

### Question.3-13

Linear regression을 위한 dataset이 다음과 같이 주어졌다.

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$$

이때, dataset은  $y = ax$ 에서부터 만들어졌다.

따라서 linear regression을 통해 predictor를 학습시킬때,

model은  $\hat{y} = \theta x$ , loss는 square error, cost는 MSE를 사용할 수 있다.

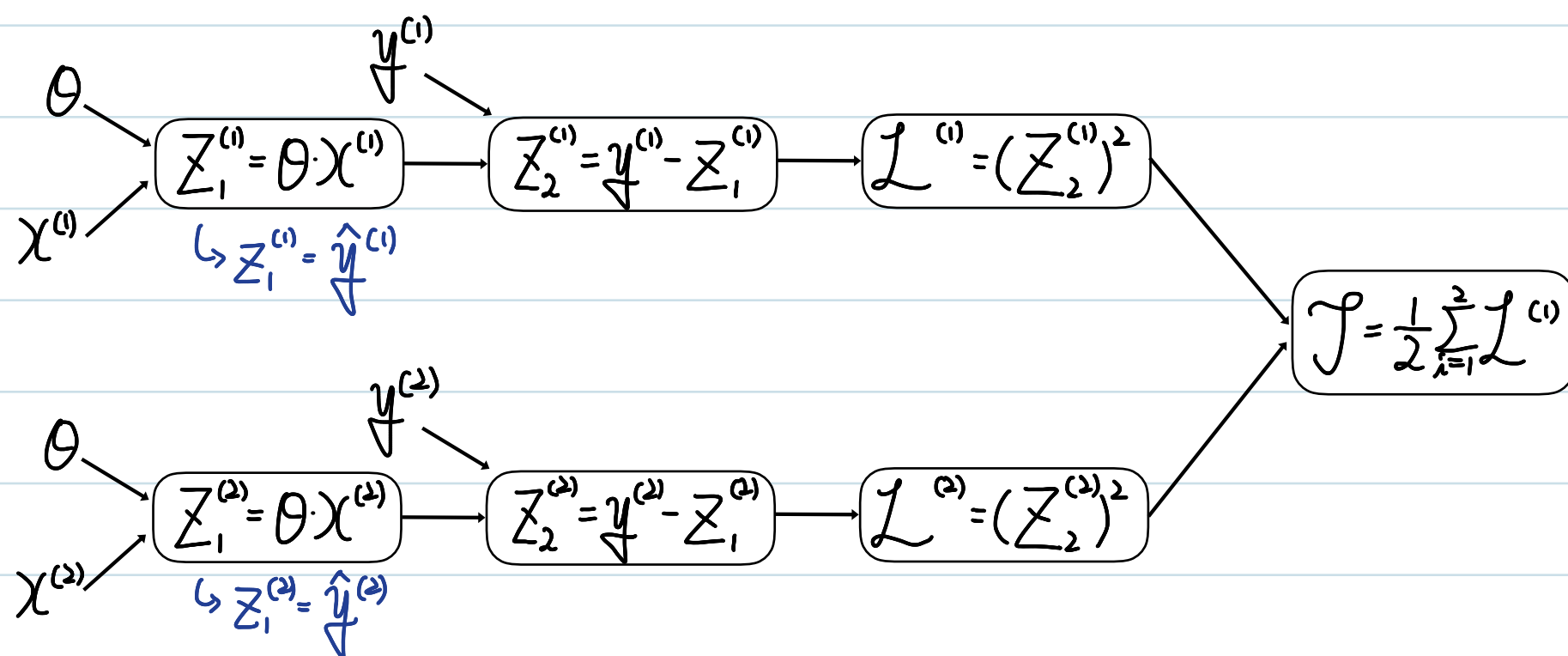
$\theta$ 를 update하기 위해 2개의 data sample를 이용할때, 1번의 iteration에 대해  $\theta$ 가 dataset을 잘 표현하는

$\theta$ 로 update되는 과정을 설명하시오.

단, forward/backward propagation을 설명하기 위해 각 연산은 basic building node들을 이용하시오.

