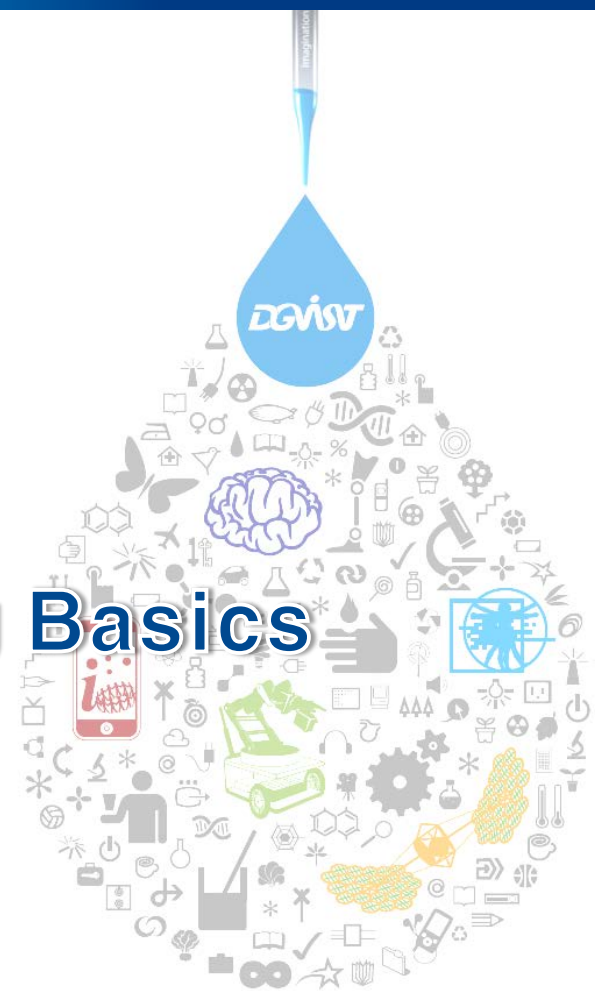# Deep Learning Seminar
# Chapter 5. Machine Learning Basics

