

1. Regression : Approximation of target function based on given data.

2. Linear Regression : Regression of "Linear form"
or Linear Combination (Linear mapping)

* Statistics / Deep Learning의 LR 차이

↳ noise를 더해준다.
[3] ↗ 공식: noise
이거를 예상!

3. Loss, Cost function

$$SSE = Loss_{(i)} = (y_i - \hat{y}_i)^2$$

↓

$$MSE = \text{cost} = \frac{1}{N} \sum_{i=1}^N Loss_{(i)} = \frac{1}{N} \sum (y_i - \hat{y}_i)^2$$

$$\theta := \theta - \alpha \nabla L(\theta) \quad * \nabla L = \frac{\partial L}{\partial \theta}$$

$$\theta := \theta - \alpha \nabla J(\theta) \quad \nabla J = \frac{\partial J}{\partial \theta}$$

4. gradient descent method

model $\hat{y} = \theta_n x_n + \theta_{n-1} x_{n-1} + \dots + \theta_1 x_1 + \theta_0$

params $(\vec{w}, b) = (\theta_n, \theta_{n-1}, \dots, \theta_1, \theta_0)$

loss $\mathcal{L}^{(i)}(\vec{\theta}) = (y^{(i)} - \hat{y}^{(i)})^2 = (y^{(i)} - (\theta_n x_n^{(i)} + \theta_{n-1} x_{n-1}^{(i)} + \dots + \theta_1 x_1^{(i)} + \theta_0))^2$

cost $J(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{(i)}$

Fast campus

gradient $\nabla_{\vec{\theta}} \mathcal{L}^{(i)}(\vec{\theta}) = \left(\frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_n}, \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_{n-1}}, \dots, \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_1}, \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_0} \right)$

gradient descent $\vec{\theta} := \vec{\theta} - \alpha \nabla_{\vec{\theta}} \mathcal{L}^{(i)}(\vec{\theta})$

$$\theta_n := \theta_n - \alpha \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_n}$$

$$\theta_{n-1} := \theta_{n-1} - \alpha \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_{n-1}}$$

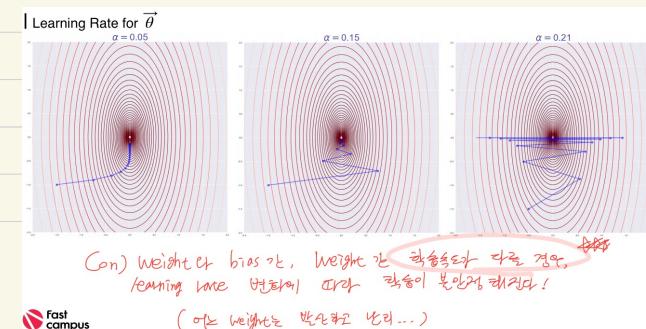
$$\vdots$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial \mathcal{L}^{(i)}(\vec{\theta})}{\partial \theta_1}$$

Fast campus
신경식 강사.

5. learning rate

$$\theta = \theta - \alpha \frac{\partial L}{\partial \theta} \Rightarrow \alpha \uparrow \text{면, 발산 위험} \uparrow \text{or 진동으로 이동할 확률} \uparrow \text{가능.}$$



• θ update의 발산 or 진동은 learning rate 때문에 input의 평균이 0이 아님!

$$\Delta \theta = \alpha \cdot \frac{\partial L}{\partial \theta}$$

$$\frac{\partial L}{\partial \theta} = -2x_{(i)}(y_{(i)} - \theta x_{(i)})$$

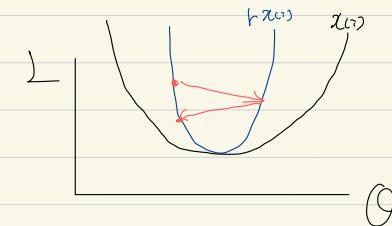
만약, $x_{(i)}$ 가 b -x_i 가 되면,

$$\frac{\partial L}{\partial \theta} = -2x_{(i)}^2(y_{(i)} - \theta x_{(i)})$$

∴ input이 b 에 증가시,

learning rate(α)의 영향은

$$\alpha \cdot 1 \rightarrow \alpha \cdot t^2 으로!$$



실제로 input이 증가시, loss function의 U 모양이
더 가팔다짐 (Convex ↑)

$$\therefore Loss_{(i)} = \theta^2 x_{(i)}^2 - 2\theta x_{(i)} y_{(i)} + y_{(i)}^2$$

• 즉, LR↓ 때도, input이 크면 $\Delta \theta$ 크거나 커서 불안정한 학습 가능하다는 것.