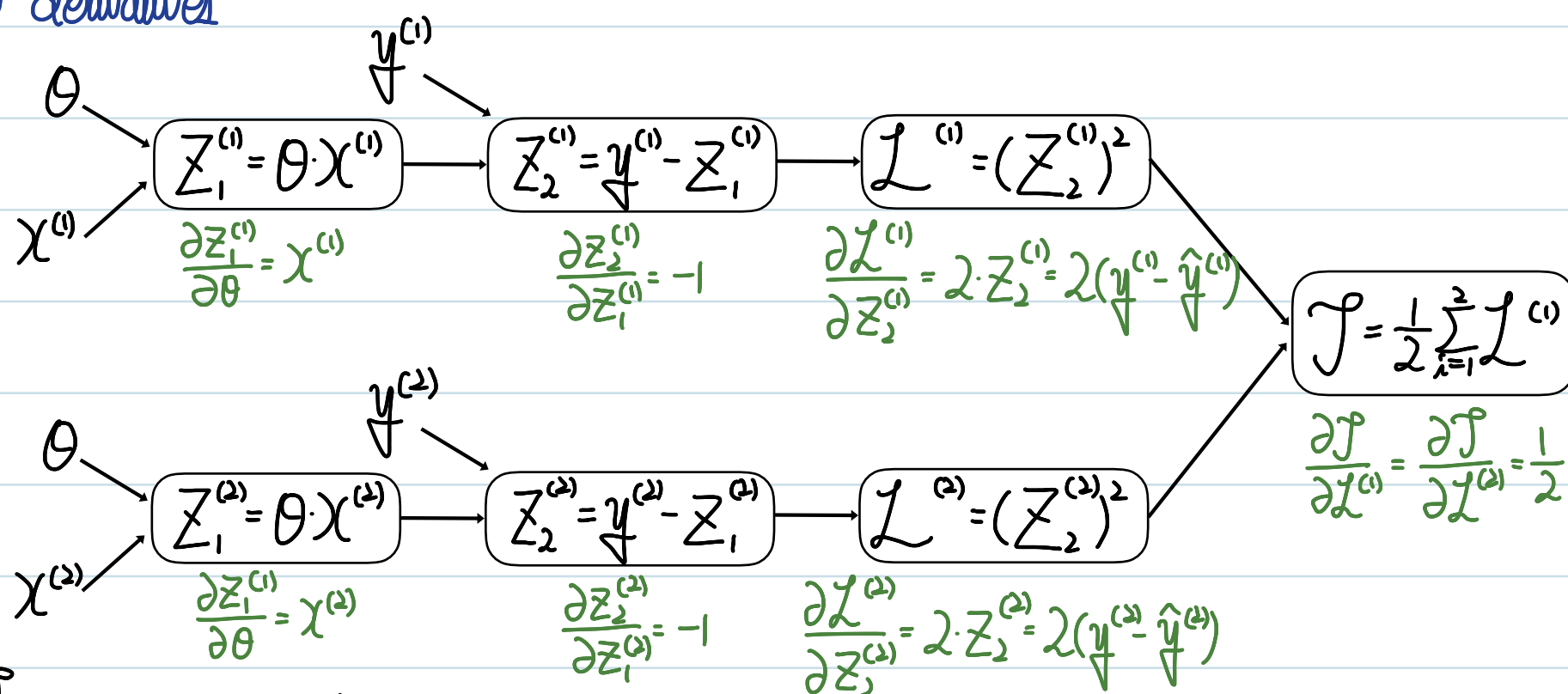
### ① model setting

위의 상황에서  $\theta$ 를 update하기 위해 2개의 data sample를 이용하기 때문에 cost에 대한 gradient descent method를 사용해야한다. 따라서 주어진 상황을 basic building node로 표현하면 다음과 같다.