2017-07-14

Jinwook Kim

# Book Information
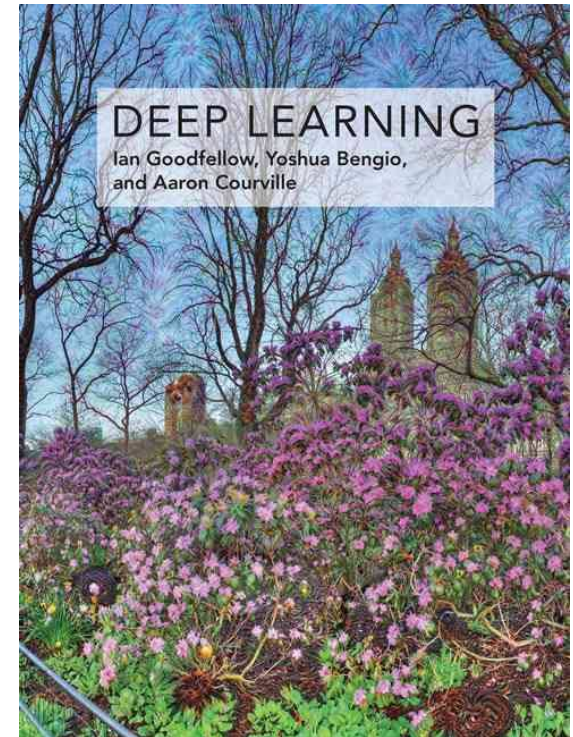
**Title: Deep learning**
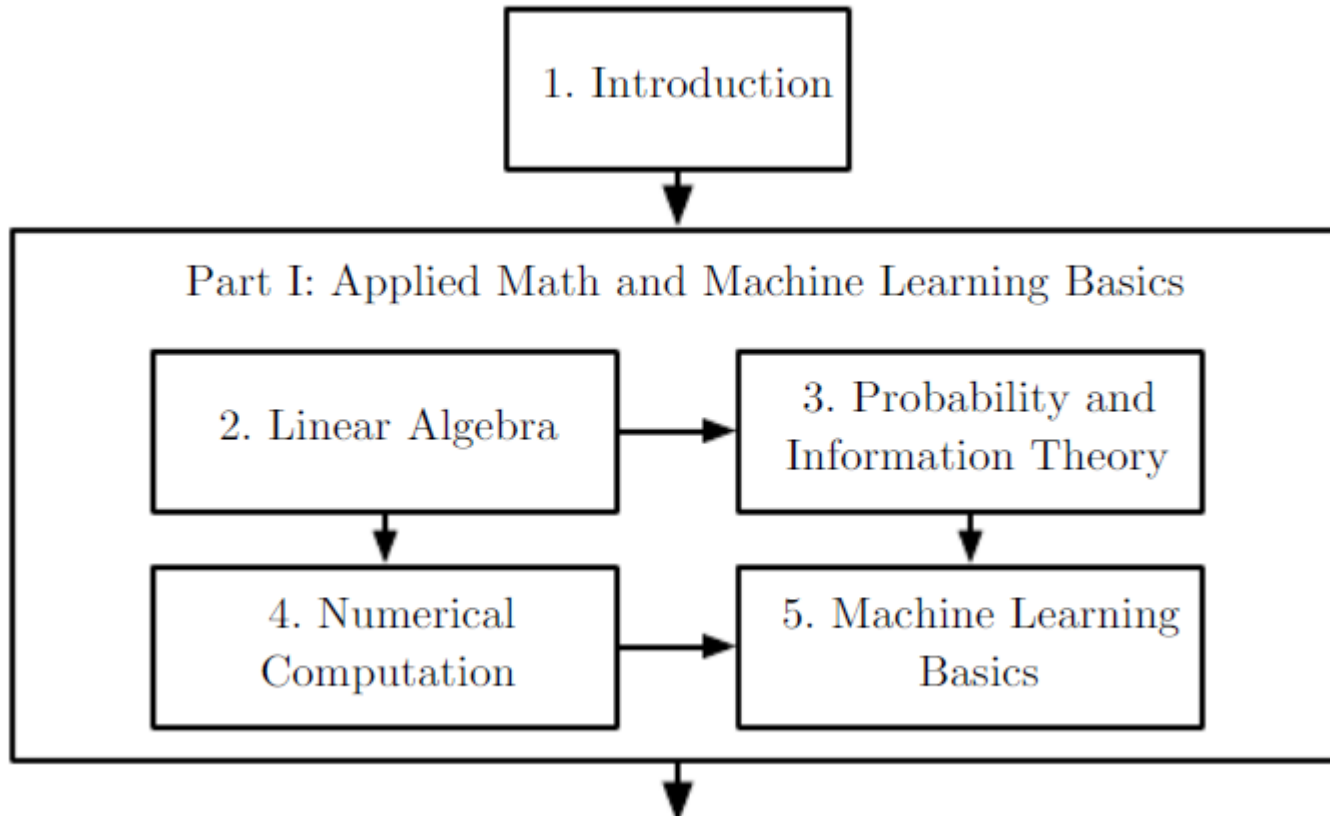
**Authors:**

- Ian Goodfellow
- Yoshua Bengio
- Aaron Courville

**Released: November 10, 2016**

**ISBN: 9780262337434**

# Chapter Organization (Part 1)

# Chapter Organization (Part 2)



Part II: Deep Networks: Modern Practices

6. Deep Feedforward Networks

7. Regularization    8. Optimization    9. CNNs    10. RNNs

11. Practical Methodology    12. Applications

InfoLab  DGIST 대구경북과학기술원

# Chapter Organization (Part 3)



Part III: Deep Learning Research

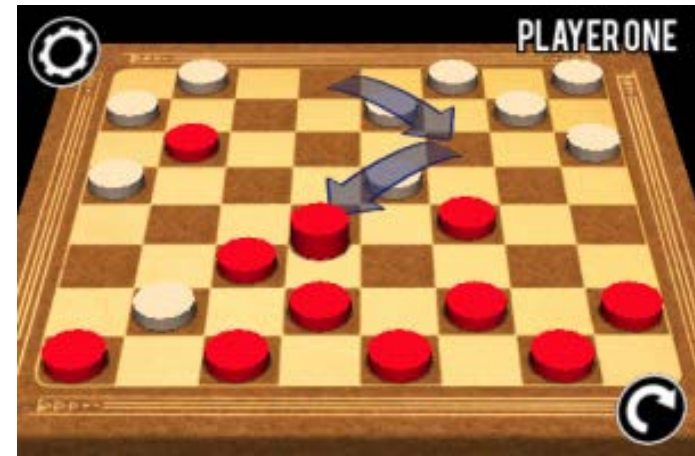| 13. Linear Factor Models | 14. Autoencoders | 15. Representation Learning |
| 16. Structured Probabilistic Models | 17. Monte Carlo Methods |
| 19. Inference | 18. Partition Function |
| 20. Deep Generative Models |

