

머신러닝 개요

Lecture 5: ML Algorithms

College of Information and Electronic Engineering

Kyung Hee Univeristy

Prof. Wonha Kim

(wonha@khu.ac.kr)

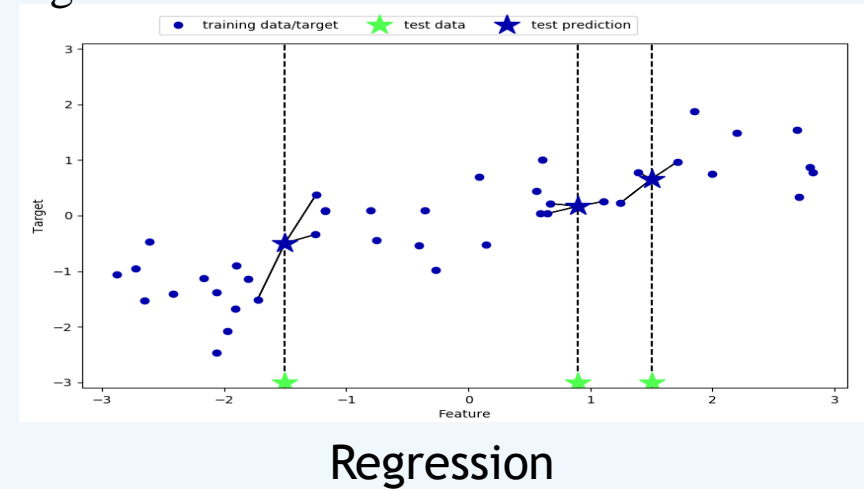
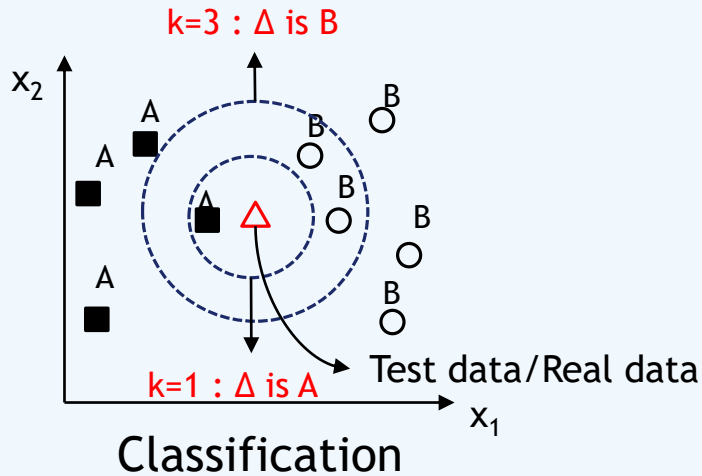
Contents

- Supervised Learning
 1. K-nearest neighbor (k-NN)
 2. Maximum A Posteriori (MAP) Estimation
 3. Regression
- Unsupervised Learning
 4. K-mean



1. k-Nearest Neighbors(k-NN) (1/3) - Supervised Learning

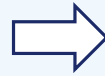
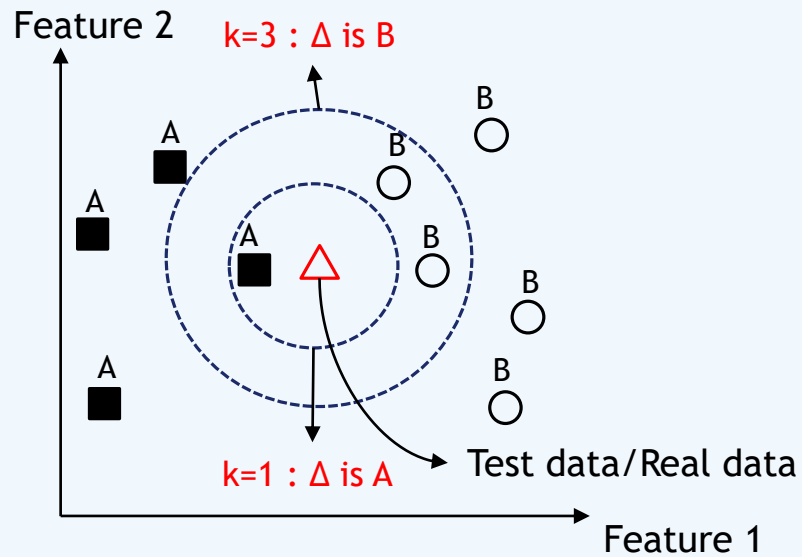
- Training set $D = \{(X_i, y_i)\}_{i=1}^n$, (X_i : Feature vector, y_i : Label/Target)
- Task : Assign the proper label to a query vector X' .
- Calculate distances from X' to each training vector X_i ($i=1, \dots, n$): $d_i(X', X_i) = ||X' - X_i||$
- Select k nearest vectors. k should be odd to prevent ties.
- For classification, assign the majority label of k nearest vectors as the label of X' .
- For regression, assign average target value of k nearest vectors as the target value of X' .



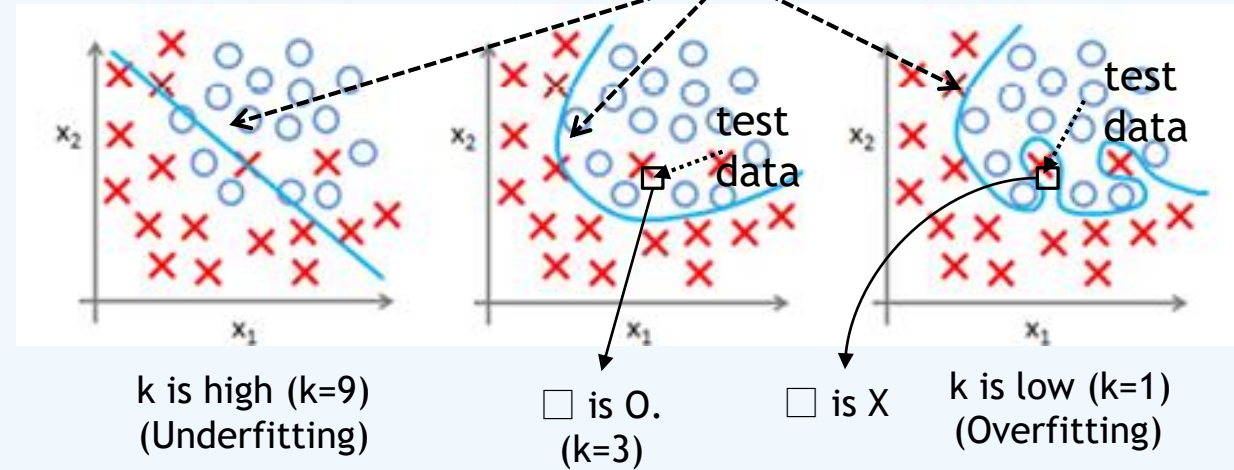
- Weight can be assigned as decreasing with distance so that nearer data contributes more.
So, $d_i^w(X', X_i) = w \cdot ||X' - X_i||$, $w \propto 1/||X' - X_i||$
- Instance-based learning, Non-parametric method. (No use of any probability model)

