

# Deep Learning Assignment 2

Saeedreza Zouashkiani

November 21, 2022

1

## 1.1

Consider the Taylor expansion of the loss function at  $\omega$ :

$$J(\omega + \delta\omega) - J(\omega) \approx \nabla_{\omega} L \delta\omega + \frac{1}{2} \delta\omega^T H \delta\omega \quad (1)$$

The first term goes to zero when the algorithm converges. So, our problem is to minimize:

$$C = \frac{1}{2} \delta\omega^T H \delta\omega \quad (2)$$

Our goal is then to set one of the weights to zero which is denoted by  $\omega_i$  to minimize the increase in error given by equation (2). Then to eliminate  $\omega_i$ :

$$\omega_i + \mathbf{e}_i^T \cdot \delta\omega = 0 \quad (3)$$

Where  $\mathbf{e}_i$  is a vector with a 1 corresponding to  $\omega_i$  and zero elsewhere. To minimize equation (2) constrained on equation (3), we use Lagrangian multipliers:

$$L = \frac{1}{2} \delta\omega^T H \delta\omega + \lambda(\omega_i + \mathbf{e}_i^T \cdot \delta\omega) \quad (4)$$

Differentiating with respect to  $\delta\omega$  and setting to zero:

$$H \delta\omega + \lambda \mathbf{e}_i = 0 \quad (5)$$

Then if we solve for  $\lambda$ :

$$\lambda = \frac{\omega_i}{(H^{-1})_{ii}} \quad (6)$$

Where  $H^{-1}$ , corresponds to the  $i^{th}$  element in the inverse Hessian matrix. Then by plugging (6) into (5)

$$\delta\omega = -\frac{\omega_i}{(H^{-1})_{ii}} H^{-1} \cdot \mathbf{e}_i \quad (7)$$

Then plugging (7) into (2):

$$C = \frac{1}{2} \frac{\omega_i^2}{(H^{-1})_{ii}} \quad (8)$$

## 1.2

If  $\mathbf{H} = \mathbf{I}$ , then  $\mathbf{H}^{-1} = \mathbf{I}$ . Therefore:

$$\delta \boldsymbol{\omega} = -\boldsymbol{\omega}_i \quad (9)$$

$$\mathbf{C} = \frac{1}{2} \boldsymbol{\omega}_i^2 \quad (10)$$

In the previous part, we needed to compute the inverse Hessian first, and compute  $\frac{1}{2} \frac{\boldsymbol{\omega}_i^2}{(\mathbf{H}^{-1})_{ii}}$  for each element and remove the one with the least increase in the loss, but in here we compute  $\frac{1}{2} \boldsymbol{\omega}_i^2$  and set  $\delta \boldsymbol{\omega} = -\boldsymbol{\omega}_i$  for that  $i$  that minimizes the loss.

## 2

### 2.1

This is solved in part 3.1

### 2.2

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{N} \|(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I}) \boldsymbol{\epsilon}\|_2^2 \right] \\ &= \frac{1}{N} \mathbb{E} [\boldsymbol{\epsilon}^T (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I})^T (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{I}) \boldsymbol{\epsilon}] \\ &= \frac{1}{N} \mathbb{E} [\boldsymbol{\epsilon}^T (\mathbf{B})^T (\mathbf{B}) \boldsymbol{\epsilon}] \end{aligned} \quad (11)$$

Since  $\mathbf{A}$  is symmetric and  $\mathbf{A}^T \mathbf{A} = \mathbf{A}^2$ , then it is idempotent. That is its eigenvalues consist of  $d$  ones and  $n - d$  zeros. Therefore using 11:

$$\mathbf{B}^T \mathbf{B} = (\mathbf{A} - \mathbf{I})^T (\mathbf{A} - \mathbf{I}) = \mathbf{A}^2 - 2\mathbf{A} + \mathbf{I} = \mathbf{I} - \mathbf{A} \quad (12)$$

