

# Chapter 5

# Machine Learning Basics

MILab Undergraduate student, TaeHyeon Kim

2023. 03. 07



# Overview

---

## Chapter 5. Machine Learning Basics

- 5.1 Learning Algorithms
- 5.2 Capacity, Overfitting and Underfitting
- 5.3 Hyperparameters and Validation Sets
- 5.4 Estimators, Bias and Variance
- 5.5 Maximum Likelihood Estimation
- 5.6 Bayesian Statistics
- 5.7 Supervised Learning Algorithms
- 5.8 Unsupervised Learning Algorithms
- 5.9 Stochastic Gradient Descent
- 5.10 Building a Machine Learning Algorithm
- 5.11 Challenges Motivating Deep Learning

# 5.1 Learning Algorithms

---

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”

Text 1 – Description of Machine Learning Algorithm

- Experience  $E$  : The type of input data (supervised or not)
- Task  $T$  : The types of training algorithm
  - Classification, Classification with missing inputs, Regression, Transcription, Machine translation, Structured output, Anomaly detection, Synthesis and sampling, Imputation of missing values, Denoising, Density estimation or probability mass function estimation
- Performance  $P$  : The efficiency of learning (Accuracy, Error rate)

# 5.1 Learning Algorithms (Example)

---

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”

Text 1 – Description of Machine Learning Algorithm

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

- Experience  $E$  : Supervised Learning ( $\mathbf{x}$  and  $\mathbf{y}$  is given.)
- Task  $T$  : Linear Regression
- Performance  $P$  : Mean Squared Error

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})})_i^2.$$

# 5.1 Learning Algorithms (Example)

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”

Text 1 – Description of Machine Learning Algorithm

$$\hat{y} = \mathbf{w}^\top \mathbf{x}$$

You can also get the  $\mathbf{w}$  value from the calculating ‘Mean Squared Error’

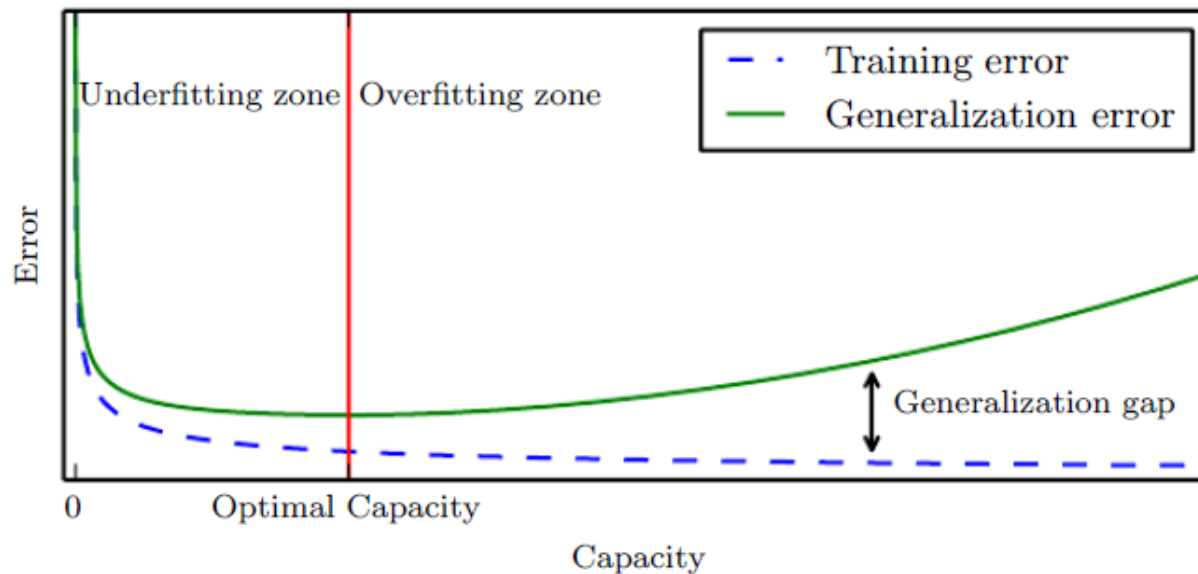
$$\begin{aligned} \nabla_{\mathbf{w}} \text{MSE}_{\text{train}} &= 0 \\ \Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 &= 0 \end{aligned} \quad \left| \quad \begin{aligned} &\Rightarrow \frac{1}{m} \nabla_{\mathbf{w}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0 & (5.8) \\ &\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0 & (5.9) \\ &\Rightarrow \nabla_{\mathbf{w}} \left( \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0 & (5.10) \\ &\Rightarrow 2\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0 & (5.11) \\ &\Rightarrow \mathbf{w} = \left( \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} & (5.12) \end{aligned}$$