# Machine Learning and AI

# Chapter 5. Machine Learning Basics

- **This chapter provides a brief course in the most important general principles**

  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

# Checkers Game: The First ML Application

- **The Samuel's Checkers-playing program appears to be the world's first self-learning program (Arthur Samuel, 1959)**

  - ➤ Over thousands of games, the program started to learn to recognize the patterns, *which patterns led to win or lose*

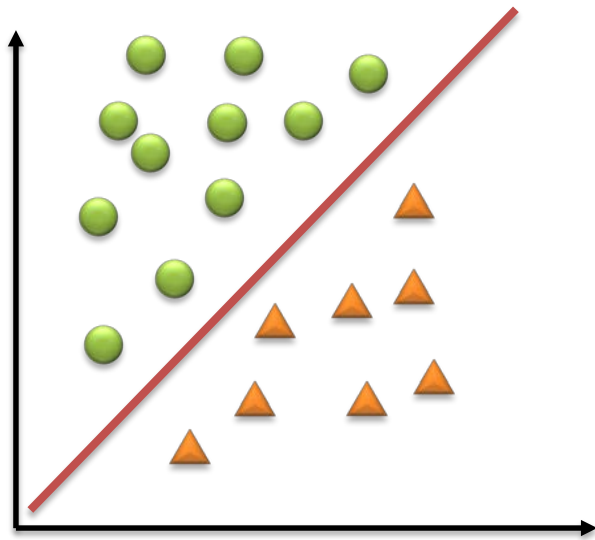- **Finally, the program played much better than Samuel himself could**

**InfoLab** DGIST 대구경북과학기술원

# Learning Algorithms

- **A field of study that gives computers the ability to learn without being *explicitly* programmed (Arthur Samuel, 1959)**

- **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 $E$ (Tom M. Michell, 1997)**

  - $E$: the experience of playing thousands of games
  - $T$: playing checkers game
  - $P$: the fraction of games it wins against human opponents

