha(1-\alpha)W_1^T W_1 + 2\alpha(1-\alpha)W_1^T W_2 - \alpha(1-\alpha)W_2 W_2^T \\
={} & -\alpha(1-\alpha)[W_1^T W_1 - 2W_1^T W_2 x + W_2^T W_2] \\
={} & -\alpha(1-\alpha)[W_1 - W_2]^T[W_1 - W_2] \\
\leq{} & 0 \equiv RHS
\end{aligned}
$$

$$(44)$$

Using a similar argument as that for Section 2.1.1.1, the loss function $\mathcal{L}_\mathcal{W}$ is shown to be convex with respect to $W$.

Therefore, the loss function $\mathcal{L}$ is a convex function with respect to W.

**2.1.1.3 Proof of Convexity of $\mathcal{L}$ to $b$**

Instead of performing a similar proof for $b$, the bias can be thought of as the $(N+1)$th dimension of the weight. Thus, $b = W_{N+1}$ can be grouped together with matrix $W$, while the vector $x$ will be expanded to have $x_{N+1}^m = 1 \ \forall m$, as shown in equation 45. Using the proof for the convexity of $\mathcal{L}$ with respect to $W$, by extension $\mathcal{L}$ is a convex function of the bias, $b$.

$$
L = \sum_{m=1}^{M} \frac{1}{2M} \sum_{n=1}^{N+1} \left(W_n^T x_n^{(m)} - y^{(m)}\right)^2 + \frac{\lambda}{2}W^T W
\tag{45}
$$

### 2.1.2 DeNormalization

Assuming $\lambda = 0$ from equation 34 we get equation 46.

$$\sum_{m=1}^{M} \frac{1}{2M} \left[ (W^T x^{(m)} + b - y^{(m)})^2 \right] \tag{46}$$

The original optimal weights and optimal bias are optimal to the non-transformed dataset. This means that the $\frac{\partial L}{\partial W} = 0$ and $\frac{\partial L}{\partial b} = 0$. More specifically, the partial gradient of the loss in equation 46 with respect to a specific weight $W_n$,

$$\frac{\partial L}{\partial W_n} = \sum_{m=1}^{M} \frac{1}{M} \Big( \sum_{i=1, i \neq n}^{N} W_i x_i^m + W_n x_n^m + b - y^m \Big) x_n^m = 0 \tag{47}$$

and for the bias, $b$

$$\frac{\partial L}{\partial b} = \sum_{m=1}^{M} \frac{1}{M} \Big( \sum_{i=1, i \neq n}^{N} W_i x_i^m + W_n x_n^m + b - y^m \Big) = 0 \tag{48}$$

A single dimension of $x$ scales by $\alpha > 1$ and shifts by $\beta > 1$. This is the same as de-normalizing the data point instead of normalizing it which normally deducts each data point by the mean and scaled by its variance. The reason for normalizing is to prevent any component from dominating the sum and to prevent the weights from learning the high bias that is not needed for prediction. As this is the opposite of normalization, it ends up training slower. However, this will not change the global minimum value as will be shown below.

Note: This note is added after submission. Realize after 1 hour after submitting that the below 3 equations is flawed as sum of products is not product of sums. Should have not brought in the sum. With a similar proof, each individual gradient would be 0. So the sum of all the gradients would be 0. Hence, proving this theorem. So shouldn't have brought in the training case. Or wait, the original equation's gradient of the sum is 0, but doesn't show that each individual gradient itself is 0. So although the individual original gradient sums to 0, each individual gradient may not necessarily be 0. Therefore, this proof, doesn't work. Sigh. Update: This equation should be

true since all the terms inside is 0, which means we can factor out the 0's and group the x term on the right :) The reason why all term inside is 0 is because it is optimal with respect to the weights. What this means is that it passes by every point. What this means is that it doesn't matter if I calculate the gradient with respect to 1, 2, ..., M training points, the gradient must be 0 since it passes by all of them. Take batch size = 1, this means that each gradient must be = 0, therefore, we are able to group the X up. Proven! :) End of note after submission =D

Expanding equation 47 to account for the sum of $M$ training cases, we find equation 49.

$$\frac{\partial L}{\partial W_n} = \frac{1}{M}(\sum_{i=1,i\neq n}^{N} W_i(x_i^1+...+x_i^M)+W_n(x_n^1+...+x_n^M)+Mb-(y^1+...+y^M))(x_n^1+...+x_n^M) = 0 \tag{49}$$

which we can rewrite in simpler terms as equation 50

$$\frac{\partial L}{\partial W_n} = \frac{1}{M}(\sum_{i=1,i\neq n}^{N} W_i\mathbf{x_i} + W_n(\mathbf{x_n}) + Mb - \mathbf{y})(\mathbf{x_n}) = 0 \tag{50}$$

Since equation 50 is equal to $0$ for any value of $\mathbf{x_n}$, we must have that the inner term, $(\sum_{i=1,i\neq n}^{N} W_i\mathbf{x_i} + W_n(\mathbf{x_n}) + Mb - \mathbf{y})$ is equal 0.

This is shown clearly in equation 51.

$$(\sum_{i=1,i\neq n}^{N} W_i\mathbf{x_i} + W_n(\mathbf{x_n}) + Mb - \mathbf{y}) = 0 \tag{51}$$

As the model changes by scaling by a positive constant of $\alpha$ and shifted by a positive constant of $\beta$ within the square term, it does not change the minimum value as it is these transformations happen within the square term. As a result, the minimum value remains at 0 which is the lowest value for any square term. Therefore, the new global minimum value of the transformed convex loss function will remain the same compare to the original loss function global minimum.

To illustrate, lets re-write the original loss function from equation 46 to account for the transformed Loss Function in equation 52.

$$L' = \sum_{m=1}^{M} \frac{1}{2M} \left[ (\sum_{i=1,i\neq n}^{N} W_i^i x_i^{m'} + W_n'(x_n^m) + Mb' - y^m)^2 \right] \tag{52}$$

where the new variables are appended with a ' to show that their values could or has changed. This can be re-written as equation 53.

$$L' = \frac{1}{2M} \left[ (\sum_{i=1,i\neq n}^{N} W_i' \mathbf{x_i} + W_n'(\mathbf{x_n'}) + Mb' - \mathbf{y})^2 \right] \tag{53}$$

Taking the gradient with respect to $W$ results in equation

54.

$$\frac{\partial L'}{\partial W_n} = \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i' \mathbf{x_i} + W_n'(\mathbf{x_n'}) + Mb' - \mathbf{y}) \mathbf{x_n'} \right]$$

$$= \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i' \mathbf{x_i} + W_n'(\alpha \mathbf{x_n} + \beta) + Mb' - \mathbf{y})(\alpha \mathbf{x_n} + \beta) \right] \tag{54}$$

To achieve the global minimum, a possible solution assignment would be to re-substitute in the original values for $W_i' = W_i \forall i \neq n$ and $W_n' = \frac{W_n}{\alpha}$ into equation 54 as shown in equation 55.

$$\frac{\partial L'}{\partial W_n} = \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i \mathbf{x_i} + \frac{W_n}{\alpha}(\alpha \mathbf{x_n} + M\beta) + Mb' - \mathbf{y})(\alpha \mathbf{x_n} + M\beta) \right]$$

$$= \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i \mathbf{x_i} + W_n \mathbf{x_n} + \frac{W_n M\beta}{\alpha} + Mb' - \mathbf{y})(\alpha \mathbf{x_n} + M\beta) \right] \tag{55}$$

$$= \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i \mathbf{x_i} + W_n \mathbf{x_n} + M(\frac{W_n \beta}{\alpha} + b') - \mathbf{y})(\alpha \mathbf{x_n} + M\beta) \right]$$

Further setting $b' = b - \frac{W_n \beta}{\alpha}$ results and using the result from equation 51, we get equation 56.

$$\frac{\partial L^{'}}{\partial W_n} = \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i \mathbf{x_i} + W_n \mathbf{x_n} + M(\frac{W_n\beta}{\alpha} + b - \frac{W_n\beta}{\alpha}) - \mathbf{y})(\alpha \mathbf{x_n} + M\beta) \right]$$

$$= \frac{1}{M} \left[ (\sum_{i=1,i\neq n}^{N} W_i \mathbf{x_i} + W_n \mathbf{x_n} + Mb - \mathbf{y})(\alpha \mathbf{x_n} + M\beta) \right] \quad (56)$$

$$= \frac{1}{M} \left[ (0)(\alpha \mathbf{x_n} + M\beta) \right]$$

$$= 0$$

Hence, from the final equation in 56, the partial gradient is 0, suggesting that this assignment results in an optimal assignment of $W^{'}$ and $b^{'}$.

As this assignment results in a similar inner term of equation 51, this suggest that the new loss function will end up being equal to the old loss function. Hence, the global minimum remains the same.

The proof that this assignment works on bias is very similar and is omitted.

From this assignment of $W^{'}$ and $b^{'}$, it shows that the weight vector $W$ will move downwards to reach $W'$ and similarly for $b$ to $b^{'}$. The optimal weights after learning will be lower compared to the optimal weights. The biased will be lower or higher or same depending on the sign of $W_n$ to converge into the biased learned from the non-transformed original data.

### 2.1.3   Regularization

The new minimum value will increase as the weight and bias will fit to the regularize model more and less to the training set. This means the loss with respect to the training set only will increase as the regularization penalizes large values of $W$. However, this regularization helps the prediction to be more robust on the validation and the test set by preventing the model from over-fitting to the training set.

