

# Statistical Machine Learning

3주차

담당: 11기 명재성

# Review

---

- Regression

$$Y_i \stackrel{ind}{\sim} N(\mu_i(\mathbf{X}_i), \sigma) \quad \text{where} \quad E[Y_i] = \mu_i(\mathbf{X}_i)$$

$$\begin{aligned}\mu_i(\mathbf{X}_i) &= \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_p X_{pi} \\ &= \boldsymbol{\beta}^T \mathbf{X}_i\end{aligned}$$

# Review

---

- Logistic Regression

$$Y_i \overset{ind}{\sim} \text{Bernoulli}(\pi_i(\mathbf{X}_i)) \quad \text{where} \quad E[Y_i] = \pi_i(\mathbf{X}_i)$$

$$\begin{aligned} \text{logit}(\pi_i(\mathbf{X}_i)) &= \log\left(\frac{\pi_i(\mathbf{X}_i)}{1 - \pi_i(\mathbf{X}_i)}\right) = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_p X_{pi} \\ &= \boldsymbol{\beta}^T \mathbf{X}_i \end{aligned}$$

# Review

---

- Estimation

$$\underset{\boldsymbol{\beta}}{\operatorname{argmin}} L[\tau(\theta), T(X)] \Leftrightarrow \underset{\boldsymbol{\beta}}{\operatorname{argmax}} L(\boldsymbol{\beta}, \sigma)$$

- Regression  $\rightarrow$  SSE
- Logistic Regression  $\rightarrow$  Cross Entropy

# Review

---

- Cross Entropy

$$CE_i = - \sum_{k=1}^C y_{ik} \log \pi_i(k)$$

$y_{ik}$  : the  $k^{th}$  value in  $y_i$

$\pi_i(k)$  : the probability for the  $i^{th}$  observation to belonging to Class k

# Review

---

- For  $C = 3$  (number of Class)

Class 1:  $y_i = (1, 0, 0)$

Class 2:  $y_i = (0, 1, 0)$

Class 3:  $y_i = (0, 0, 1)$

⇒ One-Hot encoding

$$\sum_{i=1}^n CE_i = - \sum_{k=1}^C [y_{i1} \log \pi_i(1) + y_{i2} \log \pi_i(2) + y_{i3} \log \pi_i(3)]$$

# Review

---

- IF Class 2,

Class 1:  $y_i = (1, 0, 0)$

Class 2:  $y_i = (0, 1, 0)$

Class 3:  $y_i = (0, 0, 1)$

$$CE_i = -[0 * \log \pi_i(1) + 1 * \log \pi_i(2) + 0 * \log \pi_i(3)]$$

$$= -\log \pi_i(2)$$

$\Rightarrow$  IF  $\pi_i(2) = 1$ , then  $CE_i = -\log 1 = 0$  (minimum Loss)

# Review

---

- Categorical Cross Entropy

$$CE_i = - \sum_{k=1}^C y_{ik} \log \pi_i(k)$$

- Likelihood of ?

$$Y_i \stackrel{ind}{\sim} Multi(\pi_1, \dots, \pi_k)$$



# Review

---

- For  $C = 2$  (number of Class)

$$P(Y_i = 1 | \mathbf{X}_i) = \pi_i(\mathbf{X}_i) = \frac{e^{\boldsymbol{\beta}^T \mathbf{X}_i}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_i}} = \frac{1}{1 + e^{-\boldsymbol{\beta}^T \mathbf{X}_i}} \quad (\text{sigmoid function})$$

# Review

---

- For  $C = 2$  (number of Class)

$$P(Y = 1|\mathbf{X}_i) = \pi(\mathbf{X}_i) = \frac{e^{\boldsymbol{\beta}^T \mathbf{X}_i}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_i}} = \frac{1}{1 + e^{-\boldsymbol{\beta}^T \mathbf{X}_i}} \quad (\text{sigmoid function})$$

