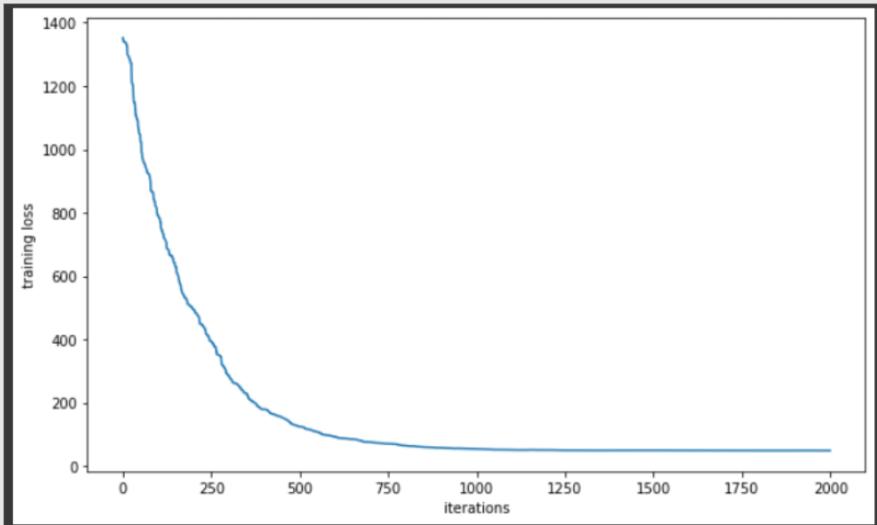


Part 1 Coding (60%)

Linear regression model :

1. Training loss over iterations



2. MSE of testing data

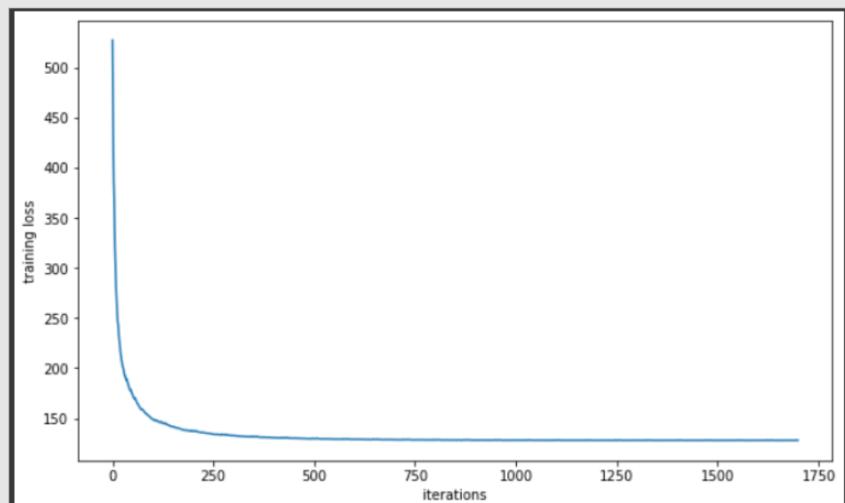
MSE: 110.45143802785354

3. weights and intercepts of the regression line

52.762954387738326 * x + -0.5343233718528416

Logistic regression model :

1. Training loss over iterations



2. Cross entropy of testing data

Cross Entropy: 47.38970895952915

3. weights and intercepts of the regression line

4.776041298265617 * x + 1.7238598687659066

Part 2 : Short answer

1. What's the difference between Gradient Descent, Mini-Batch Gradient Descent, and Stochastic Gradient Descent

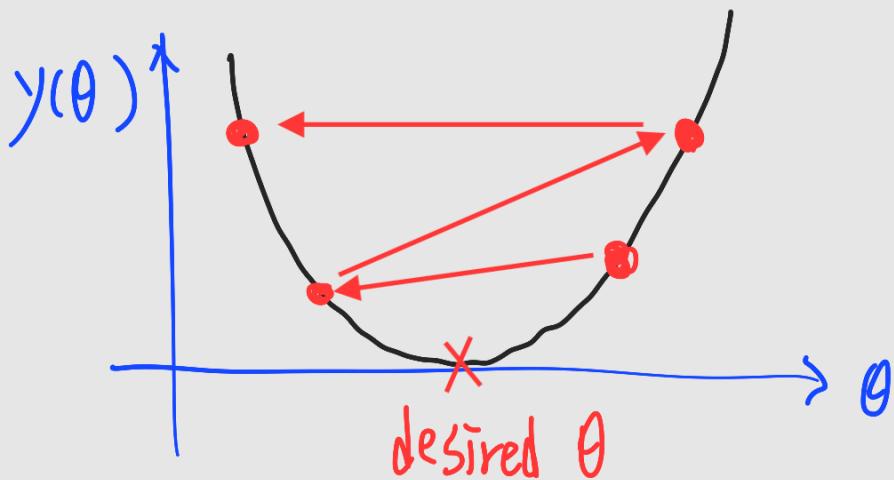
The difference between them is "how much training data is used" for one step of updating the weight. In Gradient Descent, all of the training data is seen before we update the parameter once, which is slow to converge and may result in implementation issue because the training data may be too large to put them all in the memory. For Mini Batch GD, we split training data into small batches and each of them is used for updating ^{weights} once. This batch size is prefer to be the power of 2, which is in favor of the data stored in a computer. Also we need ^{way} to make sure the batch size is small enough to put in the

memory of CPU/GPU. SGD is the case that whenever we see a (X, y) pair, we then update the parameter immediately. The pros of this method is we can derive a not-bad solution in same amount of execution time compared to GD and Mini-batch GD, while the cons is, if we plot the loss-vs-iteration trend, we will observe a escalated curve, not showing how the error decreases after each iteration smoothly.