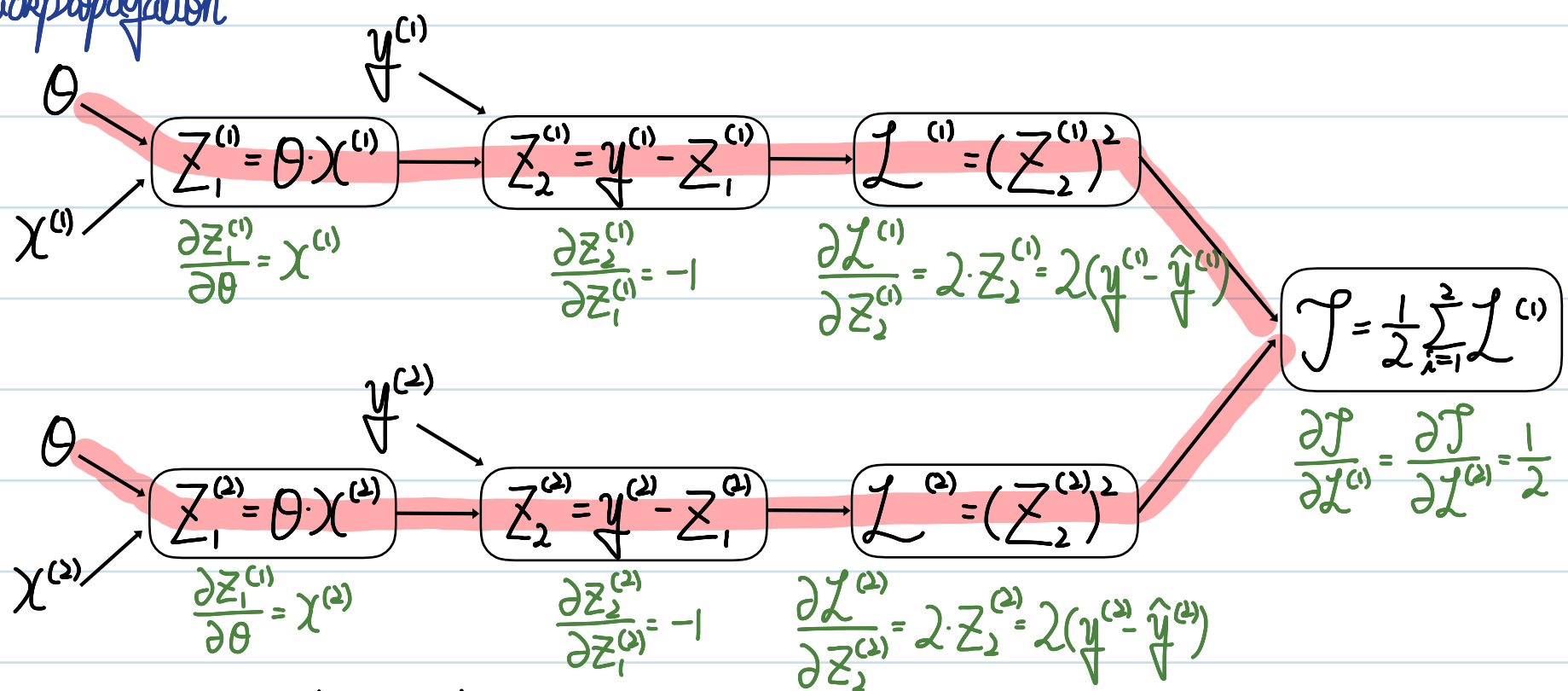
여기서  $z_1^{(1)}, z_1^{(2)}$ 는 각각  $\hat{y}^{(1)}, \hat{y}^{(2)}$  이므로  $z_2^{(1)}, z_2^{(2)}$ 는 각각  $y^{(1)} - \hat{y}^{(1)}, y^{(2)} - \hat{y}^{(2)}$ 가 된다.  
 $\theta$ 의 update에 필요한 partial derivative들을 구하면 다음과 같다.

### ② partial derivatives



따라서  $\frac{\partial J}{\partial \theta}$ 는 다음과 같이 partial derivative가 전파된다.

