

Homework 5: Gradient Calculations and Nonlinear Optimization

Introduction to Machine Learning

Shangwen Yan | N17091204 | sy2160

1.

(a)

$$\mathbf{A} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{z} = \mathbf{A}\mathbf{w}$$

$$J(\mathbf{w}) = g(\mathbf{z}) = \sum_{i=1}^n g_i(z_i), \quad g_i(z_i) = (y_i - \frac{1}{z_i})^2, \quad z_i = w_0 + \sum_{j=1}^d x_{ij}w_j$$

(b)

$$\nabla J(\mathbf{w}) = \mathbf{A}^T \nabla_{\mathbf{z}} g(\mathbf{z}), \quad \nabla_{\mathbf{z}} g(\mathbf{z}) = \begin{bmatrix} g'_1(z_1) \\ \vdots \\ g'_n(z_n) \end{bmatrix}, \quad g'_i(z_i) = \frac{2}{z_i^2} (y_i - \frac{1}{z_i})$$

(c) \mathbf{w}^k is the estimate sequence of \mathbf{w} :

$$\mathbf{w}^{k+1} = \mathbf{w}^k - a_k \nabla J(\mathbf{w}^k)$$

(d)

```
1 import numpy as np
2 def loss_function_grad(X,y,w):
3     A = np.column_stack((np.ones(n, ), X))
4     z = A.dot(w)
5     J = np.sum((y-1/z)**2)
6     dJ_dz = (2/z**2)*(y-1/z)
7     J_grad = A.T.dot(dJ_dz)
8     return J,J_grad
```

3.

(a) $\nabla_P z_i = \mathbf{x}_i \mathbf{x}_i^T$

(b) assume that $g_i(z_i) = [\frac{z_i}{y_i} - \ln(z_i)]$, $J(\mathbf{P}) = g(\mathbf{z}) = \sum_{i=1}^n g_i(z_i)$, $g'_i(z_i) = [\frac{1}{y_i} - \frac{1}{z_i}]$

$$\mathbf{X}_{n \times m} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}, \quad \nabla_{\mathbf{z}} g(\mathbf{z}) = \begin{bmatrix} g'_1(z_1) \\ \vdots \\ g'_n(z_n) \end{bmatrix}, \quad \mathbf{H}_{m \times n} = \begin{bmatrix} (\nabla_{\mathbf{z}} g(\mathbf{z}))^T \\ \vdots \\ (\nabla_{\mathbf{z}} g(\mathbf{z}))^T \end{bmatrix}$$

$$\nabla_P J(\mathbf{P})_{m \times m} = (\mathbf{X}^T \cdot \mathbf{H}) \mathbf{X}$$

A.B produce matrix C, $C_{ij} = A_{ij} B_{ij}$

(c)

```
1 def loss_function_grad(X,y,P):
2     z = []
3     for i in range(X.shape[0]):
4         zi = X[i,:].dot(P).dot(X[i,:].T)
5         z.append(zi)
6     z = np.array(z)
7     J = np.sum((z/y)-np.log(z))
8     dJ_dz = (1/y)-(1/z)
9     J_grad = (X.T*dJ_dz[None,:]).dot(X) #broadcast
10    return J,J_grad
```

(d)

```
1 def loss_function_grad(X,y,P):
2     z = X.dot(P).dot(X.T).dot(np.ones(X.shape[0]))
3     z = np.squeeze(z)
4     J = np.sum((z/y)-np.log(z))
5     dJ_dz = (1/y)-(1/z)
6     J_grad = (X.T*dJ_dz[None,:]).dot(X)
7     return J,J_grad
```

4.

(a)

$$\nabla_{\mathbf{b}} J(\mathbf{b}) = 2 \begin{bmatrix} \sum_{i=1}^n [(-e^{-a_1 x_i})(y_i - \sum_{j=1}^d b_j e^{-a_j x_i})] \\ \vdots \\ \sum_{i=1}^n [(-e^{-a_d x_i})(y_i - \sum_{j=1}^d b_j e^{-a_j x_i})] \end{bmatrix} \Big|_{\mathbf{a}=\hat{\mathbf{a}}}$$

and we use gradient descent method to update \mathbf{b} , thus \mathbf{b}^k is the estimate sequence of \mathbf{b} :

$$\mathbf{b}^{k+1} = \mathbf{b}^k - a_k \nabla J(\mathbf{b}^k)$$

after certain times of iteration or stop criterion, we can get $\hat{\mathbf{b}}$

(b)

$$J_1(\mathbf{a}) = \min_{\mathbf{b}=\hat{\mathbf{b}}} J(\mathbf{a}, \mathbf{b})$$

here $\hat{\mathbf{b}}$ is constant, and we can use the same method as (a) to compute $\nabla_{\mathbf{a}} J(\mathbf{a}, \mathbf{b})$

$$\nabla J(\mathbf{a}, \mathbf{b}) = \nabla_{\mathbf{a}} J_1(\mathbf{a}, \mathbf{b})|_{\mathbf{b}=\hat{\mathbf{b}}} = 2 \left[\begin{array}{c} \sum_{i=1}^n (b_1 x_i e^{-a_1 x_i}) (y_i - \sum_{j=1}^d b_j e^{-a_j x_i}) \\ \vdots \\ \sum_{i=1}^n b_d e^{-a_d x_i} (y_i - \sum_{j=1}^d b_j e^{-a_j x_i}) \end{array} \right]_{\mathbf{b}=\hat{\mathbf{b}}}$$