The large weights will reduce as they are regularized whereas the small weights will increase. The bias should remain the same as it is not affected by the regularized term.

The new minimum loss will increase as it is being compared to the training set labels, but is being optimize for both the training set and the regularized term.

### 2.1.4 Binary Classifiers for Multi-class Classification with D classes

Given $D > 2$ and only able to use binary classifiers. A method to solve a multi-class classification task using a number of binary classifiers would be to assign a binary classifier to each class in the $D$ classes (see Figure 5).

Each binary classifier would predict if a given test input belong to a specific class. This assumes that each test input can belong to more than 1 class as it is not constrained to belong to more than 1 class or no classes at all from this design.
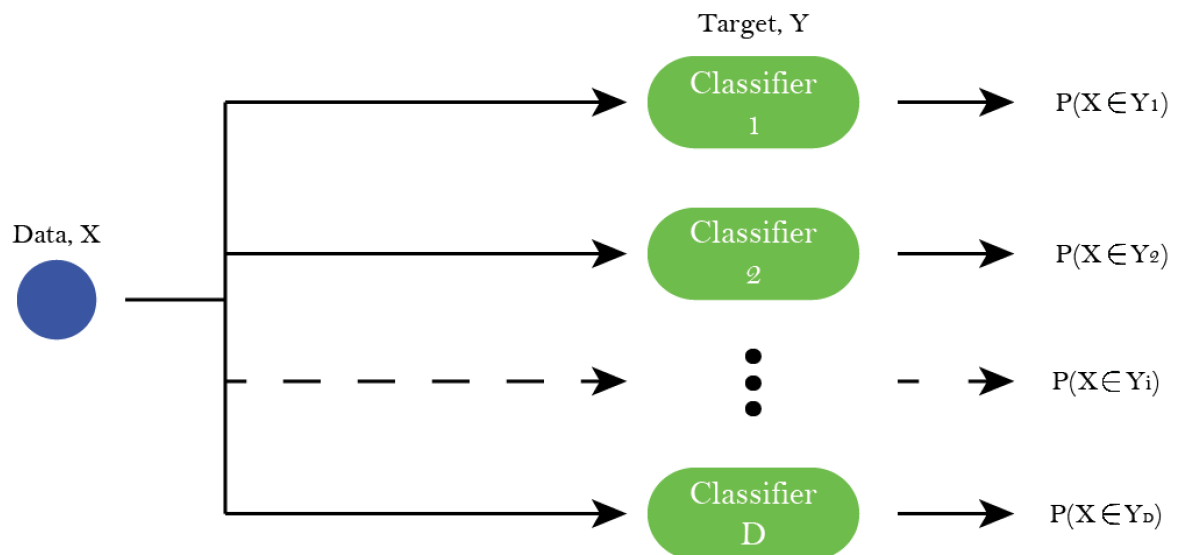


Figure 5: D binary classifiers for multi-class classification

## 2.2 Stochastic Gradient Descent

Before proceeding with our explanations for the subsequent sections, we would like to clarify on the convergence logic we selected for training our '3/5' digit classifier.

The model weights are randomly initialised. In our model, training convergence is implemented through early stopping. In our case, it is defined to be when the average validation MSE for the latest epoch is more than 0.99 of the average validation MSE for the previous epoch.

The logic is that when the average relative learnings between epochs is less than 1%, the training is terminated prematurely for shorter computation time. This makes sense since the model is not learning much from the features.

Training convergence also happens when the average MSE of one epoch is larger than its predecessor. This is done since an increase in the average validation MSE is a sign of over-fitting.

### 2.2.1 Tuning Learning Rate, $\eta$

Our programmed definition of convergence can be found in Section 2.2.

The values $\eta$ was generated based on increasing orders of 1 (see Table 4). Higher values of $\eta$ lead to fewer number of updates to reach convergence. However, from trial and error, it was discovered that there is an upper bound for $\eta$ to ensure training convergence. When the value of $\eta$ is more than 0.1 (i.e. 1, 10, etc.), the training does not converge but diverges instead as it overshoots past the minimum value.

From our experiments, The best value for $\eta$ is selected to be 0.1 as it takes the lowest number of iterations to reach the lowest convergence value. The results also summarised in Figure 6.

Table 4: Number of updates until convergence for various values of $\eta$

| $\eta$ | Number of Updates |
|---|---|
| 0.001 | 1863 |
| 0.01 | 337 |
| 0.1 | 57 |
| $> 0.3$ | N/A |



Figure 6: Graph of training mean squared error (MSE) against number of updates for the best learning rate found, $\eta = 0.1$. ($B = 50, \lambda = 1$)

### 2.2.2 Mini-batch Size

Our programmed definition of convergence can be found in Section 2.2.

From Table 5, the same pattern for $\eta$ is observed where higher values of $\eta$ lead to fewer number of updates. Meanwhile, there is an optimum value for batch size, $B$. The best mini-batch size is $B = 100$, which leads to the fewest number of updates, regardless of $\eta$. However, when the $B = 700$, early-stopping fails to occur in a reasonable time.

Please refer to Figures 7 and 8 for more information.

Table 5: Number of updates until training convergence for various values of $\eta$ and $B$

| $B$ | $\eta$ | | |
|---|---|---|---|
| | 0.001 | 0.01 | 0.1 |
| 10 | 2871 | 491 | 211 |
| 50 | 1963 | 337 | 57 |
| 100 | 1114 | 274 | 43 |
| 700 | 6001 | 6001 | 6001 |

Figure 7: Subplots of training MSE against number of updates for batch sizes, $B = 10, 50$ and learning rates, $\eta =$ **0.001**, **0.01**, **0.1**. ($\lambda = 1$ for all cases)
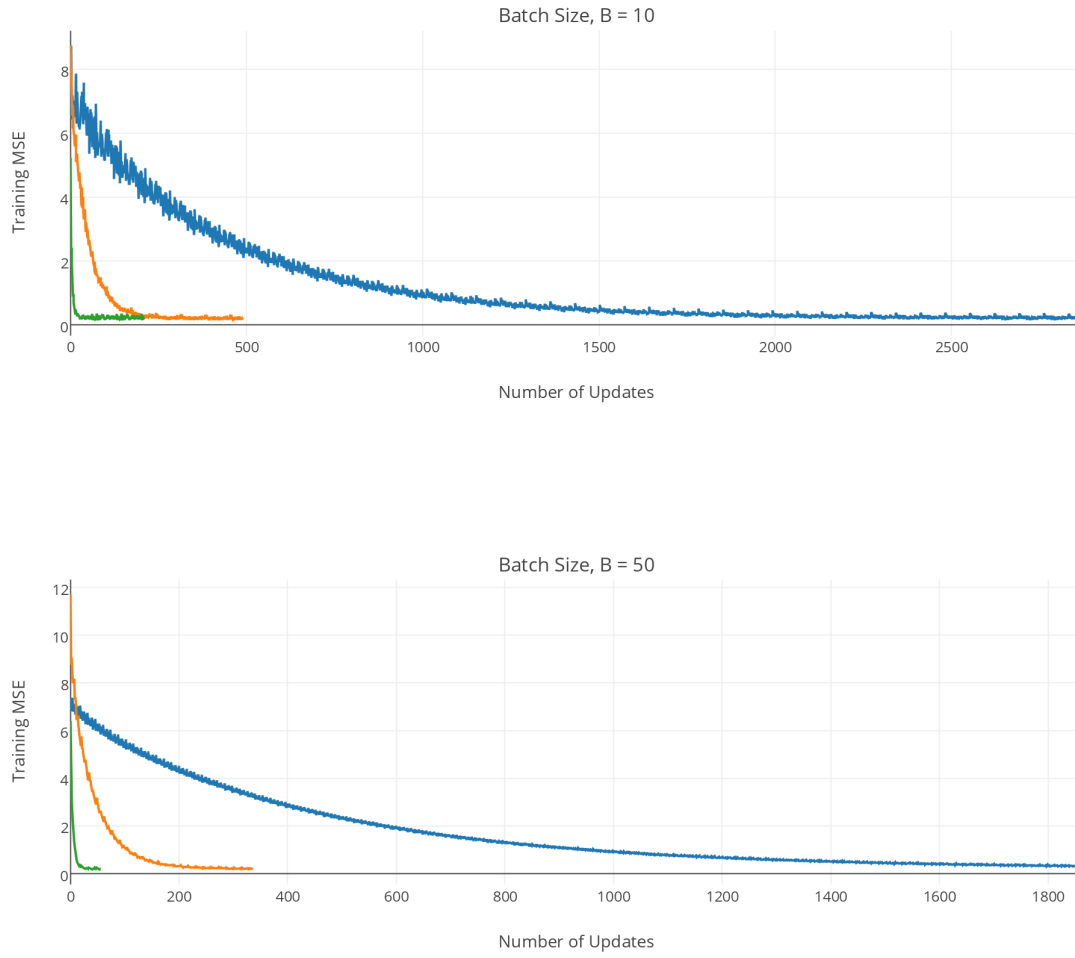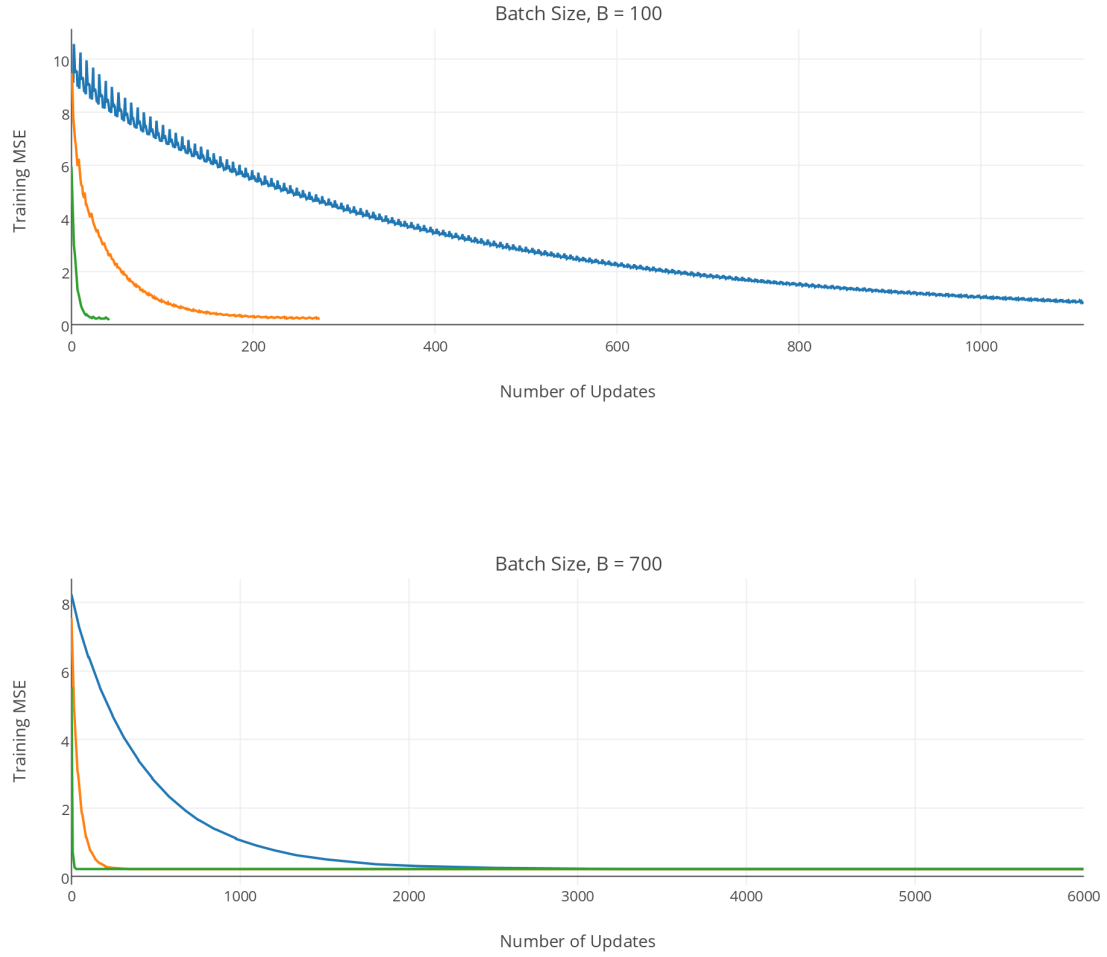
Figure 8: Subplots of training MSE against number of updates for batch sizes, $B = 100, 700$ and learning rates, $\eta = $ **0.001**, **0.01**, **0.1**. ($\lambda = 1$ for all cases)

