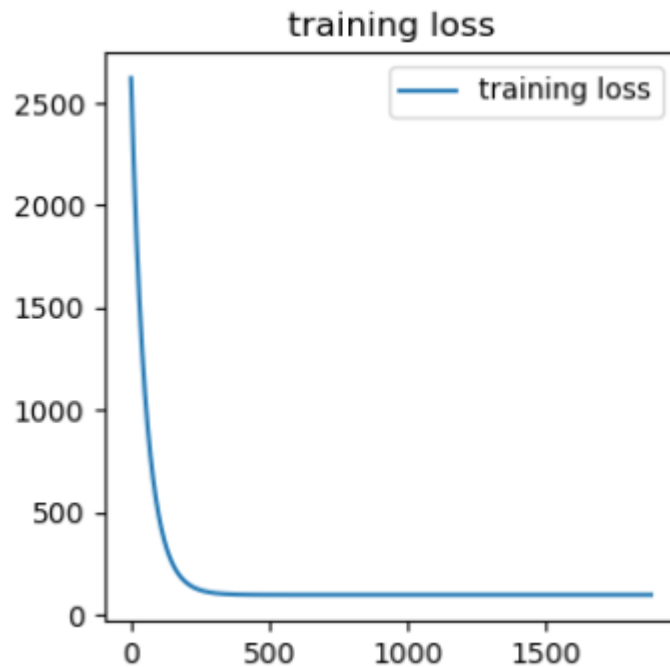# Part.1 Coding

Linear regression model
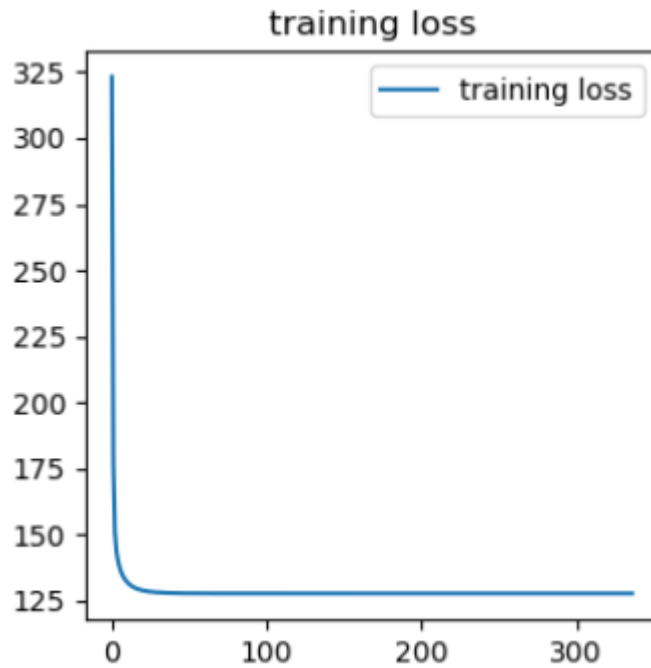
1. Learning Curve:



2. Mean square error: 110.43
3. Weight: 52.74, intersection: -0.33

```
p1:
weight: 52.74353923553911, intersection: -0.33375914007757934
mean_square_error: 110.43818874895103
```

Logistic regression model

1. Learning Curve:

training loss

2. Cross entropy error: 47.24
3. Weight: 4.87, intersection: 1.71

```
p2:
weight: 4.87690932316989, intersection: 1.7116342367476103
cross_entropy_error: 47.2476776764201
```

## Part 2. Questions

1. What's the difference between Gradient Descent, Mini-Batch Gradient Descent, and Stochastic Gradient Descent?
   - Gradient Descent: Feed in all of the training data for every iteration. May be slow for each iteration when dataset is too huge.
   - Mini-Batch Gradient Descent: Feed in a subset of training data for every iteration. Capable to deal with large dataset, but needs more iteration to cover the whole dataset.
   - Stochastic Gradient Descent: Feed in a single data to train for every iteration. Learns very fast, but can't ensure it's converging to correct direction.
2. Will different values of learning rate affect the convergence of optimization? Please explain in detail.
   - High learning rate would cause the model to converge faster, but setting to high might cause the model keep swinging around the minimum and not converging.

- Low learning rate gives the model more precision, but slower to converge. Setting too low might cause the model stuck in a local minimum, instead of global minimum.

3.

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

$$1 + \exp(-a) = \frac{1}{y}$$

$$\exp(-a) = \frac{1 - y}{y}$$

$$-a = \ln\left(\frac{1 - y}{y}\right)$$

$$a = -\ln\left(\frac{1 - y}{y}\right)$$

$$\sigma^{-1}(y) = a = \ln\left(\frac{y}{1 - y}\right)$$

4.

$$\nabla_{w_j} E = -\sum_{n=1}^{N}\sum_{k=1}^{K} \frac{\partial}{\partial a_j}\frac{\partial a_j}{\partial w_j} t_{nk}\ln y_{nk}$$

$$\frac{\partial a_j}{\partial w_j} = \phi_n$$

$$\Rightarrow -\sum_{n=1}^{N}\sum_{k=1}^{K} \frac{\partial}{\partial a_j}\frac{\partial a_j}{\partial w_j} t_{nk}\ln y_{nk} = -\sum_{n=1}^{N}\sum_{k=1}^{K} t_{nk}\phi_n \frac{\partial}{\partial a_j}\ln y_{nk}$$

$$= -\sum_{n=1}^{N}\sum_{k=1}^{K} \phi_n \frac{t_{nk}}{y_{nk}}\left(I_{kj} - y_{nj}\right)y_{nk}$$

$$= -\sum_{n=1}^{N}\sum_{k=1}^{K} \phi_n t_{nk}\left(I_{kj} - y_{nj}\right)$$

$$= -\sum_{n=1}^{N} \phi_n\left(\sum_{k=1}^{K} t_{nk}I_{kj} - \sum_{k=1}^{K} t_{nk}y_{nj}\right)$$

Note that $I_{kj}$ = {1 when k=j, 0 otherwise} and sum of $t_{nk}$ over k is 1.

So, we can get $\sum_{k=1}^{K} t_{nk}I_{kj} = t_{nj}$ and $\sum_{k=1}^{K} t_{nk}y_{nj} = y_{nj}$.

$$\Rightarrow \nabla_{w_j} E = -\sum_{n=1}^{N} \phi_n(t_{nj} - y_{nj}) = \sum_{n=1}^{N} \phi_n(y_{nj} - t_{nj})$$

$\blacksquare$