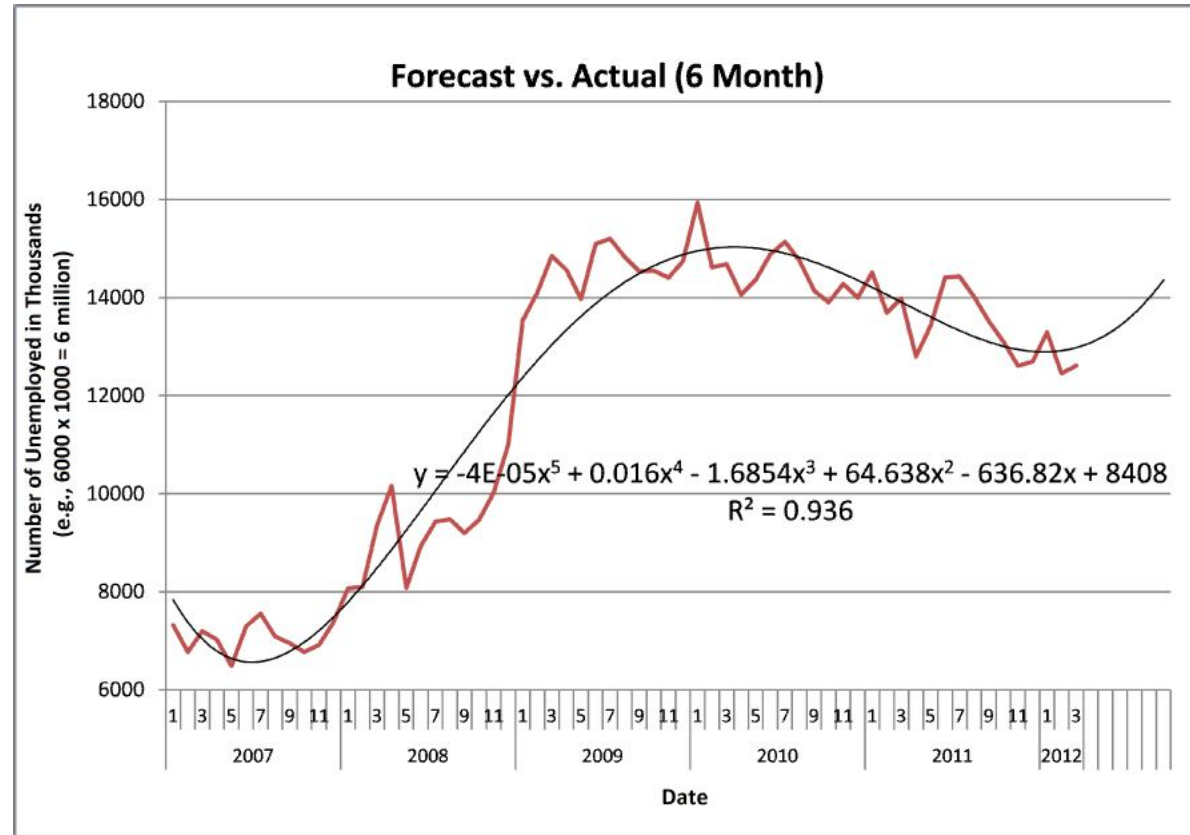  - By its definition, *Samuel's program has **learned** to play checkers*
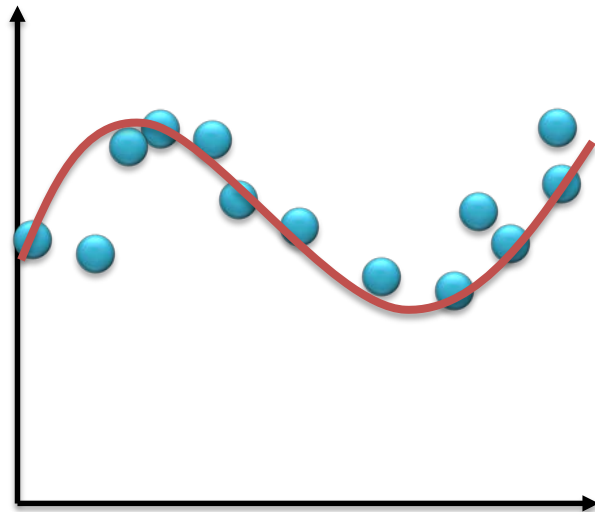
*InfoLab* *DGIST* 대구경북과학기술원

# Task, $T$

- **Classification (with missing inputs)**
- **Regression**
- **Transcription**
- **Machine translation**
- **Structured output**
- **Anomaly detection**
- **Synthesis and sampling**
- **Imputation of missing values**
- **Denoising**
- **Density estimation**
- **Probability mass function estimation**

# Classification

# Regression



Forecast vs. Actual (6 Month)

$y = -4\text{E-}05x^5 + 0.016x^4 - 1.6854x^3 + 64.638x^2 - 636.82x + 8408$

$R^2 = 0.936$

**InfoLab**  DGIST 대구경북과학기술원

# Transcription and Machine Translation





ВЫХОД В ГОРОД

ACCESS TO CITY



AVENUE DES SAPINS  AVENUE DES SAPINS  AVENUE DES SAPINS  AVENUE DES SAPINS

Avenue des Sapins

**InfoLab** DGVSV 대구경북과학기술원

# Performance Measure, $P$

- **Accuracy (error rate) for categorical data**
  - To measure the proportion of examples for which model produces the correct output

- **Average log-probability for density estimation**
  - To measure continuous-valued score for each example

- **E.g., test set for validating model with unseen data**
  1. Separate the data into training set and test set
  2. Train the model with training set
  3. Measure the model's performance with test set