### 2.2.3   Generalization

Based on validation set plot in Figure 9, the best value for $\lambda$ is picked based on the highest validation accuracy obtained. As observed from Table 6, $\lambda$ = 0.001, 0.01 and 0.1 all give the best validation accuracy of 93.0%. In this case, the best value of $\lambda$ chosen is 0.01, the middle ground between the three possible values.

As seen from the Figure 9, higher values of $\lambda$ initially increases the test set accuracy. This is because $\lambda$ prevents over-fitting of the model to the training set by penalizing large values of the weights. From a bias-variance trade-off perspective, incorporating $\lambda$ effectively reduces the model variance at the cost of a slight increase in bias. This effectively leads to an overall increase in the test set accuracy. This is true up to $\lambda = 0.1$.

When $\lambda = 1$, the validation and test set accuracy drops significantly. For this high value of $\lambda$, the weights are being penalized too much and it prevents the model form effectively learning and capturing key features in the input data. From a bias-variance trade-off perspective, the high $\lambda$ value has increased the model bias significantly to the point there is a net decrease in the model's performance.

$\lambda$ has to be tuned to the validation set instead of the training set. If $\lambda$ was tuned based on the training set, the model would further over-fitting the model to the training set. This is because $\lambda$ would have been tuned to optimise the value of the training MSE to data that was not used for training. In other words, it needs to perform well on data it has not seen before that is modeled by the validation set.

Table 6: Validation and test set accuracies for various values of weight-decay regularizer, $\lambda$

| $\lambda$ | Validation Accuracy | Test Accuracy |
|---|---|---|
| 0 | 0.930 | 0.900 |
| 0.0001 | 0.910 | 0.900 |
| 0.001 | 0.930 | 0.910 |
| 0.01 | 0.930 | 0.910 |
| 0.1 | 0.930 | 0.910 |
| 1 | 0.830 | 0.795 |



Figure 9: Graph of validation and test set accuracies against $\log_{10}(\lambda)$. ($\eta = 0.1, B = 50$)

# 3 Appendices

Dear Teaching Assistants, we implemented our code separately as we wanted to maximize our learning. The plots and code snippets pasted in this report come from two separate solutions which are both added as Appendices. Despite the different implementation, our results were very similar, indicating that our implementation should be correct.

For Stochastic Gradient Descent implementation, we started from Jimmy's posted code as a starter code and work from there for this assignment as suggested from Jimmy.

## 3.1 Entire Code 1: Chee Loong Soon's version

```python
1  {python}
2
3  # Assignment 1
4  # Optimization
5  # Early Stopping
6  # Learning rate decay
7  # Momentum
8
9  import tensorflow as tf
10 import numpy as np
11 import sys
12
13 # 1.2 Euclidean Distance Function
14 # 1.2.2 Pairwise Distances
15 # Write a vectorized Tensorflow Python function that
   ↪  implements
```

```python
16   # the pairwise squared Euclidean distance function for two
      ↪   input matrices.
17   # No Loops and makes use of Tensorflow broadcasting.
18   def PairwiseDistances(X, Z):
19       """
20       input:
21           X is a matrix of size (B x N)
22           Z is a matrix of size (C x N)
23       output:
24           D = matrix of size (B x C) containing the pairwise
      ↪   Euclidean distances
25       """
26       B = X.get_shape().as_list()[0]
27       N = X.get_shape().as_list()[1]
28       C = Z.get_shape().as_list()[0]
29       # Ensure the N dimensions are consistent
30       assert  N == Z.get_shape().as_list()[1]
31       # Reshape to make use of broadcasting in Python
32       X = tf.reshape(X, [B, 1, N])
33       Z = tf.reshape(Z, [1, C, N])
34       # The code below automatically does broadcasting
35       D = tf.reduce_sum(tf.square(tf.sub(X, Z)), 2)
36       return D
37
38   # 1.3 Making Predictions
39   # 1.3.1 Choosing nearest neighbours
40   # Write a vectorized Tensorflow Python function that takes a
      ↪   pairwise distance matrix
41   # and returns the responsibilities of the training examples to
      ↪   a new test data point.
```

```python
42  # It should not contain loops.
43  # Use tf.nn.top_k
44  def ChooseNearestNeighbours(D, K):
45      """
46      input:
47          D is a matrix of size (B x C)
48          K is the top K responsibilities for each test input
49      output:
50          topK are the value of the squared distances for the
    ↪ topK
51          indices are the index of the location of these squared
    ↪ distances
52      """
53      # Take topK of negative distances since it is the closest
         ↪ data.
54      topK, indices = tf.nn.top_k(tf.neg(D), K)
55      return topK, indices
56
57  # 1.3.2 Prediction
58  # Compute the k-NN prediction with K = {1, 3, 5, 50}
59  # For each value of K, compute and report:
60      # training MSE loss
61      # validation MSE loss
62      # test MSE loss
63  # Choose best k using validation error = 50
64  def PredictKnn(trainData , testData, trainTarget,  testTarget,
    ↪ K):
65      """
66      input:
67          trainData
```

37

```
68          testData
69          trainTarget
70          testTarget
71      output:
72          loss
73      """
74      D = PairwiseDistances(testData, trainData)
75      topK, indices = ChooseNearestNeighbours(D, K)
76      # Select the proper outputs to be averaged from the target
          ↪ values and average them
77      trainTargetSelectedAveraged =
          ↪ tf.reduce_mean(tf.gather(trainTarget, indices), 1)
78      # Calculate the loss from the actual values
79      # Divide by 2.0 since it's average over 2M instead of M
          ↪ where M = number of training data.
80      loss =
          ↪ tf.reduce_mean(tf.square(tf.sub(trainTargetSelectedAveraged,
          ↪ testTarget)))/2.0
81      return loss
82
83  # Plot the prediction function for x = [0, 11] on training
    ↪ data.
84  def PredictedValues(x, trainData, trainTarget, K):
85      """
86      Plot the predicted values
87      input:
88          x = test target to plot and predict
89      """
90      D = PairwiseDistances(x, trainData)
91      topK, indices = ChooseNearestNeighbours(D, K)
```

