



THE CHINESE UNIVERSITY OF HONG KONG, SHENZHEN

DDA3020

MACHINE LEARNING

---

## Assignment1 Report

---

*Author:*  
Shi Wenlan

*Student Number:*  
119010265

October 2, 2022

September 29, 2022

Homework due: **11:59 pm, Oct 17, 2022.** The weight of this assignment among all assignments is **13/60.**

## 1 Written Problems (50 points)

1.1 Given the following denominator layout derivatives, (15 points) 1x d

- Differentiation of a scalar function w.r.t. a vector:** If  $f(\mathbf{w})$  is a scalar function of  $d$  variables,  $\mathbf{w}$  is a  $d \times 1$  vector, then differentiation of  $f(\mathbf{w})$  w.r.t.  $\mathbf{w}$  results in a  $d \times 1$  vector

$$\frac{df(\mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix} \quad h \times d$$

- Differentiation of a vector function w.r.t. a vector:** If  $\mathbf{f}(\mathbf{w})$  is a vector function of size  $h \times 1$  and  $\mathbf{w}$  is a  $d \times 1$  vector, then differentiation of  $\mathbf{f}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  results in a  $d \times h$  vector

$$\frac{d\mathbf{f}(\mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \cdots & \frac{\partial f_h}{\partial w_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial w_d} & \cdots & \frac{\partial f_h}{\partial w_d} \end{bmatrix}$$

Please prove the following derivatives, and  $\mathbf{X}$  and  $\mathbf{y}$  are not functions of  $\mathbf{w}$ :

$$\begin{aligned} \frac{d(\mathbf{X}^\top \mathbf{w})}{d\mathbf{w}} &= \mathbf{X}, \\ \frac{d(\mathbf{y}^\top \mathbf{X}\mathbf{w})}{d\mathbf{w}} &= \mathbf{X}^\top \mathbf{y} \\ \frac{d(\mathbf{w}^\top \mathbf{X}\mathbf{w})}{d\mathbf{w}} &= (\mathbf{X} + \mathbf{X}^\top)\mathbf{w} \end{aligned}$$

Solution:

Assume  $\mathbf{X}^\top = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,d} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ X_{h,1} & X_{h,2} & \cdots & X_{h,d} \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_h \end{bmatrix} = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_d]$   $\mathbf{y}^\top = [y_1, y_2, \dots, y_d]$   
 $\mathbf{w}^\top = [w_1, w_2, \dots, w_d]$

$$\textcircled{1} \frac{d(\mathbf{X}^\top \mathbf{w})}{d(\mathbf{w})} = \begin{bmatrix} \frac{\partial(X_1 \mathbf{w})}{\partial w_1} & \frac{\partial(X_1 \mathbf{w})}{\partial w_2} & \cdots & \frac{\partial(X_1 \mathbf{w})}{\partial w_d} \\ \frac{\partial(X_2 \mathbf{w})}{\partial w_1} & \frac{\partial(X_2 \mathbf{w})}{\partial w_2} & \cdots & \frac{\partial(X_2 \mathbf{w})}{\partial w_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial(X_h \mathbf{w})}{\partial w_1} & \frac{\partial(X_h \mathbf{w})}{\partial w_2} & \cdots & \frac{\partial(X_h \mathbf{w})}{\partial w_d} \end{bmatrix} = \begin{bmatrix} X_{1,1} & X_{2,1} & \cdots & X_{h,1} \\ X_{1,2} & X_{2,2} & \cdots & X_{h,2} \\ \vdots & \vdots & \ddots & \vdots \\ X_{1,d} & X_{2,d} & \cdots & X_{h,d} \end{bmatrix} = \mathbf{X}$$

② According to ①

$$\frac{d(\mathbf{y}^\top \mathbf{X}\mathbf{w})}{d(\mathbf{w})} = \frac{d((\mathbf{y}^\top \mathbf{X})\mathbf{w})}{d(\mathbf{w})} = (\mathbf{y}^\top \mathbf{X})^\top = \mathbf{X}^\top \mathbf{y}$$