**InfoLab** DGIST 대구경북과학기술원

# Experience, $E$

- **Unsupervised learning algorithms,**
  - Experiences a dataset containing many features
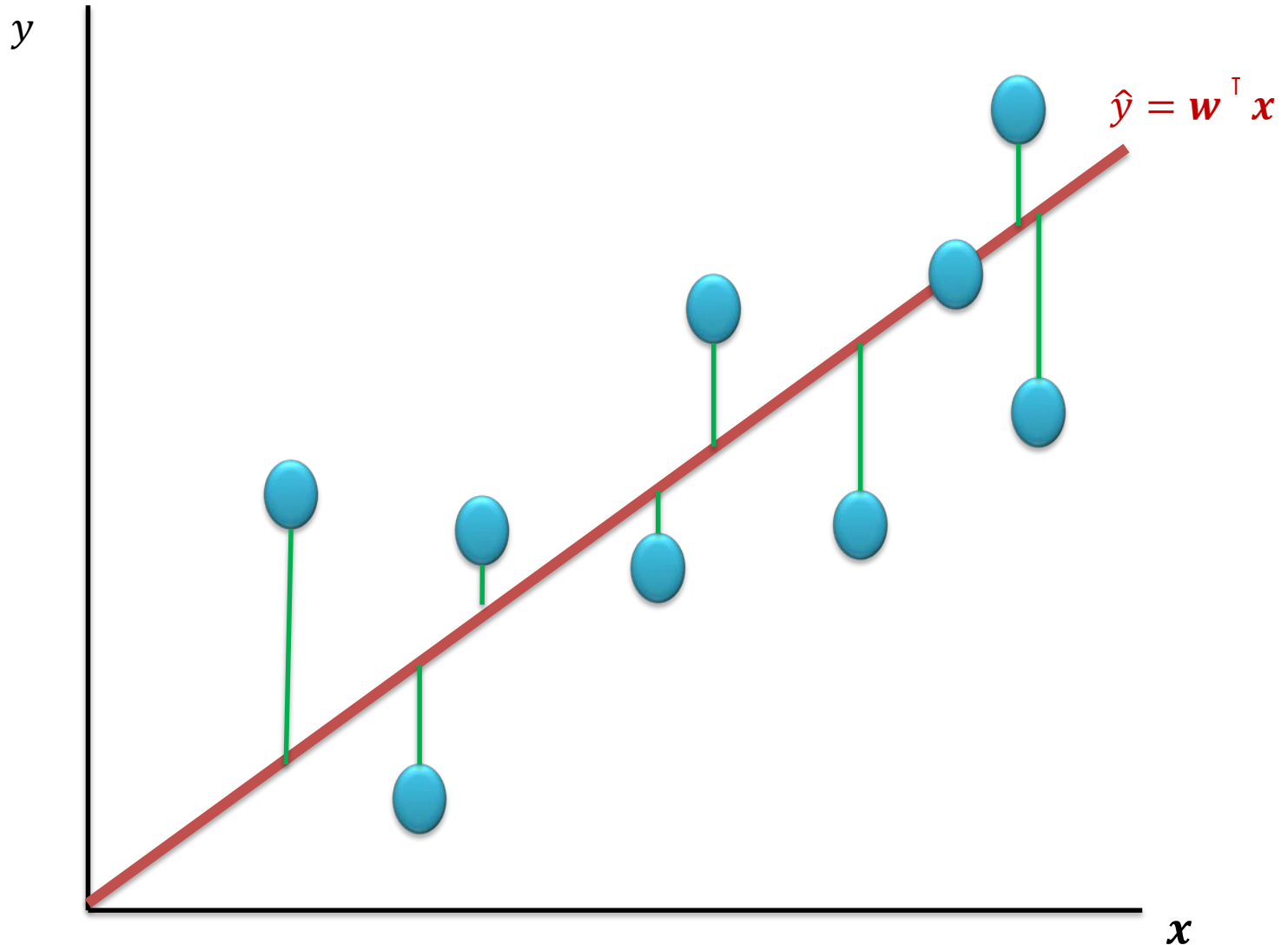  - Learns useful properties of the structure of the dataset

- **Supervised learning algorithms,**
  - Experiences a dataset associated with labels
  - Learns to predict the labels from the data

- **Reinforcement learning algorithms,**
  - Not just experience with a fixed dataset, but interact with an environment
  - Learns actions to maximize cumulative rewards

# Example: Linear Regression



$\hat{y} = \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}$

$y$

$x$

**InfoLab** DGIST 대구경북과학기술원

# Formal Definition of Linear Regression

- **Task, $T$: to predict $y$ from $x$ by outputting $\hat{y} = w^\top x$**

  $x \in \mathbb{R}^n$: input data

  $y \in \mathbb{R}$: output value ($\hat{y}$: predicted by the model)

  $w \in \mathbb{R}^n$: parameters (or weights)