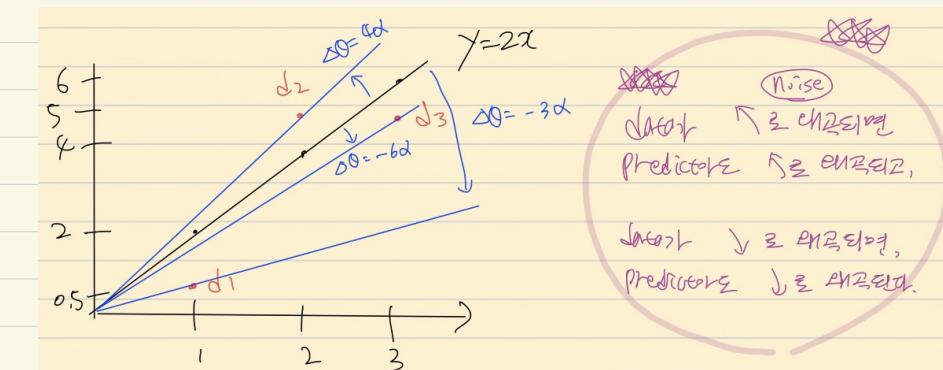
⇒ Cost function 활용! (안정적, 빠르게)

데이터를 Minibatch size(N)로 묶어, 각 데이터로

$\Delta \theta$ 크기를 $\frac{1}{n}$ 로 pooling 하는 효과! * $J(\theta) = \frac{1}{n} \sum L(\theta)$

6. Noise

Noise가 학습에 미치는 영향.



7. Learning - process

```
Total Learning Process
# dataset preparation
dataset_gen = dataset_generator()
dataset_gen.set_coefficient([5, 0])
x_data, y_data = dataset_gen.make_dataset()

# model implementation
node1 = nodes.mul_node()
node2 = nodes.minus_node()
node3 = nodes.square_node()

# hyperparameter setting
epochs = 5 # total epoch setting
lr = 0.01 # learning rate setting
th = -1 # arbitrary theta
loss_list = []
th_list = []

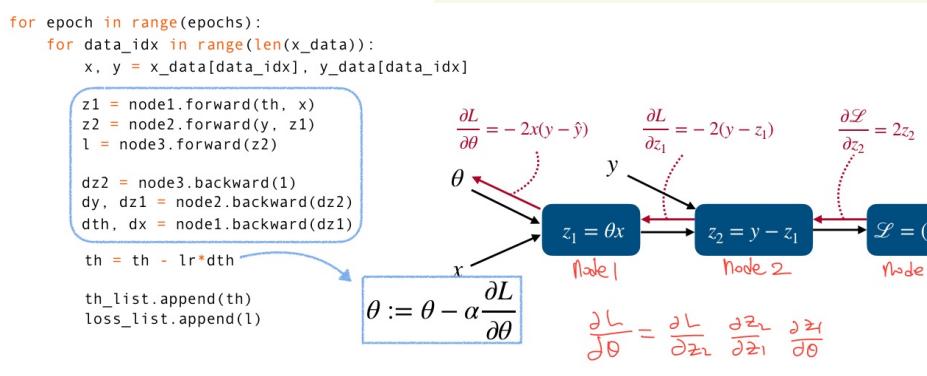
for epoch in range(epochs):
    for data_idx in range(len(x_data)):
        x, y = x_data[data_idx], y_data[data_idx]

        z1 = node1.forward(th, x)
        z2 = node2.forward(y, z1)
        l = node3.forward(z2)

        dz2 = node3.backward(1)
        dz1 = node2.backward(dz2)
        dth, dx = node1.backward(dz1)

        th = th - lr * dth
        th_list.append(th)
        loss_list.append(l)

# square error loss implementation
node2 = nodes.minus_node()
node3 = nodes.square_node()
```



8. Hyper parameter tuning

- learning rate
- initial theta
- iteration (mini batch size)
- epoch
- Coefficient.

9. Mini batch \Rightarrow Cost function 합성 연산!

Iterations/EPOCHs and Mini-batch

Total # sample: 1000
Mini-batch size: 20

Iteration per epoch: $\frac{\text{total sample}}{\text{mini-batch size}} = \frac{1000}{20} = 50$ | Epoch of 50 iterations (batch size = 20)
iterations for 5 epochs: $5 * 50 = 250$ Iterations

$\nabla L_{(i)} = -2x_{(i)}(y_{(i)} - \theta x_{(i)})$

$\nabla J = \frac{1}{n} \sum -2x_{(i)}(y_{(i)} - \theta x_{(i)})$

$\frac{\partial L}{\partial \theta} = \frac{1}{n} \sum (-2x_{(i)}(y_{(i)} - \theta x_{(i)}))$