## 5.2 Capacity, Overfitting and Underfitting

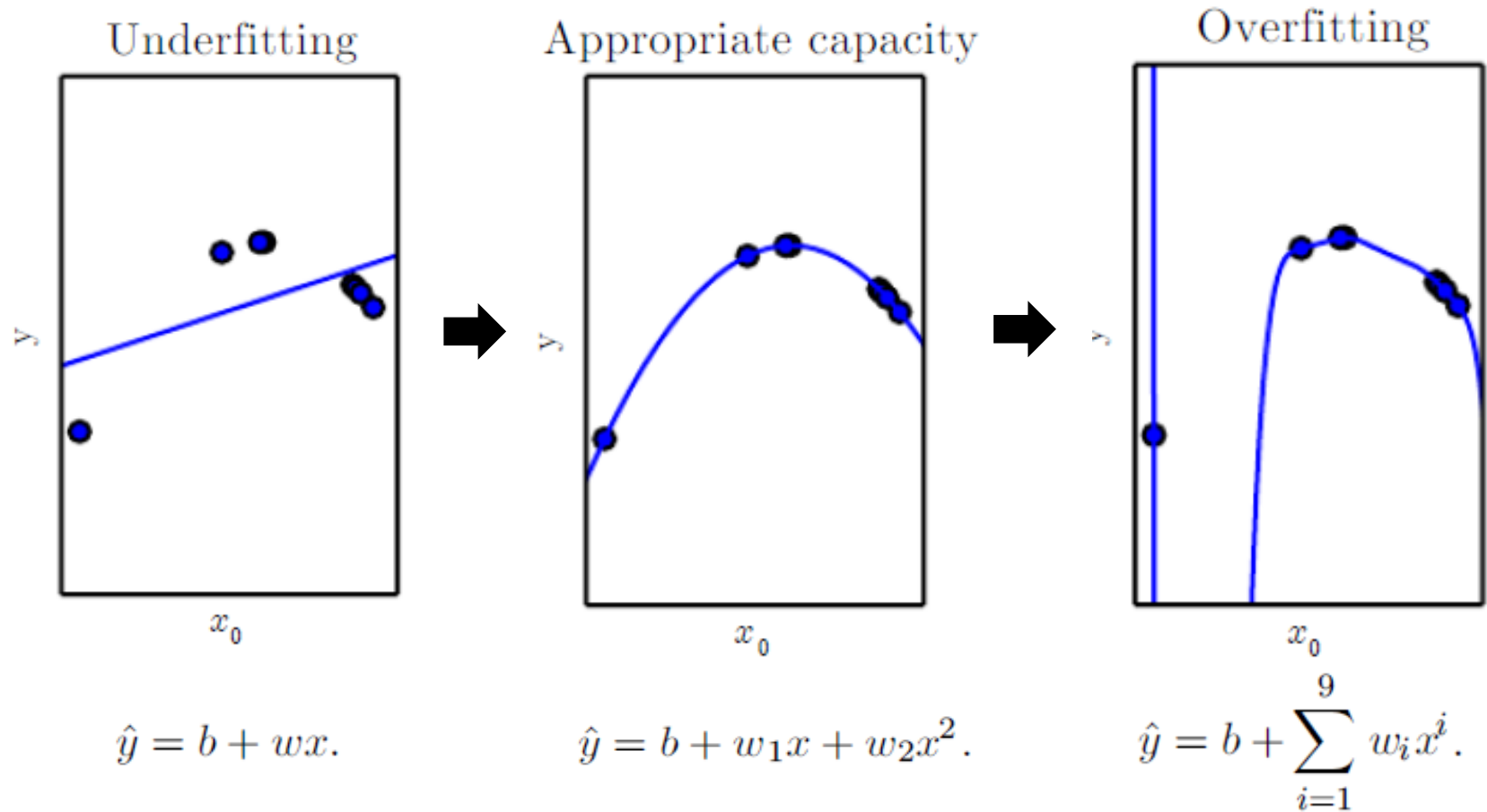
- Capacity : How could treat many variables?
- Underfitting : Training Error  $\uparrow$
- Overfitting : Generalization Error  $\uparrow$