③ In this case  $\mathbf{X}$  must be a square matrix, where  $h=d$ .

$$\mathbf{w}^\top \mathbf{X} \mathbf{w} = w_1 X_{1,1} + w_2 X_{1,2} + \cdots + w_d X_{1,d}$$

$$1 \mathbf{w}^\top \mathbf{X} \mathbf{w} = w_1 X_{1,1} + w_2 X_{1,2} + \cdots + w_d X_{1,d}$$

$$\begin{aligned} \frac{d(\mathbf{w}^\top \mathbf{X} \mathbf{w})}{d \mathbf{w}} &= \begin{bmatrix} \frac{\partial(w_1^\top \mathbf{X} \mathbf{w})}{\partial w_1} \\ \frac{\partial(w_2^\top \mathbf{X} \mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial(w_d^\top \mathbf{X} \mathbf{w})}{\partial w_d} \end{bmatrix} = \begin{bmatrix} 2w_1 X_{1,1} + w_2 X_{1,2} + \cdots + w_d X_{1,d} + X_{1,1} w_2 + \cdots + X_{d,1} w_d \\ X_{1,2} w_1 + w_1 X_{2,1} + 2w_2 X_{2,2} + w_3 X_{2,3} + \cdots + w_d X_{2,d} + X_{3,1} w_3 + \cdots + X_{d,2} w_d \\ \vdots \\ X_{1,d} w_1 + \cdots + X_{d,d} w_d + w_1 X_{d,1} + \cdots + w_{d-1} X_{d,d-1} + 2w_d X_{d,d} \end{bmatrix} \\ &= \begin{bmatrix} X_1 \mathbf{w} + \bar{X}_1 \mathbf{w} \\ X_2 \mathbf{w} + \bar{X}_2 \mathbf{w} \\ \vdots \\ X_d \mathbf{w} + \bar{X}_d \mathbf{w} \end{bmatrix} = \mathbf{X}^\top \mathbf{w} + \bar{\mathbf{X}} \mathbf{w} = (\mathbf{X} + \bar{\mathbf{X}}) \mathbf{w} \end{aligned}$$

1.2 Suppose we have training data  $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $\mathbf{y}_i \in \mathbb{R}^k$ ,  $i = 1, 2, \dots, N$ . (1) Find the closed-form solution of the following problem

$$\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N \alpha_i \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2,$$

where  $\alpha_1, \dots, \alpha_n$  are weights for different samples; (2) Show how to use gradient descent to solve the problem. (12 points)

**Solution:**

$$(1) \text{ Assume } X = [\sqrt{\alpha_1} \mathbf{x}_1^\top \quad \sqrt{\alpha_2} \mathbf{x}_2^\top \quad \dots \quad \sqrt{\alpha_N} \mathbf{x}_N^\top]^\top \quad A = [W] \quad Y = [\sqrt{\alpha_1} \mathbf{y}_1^\top \quad \sqrt{\alpha_2} \mathbf{y}_2^\top \quad \dots \quad \sqrt{\alpha_N} \mathbf{y}_N^\top]^\top \quad \text{and } (X^\top X)^{-1} \text{ exists, and } d=k \text{ or } k=1$$

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N \alpha_i \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2 &= \min_A \text{trace}((XA - Y)^\top (XA - Y)) \\ &= \min_A \text{trace}((A^\top X^\top - Y^\top)(XA - Y)) \\ &= \min_A \text{trace}(A^\top X^\top X A - A^\top X^\top Y - Y^\top X A + Y^\top Y) \\ &= \min_A \text{trace}(A^\top X^\top X A - 2Y^\top X A + Y^\top Y) \end{aligned}$$

$$\frac{\partial (A^\top X^\top X A - 2Y^\top X A + Y^\top Y)}{\partial A} = \frac{\partial (A^\top X^\top X A)}{\partial A} - 2 \frac{\partial (Y^\top X A)}{\partial A} + \frac{\partial (Y^\top Y)}{\partial A}$$