### ③ Backpropagation



여기서 각 구분에 대해 chain rule을 적용한다.

$$\frac{\partial J}{\partial z_2^{(1)}} = \frac{\partial J}{\partial L^{(1)}} \cdot \frac{\partial L^{(1)}}{\partial z_2^{(1)}} = \frac{1}{2} \cdot 2(y^{(1)} - \hat{y}^{(1)})$$

$$\frac{\partial J}{\partial z_2^{(2)}} = \frac{\partial J}{\partial L^{(2)}} \cdot \frac{\partial L^{(2)}}{\partial z_2^{(2)}} = \frac{1}{2} \cdot 2(y^{(2)} - \hat{y}^{(2)})$$

$$\frac{\partial J}{\partial z_1^{(1)}} = \frac{\partial J}{\partial z_2^{(1)}} \cdot \frac{\partial z_2^{(1)}}{\partial z_1^{(1)}} = \frac{1}{2} \cdot (-2(y^{(1)} - \hat{y}^{(1)}))$$

$$\frac{\partial J}{\partial z_1^{(2)}} = \frac{\partial J}{\partial z_2^{(2)}} \cdot \frac{\partial z_2^{(2)}}{\partial z_1^{(2)}} = \frac{1}{2} \cdot (-2(y^{(2)} - \hat{y}^{(2)}))$$

$$\frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial z_1^{(1)}} \cdot \frac{\partial z_1^{(1)}}{\partial \theta} = \frac{1}{2} (-2x^{(1)}(y^{(1)} - \hat{y}^{(1)}))$$

$$\frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial \theta} = \frac{1}{2} (-2x^{(2)}(y^{(2)} - \hat{y}^{(2)}))$$

여기 결과를 더해

$$\begin{aligned} \frac{\partial J}{\partial \theta} &= \frac{1}{2} (-2x^{(1)}(y^{(1)} - \hat{y}^{(1)})) + \frac{1}{2} (-2x^{(2)}(y^{(2)} - \hat{y}^{(2)})) \\ &= \frac{1}{2} [(-2x^{(1)}(y^{(1)} - \hat{y}^{(1)})) + (-2x^{(2)}(y^{(2)} - \hat{y}^{(2)}))] \end{aligned}$$

가 된다. 이때 각 term들은  $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ 에 대해  $\frac{\partial L^{(1)}}{\partial \theta}$ ,  $\frac{\partial L^{(2)}}{\partial \theta}$ 는

$$\frac{\partial L^{(1)}}{\partial \theta} = -2x^{(1)}(y^{(1)} - \theta \cdot x^{(1)}) \quad \frac{\partial L^{(2)}}{\partial \theta} = -2x^{(2)}(y^{(2)} - \theta \cdot x^{(2)})$$

이므로

$$\frac{\partial J}{\partial \theta} = \frac{1}{2} \left[ \frac{\partial L^{(1)}}{\partial \theta} + \frac{\partial L^{(2)}}{\partial \theta} \right]$$

이므로 즉, cost를 최소화하여  $\theta$ 를 update하려면  $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ 에 대해  $\theta$ 가 update되는 양을 평균적으로 사용한다.

### ④ gradient descent method

$\frac{\partial J}{\partial \theta}$ 와 gradient descent method를 이용하여  $\theta$ 에 대한 update를 표현하면

$$\theta := \theta - \alpha \cdot \frac{\partial J}{\partial \theta}$$

이므로 ③에서 구한  $\frac{\partial J}{\partial \theta}$ 를 대입하면

$$\begin{aligned} \theta &:= \theta + \frac{\alpha}{2} \sum_{i=1}^2 \frac{\partial L^{(i)}}{\partial \theta} \\ &= \theta + \frac{\alpha}{2} \sum_{i=1}^2 [2x^{(i)}(y^{(i)} - \hat{y}^{(i)})] \end{aligned}$$

가 된다. 이때  $\frac{\partial J}{\partial \theta}$ 는  $\frac{\partial L^{(i)}}{\partial \theta}$ 들의 평균값으로  $\frac{\partial J}{\partial \theta}$ 에 대해

Question. 3-09에서 설명한  $\theta$ 의 학습방식은 2가지 적용된다.