### < How to Solve >

1. Make the training error small.
2. Make the gap between training and test error small

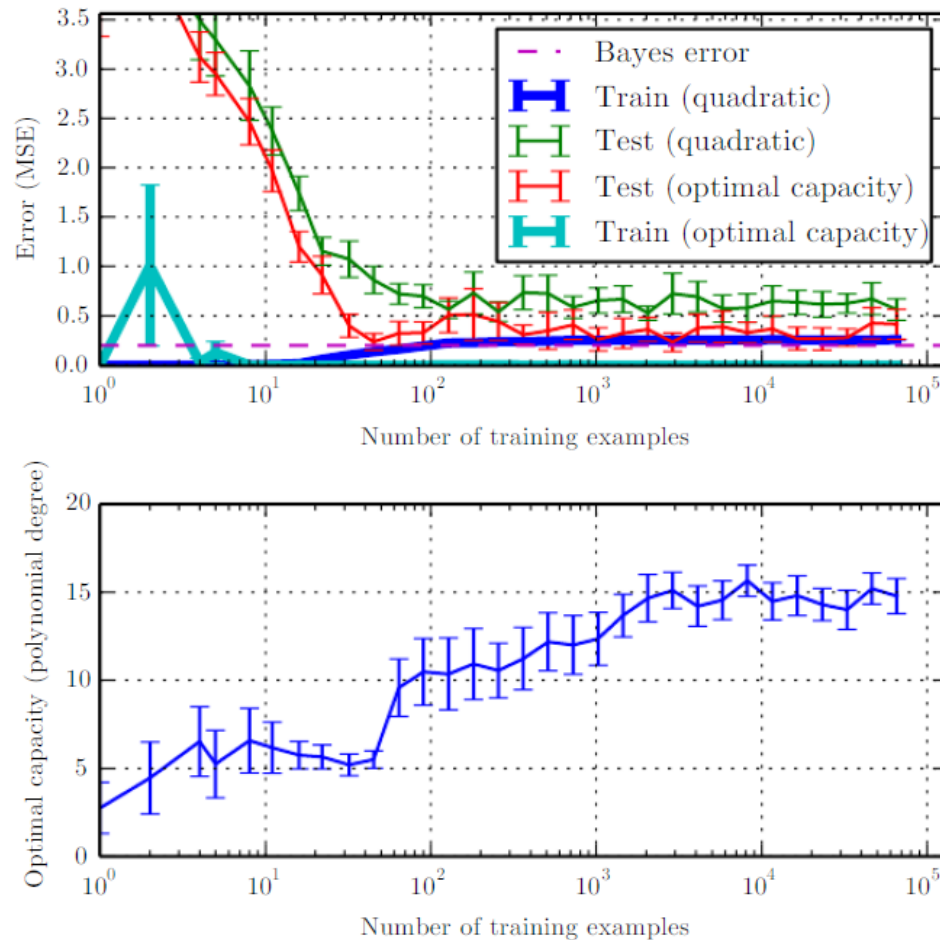


## 5.2 Capacity, Overfitting and Underfitting



## 5.2 Capacity, Overfitting and Underfitting

- The No Free Lunch Theorem





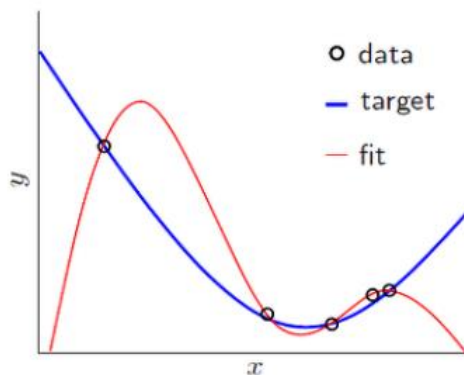
## 5.2 Capacity, Overfitting and Underfitting