$\frac{\partial J}{\partial \theta} = \frac{1}{n} \sum \frac{\partial L}{\partial \theta}$

Fast campus \rightarrow $\frac{1}{n} \sum \frac{\partial L}{\partial \theta}$

10. Quiz 2.

Review

model $\hat{y} = \theta x$

loss $L^{(i)} = (y^{(i)} - \hat{y}^{(i)})^2 = (y^{(i)} - \theta x^{(i)})^2$

cost $J(\theta) = \frac{1}{n} \sum L^{(i)} = \frac{1}{n} \sum (y^{(i)} - \theta x^{(i)})^2$

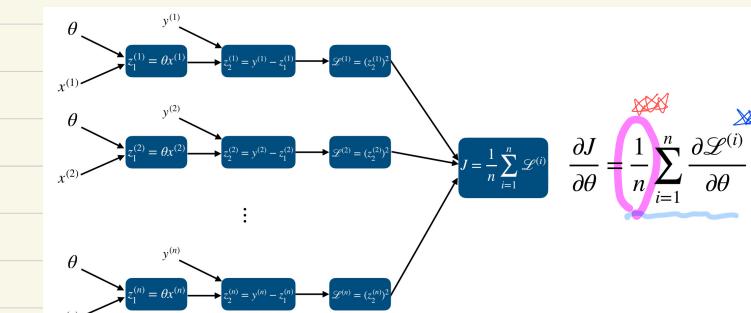
gradient $\frac{\partial J(\theta)}{\partial \theta} = -\frac{1}{n} \sum \left[2x^{(i)}(y^{(i)} - \theta x^{(i)}) \right]$

gradient descent $\theta := \theta - \alpha \frac{\partial J(\theta)}{\partial \theta}$

parameter update $\theta := \theta + \frac{\alpha}{n} \sum_{i=1}^n [2x^{(i)}(y^{(i)} - \theta x^{(i)})]$

11. mini-batch의 효과

- ① 미니 배치 Loss $L^{(i)}$ 의 θ update A) 소량의 표본 \Rightarrow 더 robust한 학습률
- ② 미니 배치의 학습률 (epoch vs iteration \downarrow)



12. Vectorization for Several Samples

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \vec{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\frac{\partial J}{\partial \vec{z}_1} = \left(\frac{-1}{n} 2z_1^{(1)}, \frac{-1}{n} 2z_1^{(2)}, \dots, \frac{-1}{n} 2z_1^{(n)} \right) \quad \frac{\partial J}{\partial \vec{z}_2} = \begin{pmatrix} -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & -1 \end{pmatrix} \quad \frac{\partial J}{\partial \vec{Z}} = \left(\frac{1}{n} 2z_2^{(1)}, \frac{1}{n} 2z_2^{(2)}, \dots, \frac{1}{n} 2z_2^{(n)} \right)$$

$$\frac{\partial J}{\partial \vec{\theta}} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right)$$

$$\frac{\partial \vec{z}_1}{\partial \theta} = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{pmatrix} \quad \frac{\partial \vec{z}_2}{\partial \theta} = \vec{y} - \vec{z}_1 \quad \frac{\partial \vec{Z}}{\partial \theta} = (\vec{z}_2)^2 \quad J = \text{mean}(\vec{Z})$$

$$\frac{\partial \vec{Z}}{\partial \vec{\theta}} = \begin{pmatrix} 2z_2^{(1)} & 0 & \dots & 0 \\ 0 & 2z_2^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 2z_2^{(n)} \end{pmatrix} \quad \frac{\partial J}{\partial \vec{Z}} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right)$$

13. Effect of Cost Node

Effect of the Cost Node

$\theta \rightarrow \vec{z}_1 = \theta \vec{x} \rightarrow \vec{z}_2 = \vec{y} - \vec{z}_1 \rightarrow \vec{Z} = (\vec{z}_2)^2 \rightarrow J = \text{mean}(\vec{Z})$

$\frac{\partial J}{\partial \vec{Z}} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right)$

$\frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial \vec{Z}} \frac{\partial \vec{Z}}{\partial \theta} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right) \begin{pmatrix} \frac{\partial Z^{(1)}(\theta)}{\partial \theta} \\ \frac{\partial Z^{(2)}(\theta)}{\partial \theta} \\ \vdots \\ \frac{\partial Z^{(n)}(\theta)}{\partial \theta} \end{pmatrix} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right) \begin{pmatrix} \frac{\partial Z^{(1)}(\theta)}{\partial \theta} \\ \frac{\partial Z^{(2)}(\theta)}{\partial \theta} \\ \vdots \\ \frac{\partial Z^{(n)}(\theta)}{\partial \theta} \end{pmatrix}$