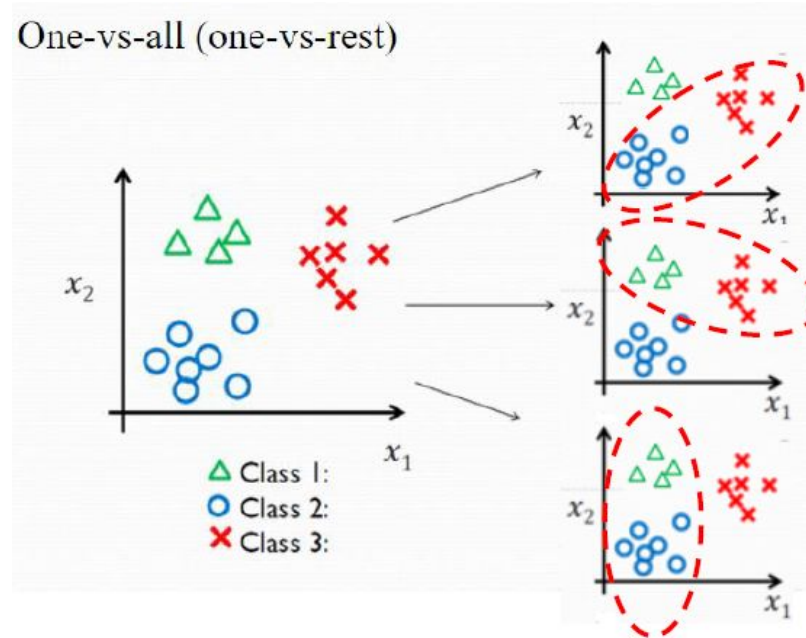
- For  $C = K > 2$  (number of Class)

$$P(Y = l|\mathbf{X}_i) = \pi_l(\mathbf{X}_i) = \frac{e^{\boldsymbol{\beta}_l^T \mathbf{X}_i}}{\sum_{c=1}^K e^{\boldsymbol{\beta}_c^T \mathbf{X}_i}} \quad (\text{softmax function})$$

# Review

- One-Vs-Rest

One-vs-all (one-vs-rest)



# Review

---

- One-Vs-Rest

```
LogisticRegression(solver='sag', multi_class='multinomial')
```

```
'ovr'
```

# Review

---

- Naïve Bayes Classifier

$$P(Y_i = k | \mathbf{X}_i) = \frac{P(\mathbf{X}_i | k)P(k)}{\sum_k P(\mathbf{X}_i | k)P(k)} \quad \text{Bayes' Theorem}$$

$$\text{where } P(\mathbf{X}_i | k) = \prod_j^p P(X_{ij} | k)$$

# Review

---

- LDA

$$P(Y_i = k | \mathbf{X}_i) = \frac{P(\mathbf{X}_i | k)P(k)}{\sum_k P(\mathbf{X}_i | k)P(k)} \quad \text{Bayes' Theorem}$$