- Use Regularization to prevent overfitting using L2 norm

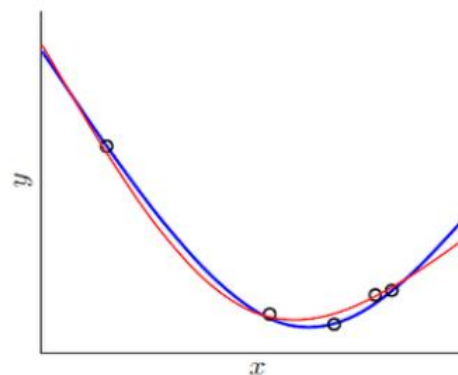
$$J(w) = \text{MSE}_{\text{train}} + \lambda w^{\top} w$$

- Regularizer :  $\Omega(w) = w^{\top} w$

- Make the value of W(weight) lower changing cost function



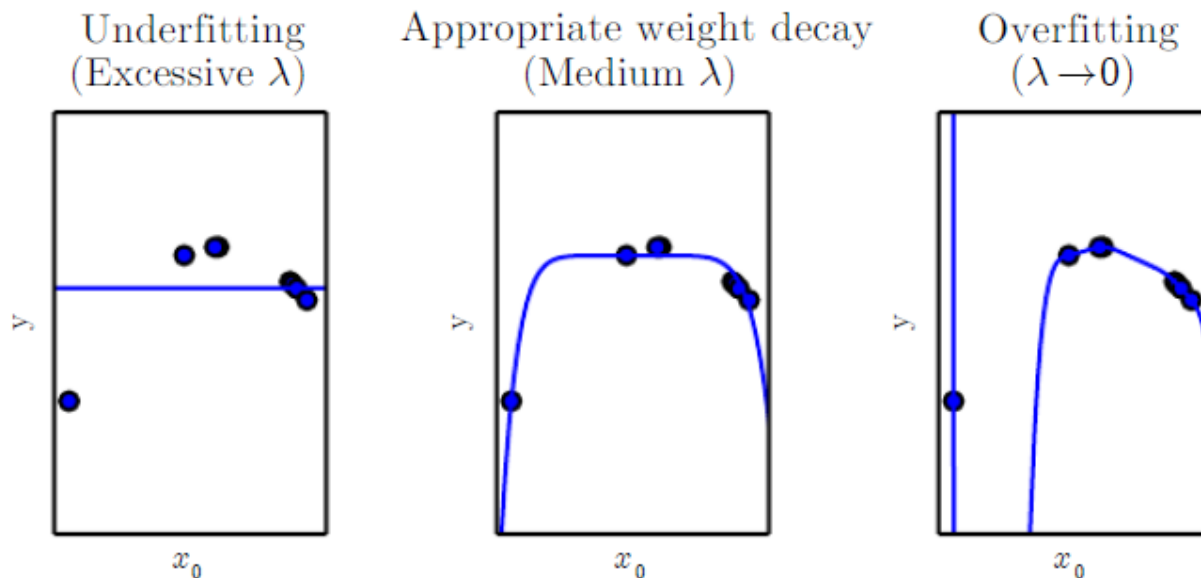
(a) without regularization



(b) with regularization

## 5.3 Hyperparameters and Validation Sets

- The definition of the ‘Model Parameter’
  - The ‘Model Parameter’ is a configuration variable that is internal to the model and whose value can be estimated from data. (ex. Avg, Stdev)
- The definition of the ‘Model Hyperparameter’
  - The ‘Model Hyperparameter’ is a configuration that is external to the model and whose value cannot be estimated from data. (ex.  $\lambda$ )



- Cross-Validation

## 5.4 Estimators, Bias and Variance

---

- Point Estimation

$$\hat{\theta}_m = g(x^{(1)}, \dots, x^{(m)}).$$

- Bias

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta,$$

- Standard Error Mean

$$\text{SE}(\hat{\mu}_m) = \sqrt{\text{Var} \left[ \frac{1}{m} \sum_{i=1}^m x^{(i)} \right]} = \frac{\sigma}{\sqrt{m}},$$

## 5.4 Estimators, Bias and Variance

---

- Get bias value of Mean and Get Variance at Bernoulli

$$P(x^{(i)}; \theta) = \theta^{x^{(i)}}(1 - \theta)^{(1-x^{(i)})}$$

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\begin{aligned} \text{bias}(\hat{\theta}_m) &= \mathbb{E}[\hat{\theta}_m] - \theta \\ &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right] - \theta \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}[x^{(i)}] - \theta \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{x^{(i)}=0}^1 \left(x^{(i)} \theta^{x^{(i)}}(1 - \theta)^{(1-x^{(i)})}\right) - \theta \\ &= \frac{1}{m} \sum_{i=1}^m (\theta) - \theta \\ &= \theta - \theta = 0 \end{aligned}$$