1. k-Nearest Neighbors (2/3) : Classification/Regression

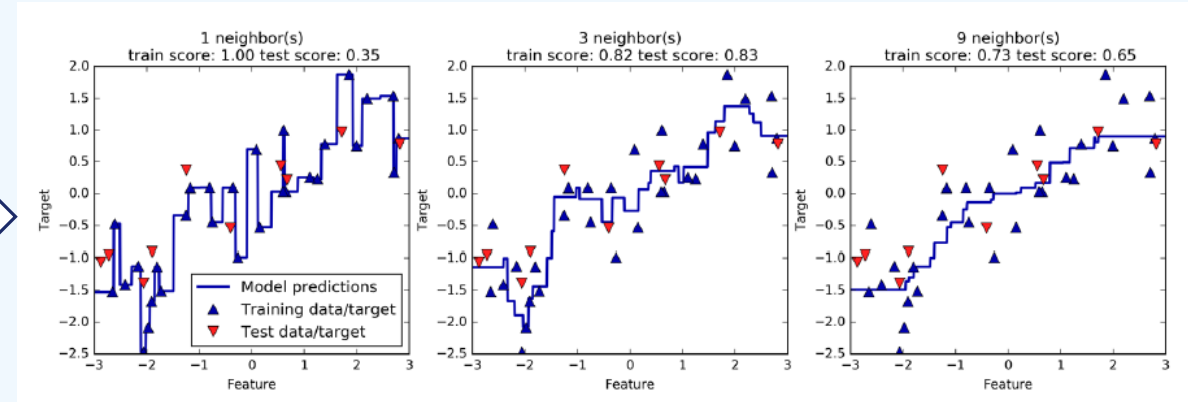
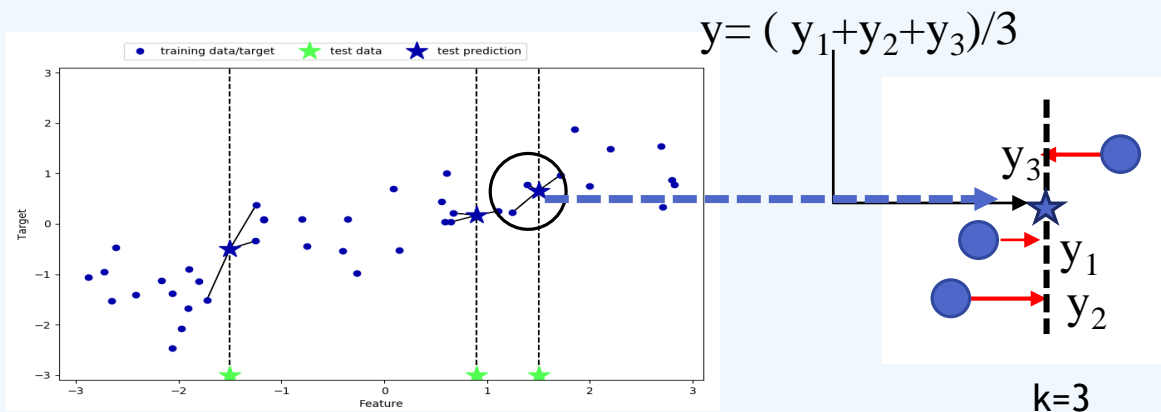
Classification



Boundary from training data



Regression



k=1, Overfitting

k=3

k=9, Underfitting



1. k-NN regression program (3/3)

```
import matplotlib.pyplot as plt
import mglearn
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsRegressor

# 데이터 셋을 만듭니다.
X, y = mglearn.datasets.make_wave(n_samples=40)

plt.plot(X, y, 'o')
plt.ylim(-3, 3)
plt.xlabel("Feature")
plt.ylabel("Target")
plt.show()

# k = 1 일 때의 예측
mglearn.plots.plot_knn_regression(n_neighbors=1)
plt.show()

# k = 3 일 때의 예측
mglearn.plots.plot_knn_regression(n_neighbors=3)
plt.show()

# wave 데이터 셋을 훈련 세트와 테스트 세트로 나눕니다.
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

# 이웃의 수를 3으로 하여 모델의 객체를 만듭니다.
reg = KNeighborsRegressor(n_neighbors=3)

# Training data와 Target을 사용하여 Model을 학습시킵니다.
reg.fit(X_train, y_train)

print("테스트 세트 예측: {}".format(reg.predict(X_test)))

print("테스트 세트 R^2: {:.2f}".format(reg.score(X_test, y_test)))
```

2. Maximum A Posteriori (MAP) Estimation (1/9) : Bayesian Rule

$$\begin{aligned} p(B|A)p(A) &= p(A \cap B) \\ p(A|B)p(B) &= p(A \cap B) \end{aligned} \Rightarrow p(A|B)p(B) = p(B|A)p(A) \Rightarrow p(A|B) = \frac{p(B|A)p(A)}{p(B)} = \frac{p(B|A)p(A)}{\sum_{\alpha \in A} p(B|A)p(A)}$$

$$P(y, x) = P(x|y)P(y) = P(x, y) = P(y|x)P(x) \Rightarrow P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

$$\bullet \quad P(y|x) = \frac{P(x|y)P(y)}{P(x)} \quad \Longleftrightarrow \quad \text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

- Posteriori Prob. $P(y | x)$ 의 추정은 대부분의 경우 x 의 공간 너무 크기 (너무 많기) 때문에 거의 불가능
- Posteriori Prob. $P(y | x)$ 는 대부분은 Bayesian Theory를 이용하여 추정함.
- Priori Prob. $P(y)$ 는 data로 부터 측정하여 구함. 즉, $P(y=i) = N_i/N = i$ 의 개수/전체 샘플 개수
- Likelihood $P(x|y)$ ($L(y|x)$)는 data의 observation를 통해서 구함. (Density Estimation method)

$$\bullet \quad \text{MAP : } x \text{가 발생 했을 때, 가장 발생 가능한 } y \text{를 추정.} \quad \Longleftrightarrow \quad y^* = \underset{y}{\operatorname{argmax}} P(y|x)$$



2. Maximum A Posteriori (MAP) (2/9): Toy Example (1/2)

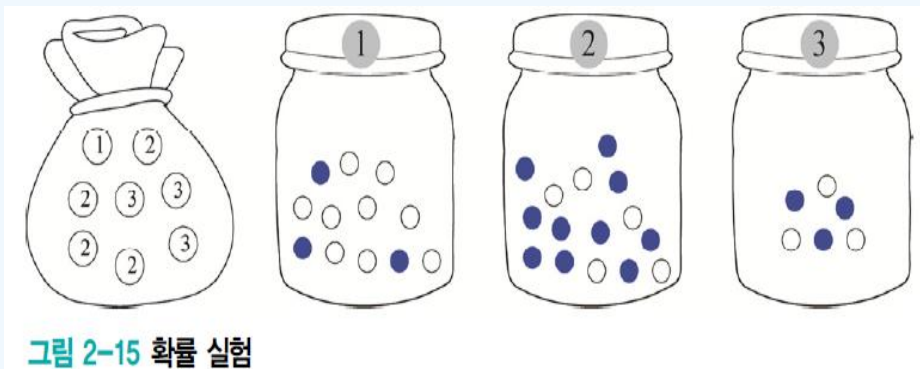


그림 2-15 확률 실험

- 주머니에서 번호를 뽑은 다음, 번호에 따라 해당 병에서 공을 뽑고 색을 관찰함
- 꺼낸 카드와 공은 꺼낸 곳에 다시 넣음.
- 번호를 y , 공의 색을 x 라는 확률변수로 표현하면 정의역은

$$y \in \{①, ②, ③\}, x \in \{\text{파랑, 하양}\}$$

- ①번 카드를 뽑을 확률은 $P(y=①)=P(①)=1/8$
- 카드는 ①번, 공은 하양일 확률은 $P(y=①, x=\text{하양})$

$$P(y=①, x=\text{하양}) = P(x=\text{하양}|y=①)P(y=①) = \frac{9}{12} \frac{1}{8} = \frac{3}{32}$$