where  $P(\mathbf{X}_i | k) \sim N_p(\boldsymbol{\mu}_k, \Sigma)$

# Review

---

- QDA

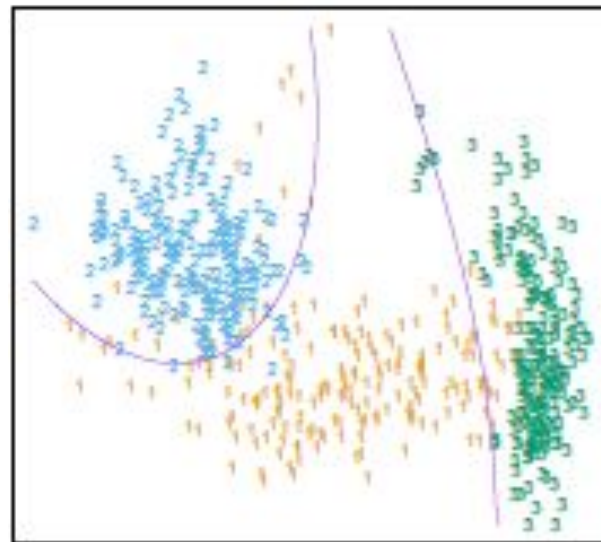
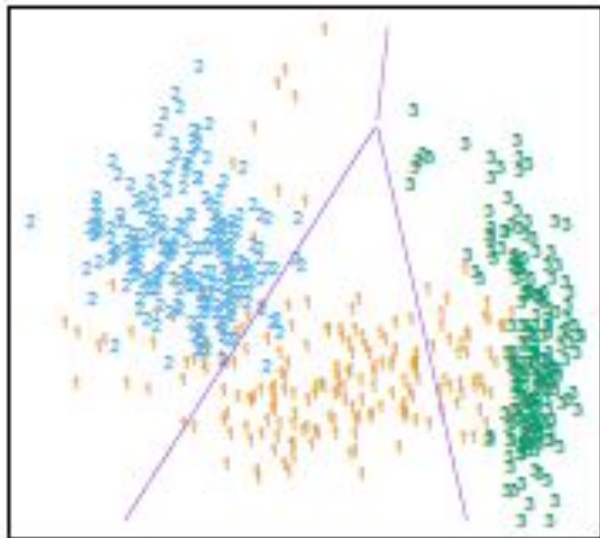
$$P(Y_i = k | \mathbf{X}_i) = \frac{P(\mathbf{X}_i | k)P(k)}{\sum_k P(\mathbf{X}_i | k)P(k)}$$

*Bayes' Theorem*

where  $P(\mathbf{X}_i | k) \sim N_p(\boldsymbol{\mu}_k, \Sigma_k)$

# Review

---





# Review

---

- 0-1 Loss

$$L[\tau(\theta), T(X)] = \sum I(Y_i \neq \hat{Y}_i)$$

- The Bayes decision rule for minimizing the loss (  $P(Y_i \neq \hat{Y}_i)$  ) is

$$\underset{k}{\operatorname{argmax}} P(Y = k|\mathbf{X})$$

# Optimization

---

- Loss function

$$L[\tau(\theta), T(X)] = L[Y, \hat{Y}]$$

$$\text{Regression} \Rightarrow L[Y, \hat{Y}] = \sum_i^n (Y_i - \hat{Y}_i)^2$$

$$\text{Classification} \Rightarrow L[Y, \hat{Y}] = - \sum_i^n \sum_k^C Y_i \log \hat{\pi}_i(k)$$

# Optimization

---

- Machine “Learning”

$$\underset{\theta}{\operatorname{argmin}} L[Y, \hat{Y}] = \hat{\theta}$$

$\Rightarrow$  *Optimization*

# Optimization

---

- Logistic Regression

$$\begin{aligned} L[Y, \hat{Y}] &= - \sum_{i=1}^n [y_i \log \pi_i + (1 - y_i) \log(1 - \pi_i)] \\ &= - \sum_{i=1}^n [y_i (\boldsymbol{\beta}^T \mathbf{x}_i) - \log(1 + \exp(\boldsymbol{\beta}^T \mathbf{x}_i))] \end{aligned}$$

# Optimization

---

- Logistic Regression

$$L[Y, \hat{Y}] = - \sum_{i=1}^n [y_i(\boldsymbol{\beta}^T \mathbf{x}_i) - \log(1 + \exp(\boldsymbol{\beta}^T \mathbf{x}_i))]$$

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} L[Y, \hat{Y}]$$

⇒ Can you solve it?