- **Experience, $E$: training set $(X^{train}, y^{train})$**

- **Performance measure, $P$:**
  **mean squared error (MSE) on $(X^{test}, y^{test})$**

$$MSE_{test} = \frac{1}{m} \sum_i (\hat{y}^{test} - y^{test})_i^2$$

$m$: number of dataset

# Find $w$ by Minimizing $MSE_{train}$

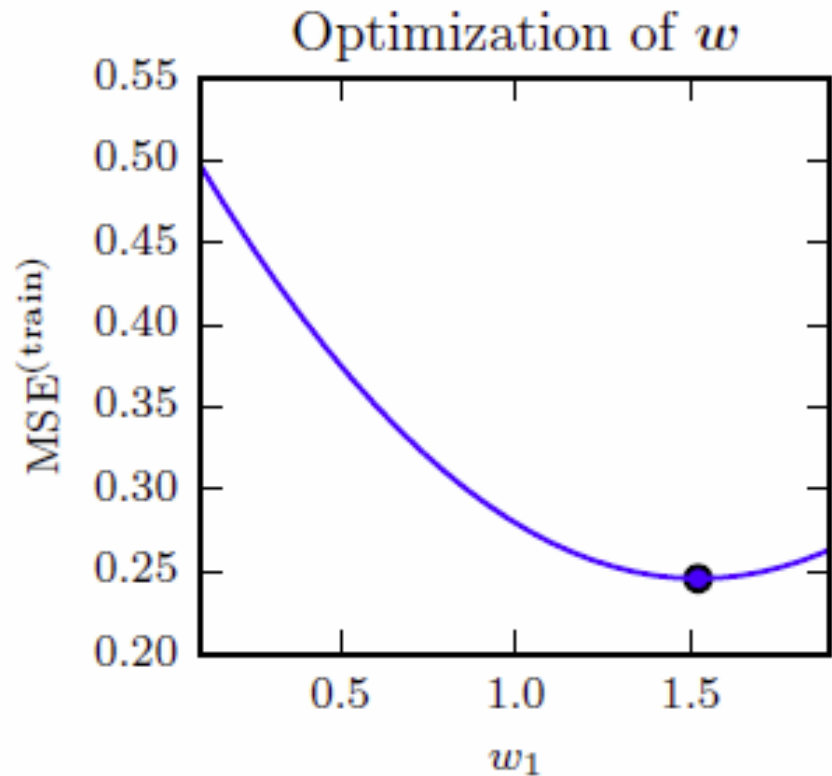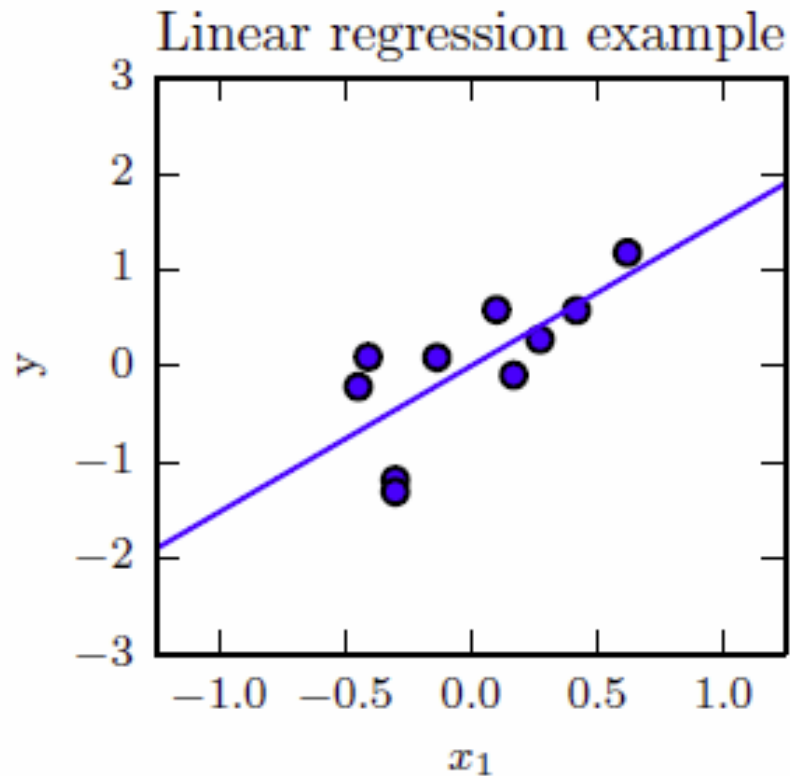$$\nabla w MSE_{train} = 0$$

$$\Rightarrow \nabla w MSE_{train} = \frac{1}{m} \sum_i (\hat{y}^{train} - y^{train})_i^2 = 0$$

$$\Rightarrow \nabla w MSE_{train} = \frac{1}{m} \sum_i (X^{train} w - y^{train})_i^2 = 0$$

