

2.

(a) Since the gradient decent dimension is randomly chosen and one dimension is only chosen once in one epoch, we have

$$E\{\Delta w(i)\} \stackrel{\text{LLN}}{=} \frac{\eta}{N} \sum_{i=1}^N \nabla_w J(w_i)$$

(b) For Batch-Gradient-Decent, we have:

$$J(w) = \sum_{i=1}^N J(w_i)$$

$$\text{So } \Delta w(i) = \eta \cdot \nabla J(w) = \eta \cdot \nabla \sum_{i=1}^N J(w_i)$$

For SGD, we have

$$E(\Delta w(i)) = (\Delta w_1 \cdot p_1 + \dots + \Delta w_N \cdot p_N) \eta$$

Since we pick with normal distribution, we have:

$$p_1 = p_2 = \dots = p_n = \frac{1}{N}$$

$$\begin{aligned} \text{So } E(\Delta w(i)) &= \frac{\eta}{N} \cdot (\Delta w_1(i) + \dots + \Delta w_N(i)) \\ &= \frac{\eta}{N} \sum_{i=1}^N \nabla J(w_i) \end{aligned}$$

Employ the linearity property of partial derivative:  $\sum \nabla f(x) = \nabla \sum f(x)$

$$\text{we have: } E(\Delta w(i)) = \frac{\eta}{N} \nabla \sum_{i=1}^N J(w_i) = \frac{1}{N} \cdot \Delta w(i), \text{ Q.E.D}$$

**Conclusion:**  $\Delta w(i)$  of BGD is  $N$  times of the  $E\{\Delta w(i)\}$  of SGD

**Explanation:** BGD considers the loss of the whole dataset and SGD considers one at one time. Since SGD data is normally random chosen, Every data has the same probability to be chosen so it is somehow equivalent with considering the whole set. Meanwhile, it's obvious that the loss of SGD should be  $\frac{1}{N}$  of the BGD for its randomness.