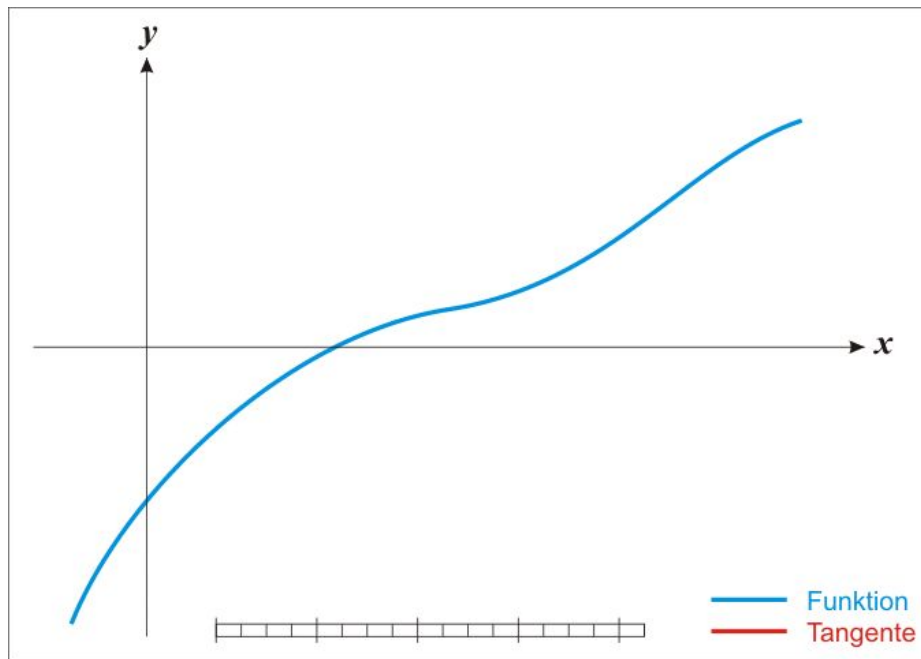
# Optimization

---

- Optimization often can be rewritten as solving equations.  
ex ) Normal equation
- Some problems do not have an explicit solution and a numerical approach should be exploited.

# Newton-Raphson Method

---



# Newton-Raphson Method

---

- Linear approximation (1<sup>st</sup> order Taylor Expansion)

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) = 0$$

$$x = x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$\Rightarrow \theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$



# Newton-Raphson Method

---

1. Initialize  $\theta^{(0)} = \theta_0$  which can be arbitrary on the domain of the function
2. Update for  $t = 0, 1, 2, 3, \dots$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

until

$$|\theta^{(t+1)} - \theta^{(t)}| < \epsilon$$

for small  $\epsilon > 0$

# Newton-Raphson Method

---

- Quadratic approximation (2<sup>nd</sup> order Taylor Expansion)

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

$$\frac{\partial}{\partial x} f(x) \approx f'(x_0) + f''(x_0)(x - x_0) = 0$$

$$x = x_0 - \frac{f'(x_0)}{f''(x_0)} \quad \Rightarrow \quad \theta^{(t+1)} = \theta^{(t)} - \frac{f'(\theta^{(t)})}{f''(\theta^{(t)})}$$

# Newton-Raphson Method

---

- Quadratic approximation (2<sup>nd</sup> order Taylor Expansion)

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + L'(\boldsymbol{\theta}_0)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T L''(\boldsymbol{\theta}_0) (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

# Newton-Raphson Method

---

- Quadratic approximation (2<sup>nd</sup> order Taylor Expansion)

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + \nabla L(\boldsymbol{\theta}_0)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}(\boldsymbol{\theta}_0) (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

where

$$\nabla L(\boldsymbol{\theta}_0) = \left. \frac{\partial}{\partial \boldsymbol{\theta}} L(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$$

$$\mathbf{H}(\boldsymbol{\theta}_0) = \left. \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} L(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$$

# Newton-Raphson Method

---

- Updating equation is

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \mathbf{H}^{-1}(\boldsymbol{\theta}^{(t)}) \nabla L(\boldsymbol{\theta}^{(t)}) \\ &= \boldsymbol{\theta}^{(t)} - \mathbf{H}^{-1}(\boldsymbol{\theta}^{(t)}) \frac{\partial}{\partial \boldsymbol{\theta}^{(t)}} L(\boldsymbol{\theta}^{(t)})\end{aligned}$$

$$cf. \quad \theta^{(t+1)} = \theta^{(t)} - \frac{f'(\theta^{(t)})}{f''(\theta^{(t)})}$$