Since  $\boldsymbol{\epsilon}^T \mathbf{B}^T \mathbf{B} \boldsymbol{\epsilon}$  is a scalar, then its trace equals itself:

$$\begin{aligned} & \frac{1}{N} \mathbb{E} [\boldsymbol{\epsilon}^T \mathbf{B}^T \mathbf{B} \boldsymbol{\epsilon}] \\ &= \frac{1}{N} \mathbb{E} [\text{Trace}(\boldsymbol{\epsilon}^T \mathbf{B}^T \mathbf{B} \boldsymbol{\epsilon})] \\ &= \frac{1}{N} \mathbb{E} [\text{Trace}(\mathbf{B}^T \mathbf{B} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T)] \\ &= \frac{1}{N} \text{Trace}(\mathbf{B}^T \mathbf{B} \boldsymbol{\Sigma}) \\ &= \frac{\sigma^2}{N} \text{Trace}(\mathbf{I} - \mathbf{A}) \\ &= \frac{N - d}{N} \sigma^2 \end{aligned} \quad (13)$$

Where the last part comes from the property of idempotent matrices. That is  $\text{Trace}(\mathbf{A}) = \text{rank}(\mathbf{X})$ .

### 2.3

As  $d$  increases and gets closer to  $N$ , the trace of  $\mathbf{I} - \mathbf{A}$  approaches zero. That is the number of independent columns of  $\mathbf{X}$  increases, so we can have a better projection matrix than before. Naturally the training error will decrease as well.

### 3

#### 3.1

We want to find:

$$\begin{aligned}
& \underset{\beta}{\operatorname{argmin}} ||Y - X\beta||_2^2 \\
&= \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta) \\
&= \underset{\beta}{\operatorname{argmin}} Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta \\
&= \underset{\beta}{\operatorname{argmin}} -2\beta^T X^T Y + \beta^T X^T X \beta
\end{aligned} \tag{14}$$

Define  $L(\beta) = 2\beta^T X^T Y + \beta^T X^T X \beta$ , and differentiate with respect to  $\beta$  and set to zero:

$$\frac{\partial L(\beta)}{\partial \beta} = -2X^T Y + 2X^T X \beta \tag{15}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \tag{16}$$

#### 3.2

Add the L2 norm of the weights to the minimization problem of the previous part:

$$\begin{aligned}
& \underset{\beta}{\operatorname{argmin}} ||Y - X\beta||_2^2 + \alpha ||\beta||_2^2 \\
&= \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta) + \alpha \beta^T \beta \\
&= \underset{\beta}{\operatorname{argmin}} Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta + \alpha \beta^T \beta \\
&= \underset{\beta}{\operatorname{argmin}} -2\beta^T X^T Y + \beta^T X^T X \beta + \alpha \beta^T \beta
\end{aligned} \tag{17}$$

Define  $L(\beta) = 2\beta^T X^T Y + \beta^T X^T X \beta + \alpha \beta^T \beta$ , and differentiate with respect to  $\beta$  and set to zero:

$$\frac{\partial L(\beta)}{\partial \beta} = -2X^T Y + 2X^T X \beta + 2\alpha \beta \tag{18}$$

$$\hat{\beta} = (X^T X + \alpha I)^{-1} X^T Y \tag{19}$$

#### 3.3

We know that both  $\beta^*$  and  $\hat{\beta}$  have an expectation value of  $\beta$ . Let  $\beta^* = Ay$  and  $\hat{\beta} = By$ . Where  $A = (X^T \Sigma^{-1} X)^{-1} \Sigma^{-1} X^T$  and  $B = (X^T X)^{-1} X^T$ .

$$\begin{aligned}
\operatorname{Var}(\beta^*) &= B \operatorname{Var}(y) B^T = X^T \Sigma^{-1} X^{-1} \\
\operatorname{Var}(\hat{\beta}) &= A \operatorname{Var}(y) A^T = (X^T X)^{-1} X^T \Sigma X (X^T X)^{-1}
\end{aligned} \tag{20}$$