Bias value about Mean

$$\begin{aligned} \text{Var}(\hat{\theta}_m) &= \text{Var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(x^{(i)}) \\ &= \frac{1}{m^2} \sum_{i=1}^m (\theta(1 - \theta)) \\ &= \frac{1}{m^2} m\theta(1 - \theta) \\ &= \frac{1}{m} \theta(1 - \theta) \end{aligned}$$

Variance

# 5.4 Estimators, Bias and Variance

- Get bias value of Mean and Get Variance at Gaussian

$$p(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}\right).$$

$$\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\begin{aligned} \text{bias}(\hat{\mu}_m) &= \mathbb{E}[\hat{\mu}_m] - \mu \\ &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right] - \mu \\ &= \left(\frac{1}{m} \sum_{i=1}^m \mathbb{E}[x^{(i)}]\right) - \mu \\ &= \left(\frac{1}{m} \sum_{i=1}^m \mu\right) - \mu \\ &= \mu - \mu = 0 \end{aligned}$$

Bias value about Mean

$$\begin{aligned} \hat{\sigma}_m^2 &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2, \\ \text{bias}(\hat{\sigma}_m^2) &= \mathbb{E}[\hat{\sigma}_m^2] - \sigma^2. \\ \mathbb{E}[\hat{\sigma}_m^2] &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2\right] \\ &= \frac{m-1}{m} \sigma^2 \\ \tilde{\sigma}_m^2 &= \frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2 \end{aligned}$$

Variance

$$\begin{aligned} \mathbb{E}[\tilde{\sigma}_m^2] &= \mathbb{E}\left[\frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2\right] \\ &= \frac{m}{m-1} \mathbb{E}[\hat{\sigma}_m^2] \\ &= \frac{m}{m-1} \left(\frac{m-1}{m} \sigma^2\right) \\ &= \sigma^2. \end{aligned}$$

Variance Estimator

# 5.4 Estimators, Bias and Variance

- Bias-Variance Trade-off

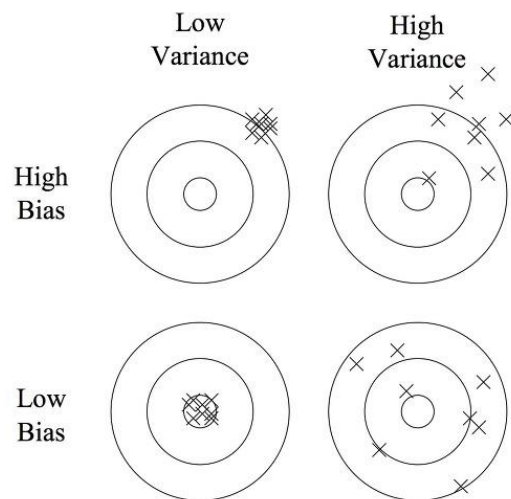
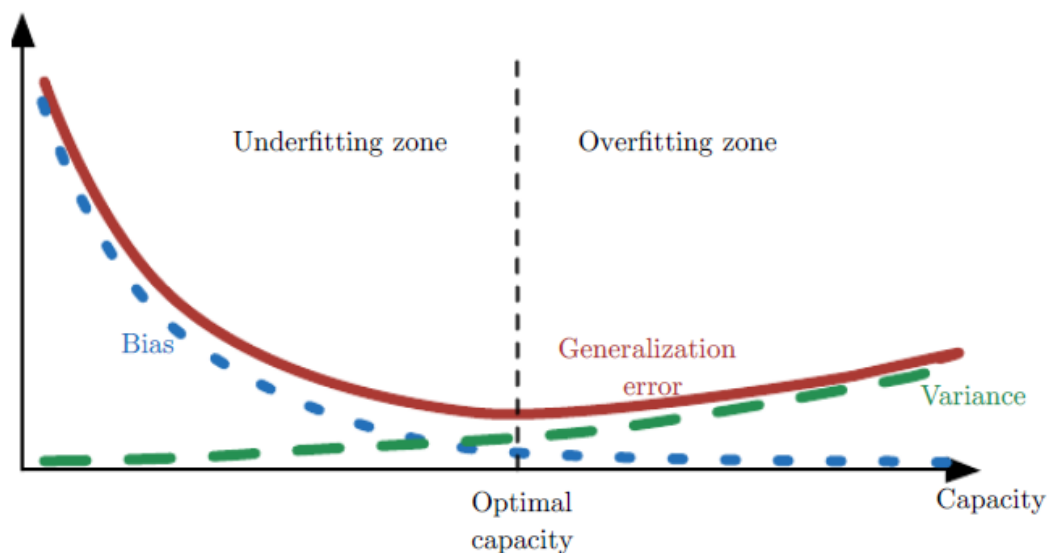