$$= (X^\top X + X^\top X)A - 2X^\top Y$$

$$= 2X^\top X A - 2X^\top Y = 0$$

$$\because X^\top X A = X^\top Y \quad \therefore \begin{bmatrix} W \\ b \end{bmatrix} = A = (X^\top X)^{-1} X^\top Y = \left[ \begin{array}{cc} \sum_{i=1}^N \alpha_i \mathbf{x}_i \mathbf{x}_i^\top & \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{x}_i^\top \\ \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{x}_i^\top & \sum_{i=1}^N \alpha_i \end{array} \right]^{-1} \left[ \begin{array}{c} \sum_{i=1}^N \alpha_i \mathbf{x}_i \mathbf{y}_i^\top \\ \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{y}_i^\top \end{array} \right]$$

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N \alpha_i \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2 &= \text{trace}((X^\top X)^{-1} X^\top X (X^\top X)^{-1} X^\top Y - 2Y^\top X (X^\top X)^{-1} X^\top Y + Y^\top Y) \\ &= \text{trace}(Y^\top X (X^\top X)^{-1} X^\top Y - 2Y^\top X (X^\top X)^{-1} X^\top Y + Y^\top Y) \\ &= \text{trace}(Y^\top Y - Y^\top X (X^\top X)^{-1} X^\top Y) \\ &= \text{trace} \left( \sum_{i=1}^N \alpha_i \mathbf{y}_i \mathbf{y}_i^\top - \left[ \sum_{i=1}^N \alpha_i \mathbf{y}_i \mathbf{x}_i^\top \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{y}_i \right] \left[ \begin{array}{cc} \sum_{i=1}^N \alpha_i \mathbf{x}_i \mathbf{x}_i^\top & \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{x}_i^\top \\ \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{x}_i^\top & \sum_{i=1}^N \alpha_i \end{array} \right]^{-1} \left[ \begin{array}{c} \sum_{i=1}^N \alpha_i \mathbf{x}_i \mathbf{y}_i^\top \\ \sum_{i=1}^N \sqrt{\alpha_i} \mathbf{y}_i^\top \end{array} \right] \right) \end{aligned}$$

(2) given a start point  $\mathbf{x}$  in the domain of  $A$   
repeat

$$1. \Delta \mathbf{x} := -X^\top (XA - Y)$$

$$2. \text{Choose the step size } t = \arg \min_{t>0} \mathbf{x}^\top (\mathbf{x} + t\Delta \mathbf{x})$$

$$3. \text{Update } \mathbf{x} = \mathbf{x} + t\Delta \mathbf{x}$$

until stopping criterion (e.g.  $\|X^\top (XA - Y)\|_2 \leq \epsilon$ ) is satisfied.

$$\begin{bmatrix} W \\ b \end{bmatrix} = \mathbf{x}$$

1.3 Prove that: (1)  $f(x) = x^4$  is convex; (2)  $f(x) = |x|$  is convex; (3)  $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|^2$  is convex, where  $\mathbf{A}$  is a matrix. (12 points)

*Solution:*

(1)  $\because$  Domain of  $f(x)$  is  $\mathbb{R}$ , which is convex

$$\therefore f''(x) = 12x^2 \geq 0$$

$\therefore$  According to SDC,  $f(x)$  is convex

(2)  $\because$  Domain of  $f(x)$  is  $\mathbb{R}$ , which is convex

$$\because |\theta x + (1-\theta)y| \leq \theta|x| + (1-\theta)|y|$$

$\therefore$  According to definition,  $f(x)$  is convex

(3)  $\because$  Domain of  $f(\mathbf{x})$  is  $\mathbb{R}^n$ , which is convex

$$f(\mathbf{x}) = (\mathbf{Ax} - \mathbf{b})^T(\mathbf{Ax} - \mathbf{b})$$

$$= (\mathbf{x}^T \mathbf{A}^T - \mathbf{b}^T)(\mathbf{A}\mathbf{x} - \mathbf{b})$$

$$= \mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} - 2\mathbf{b}^T \mathbf{A}\mathbf{x} + \mathbf{b}^T \mathbf{b}$$