Cost의 역할
forward \Rightarrow 학습률 구하기
backward $\Rightarrow \frac{1}{n}$ 배치 위한 vectorization

Back propagation (미니 배치)
Cost Node는 $\frac{1}{n}$ 배치의 역할

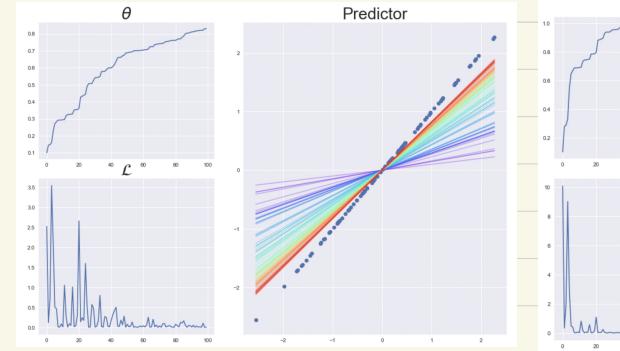
Cost Node는 $\frac{1}{n}$ 배치의 역할
Update A), 미니 배치의 역할

* forward propagation 역할,
* mini-batch의 역할
* backward propagation 역할

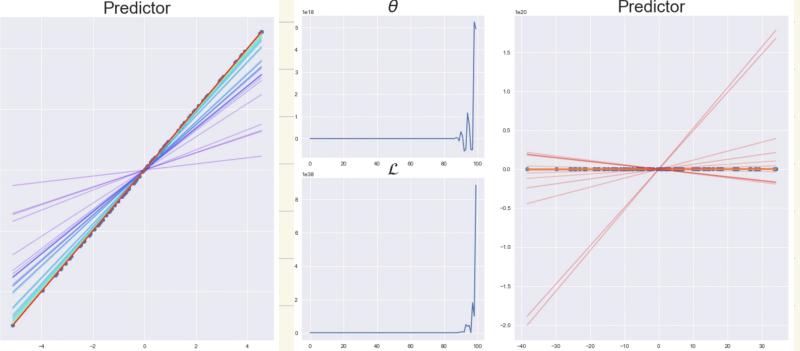
17. Learning with Various Standard Deviation

- 데이터의 표준편차가 클수록, 학습률 확률의 보통이 넓어짐. (심지어 6이 크면, 0이 빨라짐!)
- 따라서, 표준편차가 큰 데이터를 대상으로 할 때, Learning rate를 ↓↓ 할 필요가 있음.

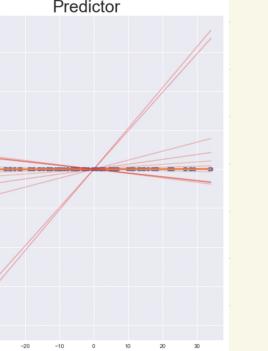
$\sigma = 1$



$\sigma = 2$



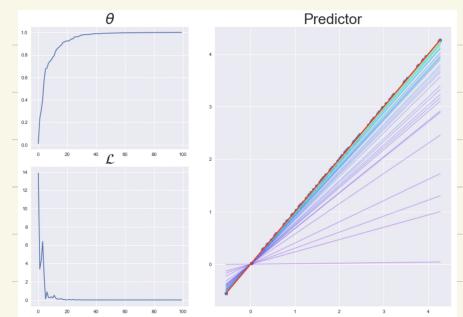
$\sigma = 15$



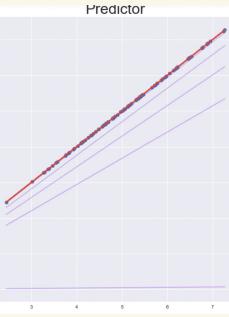
18. Learning with various mean

- 데이터의 mean의 절대값이 클수록 θ 확률 진동이 커지는 모습!
- std가 마찬가지로, 데이터의 mean (mean이) 클수록 learning rate를 작게 잡을 필요가 있음!

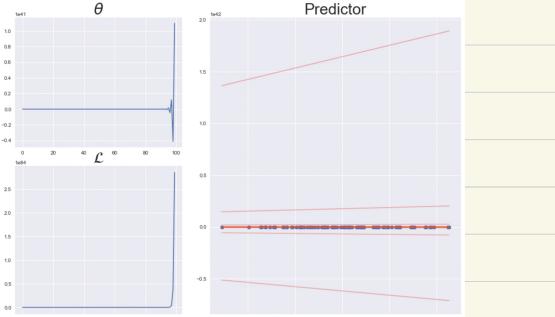
$\mu = 2$



$\mu = 5$



$\mu = 15$



※ 음수도 마찬가지 모습 (-15에서 0까지만)