Image from Quora

$$\text{MSE} = \mathbb{E} \left[ (\hat{\theta}_m - \theta)^2 \right] = \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m)$$



- Consistency

- When training data set be greater, point estimates will go to original.

$$\text{plim}_{m \rightarrow \infty} \hat{\theta}_m = \theta$$

$$P(|\hat{\theta}_m - \theta| > \epsilon) \rightarrow 0$$

# 5.5 Maximum Likelihood Estimation

---

- Maximum Likelihood Estimation
  - Estimator could be analyzed using by **bias** and **variance**
  - Also, we could know how to extract some function for good estimator
  - One of the general principle is the **Maximum Likelihood Estimation**

$$\theta_{\text{ML}} = \arg \max_{\theta} p_{\text{model}}(\mathbb{X}; \theta),$$

$$= \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(x^{(i)}; \theta).$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(x^{(i)}; \theta).$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta).$$

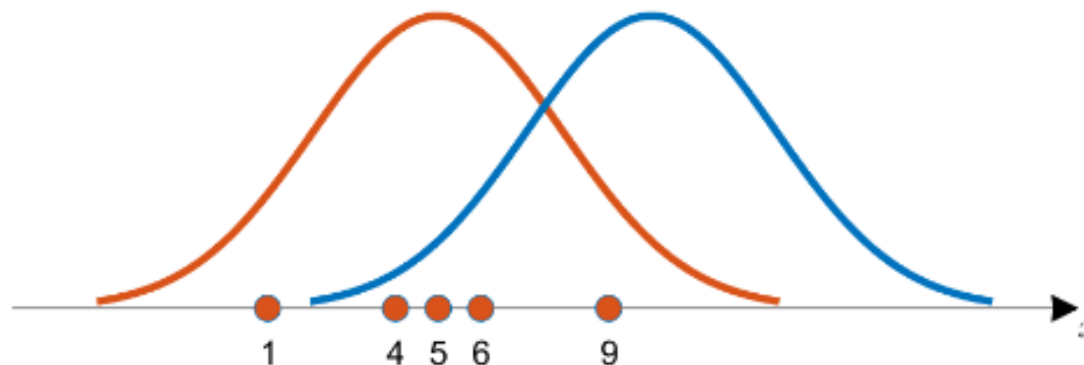
$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})].$$

$$- \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})],$$

# 5.5 Maximum Likelihood Estimation

---

- Maximum Likelihood Estimation
  - Estimator could be analyzed using by **bias** and **variance**
  - Also, we could know how to extract some function for good estimator
  - One of the general principle is the **Maximum Likelihood Estimation**



Which distribution represents the points as well? Red or Blue?



# 5.6 Bayesian Statistics

---

- Bayesian Statistics
  - In Frequentist Statistics, True Parameter  $\theta$  is fixed and unknown, and usually estimate some data using variable the hat of  $\theta$ .
  - But, in Bayesian Statistics, the probability means certainty. True parameter is uncertain, and it is used to random variable.
  - Prior Probability Distribution is represented to  $p(\theta)$ .

$$p(\theta \mid x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} \mid \theta)p(\theta)}{p(x^{(1)}, \dots, x^{(m)})}$$

$$p(x^{(m+1)} \mid x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} \mid \theta)p(\theta \mid x^{(1)}, \dots, x^{(m)})d\theta$$

## 5.6 Bayesian Statistics

---

- Bayesian Statistics
  - In contrast to MLE estimating point using  $\theta$ , predict **using all over the distributions**
  - Prior Distribution : in ‘roll the dice’ case

$$p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(x^{(1)}, \dots, x^{(m)})}$$

$$p(x^{(m+1)} \mid x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)})d\boldsymbol{\theta}$$

## 5.6 Bayesian Statistics

---

- Maximum a Posteriori (MAP) Estimation

$$p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(x^{(1)}, \dots, x^{(m)})}$$

$$p(x^{(m+1)} \mid x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)})d\boldsymbol{\theta}$$

$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} \log p(\boldsymbol{x} \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

Use only one estimation using posterior distribution, and the others are not.