# Newton-Raphson Method

---

$$\text{Loss}[\boldsymbol{\beta}] = - \sum_{i=1}^n [y_i(\boldsymbol{\beta}^T \mathbf{X}_i) - \log(1 + \exp(\boldsymbol{\beta}^T \mathbf{X}_i))]$$

$$\nabla L(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} L(\boldsymbol{\beta}) = - \sum_{i=1}^n \left[ y_i \mathbf{X}_i - \frac{\exp(\boldsymbol{\beta}^T \mathbf{X}_i)}{1 + \exp(\boldsymbol{\beta}^T \mathbf{X}_i)} \mathbf{X}_i \right]$$

$$\mathbf{H}(\boldsymbol{\beta}) = \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[ \left( \frac{\exp(\boldsymbol{\beta}^T \mathbf{X}_i)}{1 + \exp(\boldsymbol{\beta}^T \mathbf{X}_i)} \right) \left( \frac{1}{1 + \exp(\boldsymbol{\beta}^T \mathbf{X}_i)} \right) \mathbf{X}_i \mathbf{X}_i^T \right]$$

# Newton-Raphson Method

---

$$\text{Loss}[\boldsymbol{\beta}] = - \sum_{i=1}^n [y_i(\boldsymbol{\beta}^T \mathbf{x}_i) - \log(1 + \exp(\boldsymbol{\beta}^T \mathbf{x}_i))]$$

Update

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \mathbf{H}^{-1}(\boldsymbol{\beta}^{(t)}) \nabla L(\boldsymbol{\beta}^{(t)})$$

until

$$\|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^t\| < \epsilon \quad \text{for small } \epsilon > 0$$

# Newton-Raphson Method

Solvers					
Penalties	'liblinear'	'lbfgs'	'newton-cg'	'sag'	'saga'
Multinomial + L2 penalty	no	yes	yes	yes	yes
OVR + L2 penalty	yes	yes	yes	yes	yes
Multinomial + L1 penalty	no	no	no	no	yes
OVR + L1 penalty	yes	no	no	no	yes
Behaviors					
Penalize the intercept (bad)	yes	no	no	no	no
Faster for large datasets	no	no	no	yes	yes
Robust to unscaled datasets	yes	yes	yes	no	no



# Gradient Descent

---

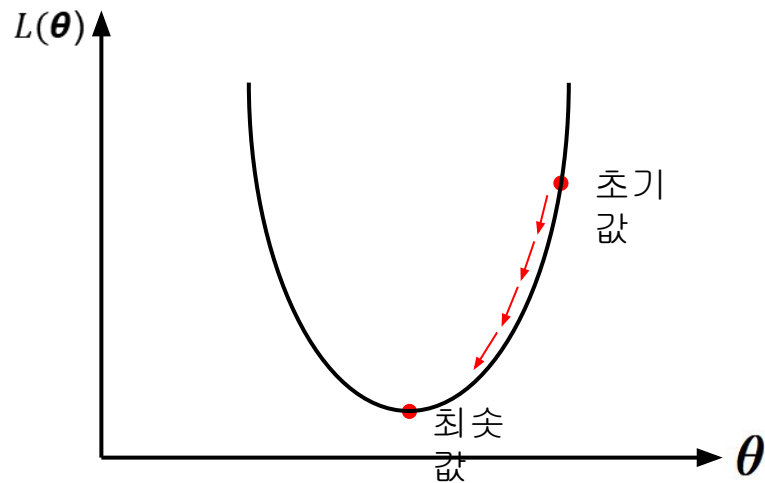
- Newton-Raphson is expensive to compute due to the computation of the inverse of Hessian.