38

```python
92      predictedValues = tf.reduce_mean(tf.gather(trainTarget,
         ↪ indices), 1)
93      return predictedValues

94

95  # 1.4 Soft-Knn & Gaussian Processes
96  # 1.4.1.1 Soft Decisions
97  # Write a Tensorflow python program based on the soft k-NN
     ↪ model to compute
98  # predictions on the data1D.npy dataset.
99  # Set lambda = 100 NOT 10 as given in assignment handout
100  # and plot the test-set prediction of the model.

101

102  # Predict values using soft decision
103  def PredictedValuesSoftDecision(x, trainData, trainTarget):
104      hyperParam = 100
105      D1 = PairwiseDistances(x, trainData)
106      K1 =  tf.exp(-hyperParam*D1)
107      sum1 = tf.reduce_sum(tf.transpose(K1), axis=0)
108      N = sum1.get_shape().as_list()[0]
109      sum1 = tf.reshape(sum1, [N,1])
110      rStar = tf.div(K1, sum1)
111      predictedValues = tf.matmul(rStar,trainTarget)
112      return predictedValues

113

114  # Predict values using Gaussian
115  # 1.4.1.1 Gaussian Processes
116  def PredictedValuesGaussianProcesses(x, trainData,
     ↪ trainTarget):
117      hyperParam = 100
118      D1 = PairwiseDistances(x, trainData)
```

39

```python
119     K1 =  tf.exp(-hyperParam*D1)

120     D2 = PairwiseDistances(trainData, trainData)

121     K2 =  tf.matrix_inverse(tf.exp(-hyperParam*D2))

122     rStar = tf.matmul(K1, K2)

123     predictedValues = tf.matmul(rStar,trainTarget)

124     return predictedValues

125

126 # Comment on the difference you observe between two programs

127 # Gaussian has higher loss.

128

129 # 2 Linear and Logistic Regression

130 # 2.2 Stochastic Gradient Descent

131 # Implement linear regression and stochastic gradient descent
    ↪  algorithm

132 # with mini-batch size B = 50.

133 def buildGraph(learningRate, weightDecayCoeff):

134     # Variable creation

135     W = tf.Variable(tf.truncated_normal(shape=[64, 1],
        ↪  stddev=0.5), name='weights')

136     b = tf.Variable(0.0, name='biases')

137     X = tf.placeholder(tf.float32, [None, 64], name='input_x')

138     y_target = tf.placeholder(tf.float32, [None,1],
        ↪  name='target_y')

139     weightDecay =
        ↪  tf.div(tf.constant(weightDecayCoeff),tf.constant(2.0))

140     # Graph definition

141     y_predicted = tf.matmul(X,W) + b

142     # Error definition

143     meanSquaredError =
        ↪  tf.reduce_mean(tf.reduce_mean(tf.square(y_predicted -
        ↪  y_target),
```

```
144                                               ↪ reduction_indices=1,

145                                               ↪ name='squared_error')

146                              name='mean_squared_error')

147     weightDecayMeanSquareError =
          ↪ tf.reduce_mean(tf.reduce_mean(tf.square(weightDecay)))

148     weightDecayTerm = tf.multiply(weightDecay,
          ↪ weightDecayMeanSquareError)

149     meanSquaredError =
          ↪ tf.add(meanSquaredError,weightDecayTerm)

150

151     # Training mechanism

152     optimizer =
          ↪ tf.train.GradientDescentOptimizer(learning_rate =
          ↪ learningRate)

153     train = optimizer.minimize(loss=meanSquaredError)

154     return W, b, X, y_target, y_predicted, meanSquaredError,
          ↪ train

155

156

157 def ShuffleBatches(trainData, trainTarget):

158     rngState = np.random.get_state()

159     np.random.shuffle(trainData)

160     np.random.set_state(rngState)

161     np.random.shuffle(trainTarget)

162     return trainData, trainTarget

163

164 def LinearRegression(trainData, trainTarget, validData,
      ↪ validTarget, testData, testTarget):
```

41

```python
165     figureCount = 30
166     # 2.2.3 Generalization (done by partner)
167     # Run SGD with B = 50 and use validation performance to
          ↪ choose best weight decay coefficient
168     # from weightDecay = {0., 0.0001, 0.001, 0.01, 0.1, 1.}
169     # Plot weightDecay vs test set accuracy. (Done by partner)
170     weightDecayTrials= [0.0, 0.0001, 0.0001, 0.01, 0.1, 1.0]
171     # Plot total loss function vs number of updates for the
          ↪ best learning rate found
172     learningRateTrials = [0.1, 0.01, 0.001]
173     # 2.2.2 Effect of the mini-batch size
174     # Run with Batch Size, B = {10, 5, 100, 700} and tune the
          ↪ learning rate separately for each mini-batch size.
175     # Plot  the total loss function vs the number of updates
          ↪ for each mini-batch size.
176     miniBatchSizeTrials = [10, 50, 100, 700]
177     learningRate = 0.01
178     miniBatchSize = 10
179     weightDecayCoeff = 1.0
180     # for weightDecayCoeff in weightDecayTrials:
181     for miniBatchSize in miniBatchSizeTrials:
182         for learningRate in learningRateTrials:
183             # Build computation graph
184             W, b, X, y_target, y_predicted, meanSquaredError,
                  ↪ train = buildGraph(learningRate,
                  ↪ weightDecayCoeff)
185             # Initialize session
186             init = tf.global_variables_initializer()
187             sess = tf.InteractiveSession()
188             sess.run(init)
```

```python
            initialW = sess.run(W)
            initialb = sess.run(b)

            # print "Initial weights: %s, initial bias: %.2f",
            ↪ initialW, initialb
            # Training model
            numEpoch = 200
            currEpoch = 0
            wList = []

            xAxis = []
            yTrainErr = []
            yValidErr = []
            yTestErr = []
            numUpdate = 0
            step = 0
            errTrain = -1
            errValid = -1
            errTest = -1
            while currEpoch <= numEpoch:
                # Shuffle the batches and return
                trainData, trainTarget =
                ↪ ShuffleBatches(trainData, trainTarget)
                step = 0
                # Full batch
                while step*miniBatchSize < 700:
                    _, errTrain, currentW, currentb, yhat =
                    ↪ sess.run([train, meanSquaredError, W,
                    ↪ b, y_predicted], feed_dict={X:
                    ↪ trainData[step*miniBatchSize:(step+1)*miniBatchSi
                    ↪ y_target:
                    ↪ trainTarget[step*miniBatchSize:(step+1)*miniBatch
```

```python
214              wList.append(currentW)
215              #if not (step*miniBatchSize % 50):
216              #    print "Iter: %3d, MSE-train: %4.2f,
                 ↪ weights: %s, bias: %.2f", step, err,
                 ↪ currentW.T, currentb
217              step = step + 1
218              xAxis.append(numUpdate)
219              numUpdate += 1
220              yTrainErr.append(errTrain)
221              errValid = sess.run(meanSquaredError,
                 ↪ feed_dict={X: validData, y_target:
                 ↪ validTarget})
222              errTest = sess.run(meanSquaredError,
                 ↪ feed_dict={X: testData, y_target:
                 ↪ testTarget})
223              yValidErr.append(errValid)
224              yTestErr.append(errTest)
225          # Testing model
226          # TO know what is being run
227          currEpoch += 1
228      print "LearningRate: " , learningRate, " Mini
             ↪ batch Size: ", miniBatchSize
229      print "Iter: ", numUpdate
230      print "Final Train MSE: ", errTrain
231      print "Final Valid MSE: ", errValid
232      print "Final Test MSE: ", errTest
233      import matplotlib.pyplot as plt
234      plt.figure(figureCount)
235      figureCount = figureCount + 1
236      plt.plot(np.array(xAxis), np.array(yTrainErr))
```

```python
237            plt.savefig("TrainLossLearnRate" +
             ↪  str(learningRate) + "Batch" +
             ↪  str(miniBatchSize) + '.png')

238

239            plt.figure(figureCount)
240            figureCount = figureCount + 1
241            plt.plot(np.array(xAxis), np.array(yValidErr))
242            plt.savefig("ValidLossLearnRate" +
             ↪  str(learningRate) + "Batch" +
             ↪  str(miniBatchSize) + '.png')
243            plt.figure(figureCount)
244            figureCount = figureCount + 1
245            plt.plot(np.array(xAxis), np.array(yTestErr))
246            plt.savefig("TestLossLearnRate" +
             ↪  str(learningRate) + "Batch" +
             ↪  str(miniBatchSize) + '.png')
247        return

248

249    def SortData(inputVal, outputVal):
250        """
251        This sorts a given test set by the dataValue before
         ↪  plotting it.
252        """
253        p = np.argsort(inputVal, axis=0)
254        inputVal = np.array(inputVal)[p]
255        outputVal = np.array(outputVal)[p]
256        inputVal = inputVal[:, :,0]
257        outputVal = outputVal[:, :,0]
258        return inputVal, outputVal

259
```

```python
if __name__ == "__main__":
    print 'helloworld'
    N = 2 # number of dimensions
    B = 3 # number of test inputs (To get the predictions for
     ↪ all these inputs
    C = 2 # number of training inputs (Pick closest k from
     ↪ this C)
    X = tf.constant([1, 2, 3, 4, 5, 6], shape=[3, 2])
    Z = tf.constant([21, 22, 31, 32], shape=[2, 2])
    # Need to put seed so random_uniform doesn't generate new
     ↪ random values
    # each time you evaluate when you print, so then the
     ↪ values would be
    # inconsistent as to what you would have used or checked
    #X = tf.random_uniform([B, N], seed=111)*30
    #Z = tf.random_uniform([C, N], seed=112)*30
    D = PairwiseDistances(X, Z)
    K = 1 # number of nearest neighbours
    # You calculate all the pairwise distances between each
     ↪ test input
    # and existing training input
    topK, indices = ChooseNearestNeighbours(D, K)
    # Prediction
    #for K in [1, 3, 5, 50]:
    for K in [1]:
        np.random.seed(521)
        Data = np.linspace(1.0 , 10.0 , num =100) [:,
         ↪ np.newaxis]
        Target = np.sin( Data ) + 0.1 * np.power( Data , 2) +
         ↪ 0.5 * np.random.randn(100 , 1)
```