# 5.7 Supervised Learning Algorithms

---

- Supervised Learning Algorithms
  - The learning algorithm means that learns the relationship between the input and output data
  - Answer Label is contained in train data set
- Probabilistic Supervised Learning



# 5.8 Unsupervised Learning Algorithms

---

- Unsupervised Learning Algorithms
  - Learning without answer label

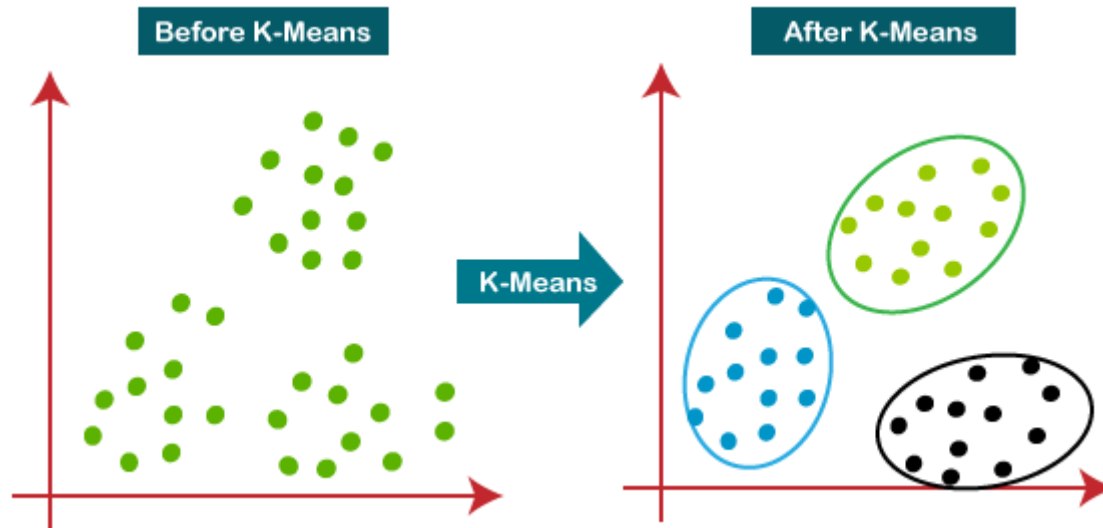
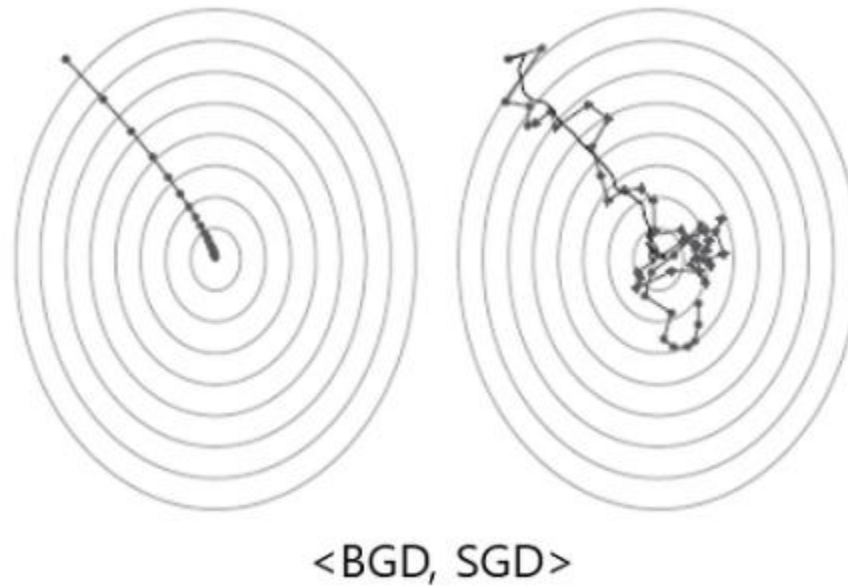