- 하얀 공이 뽑힐 확률

$$\begin{aligned} P(\text{하양}) &= P(\text{하양}|①)P(①) + P(\text{하양}|②)P(②) + P(\text{하양}|③)P(③) \\ &= \frac{9}{12} \frac{1}{8} + \frac{5}{15} \frac{4}{8} + \frac{3}{6} \frac{3}{8} = \frac{43}{96} \end{aligned}$$

- 합 규칙 (Marginalization)

$$P(x) = \sum_y P(y, x) = \sum_y P(x|y)P(y)$$

2. Maximum A Posteriori (MAP) (3/9): Toy Example (2/2)

- 하얀 공이 나왔을 때 ($x=\text{하양}$), 어느 병에서 나왔는지 추정하시오. (가장 발생 가능한 $y=?$)

$$\hat{y} = \operatorname{argmax}_y P(y|x = \text{하양}) = \operatorname{argmax}_y \frac{P(x = \text{하양}|y)P(y)}{P(x = \text{하양})}$$

$$P(\text{①}|\text{하양}) = \frac{P(\text{하양}|\text{①})P(\text{①})}{P(\text{하양})} = \frac{\frac{9}{12} \cdot \frac{1}{8}}{\frac{43}{96}} = \frac{9}{43}$$

$$P(\text{②}|\text{하양}) = \frac{P(\text{하양}|\text{②})P(\text{②})}{P(\text{하양})} = \frac{\frac{5}{15} \cdot \frac{4}{8}}{\frac{43}{96}} = \frac{16}{43}$$

$$P(\text{③}|\text{하양}) = \frac{P(\text{하양}|\text{③})P(\text{③})}{P(\text{하양})} = \frac{\frac{3}{6} \cdot \frac{3}{8}}{\frac{43}{96}} = \frac{18}{43}$$

⇒ ③번 병일 확률이 가장 높음

- Bayesian rule의 Decision에 적용

$$\begin{array}{ccc} \text{Posteriori} & \text{Likelihood} & \text{Priori} \\ \text{Prob.} & & \text{Prob.} \\ \overbrace{P(y|x)} & = \frac{\overbrace{P(x|y)} \overbrace{P(y)}}{\underbrace{P(x)}} & \\ & \text{Evidence} & \end{array}$$



2. Maximum A Posteriori (MAP) (4/9) : Estimation



- Jang, Hana (A pro golfer)
 - **Evidence** : We *observe* that 1 of 10 players makes par at a hole. ($P(x)$): $P(\text{par}) = 1/10$.
 - **Prior** : Among 20 players, Jang hits one time. The prob. that Jan hits. ($P(y)$): $P(\text{Jang})=1/20$.
 - **Likelihood** : Data shows that Jan Hana makes 6 pars at 7 hits. ($P(x/y)$) : $P(\text{par}/\text{Jang}) = 6/7$

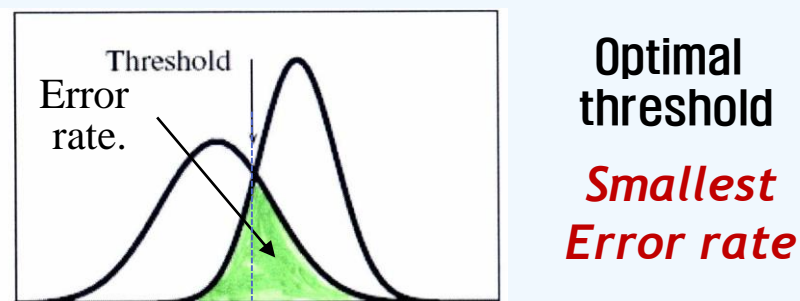
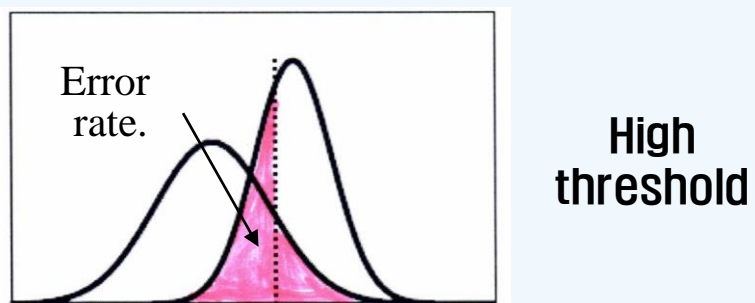
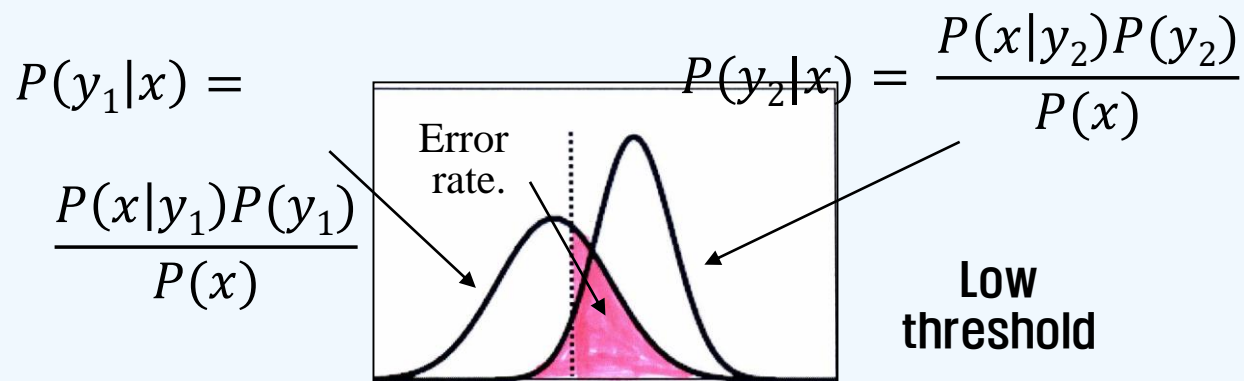


- Prof. Kim, Wonha
 - **Evidence** : We *observe* that 1 of 10 players make par at a hole. : $P(\text{par}) = 1/10$.
 - **Prior** : Among 20 players, Prof. Kim played 7 times. The prob. of Kim's hitting. : $P(\text{Kim})=7/20$.
 - **Likelihood** : *Data* shows that Kim makes the one par at 6 hits. : $P(\text{par}/\text{Kim}) = 1/6$.

