

**Instruction**

- We suggest you type your answer using L<sup>A</sup>T<sub>E</sub>X or Microsoft Word.
- This homework contains both theoretical and coding parts. For questions with [coding] mark, you need to write a small program to answer the question. We recommend you use Python, R or MATLAB to do the problem. Append your code after the answer.
- Submit a single PDF of your write-up on bCourses. Make sure your answer is legible.

**1. (MLE of multivariate normal distribution)**

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  be  $n$  independent samples drawn from the following multivariate normal distributions:

- (a)  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with mean  $\boldsymbol{\mu}$  and diagonal covariance matrix  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_{d \times d}$
- (b)  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with mean  $\boldsymbol{\mu}$  and diagonal covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_d^2 \end{bmatrix}$$

Compute the maximum likelihood estimator for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  and compute the Fisher information matrix for part (a).

**2. (Weighted least squares)**

Consider a dataset  $(\mathbf{x}_i, t_i)_{i=1, \dots, n}$  which each data is associated with a weighting factor  $\beta_i$ , so the sum-of-squares error function is

$$f(\mathbf{w}) = \sum_{i=1}^n \beta_i (t_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

Write  $f(\mathbf{w})$  in matrix form using

$\mathbf{X}$ :  $n \times d$  design matrix

$\mathbf{t}$ :  $n \times 1$  target vector

$\mathbf{B}$ :  $n \times n$  diagonal matrix with  $\beta_i$  on the diagonal

and find an expression for the solution  $\mathbf{w}^*$ .

### 3. (Linear regression on MNIST dataset)

[Coding] In this problem we will apply the least squares linear regression model on the hand-written digit recognition MNIST dataset.

The dataset can be downloaded from bCourses. It contains four files: 60000 training images, 60000 training labels, 10000 testing images and 10000 testing labels. Images are originally 28 pixel by 28 pixel grey scale image of a hand-written digit from 0 to 9. The actual data  $\mathbf{x}_i$  is flattened into a  $28 \times 28 = 784$  dimensional vector. Label is a single number shows the true digit of the corresponding image represents.

To train a least squares linear regression model, we solve

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \sum_{i=1}^n \|\mathbf{W}^\top \mathbf{x}_i - \mathbf{y}_i\|_2^2 \quad (1)$$

where

- $n$ : Number of samples (= 60000)
- $k$ : Number of classes (= 10)
- $\mathbf{x}_i$ :  $784 \times 1$  vector denotes a single data point
- $\mathbf{y}_i$ :  $10 \times 1$  vector denotes the one-hot encoding of the label, i.e. map  $j = \{0, 1, \dots, 9\}$  to the standard basis vector  $\mathbf{e}_j$ .
- $\mathbf{W}$ :  $784 \times 10$  matrix denotes the regression coefficients

By convention the above optimization problem usually is written as

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 \quad (2)$$

where

- $\mathbf{X}$ :  $60000 \times 784$  matrix where each row represents a single data  $\mathbf{x}_i^\top$
- $\mathbf{Y}$ :  $60000 \times 10$  matrix where each row represents a single encoded label  $\mathbf{y}_i^\top$
- $\|\cdot\|_F$ : denote the Frobenius norm [wiki]

Usually machine learning package takes  $\mathbf{X}$  and  $\mathbf{Y}$  as input.

In this problem, you need to use machine learning package to perform the following tasks:

- Read the data
- One-hot encode the labels to  $\mathbf{y}_i$
- Train the model using least squares regression
- Make predictions on both training and testing data
- Decode the predicted  $\mathbf{y}_i$  to labels

- Report the training and testing accuracy
- Find the first mis-classified data in the testing set and visualize the data as an image. What is the predicted digit using linear regression? What is your own prediction before seeing the true label? Are you doing better than the linear regression model?

Skeleton code using Python with `scikit-learn` package is provided on the next page.

```

1 from mnist import MNIST
  import numpy as np
3 import matplotlib.pyplot as plt
  import sklearn.metrics as metrics
5 from sklearn import linear_model
  from sklearn.preprocessing import OneHotEncoder
7
9 def load_dataset():
    mndata = MNIST('./data/')
11    X_train, labels_train = map(np.array, mndata.load_training())
    X_test, labels_test = map(np.array, mndata.load_testing())
13    X_train = X_train/255.0
    X_test = X_test/255.0
15
    X_train = np.asarray(X_train).reshape(60000, 784)
17    labels_train = np.asarray(labels_train).reshape(60000, 1)
    X_test = np.asarray(X_test).reshape(10000, 784)
19    labels_test = np.asarray(labels_test).reshape(10000, 1)
21
    return X_train, labels_train, X_test, labels_test
23 def plot(x):
    image = x.reshape(28, 28)
25    plt.imshow(image, cmap='grey')
    plt.axis('off')
27
    plt.show()
29
30 if __name__ == "__main__":
31
    ##### Read data #####
33    X_train, labels_train, X_test, labels_test = load_dataset()
35
    ##### One hot encode #####
    # (Use sklearn.preprocessing.OneHotEncoder)
37
    ##### Train #####
39    # (Use sklearn.linear_model.LinearRegression)
41
    ##### Predict #####
    # (Use sklearn.linear_model.LinearRegression)
43
    ##### Decode #####
45    # (Use numpy.argmax)
47
    ##### Print accuracy #####
    print("Train accuracy: {}".format(metrics.accuracy_score(labels_train, pred_labels_train)))
49    print("Test accuracy: {}".format(metrics.accuracy_score(labels_test, pred_labels_test)))
51
    ##### Plot first mis-classified data #####
    # (Use the provided plot(x) function)

```

code/IEOR290\_HW3\_P3\_skeleton.py