46

```
283        randIdx = np.arange(100)
284        np.random.shuffle(randIdx)
285        # data1D.npy
286        trainData, trainTarget  = Data[randIdx[:5]],
           ↪ Target[randIdx[:5]]
287        trainData, trainTarget  = Data[randIdx[:80]],
           ↪ Target[randIdx[:80]]
288        validData, validTarget = Data[randIdx[80:90]],
           ↪ Target[randIdx[80:90]]
289        testData, testTarget = Data[randIdx[90:93]],
           ↪ Target[randIdx[90:93]]
290        testData, testTarget = Data[randIdx[90:100]],
           ↪ Target[randIdx[90:100]]
291
292        #trainData, trainTarget = SortData(trainData,
           ↪ trainTarget)
293        #validData, validTarget = SortData(validData,
           ↪ validTarget)
294        testData, testTarget = SortData(testData, testTarget)
295
296
297        # Convert to tensors from numpy
298        trainData = tf.pack(trainData)
299        validData = tf.pack(validData)
300        testData = tf.pack(testData)
301        trainTarget = tf.pack(trainTarget)
302        validtarget = tf.pack(validTarget)
303        testTarget = tf.pack(testTarget)
304        trainMseLoss = PredictKnn(trainData, trainData,
           ↪ trainTarget, trainTarget, K)
```

47

```python
305         validationMseLoss = PredictKnn(trainData, validData,
            ↪ trainTarget, validTarget, K)
306         testMseLoss = PredictKnn(trainData, testData,
            ↪ trainTarget, testTarget, K)
307         init = tf.global_variables_initializer()
308         '''
309         with tf.Session() as sess:
310             sess.run(init)
311             print 'K ' + str(K)
312             print 'trainMseLoss'
313             print sess.run(trainMseLoss)
314             print 'validationMseLoss'
315             print sess.run(validationMseLoss)
316             print 'testMseLoss'
317             print sess.run(testMseLoss)
318         '''
319         # Plot the prediction for the x below
320         x = np.linspace(0.0, 11.0, num=1000)[:, np.newaxis]
321         xTensor = tf.pack(x)
322         predictedValuesKnn = PredictedValues(xTensor,
            ↪ trainData, trainTarget, K)
323         predictedValuesSoft =
            ↪ PredictedValuesSoftDecision(testData, trainData,
            ↪ trainTarget)
324         predictedValuesGaussian =
            ↪ PredictedValuesGaussianProcesses(testData,
            ↪ trainData, trainTarget)
325         lossSoft =
            ↪ tf.reduce_mean(tf.square(tf.sub(predictedValuesSoft,
            ↪ testTarget)))/2.0
```

48

```python
326        lossGaussian =
       ↪ tf.reduce_mean(tf.square(tf.sub(predictedValuesGaussian,
       ↪ testTarget)))/2.0
327    import matplotlib.pyplot as plt
328    plt.figure(0)
329    init = tf.global_variables_initializer()
330    with tf.Session() as sess:
331        sess.run(init)
332        plt.figure(K+100)
333        plt.scatter(sess.run(trainData),
            ↪ sess.run(trainTarget))
334        plt.plot(sess.run(xTensor),
            ↪ sess.run(predictedValuesKnn))
335        fileName = str("KNN") + str(K) +
            ↪ str("trainingGraph.png")
336        plt.savefig(fileName)
337
338        # Plot for SoftDecision
339        plt.figure(K+101)
340        plt.title("Soft Decision kNN on Test Set, MSE = "
            ↪ + str(sess.run(lossSoft)))
341        plt.xlabel("Data Value")
342        plt.ylabel("Target Value")
343        plt.scatter(sess.run(testData),
            ↪ sess.run(testTarget), label= "testValue")
344        plt.plot(sess.run(testData),
            ↪ sess.run(predictedValuesSoft), label =
            ↪ "predicted")
345        plt.legend()
346        fileName = str("SoftDecision.png")
```

```
347            plt.savefig(fileName)
348            print 'SoftDecisionLoss'
349            print sess.run(lossSoft)
350
351            # Plot for Gaussian
352            plt.figure(K+102)
353            plt.title("Gaussian Process Regression on Test
                ↪ Set, MSE = " + str(sess.run(lossGaussian)))
354            plt.xlabel("Data Value")
355            plt.ylabel("Target Value")
356            plt.scatter(sess.run(testData),
                ↪ sess.run(testTarget), label = "testValue")
357            plt.plot(sess.run(testData),
                ↪ sess.run(predictedValuesGaussian), label =
                ↪ "predicted")
358            plt.legend()
359            fileName = str("ConditionalGaussian.png")
360            plt.savefig(fileName)
361            print 'ConditionalGaussianLoss'
362            print sess.run(lossGaussian)
363    # Part 2
364    with np.load ("tinymnist.npz") as data :
365        trainData, trainTarget = data ["x"], data["y"]
366        validData, validTarget = data ["x_valid"], data
                ↪ ["y_valid"]
367        testData, testTarget = data ["x_test"], data
                ↪ ["y_test"]
368        LinearRegression(trainData, trainTarget,validData,
                ↪ validTarget, testData, testTarget)
```