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \mathbf{H}^{-1}(\boldsymbol{\theta}^{(t)}) \nabla L(\boldsymbol{\theta}^{(t)})$$

$$\Rightarrow \boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta^{(t)} \nabla L(\boldsymbol{\theta}^{(t)})$$

# Gradient Descent

---



# Gradient Descent

---

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta^{(t)} \nabla L(\boldsymbol{\theta}^{(t)}) \quad \text{or} \quad \boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \nabla L(\boldsymbol{\theta}^{(t)})$$

```
t0, t1 = 5, 50 # learning schedule hyperparameters
```

```
def learning_schedule(t):  
    return t0 / (t + t1)
```

```
eta = learning_schedule(epoch * m + i)  
theta = theta - eta * gradients
```

# Batch Gradient Descent

---

- Regression → SSE

$$\nabla L(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} L(\boldsymbol{\beta}) = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- Logistic Regression → Cross Entropy

$$\nabla L(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} L(\boldsymbol{\beta}) = - \sum_{i=1}^n \left[ y_i \mathbf{x}_i - \frac{\exp(\boldsymbol{\beta}^T \mathbf{x}_i)}{1 + \exp(\boldsymbol{\beta}^T \mathbf{x}_i)} \mathbf{x}_i \right]$$

# Steepest Descent

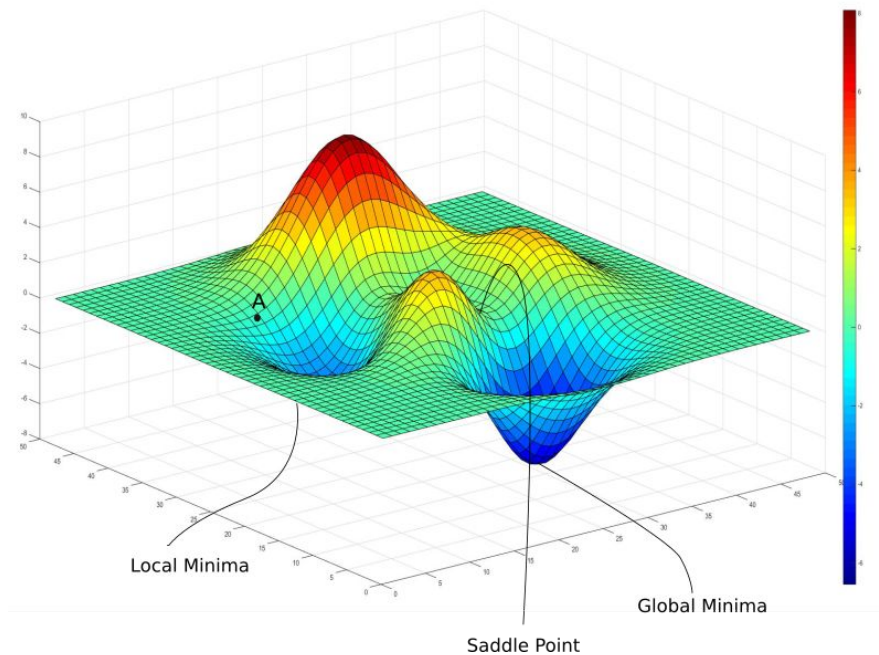
---

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \frac{\partial}{\partial \boldsymbol{\theta}^{(t)}} L(\boldsymbol{\theta}^{(t)}) \quad \Leftrightarrow \quad (\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}) \propto - \frac{\partial L(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}^{(t)}}$$

$$\frac{\partial L(\boldsymbol{\theta}^{(t+1)})}{\partial \eta} = \left[ \frac{\partial L(\boldsymbol{\theta}^{(t+1)})}{\partial \boldsymbol{\theta}^{(t+1)}} \right]^T \frac{\partial \boldsymbol{\theta}^{(t+1)}}{\partial \eta} = - \left[ \frac{\partial L(\boldsymbol{\theta}^{(t+1)})}{\partial \boldsymbol{\theta}^{(t+1)}} \right]^T \frac{\partial L(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}^{(t)}} \stackrel{\text{set}}{=} 0$$

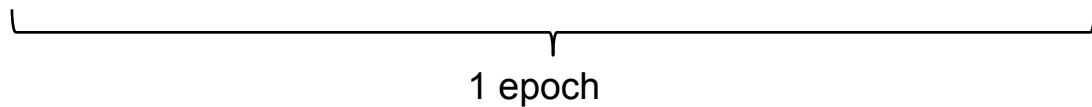
$\Rightarrow \boldsymbol{\theta}^{(t+2)}$  and  $\boldsymbol{\theta}^{(t+1)}$  are orthogonal