Image from Javatpoint

## 5.9 Stochastic Gradient Descent

---

- Stochastic Gradient Descent
  - Use the Mini-Batch not using all Batch
  - The one step of BGD takes a long time
  - SGD is faster than BGD normally, but it has uncertainty more than BGD



# 5.10 Building a Machine Learning Algorithm

---

- Building a Machine Learning Algorithm
  1. Link the dataset and cost function
  2. Decide the model and the type of optimization

In Linear Regression,

$$J(\mathbf{w}, b) = -\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} \log p_{model}(y \mid \mathbf{x})$$

$$J(\mathbf{w}, b) = \lambda \|\mathbf{w}\|_2^2 - \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} \log p_{model}(y \mid \mathbf{x})$$

$$J(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} \|\mathbf{x} - r(\mathbf{x}; \mathbf{w})\|_2^2$$

# 5.11 Challenges Motivating Deep Learning

---

- Challenges Motivating Deep Learning
  - The Curse of Dimensionality





# 5.11 Challenges Motivating Deep Learning

- Challenges Motivating Deep Learning

$$f^*(\mathbf{x}) \approx f^*(\mathbf{x} + \epsilon)$$

- Local Constancy and Smoothness Regularization
- Manifold Learning

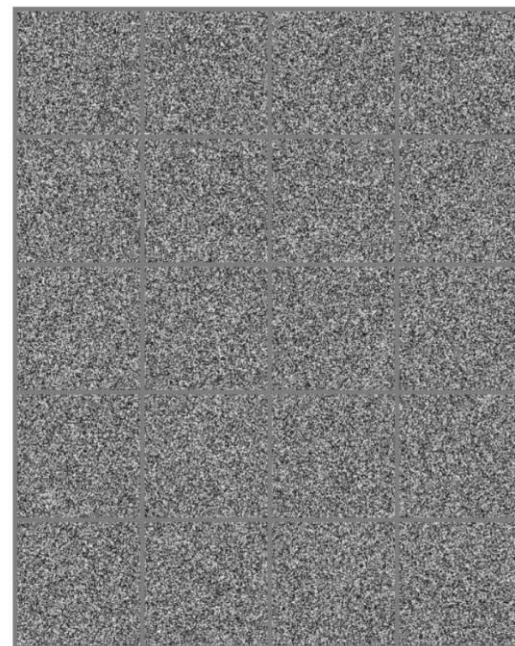
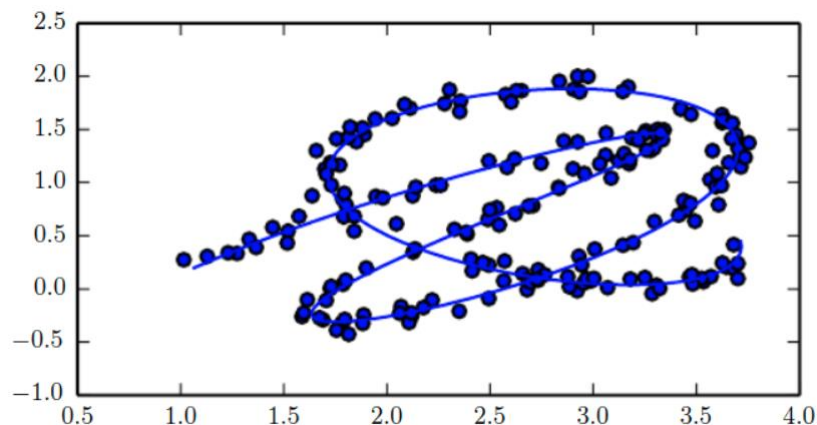


Figure 5.12: Sampling images uniformly at random (by randomly picking each pixel according to a uniform distribution) gives rise to noisy images. Although there is a nonzero probability of generating an image of a face or of any other object frequently encountered in AI applications, we never actually observe this happening in practice. This suggests that the images encountered in AI applications occupy a negligible proportion of the volume of image space.

The End.  
Thank you for watching!

MILab Undergraduate student, TaeHyeon Kim

2023. 03. 07



+ ) Set rules on our study  
(for example, contents, ...)

+ ) How about to put PDFs together?