## 3.2 Entire Code 2: FuYuan Tee's version

### 3.2.1 Question 1: k-Nearest Neighbour

```python
import tensorflow as tf
import numpy as np

import matplotlib.pyplot as plt

import plotly.plotly as py
import plotly.graph_objs as go
from plotly import tools

import plotly.offline as pyo
pyo.init_notebook_mode(connected=True)

# Generate data
np.set_printoptions(precision=3)
np.random.seed(521)

# Generating data
Data = np.linspace(1.0, 10.0, num=100) [:, np.newaxis]
Target = np.sin(Data) + 0.1 * np.power(Data, 2) \
    + 0.5 * np.random.randn(100, 1)

# Generating a random index
randIdx = np.arange(100)
np.random.shuffle(randIdx)

# Partitioning 100 datapoints into training, validation and
```

```
27  # test sets consisting of 80, 10 and 10 points respectively.
28  trainData, trainTarget = Data[randIdx[:80]],
    ↪  Target[randIdx[:80]]
29  validData, validTarget = Data[randIdx[80:90]],
    ↪  Target[randIdx[80:90]]
30  testData, testTarget = Data[randIdx[90:]],
    ↪  Target[randIdx[90:]]
31
32  # Defining TensorFlow Variables
33
34  # Calculates pairwise squared Euclidean distance between
    ↪  matrices X and Z
35  # X is the input matrix which houses data points to be
    ↪  calculated,
36  # which are calculated against reference datapoints in Z
37  def euclideanDistance(X, Z):
38      if tf.TensorShape.num_elements(X.get_shape()) == 3:
39          return tf.reduce_sum(tf.square(X - tf.transpose(Z)),
            ↪  axis=2) # Sums feature deviations for each
            ↪  variable
40      else:
41          return tf.square(X - tf.transpose(Z),
            ↪  name='Euclidean_Distance_Matrix')
42
43  # Produces a BxC responsibility vector
44  def responsibilityVector(D, k):
45      def flatten(tensor):
46          return tf.reshape(tensor, [-1])
47
48      # Obtains indices of top-K points
```

```python
49      _, idx = tf.nn.top_k(tf.transpose(-D), k)

50

51      # Creates a step sequence for M values repeated k times
52      M = tf.shape(D)[1]
53      step_seq =
        ↪ flatten(tf.transpose(tf.reshape(tf.tile(tf.range(0,
        ↪ M), [k]), [k, M]))

54

55      # Form new index key compatible for subsequent
        ↪ sparse_to_dense op
56      sparse_idx = tf.pack([step_seq, flatten(idx)], axis=1,
        ↪ name='Sparse_Indices')

57

58      # Forms dense tensor
59      return tf.sparse_to_dense(tf.cast(sparse_idx, tf.int32), \

60
                                ↪ tf.cast(tf.shape(tf.transpose(D)),
                                ↪ tf.int32), \

61
                                ↪ tf.fill([tf.shape(sparse_idx)[0]],
                                ↪ tf.divide(1.0, tf.cast(k,
                                ↪ tf.float32))), \
62                              validate_indices=False, \
63                              name='Responsibility_Vector')

64

65  # Function that creates a TensorFlow model
66  def buildGraph(k_):
67      k = tf.constant(k_, name='k') # Hyperparameter
68      X = tf.placeholder(tf.float32, shape=[None, None],
        ↪ name='Training_Data')
```

```python
69      Y = tf.placeholder(tf.float32, shape=[None, None],
         ↪  name='Training_Target')
70      Z = tf.placeholder(tf.float32, shape=[None, None],
         ↪  name='Input_Data')
71      T = tf.placeholder(tf.float32, shape=[None, None],
         ↪  name='Input_Target')
72
73      D = euclideanDistance(X, Z)
74      R = responsibilityVector(D, k)
75
76      Y_hat = tf.matmul(R, Y, name='Y_hat')
77      MSE = tf.divide(tf.reduce_sum(tf.square(T - Y_hat)), \
78                      tf.scalar_mul(2, tf.cast(tf.shape(Z)[0],
                         ↪  tf.float32)), \
79                      name='Mean_Squared_Error') # Half of
                         ↪  theorectical MSE
80
81      return X, Y, Z, T, Y_hat, MSE
82
83  # Function used to adapt 'kNN' function based on type of
     ↪  dataset used for calculation
84  # Mode consists of either ['train', 'validation', 'test',
     ↪  'full']
85  def selectDataPartition(mode):
86
87      inputData = trainData
88      inputTarget = trainTarget
89
90      if mode == 'train':
91          dataSource = trainData
```

```python
92              targetSource = trainTarget
93          elif mode == 'validation':
94              dataSource = validData
95              targetSource = validTarget
96          elif mode == 'test':
97              dataSource = testData
98              targetSource = testTarget
99          elif mode == 'full':
100             dataSource = np.linspace(0.0, 11.0, num=1000) [:,
              ↪ np.newaxis]
101             targetSource = np.sin(dataSource) + 0.1 *
              ↪ np.power(dataSource, 2) \
102                            + 0.5 * np.random.randn(1000, 1)
103
104             inputData = trainData
105             inputTarget = trainTarget
106
107         return inputData, inputTarget, dataSource, targetSource
108
109     def kNN(k, mode):
110         X, Y, Z, T, Y_hat, MSE = buildGraph(k)
111
112         with tf.Session() as sess:
113     #         print "k: %d, mode: %s" % (k, mode)
114             inputData, inputTarget, dataSource, targetSource =
              ↪ selectDataPartition(mode)
115
116             error, y_pred = sess.run([MSE, Y_hat], \
117                                       feed_dict={X: inputData, Y:
                                      ↪ inputTarget, \
```

```python
118                                             Z: dataSource, T:
                                            ↪ targetSource
119                                         })
120
121     return error, np.transpose(np.append(inputData,
        ↪ inputTarget, axis=1)),
        ↪ np.transpose(np.append(dataSource, y_pred, axis=1))
122
123 # Main Function
124 k_list = [1, 3, 5, 50]
125 kNN_mode = ['train', 'validation', 'test', 'full']
126
127 MSE_list = []
128 targetSeries = []
129 predictionSeries = []
130
131 # Performs kNN based on various calculation modes
132 for i in range(4):
133     for j in range(4):
134         MSE, target, prediction = kNN(k_list[i], kNN_mode[j])
135         MSE_list.append(MSE)
136         targetSeries.append(target)
137         predictionSeries.append(prediction)
138 MSE_list = np.reshape(MSE_list, (4, 4))
139
140 # Generate interactive Plotly graph
141 def generateVisualisation(targetSeries, predictionSeries,
    ↪ k_list):
142     subplotTitleString = []
143     for i in range(len(k_list)):
```

```
144        subplotTitleString.append('k = %s' % str(k_list[i]))

145

146    fig = tools.make_subplots(rows=2, cols=2,
        ↪ subplot_titles=(subplotTitleString))

147

148    for i in range(len(k_list)):
149        traceData = go.Scatter(
150            x = targetSeries[i * 4][0],
151            y = targetSeries[i * 4][1],
152            marker = {'color': 'blue',
153                      'symbol': 200},
154            mode = 'markers',
155            name = 'data_k=' + str(k_list[i])
156        )
157        tracePred = go.Scatter(
158            x = predictionSeries[4 * (i + 1) - 1][0],
159            y = predictionSeries[4 * (i + 1) - 1][1],
160            marker = {'color': 'green'},
161            mode = 'lines',
162            name = 'prediction_k=' + str(k_list[i])
163        )

164

165        fig.append_trace(traceData, i / 2 + 1, i % 2 + 1)
166        fig.append_trace(tracePred, i / 2 + 1, i % 2 + 1)

167

168        fig['layout']['xaxis'+str(i+1)].update(title='x')
169        fig['layout']['yaxis'+str(i+1)].update(title='y')

170

171    fig['layout'].update(height=900, width=950, title='k-NN
        ↪ Regression on data1D', showlegend=False)
```

```python
172     return py.iplot(fig, filename='A1Q1_kNN_subplot2x2')
173
174 # Output summary table
175 print 'Mean Squared Error Summary:'
176 for i in range(5):
177     if i == 0:
178         print "%3s %10s %10s %6s" % ('k', 'Training',
            ↪ 'Validation', 'Test')
179     else:
180         print "%3d %10.3f %10.3f %6.3f" % (k_list[i - 1],
            ↪ MSE_list[i - 1][0], MSE_list[i - 1][1], MSE_list[i
            ↪ - 1][2])
181
182 print "\n\n\n"
183 kNN_visuals = generateVisualisation(targetSeries,
    ↪ predictionSeries, k_list)
184 kNN_visuals
```

### 3.2.2 Question 2.2: Stochastic Gradient Descent

```python
1 {python}
2 # Import relevant packages
3 import tensorflow as tf
4 import numpy as np
5 import math
6
7 import time
8
9 # Non-interactive plotting
```

```python
10  import matplotlib.pyplot as plt
11  from IPython import display
12
13  # Interactive plotting
14  from plotly import tools
15  import plotly.plotly as py
16  import plotly.graph_objs as go
17  import plotly.offline as pyo
18  from plotly.offline import download_plotlyjs
19
20  # Configure environment
21  np.set_printoptions(precision=3)
22  np.random.seed(521)
23
24  # Activate Plotly Offline for Jupyter
25  pyo.init_notebook_mode(connected=True)
26
27  # Load Tiny MNIST dataset
28  with np.load ("tinymnist.npz") as data:
29      trainData, trainTarget = data ["x"], data["y"]
30      validData, validTarget = data ["x_valid"], data
          ↪ ["y_valid"]
31      testData, testTarget = data ["x_test"], data ["y_test"]
32
33  # Create Tensorflow Graph
34  def buildGraph(eta, lambda_):
35      # Model inputs
36      X = tf.placeholder(tf.float32, shape=[None, None],
          ↪ name='Input')
37      Y = tf.placeholder(tf.float32, shape=[None, None],
          ↪ name='Target')
```

```python
38
39      # Model variables
40      W = tf.Variable(tf.truncated_normal(shape=[64, 1],
        ↪  stddev=0.5), name='Weights')
41      b = tf.Variable(0.0, name='Biases')
42
43      # Model parameters
44      eta = tf.constant(eta, name='Learning_Rate')
45      lambda_ = tf.constant(lambda_, name='L2_Regularizer')
46
47      # Predicted target
48      Y_hat = tf.matmul(X, W)
49
50      # Mean squared error
51      MSE = tf.scalar_mul(tf.divide(1.0, tf.cast(tf.shape(X)[0],
        ↪  tf.float32)), \
52                          tf.reduce_sum(tf.square(Y_hat - Y)))
                            ↪  \
53          + tf.scalar_mul(tf.divide(tf.cast(lambda_,
            ↪  tf.float32), 2.0), tf.matmul(tf.transpose(W),
            ↪  W))
54
55      # Basic accuracy definition (n_correct / n_total)
56      Y_hat_thresholded = tf.cast(tf.greater_equal(Y_hat, 0.5),
        ↪  tf.float32)
57      accuracy =
        ↪  tf.divide(tf.reduce_sum(tf.cast(tf.equal(Y_hat_thresholded,
        ↪  Y), tf.float64)), \
58                          tf.cast(tf.shape(X)[0], tf.float64))
59
```

```python
60      # Basic gradient descent optimizer
61      optimizer =
         ↪ tf.train.GradientDescentOptimizer(eta).minimize(MSE)
62
63      return W, b, X, Y, Y_hat, MSE, accuracy, optimizer
```