- We know that either Kim or Jang makes par. Estimate who is he/she?
 - Prob. of Kim's hitting : $P(\text{Kim}/\text{Par}) = P(\text{Par}/\text{Kim}) * P(\text{Kim})/P(\text{Par}) = (1/6)*(7/20)/(1/10) = 7/12$.
 - Prob. of Jang's hitting : $P(\text{Jang}/\text{Par}) = P(\text{Par}/\text{Jang}) * P(\text{Jang})/P(\text{Par}) = (6/7)*(1/20)/(1/10) = 6/14$.
 - $\Rightarrow P(\text{Par}/\text{Jang}) > P(\text{Par}/\text{Kim})$: *Definitely, Jang plays much better than Prof. Kim.*
 - $\Rightarrow P(\text{Kim}) > P(\text{Jang})$: *Unfair. Prof. Kim plays 7 times more than Jang.*
 - $\Rightarrow P(\text{Kim}/\text{Par}) > P(\text{Jang}/\text{Par})$: *We must estimate that Prof. Kim makes the par.*

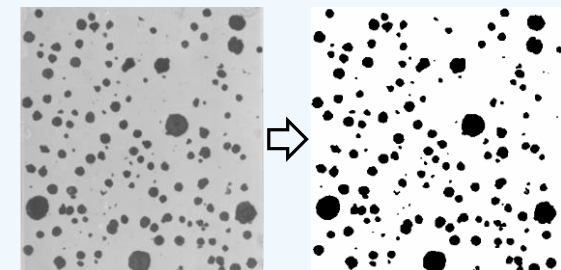
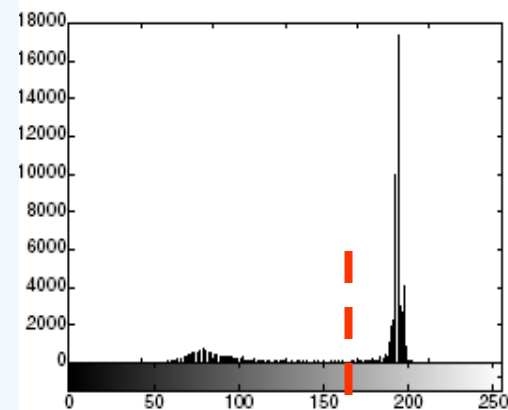
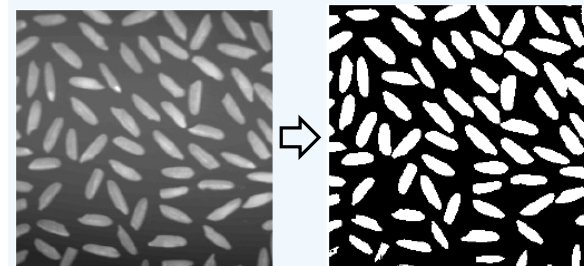
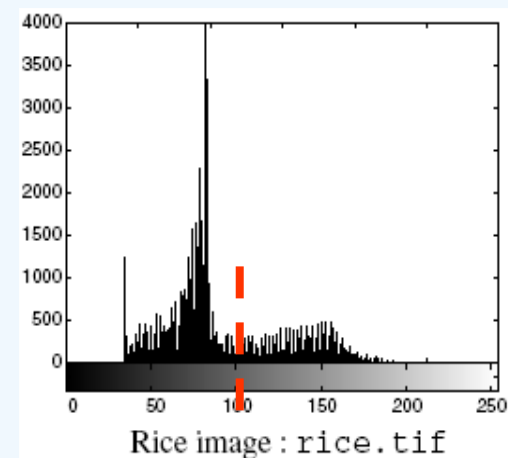
2. Maximum A Posteriori (MAP) (5/9) : Segmentation

- Optimal Decision (Classification)



- EX: Segmentation

- 두 물체의 밝기에 따른 확률 분포
- 두 물체를 구분하는 최적의 밝기 threshold 설정하기



2. MAP Estimation (6/9): Bayes Decision Rule (BDR)

- Bayesian Decision for two classes.

- For observed data X , there are two classes, i.e, $y=0$, or $y=1$.
- We decide y^* in following way :

$$y^*(x) = \begin{cases} 0, & \text{if } P_{y|X}(0|X) \geq P_{y|X}(1|X) \\ 1, & \text{if } P_{y|X}(1|X) \geq P_{y|X}(0|X) \end{cases} \equiv \underset{y}{\operatorname{argmax}} P_{y|X}(y|X) = \underset{y}{\operatorname{argmax}} P_{X|y}(X|y) P_y(y)$$

- Bayesian Decision for n classes.

- For observed data X , there are n classes, i.e, $y=1, \dots, n$.
- $y^*=i$ if $P_{y|X}(y = i|X) > P_{y|X}(y = j|X)$ for $j \neq i$

$$\Rightarrow y^* = \underset{y}{\operatorname{argmax}} P_{y|X}(y|X) = \underset{y}{\operatorname{argmax}} P_{X|y}(X|y) P_y(y)$$



2. MAP Estimation (7/9) :Decision

- Maximum Likelihood Estimation (MLE)

$$y^*(x) = \operatorname{argmax}_y P_{X|y}(x|y)$$

- Parameter w is known and the optimal w should be founded .

- Maximum A Posteriori (MAP) Estimation .

$$y^*(x) = \operatorname{argmax}_y P_{y|X}(y|X) \quad : \text{Not much used}$$

$$y^*(x) = \operatorname{argmax}_y P_{X|y}(X|y) P_y(y)$$

$$y^*(x) = \operatorname{argmax}_y [\log(P_{X|y}(X|y)) + \log(P_y(y))] : \text{Most frequently used}$$

- Data X are observed.
- Parameters y (or prior knowledge) are unknown and probabilistically presumed.
- If $P_y(w)$ is uniform, MLE =MAP.



2 MAP Estimation (8/9) : Classification for Gaussian

- MAP for Multivariate Gaussian Classifier (Model based learning)

- The pdf of class i is a N-dimensional Gaussian with mean μ_i and covariance

$$P_{x|y}(X|y = i) = \frac{1}{\sqrt{(2\pi)^N |\Sigma_i|}} \exp\left\{-\frac{1}{2} (X - \mu_i)^T \Sigma_i^{-1} (X - \mu_i)\right\}$$

-MAP Classifier

$$i^*(x) = \underset{i}{\operatorname{argmax}} \left[-\frac{1}{2} (X - \mu_i)^T \Sigma_i^{-1} (X - \mu_i) - \frac{1}{2} \log\{(2\pi)^N |\Sigma_i|\} + \log P_y(y = i) \right]$$



- Compare with K-NN.
 - K-NN : Non-model
 - Low complexity

- Design one dimensional Gaussian K-classifiers.

-Collect i^{th} class dataset : $D^i = \{x_1^i, x_2^i, \dots, x_{n^i}^i\}$ (data size for i^{th} class is n^i .) where $i=1, \dots, k$.