$$\nabla f(\mathbf{x}) = \frac{\partial(\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} - 2\mathbf{b}^T \mathbf{A}\mathbf{x} + \mathbf{b}^T \mathbf{b})}{\partial \mathbf{x}} = 2\mathbf{A}^T \mathbf{A}\mathbf{x} - 2\mathbf{A}^T \mathbf{b}$$

$$\therefore \nabla^2 f(\mathbf{x}) = \frac{\partial(2\mathbf{A}^T \mathbf{A}\mathbf{x} - 2\mathbf{A}^T \mathbf{b})}{\partial \mathbf{x}} = 2\mathbf{A}^T \mathbf{A} \succeq 0$$

$\therefore$  According to SDC,  $f(\mathbf{x})$  is convex.

1.4 Suppose  $x_1, x_2, \dots, x_N$  are drawn from  $\mathcal{N}(\mu, \sigma^2)$ . Show that the MLE (maximum likelihood estimation) of  $\sigma^2$  is  $\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$ . (11 points)

*Solution:*

$$\text{Assume } X = \{x_i\}_{i=1}^N$$

$$\sigma_{ML}^2 = \arg \max_{\sigma^2} \log L(\sigma^2, \mu | X)$$

$$\begin{aligned} f(\sigma^2) &= \log L(\sigma^2 | X) = \log \left( \prod_{i=1}^N p(x_i | \mu, \sigma^2) \right) \\ &= \sum_{i=1}^N \log N(\mu, \sigma^2) \\ &= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(x_i - \mu)^2}{2\sigma^2} \right) \\ &= -N \log \sqrt{2\pi\sigma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 \end{aligned}$$

$$\begin{aligned} f'(\sigma^2) &= -N \frac{1}{2\sigma^2 \sigma^2} \cdot \sqrt{2\pi} \cdot \frac{1}{2\sigma^2} + \sum_{i=1}^N (x_i - \mu)^2 \frac{1}{2\sigma^4} \\ &= -\frac{N}{2\sigma^2} + \sum_{i=1}^N (x_i - \mu)^2 \frac{1}{2\sigma^4} = 0 \end{aligned}$$

$$\therefore \sum_{i=1}^N (x_i - \mu)^2 \frac{1}{\sigma^2} = N$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$\therefore \sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2$$

**2.1.** The CSV or XLS file contains a dataset for regression. There are 7750 samples with 25 features (described in the doc file). This data is for the purpose of bias correction of next-day maximum and minimum air temperatures forecast of the LDAPS model operated by the Korea Meteorological Administration over Seoul, South Korea. This data consists of summer data from 2013 to 2017. The input data is largely composed of the LDAPS model's next-day forecast data, in-situ maximum and minimum temperatures of present-day, and geographic auxiliary variables. There are two outputs (i.e. next-day maximum and minimum air temperatures) in this data. Hindcast validation was conducted for the period from 2015 to 2017.

You need to delete the first two attributes (station and date), and use attributes 3-23 to predict attributes 24 and 25. Randomly split the data into two parts, one contains 80% of the samples and the other contains 20% of the samples. Use the first part as training data and train a linear regression model and make prediction on the second part. Report the training error and testing error in terms of RMSE.

Repeat the splitting, training, and testing for 10 times. Use a loop and print the RMSEs in each trial.

**Solution:**

Read in data : `np.loadtxt(file name)`

Clean data : `np.delete( rows with NaN )`

Shuffle data : `np.random.shuffle(data)`

Split data : `train data : test data = 4 : 1`

$X = [data[2:23]]$   $w = [w, b]$  because the whole training

$Y = [data[23:]]$  ↓ process is too slow

Train: Initial value of  $w$ : get from a previous training process

(initial: `np.ones()`, after 295000 iterations)

Objective function:  $RMSE = \sqrt{\frac{1}{n} \text{trace}((Xw - Y)^T(Xw - Y))}$

Gradient :  $X^T(Xw - Y)$  and take negative

Stepsize : get from Golden Section Method iteration range from (0,2)