2. Will different values of learning rate affect the convergence of optimization? Please explain in detail.

Yes, it will. If the learning rate (η) is too small, we may need to spend lots of updating round to get the result, which is not desired if the dataset is large. On the other side, if η is too big, we may step across the point of local minimum and lead to divergent behavior. So it is simply a

eX:



kind of trade-off problem.

3. Show that the logistic sigmoid function (eq. 1) satisfies the property $\sigma(-a) = 1 - \sigma(a)$ and that its inverse is given by $\sigma^{-1}(y) = \ln\{y/(1 - y)\}$.

$$1 - \sigma(a) = 1 - \frac{1}{1 + e^{-a}} = \frac{1 + e^{-a} - 1}{1 + e^{-a}} = \frac{e^{-a}}{1 + e^{-a}}$$

$$= \frac{1}{\frac{e^a}{e^{-a}} + 1} = \frac{1}{1 + e^{-(a)}} = \sigma(-a)$$

$$\therefore f(f^{-1}(x)) = x$$

$$\begin{aligned} & \therefore \sigma(\ln\left(\frac{y}{1-y}\right)) = \frac{1}{1 + \frac{1-y}{y}} \\ & = \frac{1}{1 + e^{-\ln\left(\frac{y}{1-y}\right)}} = \frac{y}{y+1-y} = y \quad \# \end{aligned}$$

4.

int(第n筆資料是第k類)

$$\begin{aligned}
 \textcircled{1} \quad E = E(w_1, w_2, \dots, w_K) &= -\ln P(T | w_1, w_2, \dots, w_K) \\
 &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \ln P(C_k | \phi_n) \\
 &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \ln y_{nk} \rightarrow \text{第n筆資料是第k類的 soft max 結果}
 \end{aligned}$$

$$\textcircled{2} \quad P(C_k | \phi_n) = y_k(\phi_n) = \frac{e^{a_{nk}}}{\sum_j e^{a_{nj}}}$$

$$\Rightarrow \frac{\partial y_{hk}}{\partial a_{nj}} = \frac{e^{a_{hk}} e^{a_{nj}} - e^{a_{hk}} \cdot e^{a_{hk}}}{\left[\sum_j e^{a_{nj}} \right]^2} = y_{hk} - y_{hk}^2$$

(when $j=k$)

$$= y_{hk} (1 - y_{hk})$$

$$\Rightarrow \frac{\partial y_{hk}}{\partial a_{nj}} \underset{(j \neq k)}{=} \frac{-e^{a_{hk}} \cdot e^{a_{nj}}}{\left[\sum_j e^{a_{nj}} \right]^2} = -y_{hk} \cdot y_{nj}$$

$$\Rightarrow \frac{\partial y_{hk}}{\partial a_{nj}} = y_{hk} (I_{kj} - y_{nj}), \quad I_{kj} = \begin{cases} =1, & \text{if } k=j \\ =0, & \text{if } k \neq j \end{cases}$$

(3)

$$a_{nj} = w_j^T \phi_n$$

$$\Rightarrow \frac{\partial a_{nj}}{\partial w_j} = \phi_n$$

Goal

By ① + ② + ③, $(E \leftarrow y \leftarrow a \leftarrow (w))$

(4)

$$\frac{\partial E}{\partial a_{nj}} = \sum_{k=1}^K \frac{\partial E}{\partial y_{nk}} \times \frac{\partial y_{nk}}{\partial a_{nj}}$$

$$= - \sum_{k=1}^K \frac{t_{nk}}{g'_{nk}} y_{nk} (I_{kj} - y_{nj})$$

$$= - \sum_{k=1}^K t_{nk} (I_{kj} - y_{nj})$$

$$= -t_{nj} + \sum_{k=1}^K t_{nk} y_{nj} = y_{nj} - t_{nj}$$

\sim
 $j=k$ 時
那項

$$\sum_{k=1}^K t_{nk} = 1$$

Why this summation?

$$\begin{aligned}
 & y_{n1} - a_{nj} \\
 & y_{n2} - a_{nj} \\
 & \vdots \\
 & y_{nk} - a_{nj} \\
 & \vdots \\
 & y_{nK} - a_{nj}
 \end{aligned}$$

By ④

$$\frac{\partial E}{\partial w_j} = \boxed{\sum_{n=1}^N} \frac{\partial E}{\partial a_{nj}} \times \frac{\partial a_{nj}}{\partial w_j}$$

$$= \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n \#$$

Why this summation:

$$E \leftarrow \begin{array}{l} a_{1j} - w_j \\ a_{2j} - w_j \\ \vdots \\ a_{nj} - w_j \\ \vdots \\ a_{Nj} - w_j \end{array}$$