If  $XF = \Sigma X$ , then:

$$(X^T \Sigma^{-1} X)^{-1} = F(X^T X)^{-1} \tag{21}$$

It is obvious from 21 that  $F$  is non-singular. Using (21), we can write:

$$\begin{aligned}
\operatorname{Var}(\beta^*) &= F(X^T X)^{-1} \\
\operatorname{Var}(\hat{\beta}) &= (X^T X)^{-1} X^T \Sigma X (X^T X)^{-1} = (X^T X)^{-1} X^T X F (X^T X)^{-1} = F(X^T X)^{-1}
\end{aligned} \tag{22}$$

Therefore  $\boldsymbol{\beta}^*$  and  $\hat{\boldsymbol{\beta}}$  estimate the same quantity. On the other hand, to prove when the estimators are equal, we need to show that:

$$\begin{aligned}
\boldsymbol{\beta}^* &= \hat{\boldsymbol{\beta}} \\
(\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\
\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} &= \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\
\mathbf{X}^T (\mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} - \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) &= \mathbf{0} \\
\mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} &= \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\
\mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} &= \boldsymbol{\Sigma} \mathbf{X} \\
\mathbf{X} \mathbf{F} &= \boldsymbol{\Sigma} \mathbf{X}
\end{aligned} \tag{23}$$

### 3.4

Extend the error function without the L-1 norm:

$$\begin{aligned}
L(\boldsymbol{\beta}, \lambda_1) &= \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_2^2 = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda_1 \boldsymbol{\beta}^T \boldsymbol{\beta} \\
&= \mathbf{Y}^T \mathbf{Y} - 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{Y} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} + \lambda_1 \boldsymbol{\beta}^T \boldsymbol{\beta}
\end{aligned} \tag{24}$$

Define the augmented  $\hat{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda_1} \mathbf{I} \end{bmatrix}$ , and  $\hat{\mathbf{Y}} = \begin{bmatrix} \mathbf{Y} \\ \mathbf{0} \end{bmatrix}$ . Then equation (24) can be written as:

$$L(\boldsymbol{\beta}, \lambda_1) = \|\hat{\mathbf{Y}} - \hat{\mathbf{X}}\boldsymbol{\beta}\|_2^2 \tag{25}$$

Therefore the loss function  $L(\boldsymbol{\beta}, \lambda_1, \lambda_2)$  will become:

$$L(\boldsymbol{\beta}, \lambda_1, \lambda_2) = \|\hat{\mathbf{Y}} - \hat{\mathbf{X}}\boldsymbol{\beta}\|_2^2 \tag{26}$$

## 4

### 4.1

Let  $J_D = (y_d - O_D)^2$  and  $J_N = (y_d - O_N)^2$  denote the loss function of the Dropout and without Dropout network, respectively, and  $O_D = \sum_{k=1}^n \delta_k \omega_k x_k$

$$\begin{aligned}
\frac{\partial J_N}{\partial \omega_i} &= -2(y_d - O_N) \frac{\partial O_N}{\partial \omega_i} = -2(y_d - O_N) x_i = -2y_d x_i + 2 \sum_{k=1}^n \omega_k x_k x_i \\
&= -2y_d x_i + 2\omega_i x_i^2 + 2 \sum_{k=1, k \neq i}^n \omega_k x_k x_i
\end{aligned} \tag{27}$$

$$\begin{aligned}
\frac{\partial J_D}{\partial \omega_i} &= -2(y_d - O_D) \frac{\partial O_D}{\partial \omega_i} = -2(y_d - O_D) \delta_i x_i \\
&= -2y_d \delta_i x_i + 2\omega_i^2 \delta_i^2 x_i^2 + 2 \sum_{k=1, k \neq i}^n \omega_k \delta_k \delta_i x_k x_i
\end{aligned} \tag{28}$$

Since equation (12) is probabilistic, we take the expectation:

$$\begin{aligned}
\mathbb{E}[\frac{\partial J_D}{\partial \omega_i}] &= -2y_d x_i + 2\omega_i x_i^2 + 2\sigma^2 \omega_i^2 x_i^2 + 2 \sum_{k=1, k \neq i}^n \omega_k x_k x_i \\
&= \frac{\partial J_N}{\partial \omega_i} + 2\sigma^2 \omega_i^2 x_i^2
\end{aligned} \tag{29}$$

## 4.2

The overall regularized loss function can be written as:

$$J_R = (y_d - \sum_{k=1}^n \delta_k \omega_k x_k)^2 + \sigma^2 \sum_{k=1}^n w_k^2 x_k^2 \quad (30)$$

It means that the network is regularized by the square of multiplication of weights by inputs with a regularization factor of  $\sigma^2$ .