-Estimate the Gaussian parameters:

$$\hat{u}_i = \frac{1}{n^i} \sum_{j=1}^{n^i} x_j^i, \quad \hat{\Sigma}_i = \frac{1}{n^i} \sum_{j=1}^{n^i} (x_j^i - \hat{u}_i)^2, \quad \hat{P}_w(y=i) = \frac{n^i}{S} \quad \text{where } S \text{ is the total sample size.}$$

-MAP Classifier

$$i^*(x) = \underset{y}{\operatorname{argmax}} \left[-\frac{1}{2} (x - \hat{u}_i)^T \hat{\Sigma}_i^{-1} (x - \hat{u}_i) - \frac{1}{2} \log\{2\pi |\hat{\Sigma}_i|\} + \log \hat{P}_y(y = i) \right]$$



2. Maximum A Posteriori (MAP) (9/9): Training

- Data x 가 주어졌을 때, x 를 발생 시킬 확률을 최대로 하는 매개변수 θ 를 추정하는 문제

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\mathbb{X}|\theta)$$

- 수학적 계산 편의를 위해 Log 함수를 사용 할 수 있음.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(\mathbb{X}|\theta)$$

- $\mathbb{X} = [x_1, x_2, \dots, x_n]$ 에서 x_i ($i=1, \dots, n$) 가 independent and identical distribution (iid) 일 때 .

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(\mathbb{X}|\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P(x_i|\theta)$$

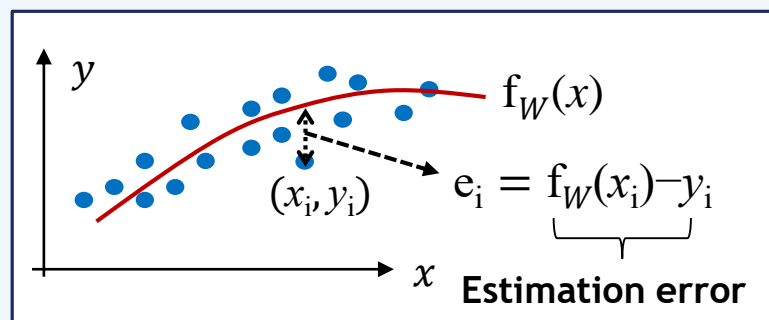
3 Machine Learning Model (1/7): Regression

- Feature Data set : $\mathbf{X} = [x_1, x_2 \cdots x_n]^T$, Target value set : $\mathbf{Y} = [y_1, y_2 \cdots y_n]^T$
- Polynomial predicting the target value : $f_w(x) \Rightarrow$ For an input feature x_i , $f_w(x_i)$ predicts its target get value y_i .

$$f_W(x_i) = w_0 + w_1x_i + w_2x_i^2 \cdots + w_px_i^p = \sum_{k=0}^p w_k x_i^k = \mathbf{x}_i^T \cdot \mathbf{W} \text{ where } \mathbf{x}_i = [1, x_i, x_i^2, \cdots, x_i^p]^T$$

$$\mathbf{W} = [w_0, w_1, w_2, \cdots, w_p]^T \quad \text{Bias}$$

- Objective (Cost) function :



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n |e_i|^2 = \frac{1}{n} \sum_{i=1}^n (f_W(x_i) - y_i)^2$$

(Mean Square Error (MSE))

$$= (\mathbf{X}_n \cdot \mathbf{W} - \mathbf{Y})^T (\mathbf{X}_n \cdot \mathbf{W} - \mathbf{Y})$$

where $\mathbf{X}_n = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} = \begin{bmatrix} 1, x_1, x_1^2, \cdots, x_1^p \\ \vdots \\ 1, x_n, x_n^2, \cdots, x_n^p \end{bmatrix}$

- ML Training : Decide parameters so that the parameter set \mathbf{W} minimizes MSE $J(\mathbf{W})$.

$$\hat{\mathbf{W}} = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$



$$J(\mathbf{W}) = (\mathbf{X}_n \cdot \mathbf{W} - \mathbf{Y})^T (\mathbf{X}_n \cdot \mathbf{W} - \mathbf{Y}) = \mathbf{W}^T \mathbf{X}_n^T \mathbf{X}_n \mathbf{W} - 2(\mathbf{X}_n \mathbf{W})^T \mathbf{Y} + \mathbf{Y}^T \mathbf{Y}$$

$$\frac{1}{\partial \mathbf{W}} J(\mathbf{W}) = 2\mathbf{X}_n^T \mathbf{X}_n \mathbf{W} - 2\mathbf{X}_n^T \mathbf{Y} = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{W}} = (\mathbf{X}_n^T \cdot \mathbf{X}_n)^{-1} \cdot \mathbf{X}_n^T \cdot \mathbf{Y}$$

4. Regression by curve fitting (2/7): Exercise

- Feature vectors and target values are

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 2 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 2 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

* Note that feature values are {1,1,2,2,3}.

- Linear Regression Model is $y = b + ax$.
- Optimal weights.

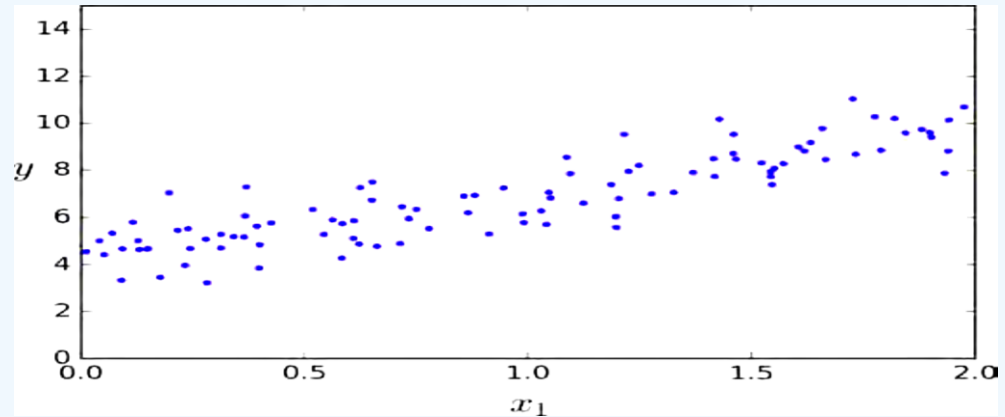
$$[b \ a]^T = (X_n^T \cdot X_n)^{-1} \cdot X_n^T \cdot y = [-0.5 \ 1.5]^T$$

$$y = -0.5 + 1.5x$$

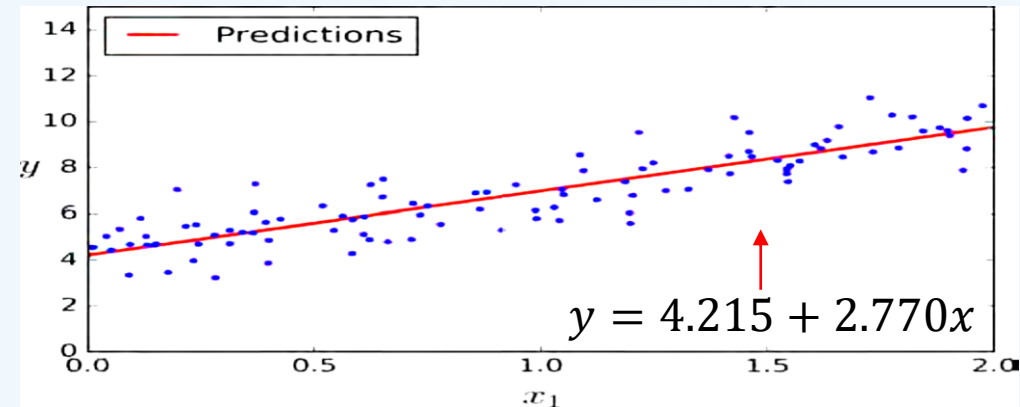
- Random data.