### 3.2.3 Question 2.2.1: Tuning the learning rate

```python
1   {python}
2   # Tune learning rate
3   MAX_ITER = 2000
4   def tuneLearningRate(etaList, batchSize=50, lambda_=1):
5       # Returns the i-th batch of training data and targets
6       # Generates a new, reshuffled batch once all previous
         ↪ batches are fed
7       def getNextTrainingBatch(currentIter):
8           currentBatchNum = currentIter % (trainData.shape[0] /
             ↪ batchSize)
9           if currentBatchNum == 0:
10              np.random.shuffle(randIdx)
11          # print 'Iteration: %4d, BatchCap: %2d, BatchNum: %2d'
             ↪ % (currentIter, trainData.shape[0] / batchSize,
             ↪ currentBatchNum)
12          lowerBoundIdx = currentBatchNum * batchSize
13          upperBoundIdx = (currentBatchNum + 1) * batchSize
14          return trainData[lowerBoundIdx:upperBoundIdx],
             ↪ trainTarget[lowerBoundIdx:upperBoundIdx]
15
16      # Generate updated plots for training and validation MSE
```

```python
17    def plotMSEGraph(MSEList, param):
18        label = '$\eta$ = ' + str(param)
19        label_classification = ['train.', 'valid.']
20
21        display.clear_output(wait=True)
22        plt.figure(figsize=(8,5), dpi=200)
23
24        for i, MSE in enumerate(MSEList):
25            plt.plot(range(len(MSE)), MSE, '.', markersize=3,
                 ↪ label=label+' '+label_classification[i])
26
27        plt.axis([0, MAX_ITER, 0, np.amax(MSEList)])
28        plt.legend()
29        plt.show()
30
31    # Calculates the ratio between the n-th average epoch MSE
       ↪ and the (n-1)-th average epoch MSE
32    def ratioAverageEpochMSE(currentValidMSE):
33        averageN =
           ↪ np.average(currentValidMSE[-(np.arange(epochSize -
           ↪ 1) + 1)])
34        averageNlessOne =
           ↪ np.average(currentValidMSE[-(np.arange(epochSize -
           ↪ 1) + epochSize)])
35        return averageN / averageNlessOne
36
37    # Returns True if the average epoch validation MSE is at
       ↪ least 99% of the previous epoch average.
38    # i.e. Returns True if the average learnings between epoch
       ↪ is less than +1%
```

```python
39      # Otherwise, returns False
40      def shouldStopEarly(currentValidMSE):
41          if currentValidMSE.shape[0] < 2 * epochSize:
42              return False
43          return True if (ratioAverageEpochMSE(currentValidMSE)
             ↪ >= 0.99) else False
44
45      summaryList = []
46      randIdx = np.arange(trainData.shape[0])
47      epochSize = trainData.shape[0] / batchSize
48
49      for eta in etaList:
50          W, b, X, Y, Y_hat, MSE, accuracy, optimizer =
             ↪ buildGraph(eta, lambda_)
51
52          with tf.Session() as sess:
53              tf.global_variables_initializer().run()
54
55              # Creates blank training and validation MSE arrays
                 ↪ for the Session
56              currentTrainMSE = np.array([])[:, np.newaxis]
57              currentValidMSE = np.array([])[:, np.newaxis]
58
59              # Runs update
60              currentIter = 0
61              while currentIter <= MAX_ITER:
62                  inputData, inputTarget =
                     ↪ getNextTrainingBatch(currentIter)
63
64                  _, trainError = sess.run([optimizer, MSE],
                     ↪ feed_dict={X: inputData, Y: inputTarget})
```

```python
          validError = sess.run([MSE], feed_dict={X:
          ↪ validData, Y: validTarget})


          currentTrainMSE = np.append(currentTrainMSE,
          ↪ trainError)
          currentValidMSE = np.append(currentValidMSE,
          ↪ validError)


          # Update graph of training and validation MSE
          ↪ arrays
          if (currentIter < 3) or (currentIter % 500 ==
          ↪ 0):
              plotMSEGraph([currentTrainMSE,
              ↪ currentValidMSE], eta)


          # At every epoch, check for early stopping
          ↪ possibilty. If so, breaks from while loop
          if currentIter % epochSize == 0:
              if shouldStopEarly(currentValidMSE):
                  break


          currentIter += 1


   # Save session results as dictionary and appends to
   ↪ MSEsummaryList
   summaryList.append(
       {
           'eta': eta,
           'B': batchSize,
           'lambda': lambda_,
```

64

```python
                   'numIter': currentIter + 1,
                   'epoch': float(currentIter + 1) /
                   ↪ (trainData.shape[0] / batchSize),
                   'trainMSE': currentTrainMSE,
                   'validMSE': currentValidMSE,
               }
           )


   return summaryList



# Main Function
etaList = [0.001, 0.01, 0.1]
tunedEtaSummary = tuneLearningRate(etaList)


# Output summary table
for summary in tunedEtaSummary:
   print 'eta: %.3f, numIter: %d, validMSE: %.3f' %
       ↪ (summary['eta'], summary['numIter'],
       ↪ summary['validMSE'][-1])


# Produce interactive graph for best learning rate
def etaIGraph(tunedEtaSummary):
   # Create plot for each summary
   traceList = []
   for summary in tunedEtaSummary:
       traceList.append(
           go.Scatter(
               x = range(summary['numIter'] + 1),
               y = summary['trainMSE'],
```

```python
114                 name = '$\\eta = ' + str(summary['eta']) + '$'
115            )
116        )
117    data = go.Data(traceList)
118
119    # Create figure layout
120    layout = go.Layout(
121        title = '$\\textit{Training performance for various
            ↪ learning rates, } \\eta$',
122        xaxis = {'title': 'Number of Updates'},
123        yaxis = {'title': 'Training MSE'},
124    )
125
126    figure = go.Figure(data=data, layout=layout)
127    return py.iplot(figure, filename='A1Q2.1_bestEtaGraph')
128 fig2_1 = etaIGraph(tunedEtaSummary)
129 fig2_1
```

### 3.2.4 Question 2.2.2: Effect of the mini-batch size

```python
1  {python}
2  # Tune the mini-batch size
3  MAX_ITER = 6000
4  def tuneBatchSize(etaList, batchSizeList, lambda_=1):
5      # Returns the i-th batch of training data and targets
6      # Generates a new, reshuffled batch once all previous
          ↪ batches are fed
7      def getNextTrainingBatch(currentIter):
8          currentBatchNum = currentIter % (trainData.shape[0] /
              ↪ batchSize)
```

66

```
9          if currentBatchNum == 0:
10             np.random.shuffle(randIdx)
11         # print currentBatchNum + 1
12         lowerBoundIdx = currentBatchNum * batchSize
13         upperBoundIdx = (currentBatchNum + 1) * batchSize
14         return trainData[lowerBoundIdx:upperBoundIdx],
         ↪ trainTarget[lowerBoundIdx:upperBoundIdx]

15

16     # Generate updated plots for training and validation MSE
17     def plotMSEGraph(MSEList, param):
18         label = '$B$ = ' + str(param[0]) + ', $\eta$: ' +
         ↪ str(param[1])
19         label_classification = ['train.', 'valid.']

20

21         display.clear_output(wait=True)
22         plt.figure(figsize=(8,5), dpi=200)

23

24         for i, MSE in enumerate(MSEList):
25             plt.plot(range(len(MSE)), MSE, '.', markersize=3,
             ↪ label=label+'\n'+label_classification[i])

26

27         plt.axis([0, MAX_ITER, 0, np.amax(MSEList)])
28         plt.legend()
29         plt.show()

30

31     # Calculates the ratio between the n-th average epoch MSE
     ↪ and the (n-1)-th average epoch MSE
32     def ratioAverageEpochMSE(currentValidMSE):
33         averageN =
         ↪ np.average(currentValidMSE[-(np.arange(epochSize -
         ↪ 1) + 1)])
```

```python
34            averageNlessOne =
              ↪  np.average(currentValidMSE[-(np.arange(epochSize -
              ↪  1) + epochSize)])
35            return averageN / averageNlessOne
36
37        # Returns True if the average epoch validation MSE is at
          ↪  least 99% of the previous epoch average.
38        # i.e. Returns True if the average learnings between epoch
          ↪  is less than +1%
39        # Otherwise, returns False
40        def shouldStopEarly(currentValidMSE):
41            if currentValidMSE.shape[0] < 2 * epochSize:
42                return False
43            return True if (ratioAverageEpochMSE(currentValidMSE)
              ↪  >= 0.99) else False
44
45        summaryList = []
46        randIdx = np.arange(trainData.shape[0])
47
48        for batchSize in batchSizeList:
49            epochSize = trainData.shape[0] / batchSize
50            batchSummary = []
51            for eta in etaList:
52                W, b, X, Y, Y_hat, MSE, accuracy, optimizer =
                  ↪  buildGraph(eta, lambda_)
53
54                with tf.Session() as sess:
55                    tf.global_variables_initializer().run()
56
57                    # Creates blank training and validation MSE
                      ↪  arrays for the Session
```

```python
currentTrainMSE = np.array([])[:, np.newaxis]
currentValidMSE = np.array([])[:, np.newaxis]


# Runs update
currentIter = 0
while currentIter <= MAX_ITER:
    inputData, inputTarget =
    ↪ getNextTrainingBatch(currentIter)


    _, trainError = sess.run([optimizer, MSE],
    ↪ feed_dict={X: inputData, Y:
    ↪ inputTarget})
    validError = sess.run([MSE], feed_dict={X:
    ↪ validData, Y: validTarget})


    currentTrainMSE =
    ↪ np.append(currentTrainMSE, trainError)
    currentValidMSE =
    ↪ np.append(currentValidMSE, validError)


    # Update graph of training and validation
    ↪ MSE arrays
    if (currentIter < 3) or (currentIter % 500
    ↪ == 0):
        plotMSEGraph([currentTrainMSE,
        ↪ currentValidMSE], [batchSize,
        ↪ eta])


    # At every epoch, check for early stopping
    ↪ possibilty. If so, breaks from while
    ↪ loop
```

```python
77                     if currentIter % epochSize == 0:
78                         if shouldStopEarly(currentValidMSE):
79                             break
80
81                     currentIter += 1
82
83          # Save session results as dictionary and appends
             ↪  to MSEsummaryList
84          batchSummary.append(
85              {
86                  'eta': eta,
87                  'B': batchSize,
88                  'lambda': lambda_,
89                  'numIter': currentIter + 1,
90                  'epoch': float(currentIter + 1) /
                     ↪  (trainData.shape[0] / batchSize),
91                  'trainMSE': currentTrainMSE,
92                  'validMSE': currentValidMSE,
93              }
94          )
95      summaryList.append(batchSummary)
96
97  return summaryList
98
99  # Main function
100 etaList = [0.001, 0.01, 0.1]
101 batchSizeList = [10, 50, 100, 700]
102 tunedBatchSizeSummary = tuneBatchSize(etaList, batchSizeList)
103
104 # Output summary table:
```

```python
105  for batchSummary in tunedBatchSizeSummary:
106      for summary in batchSummary:
107          print 'B: %5d, eta: %5.3f, numIter: %5d, validMSE:
         ↪ %3.3f' % \
108              (summary['B'], summary['eta'], summary['numIter'],
             ↪ summary['validMSE'][-1])
109
110  # Generate interactive Plotly plot
111  def batchSizeIGraphSubplot(tunedBatchSizeSummary):
112
113      # Define subplot title
114      subplotTitle = []
115      for batchSummary in tunedBatchSizeSummary:
116          subplotTitle.append('Batch Size, B  = ' +
             ↪ str(batchSummary[0]['B']))
117
118      # Define subplot figure
119      figure = tools.make_subplots(rows=4, cols=1,
          ↪ subplot_titles=(subplotTitle))
120
121      # Define color list
122      colorList = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728',
          ↪ '#9467bd', '#8c564b']
123
124      # Create plot for each summary
125      for i, batchSummary in enumerate(tunedBatchSizeSummary):
126          traceList = []
127          for j, summary in enumerate(batchSummary):
128              trace = go.Scatter(
129                  x = range(summary['numIter'] + 1),
```

```
130                    y = summary['trainMSE'],
131                    marker = {'color': colorList[j]},
132                    name = '$B=' + str(summary['B']) + ', \\eta='
                      ↪  + str(summary['eta']) + '$'
133                )
134            figure.append_trace(trace, i + 1, 1)
135
           ↪  figure['layout']['xaxis'+str(i+1)].update(title='Number
           ↪  of Updates')
136
           ↪  figure['layout']['yaxis'+str(i+1)].update(title='Training
           ↪  MSE')
137
138    # Create figure layout
139    figure['layout'].update(
140        height = 1800,
141        title = '$\\textit{Model training performance for
              ↪  various batch size, } B' + \
142            '\\textit{, and learning rate, } \\eta$',
143        showlegend = False
144    )
145
146    return py.iplot(figure,
          ↪  filename='A1Q2.2_batch_size_subplot2x2')
147 fig2_2_visual = batchSizeIGraphSubplot(tunedBatchSizeSummary)
148 fig2_2_visual
```