Stop condition:  $\text{norm(Gradient)} < 0.001$  ←

or because the whole training process is too slow

decrease of objective value  $< 0.001$

Test: Calculate RMSE on test data set

**Result snapshot:**

```
Trail 0
train RMSE = 6.574515200326087
test RMSE = 7.289944883147615
Trail 1
train RMSE = 6.540215519510085
test RMSE = 7.412679830282912
Trail 2
train RMSE = 6.752142677709038
test RMSE = 6.608645543493965
Trail 3
train RMSE = 6.64360268964265
test RMSE = 7.02734407988645
Trail 4
train RMSE = 6.6777866935260155
test RMSE = 6.9043236270052954
Trail 5
train RMSE = 6.699952345247394
test RMSE = 6.451579486781242
Trail 6
train RMSE = 6.611091036144075
test RMSE = 6.899438956434362
Trail 7
train RMSE = 6.499482804076262
test RMSE = 7.537555680283795
Trail 8
train RMSE = 6.770978396204615
test RMSE = 6.531154433978679
Trail 9
train RMSE = 6.779375874598276
test RMSE = 6.496120347208102
```

**2.2.** The classification data file contains the iris dataset. This is perhaps the best known database to be found in the pattern recognition literature. Fisher's paper is a classic in the field and is referenced frequently to this day. (See Duda Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other. There are five attributes: 1. sepal length in cm; 2. sepal width in cm; 3. petal length in cm; 4. petal width in cm; 5. class: – Iris Setosa – Iris Versicolour – Iris Virginica. You need to use the first attributes to predict the last attribute, namely, classifying the data into three classes. Randomly split the data into two parts, one contains 80% of the samples and the other contains 20% of the samples. Use the first part as training data and train a linear model and make classification on the second part. Report the training error and testing error in terms of classification error rate (number of miss-classified samples divided by number of all samples)

Repeat the splitting, training, and testing for 10 times. Use a loop and print the classification errors in each trial.

**Solution:**

Read in data : `xlrd.open_workbook(file name).sheets()[0]`

Clean data : replace string tag to array like [1 0 0]

Shuffle data : `np.random.shuffle(data)`

Split data : train data : test data = 4 : 1

$X = [data[4::1]]$   $w = [w, b]$  because the whole training

$Y = [data[4::1]]$  ↓ process is too slow

Train : Initial value of  $w$ : get from a previous training process

(initial : `np.ones()`, after 5700 iterations)

Objective function :  $RMSE = \sqrt{\frac{1}{n} \text{trace}((Xw - Y)^T(Xw - Y))}$

Gradient :  $X^T(Xw - Y)$  and take negative

Stepsize : get from Golden Section Method iteration range from (0,2)

Stop condition :  $\text{norm(Gradient)} < 0.000001$

or

decrease of objective value < 0.000001

Test : Classify the sample to the class with the highest value : e.g.  $Xw = [0.1 0.7 0.2] \rightarrow [0 1 0]$

Calculate the error rate on both data sets

**Result snapshot:**

```

Trail 0
train classification error rate = 0.125
test classification error rate = 0.2333333333333328
Trail 1
train classification error rate = 0.1583333333333333
test classification error rate = 0.0999999999999998
Trail 2
train classification error rate = 0.1333333333333333
test classification error rate = 0.1999999999999996
Trail 3
train classification error rate = 0.1500000000000002
test classification error rate = 0.1666666666666663
Trail 4
train classification error rate = 0.1750000000000004
test classification error rate = 0.0999999999999998
Trail 5
train classification error rate = 0.14166666666666672
test classification error rate = 0.1999999999999996
Trail 6
train classification error rate = 0.14166666666666672
test classification error rate = 0.1999999999999996
Trail 7
train classification error rate = 0.1583333333333333
test classification error rate = 0.0666666666666665
Trail 8
train classification error rate = 0.14166666666666672
test classification error rate = 0.0999999999999998
Trail 9
train classification error rate = 0.14166666666666672
test classification error rate = 0.1999999999999996

```