$$x = 2 \times \text{random_noise}$$

$$y = 4 + 3 \times x + \text{random_noise}$$



$$[b \ a]^T = (X_n^T \cdot X_n)^{-1} \cdot X_n^T \cdot y = [4.215 \ 2.770]^T$$



4. Regression by curve fitting (3/7) : Polynomial Regression

- Multi dimension / Polynomial Regression Model

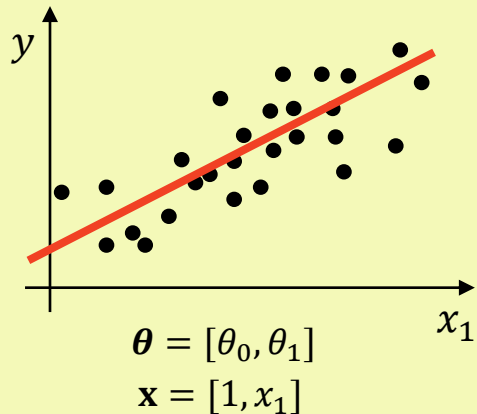
- d-features: Training Set $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ where $\mathbf{x}_i = (\overset{\text{Bias term}}{\underset{\text{Bias term}}{1}} x_{i1}, x_{i2}, \dots, x_{id})^T$

$$\hat{y} = \theta_0 + \theta_1 x_{11} + \theta_2 x_{21}^2 + \dots + \theta_p x_{n1}^p + \theta_{p+1} x_{12} + \dots + \theta_{2p} x_{n2}^p + \dots + \theta_{(d-1)p} x_{1d} + \dots + \theta_{dp} x_{nd}^p$$

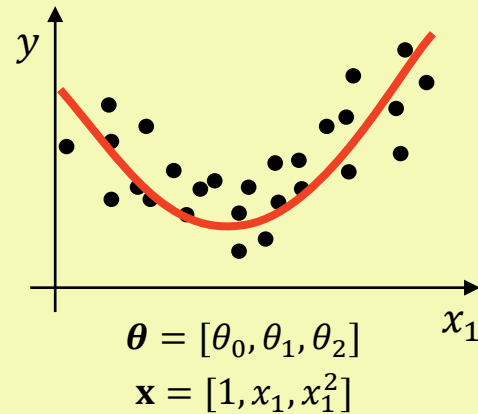
$$= \begin{bmatrix} 1 & x_{11} & x_{21}^2 & \dots & x_{n1}^p & x_{21} & x_{21}^2 & \dots & x_{2n}^p & \dots & x_{1d} & \dots & x_{nd}^p \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \vdots \\ \theta_{dp} \end{bmatrix} = \underline{J_{\theta}(\mathbf{x}) = \theta^T \cdot \mathbf{x}}$$

=> Optimal solution also is same with one-dimension model.

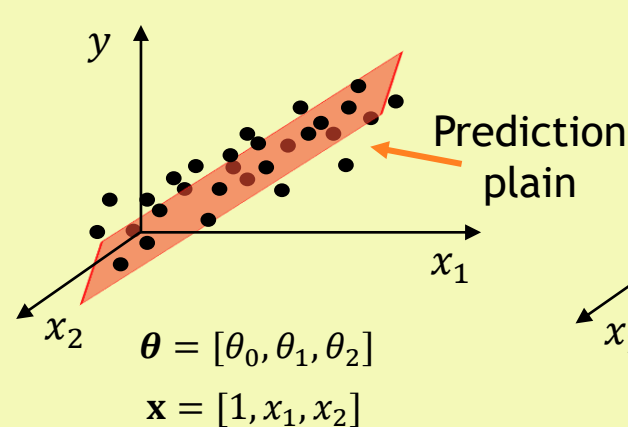
• Single feature
/ Linear



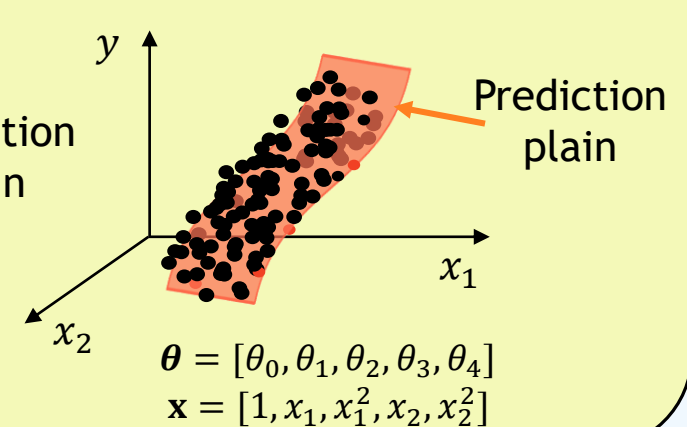
• Single feature
/ Non-Linear



• Multiple features
/ Linear

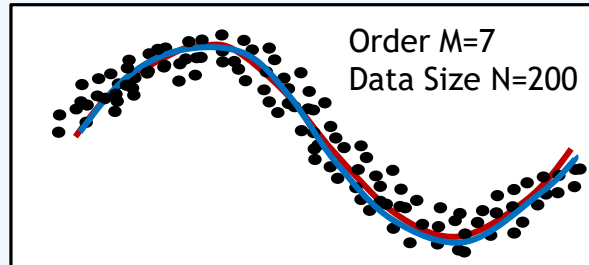
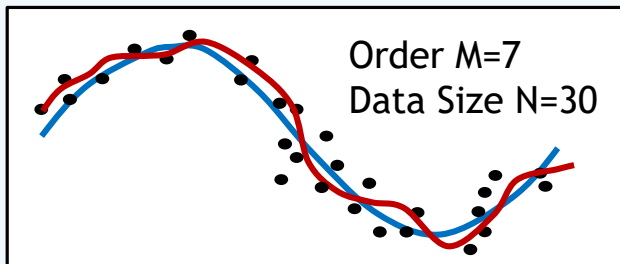