### 3.2.5 Question 2.2.3: Generalization

```python
{python}
# Tune weight decay regularizer
MAX_ITER = 2000
def tuneLambda(lambdaList, eta=0.1, batchSize=50):
    # Returns the i-th batch of training data and targets
    # Generates a new, reshuffled batch once all previous
    #  ↪ batches are fed
    def getNextTrainingBatch(currentIter):
        currentBatchNum = currentIter % (trainData.shape[0] /
            ↪ batchSize)
        if currentBatchNum == 0:
            np.random.shuffle(randIdx)
        lowerBoundIdx = currentBatchNum * batchSize
        upperBoundIdx = (currentBatchNum + 1) * batchSize
        return trainData[lowerBoundIdx:upperBoundIdx],
            ↪ trainTarget[lowerBoundIdx:upperBoundIdx]

    # Generate updated plots for training and validation MSE
    def plotMSEGraph(MSEList, param):
        label = '$\lambda$ = ' + str(param)
        label_classification = ['train.', 'valid.']

        display.clear_output(wait=True)
        plt.figure(figsize=(8,5), dpi=200)

        for i, MSE in enumerate(MSEList):
            plt.plot(range(len(MSE)), MSE, '-', label=label+'
                ↪ '+label_classification[i])
```

73

```python
25
26          plt.axis([0, MAX_ITER, 0, np.amax(MSEList)])
27          plt.legend()
28          plt.show()
29
30      # Calculates the ratio between the n-th average epoch MSE
         ↪  and the (n-1)-th average epoch MSE
31      def ratioAverageEpochMSE(currentValidMSE):
32          averageN =
             ↪  np.average(currentValidMSE[-(np.arange(epochSize -
             ↪  1) + 1)])
33          averageNlessOne =
             ↪  np.average(currentValidMSE[-(np.arange(epochSize -
             ↪  1) + epochSize)])
34          return averageN / averageNlessOne
35
36      # Returns True if the average epoch validation MSE is at
         ↪  least 99% of the previous epoch average.
37      # i.e. Returns True if the average learnings between epoch
         ↪  is less than +1%
38      # Otherwise, returns False
39      def shouldStopEarly(currentValidMSE):
40          if currentValidMSE.shape[0] < 2 * epochSize:
41              return False
42          return True if (ratioAverageEpochMSE(currentValidMSE)
             ↪  >= 0.99) else False
43
44      summaryList = []
45      randIdx = np.arange(trainData.shape[0])
46      epochSize = trainData.shape[0] / batchSize
```

74

```python
47
48      for lambda_ in lambdaList:
49          W, b, X, Y, Y_hat, MSE, accuracy, optimizer =
            ↪ buildGraph(eta, lambda_)
50
51          with tf.Session() as sess:
52              tf.global_variables_initializer().run()
53
54              # Creates blank training and validation MSE arrays
                ↪ for the Session
55              currentTrainMSE = np.array([])[:, np.newaxis]
56              currentValidMSE = np.array([])[:, np.newaxis]
57
58              # Runs update
59              currentIter = 0
60              while currentIter <= MAX_ITER:
61                  inputData, inputTarget =
                    ↪ getNextTrainingBatch(currentIter)
62
63                  _, trainError = sess.run([optimizer, MSE],
                    ↪ feed_dict={X: inputData, Y: inputTarget})
64                  validError = sess.run([MSE], feed_dict={X:
                    ↪ validData, Y: validTarget})
65
66                  currentTrainMSE = np.append(currentTrainMSE,
                    ↪ trainError)
67                  currentValidMSE = np.append(currentValidMSE,
                    ↪ validError)
68
69                  # Update graph of training and validation MSE
                    ↪ arrays
```

```python
            if (currentIter < 3) or (currentIter % 500 ==
            ↪ 0):
                plotMSEGraph([currentTrainMSE,
                ↪ currentValidMSE], lambda_)


            # At every epoch, check for early stopping
            ↪ possibilty. If so, breaks from while loop
            if currentIter % epochSize == 0:
                if shouldStopEarly(currentValidMSE):
                    break

            currentIter += 1

        # Compute validation and test accuracy
        validAccuracy = sess.run(accuracy, feed_dict={X:
        ↪ validData, Y: validTarget})
        testAccuracy = sess.run(accuracy, feed_dict={X:
        ↪ testData, Y: testTarget})

    # Save session results as dictionary and appends to
    ↪ MSEsummaryList
    summaryList.append(
        {
            'eta': eta,
            'B': batchSize,
            'lambda': lambda_,
            'numIter': currentIter + 1,
            'epoch': float(currentIter + 1) /
            ↪ (trainData.shape[0] / batchSize),
            'trainMSE': currentTrainMSE,
```

```python
                    'validMSE': currentValidMSE,
                    'validAccuracy': validAccuracy,
                    'testAccuracy': testAccuracy
                }
            )

    return summaryList

# Main Function
lambdaList = [0.0, 0.0001, 0.001, 0.01, 0.1, 1.0]
tunedLambdaSummary = tuneLambda(lambdaList)


# Output summary table
for summary in tunedLambdaSummary:
    print 'lambda: %5.4f, numIter: %5d, validMSE: %5.3f,
      ↪ validAcc: %3.3f, testAcc: %3.3f' % \
        (summary['lambda'], summary['numIter'],
          ↪ summary['validMSE'][-1], summary['validAccuracy'],
          ↪ summary['testAccuracy'])


# Produce interactive Plotly graph
def lambdaIGraph(tunedLambdaSummary):
    # Create plot for each summary
    trace1 = go.Scatter(
        x = [np.log10(summary['lambda'] + 1e-5) for summary in
          ↪ tunedLambdaSummary],
        y = [summary['validAccuracy'] for summary in
          ↪ tunedLambdaSummary],
        name = 'Validation set accuracy'
    )
```

```python
trace2 = go.Scatter(
    x = [np.log10(summary['lambda'] + 1e-5) for summary in
     ↪ tunedLambdaSummary],
    y = [summary['testAccuracy'] for summary in
     ↪ tunedLambdaSummary],
    name = 'Test set accuracy'
)


data = go.Data([trace1, trace2])


# Create figure layout
layout = go.Layout(
    title = '$\\textit{Validation and Test set accuracy
     ↪ vs. } \\lambda$',
    xaxis = {'title': '$\\log_{10}(\\lambda)$'},
    yaxis = {'title': 'Model Accuracy'},
    annotations = [
        dict(
            text = '$\\textit{Used to represent }
             ↪ \\log_{10}(\\lambda=0)$',
            x = -5,
            y = 0.90,
        )
    ]
)


figure = go.Figure(data=data, layout=layout)
return py.iplot(figure,
 ↪ filename='A1Q2.3_accuracyVsLambda')
```

```
143  fig2_3 = lambdaIGraph(tunedLambdaSummary)
144  fig2_3
```