# Stochastic Gradient Descent



# Stochastic Gradient Descent

---



- Stochastic (Randomness) → Shuffle the data

# Stochastic Gradient Descent

---

```
▶ from sklearn.linear_model import SGDRegressor

sgd_reg = SGDRegressor(loss='squared_loss',
                       max_iter=50, tol=-np.infty, penalty=None, eta0=0.1, random_state=42)
sgd_reg.fit(X, y.ravel())

sgd_reg.intercept_, sgd_reg.coef_

↳ (array([4.16782089]), array([2.72603052]))
```

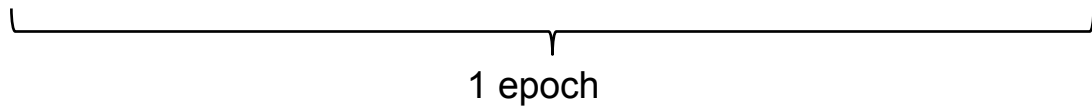
```
from sklearn.linear_model import SGDClassifier

sdg_cls = SGDClassifier(loss='log',
                        max_iter=50, tol=-np.infty, penalty=None, eta0=0.1, random_state=42))
```



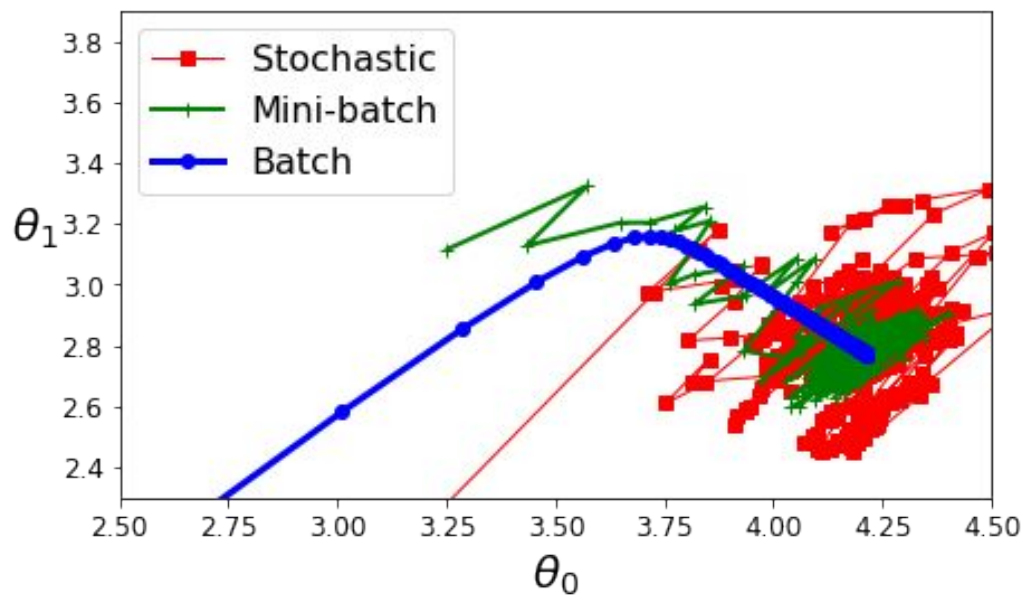
# Mini-Batch Gradient Descent

---



- $k$  batches have  $p$  data  $\rightarrow n = k \times p$

# Mini-Batch Gradient Descent



# Gradient Descent

---

- What about  $\eta^{(t)}$  ?       $\rightarrow$       in Deep Learning

# reference

자료

19-2 STAT424 통계적 머신러닝 - 박유성 교수님

교재

파이썬을 이용한 통계적 머신러닝 (2020) – 박유성

ISLR (2013) - G. James, D. Witten, T. Hastie, R. Tibshirani

The elements of Statistical Learning (2001) - J. Friedman, T. Hastie, R. Tibshirani

Hands on Machine Learning (2017) - Aurelien Geron