• Multiple features
/ Non-Linear



4. Regression(4/7) : Regularization

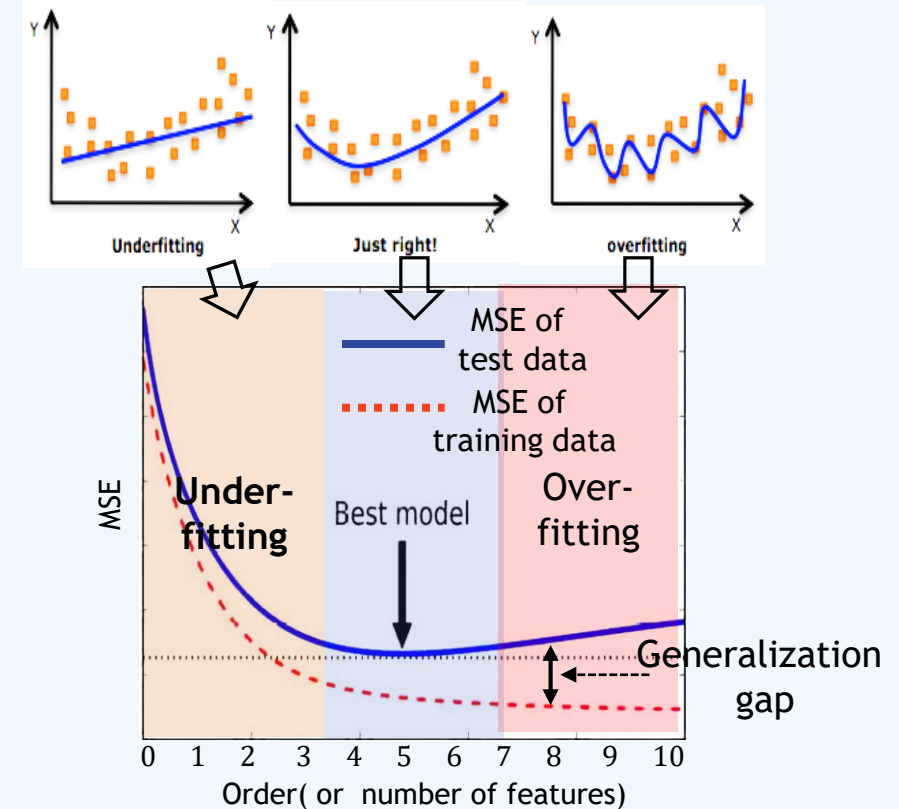
- Low order or Small features : Underfitting
 - A regression model with lower order (or small features) cannot fit to the training and test data.
 - Biased error
- Higher Order or Many features : Overfitting
 - A regression model with higher order (or many features) can reduce MSE for training data, but increase MSE for test data
 - Variance error
- The more data generally produce the more appropriate regression estimator.



However, increasing data is difficult in real applications



Regularization



4. Regression(5/7) : Regularization

• Tikhonov Regularization

$$\underbrace{J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{R(\Theta)}_{\text{규제 항}}$$

• L2 norm regulation :

- For online learning

$$J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y}) = J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \|\Theta\|_2^2$$

$$\nabla J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y}) = \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta$$

$$\Theta = \Theta - \rho \nabla J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y})$$

$$= \Theta - \rho(\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta)$$

$$= (1 - 2\rho\lambda)\Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y})$$



$$\Theta = (1 - 2\rho\lambda)\Theta - \rho \nabla J$$

• Smooth function : degree of inverse Θ smooth

- 차수가 높을 수록 커져서 penalty가 강화됨
- 즉, 4 함수는 12차 함수 보다 smooth 정도가 적음
- Data와 무관하고 model의 형태에 따라서 다름

- For batch learning

$$J_{regularized}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2$$

ox문제 나옴

도중에 추가하면 절대 안됨

$$\frac{\partial J_{regularized}}{\partial \mathbf{w}} = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0}$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

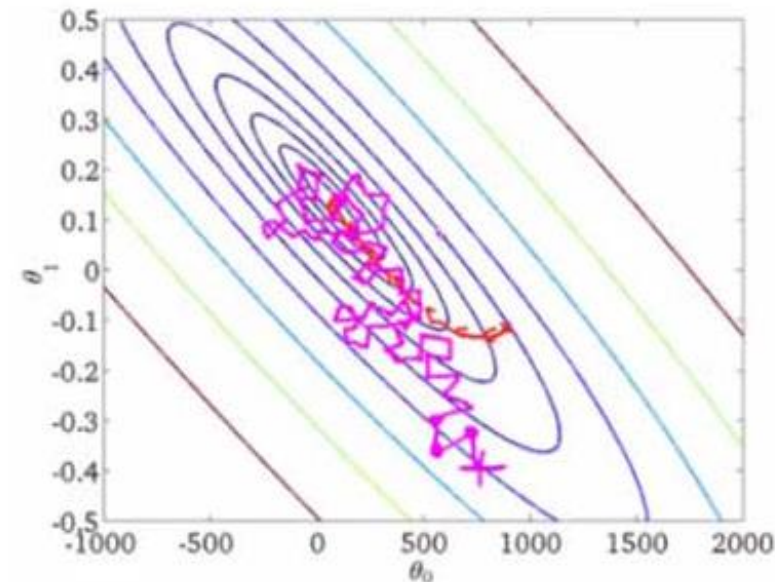
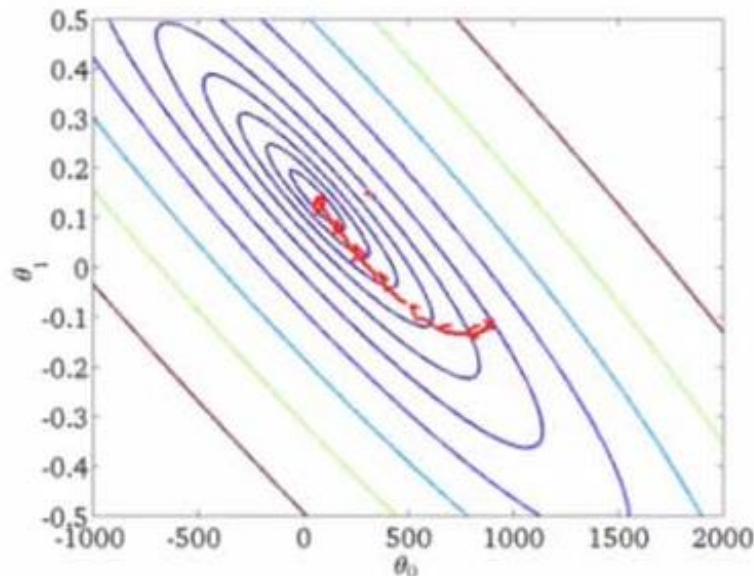


$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$



4. Regression (6/7): On-line vs Batch regression

Batch Gradient Descent	Stochastic Gradient Descent
<ul style="list-style-type: none">• Gently converges to the (local) minimum• Very slow• Intractable for datasets that don't fit in memory• No online learning	<ul style="list-style-type: none">• Faster• Online learning• Heavy fluctuation• Capability to jump to new (potentially better local minima)• Complicated convergence (overshooting)



4. Regression(7/7) : Example

예제 5-1 리지 회귀

훈련집합 $\mathbb{X} = \{\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 3 \\ 3 \end{pmatrix}\}$, $\mathbb{Y} = \{y_1 = 3.0, y_2 = 7.0, y_3 = 8.8\}$ 이 주어졌다고 가정하자. 특징 벡터가 2차원이므로 $d=2$ 이고 샘플이 3개이므로 $n=3$ 이다. 훈련집합으로 설계행렬 \mathbf{X} 와 레이블 행렬 \mathbf{y} 를 다음과 같이 쓸 수 있다.

$$\mathbf{X} = \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix}$$

이 값들을 식 (5.29)에 대입하여 다음과 같이 $\hat{\mathbf{w}}$ 을 구할 수 있다. 이때 $\lambda = 0.25$ 라 가정하자.

$$\hat{\mathbf{w}} = \left(\begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix} + \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix} = \begin{pmatrix} 1.4916 \\ 1.3607 \end{pmatrix}$$

따라서 하이퍼 평면은 $y = 1.4916x_1 + 1.3607x_2$ 이다. 새로운 샘플로 $\mathbf{x} = (5 \ 4)^T$ 가 입력되면 식 (5.30)을 이용하여 12.9009를 예측한다.