## 5

Let the cost function be:

$$L = \frac{1}{2} \boldsymbol{\omega}^T \mathbf{H} \boldsymbol{\omega} \quad (31)$$

Computing the first and second order derivatives:

$$\nabla_{\boldsymbol{\omega}} L = \frac{\partial L}{\partial \boldsymbol{\omega}} = \mathbf{H} \boldsymbol{\omega} \quad (32)$$

$$\nabla_{\boldsymbol{\omega}}^2 L = \frac{\partial^2 L}{\partial \boldsymbol{\omega}^2} = \mathbf{H} \quad (33)$$

### 5.1

$$\boldsymbol{\omega}^{(t+1)} = \boldsymbol{\omega}^{(t)} - \epsilon \mathbf{H} \boldsymbol{\omega}^{(t)} = \mathbf{Q}(\mathbf{I} - \epsilon \boldsymbol{\Lambda}) \mathbf{Q}^T \boldsymbol{\omega}^{(t)} \quad (34)$$

### 5.2

$$\boldsymbol{\omega}^{(t+1)} = \mathbf{Q}(\mathbf{I} - \epsilon \boldsymbol{\Lambda}) \mathbf{Q}^T \boldsymbol{\omega}^{(0)} \quad (35)$$

### 5.3

$|\mathbf{I} - \epsilon \boldsymbol{\Lambda}| < 1$  must be satisfied. Therefore:

$$\epsilon < \frac{2}{\lambda_{max}} \quad (36)$$

### 5.4

Because of 33, the newtons method will be:

$$\boldsymbol{\omega}^{(t+1)} = (1 - \epsilon) \boldsymbol{\omega}^{(t)} \quad (37)$$

Then by choosing  $\epsilon = 1$ , the newton method will converge in one step.

### 5.5

Because in Newton's method, calculating the inverse of the Hessian matrix is expensive, gradient descent is used instead.

## 6

### 6.1

It will learn a similarity measure between 2 inputs, since the inputs are both mapped to a latent space and the distance between the 2 latent vectors is minimized. This can be used for images and text, for example.

### 6.2

$$\begin{aligned}
 \mathbf{z}_1 &= \mathbf{W}\mathbf{x}_1 + \mathbf{b} \\
 \mathbf{z}_2 &= \mathbf{W}\mathbf{x}_2 + \mathbf{b} \\
 \Delta_1 &= \frac{\partial J}{\partial \mathbf{h}_1} = 2\mathbf{h}_1 - \mathbf{h}_2 \\
 \Delta_2 &= \frac{\partial J}{\partial \mathbf{h}_2} = 2\mathbf{h}_2 - \mathbf{h}_1 \\
 \Delta_3 &= \frac{\partial J}{\partial \mathbf{z}_1} = \Delta_1 \odot (1 - \mathbf{h}_1^2) \\
 \Delta_4 &= \frac{\partial J}{\partial \mathbf{z}_2} = \Delta_2 \odot (1 - \mathbf{h}_2^2) \\
 \frac{\partial J}{\partial \mathbf{W}} &= \Delta_3 \mathbf{x}_1^T + \Delta_4 \mathbf{x}_2^T + 2\mathbf{W} \\
 \frac{\partial J}{\partial \mathbf{b}} &= \Delta_3 + \Delta_4
 \end{aligned} \tag{38}$$

Then the gradient descent update rule for  $m$  batch size is:

$$\begin{aligned}
 \mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} - \epsilon \frac{1}{m} \sum_{i=1}^m \frac{\partial J_i}{\partial \mathbf{W}} \\
 \mathbf{b}^{(t+1)} &= \mathbf{b}^{(t)} - \epsilon \frac{1}{m} \sum_{i=1}^m \frac{\partial J_i}{\partial \mathbf{b}}
 \end{aligned} \tag{39}$$