4. K-mean (1/3) : Unsupervised Learning

- For d-features (dimension) and n- observation (or data) set $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$, partition the n observation into $k(\leq n)$ sets $\mathbf{S} = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares(WCSS) (Wikipedia)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 \quad \boldsymbol{\mu}_i : \text{ is the mean of elements in } S_i.$$

- Algorithms

- : Give an initial set of k means $\mathbf{u}_1^{(1)}, \dots, \mathbf{u}_k^{(1)}$

- : Do iteration until $\|\mathbf{u}_i^{(t+1)} - \mathbf{u}_i^{(t)}\| > \varepsilon$

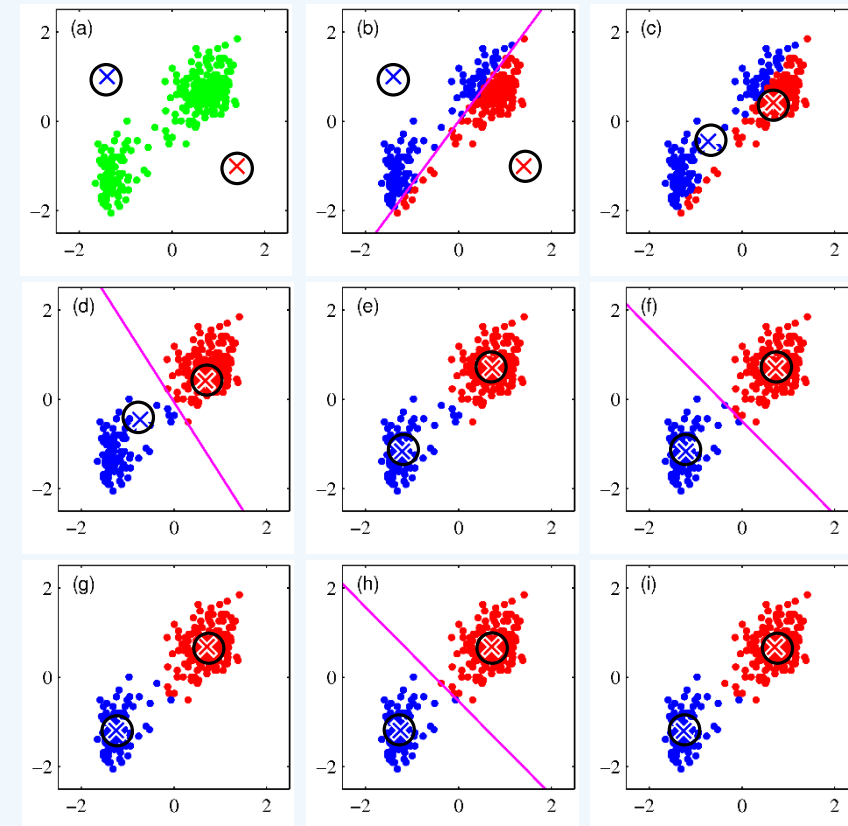
- (1) Assign each observation to the cluster whose mean has the least distance.

$$i^* = \underset{1 \leq i \leq k}{\operatorname{argmin}} \|\mathbf{x}_j - \mathbf{u}_i\| \Rightarrow \mathbf{x}_j \rightarrow S_{i^*} \quad (j = 1, \dots, N)$$

- (2) Update the means to be the centroids of observations in the new clusters

$$\mathbf{u}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j \quad |S_i^{(t)}| \text{ is the number of data in class } S_i^{(t)}$$

K(=2)-mean



4. K-mean Clustering (2/3) : Constructing clusters

- Mean Splitting (for deciding proper k).

1. Set $k=1$, calculate $u^{(i)}$.
2. Splitting the mean for $k=2$ by perturbing randomly.

$$u_1^{(2)} \leftarrow u^{(1)}, \quad u_2^{(2)} \leftarrow (1 + \epsilon)u^{(1)} \quad (\epsilon \ll 1)$$

3. Run $k=2$ means.
4. Deleting empty clusters
 - Check the number of elements in each cluster.
 - Delete the cluster having too small elements
 - Split the cluster having the most elements.

4. Splitting the means for $k=4$.

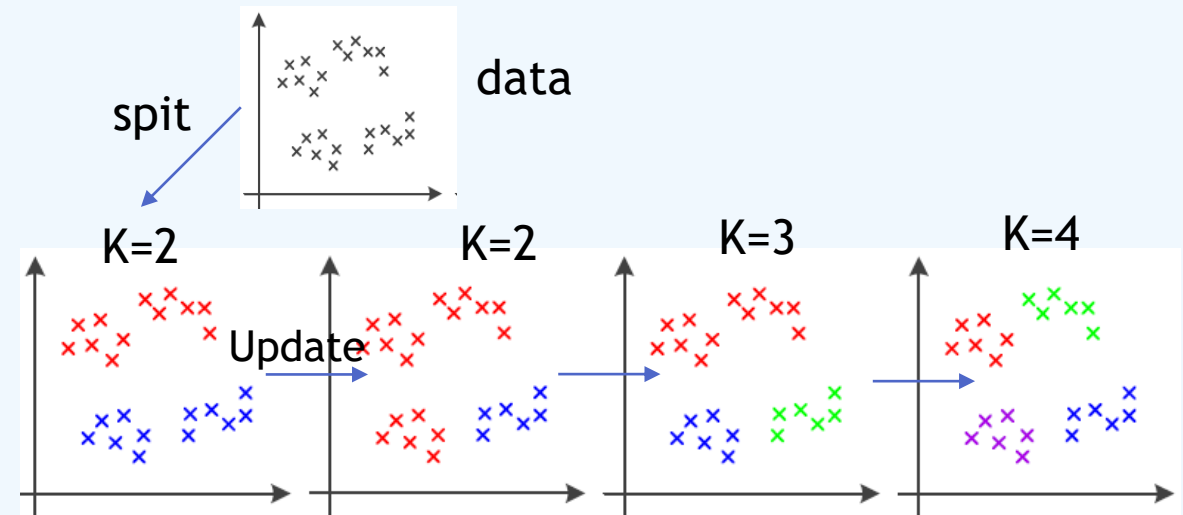
$$u_1^{(4)} \leftarrow u_1^{(2)}, \quad u_2^{(4)} \leftarrow (1 + \epsilon)u_1^{(2)}$$

$$u_3^{(4)} \leftarrow u_2^{(2)}, \quad u_4^{(4)} \leftarrow (1 + \epsilon)u_2^{(2)}$$

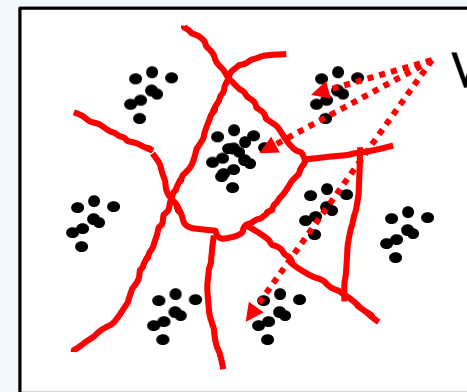
5. Repeat splitting until finding the proper k .

(* Random initialization of means are not good.)

- Clustering with $k=4$



- Voronoi Cell



Voronoi Cells: $S = \{S_1, S_2, \dots, S_k\}$

4. K-mean Clustering (3/3) : Application

- Image segmentation
 - The Image Data form is $I(x,y) = [R, G, B]$.
 - K-mean clustering let neighbored and similar color pixels be a set.