…

$$\Rightarrow w = (X^{train^\top} X^{train})^{-1} X^{train^\top} y^{train}$$
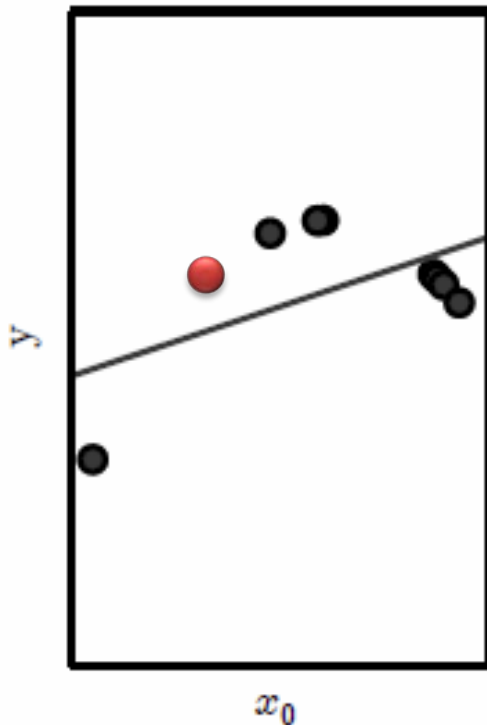
# Linear Regression Problem

# Generalization

- **In ML, generalization is the ability to perform well on previously *unobserved inputs***
  1. Making the training error small
  2. Making the gap between training and test error small

- ***Underfitting* occurs when the model is not able to obtain a sufficiently low error value on training set**

- ***Overfitting* occurs when the gap between the training error and the test error is too large**
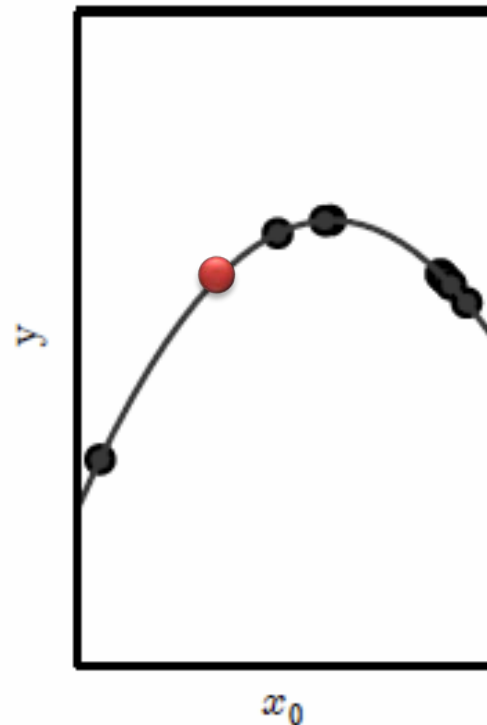
**InfoLab** DGIST 대구경북과학기술원

# Controlling a Model with its Capacity



Underfitting · Appropriate capacity · Overfitting

Linear (degree-1)    Quadratic (degree-2)    Polynomial (degree-9)

**InfoLab** *DGIST* 대구경북과학기술원

# Relationship between Capacity and Error

InfoLab  DGIST 대구경북과학기술원

# Training Set Size and Generalization



Polynomial (degree-9)

Christopher, M. Bishop. *PATTERN RECOGNITION AND MACHINE LEARNING*. Springer-Verlag New York, Chapter 1.

# Regularization

- **Regularization is any modification to prevent overfitting**

- **Regularization is able to *control the performance* of an algorithm**
  - ➤ Intended to reduce test error but not training error

- **Example: weight decay for regression problem**

$$J(w) = MSE_{train} + \lambda w^{\mathsf{T}} w$$

  - $J(\cdot)$: cost function **to be minimized** on training
  - $\lambda$: control factor of the preference for smaller weight ($\lambda \geq 0$)
  - ➤ Trades-off between fitting the training data and being small $w$

# 9-degree Regression Model with Weight Decay



Underfitting (Excessive $\lambda$)

Appropriate weight decay (Medium $\lambda$)

Overfitting ($\lambda \to 0$)

Large $w$ over-fits to training data

InfoLab  DGIST 대구경북과학기술원

# Hyperparameters vs. parameters

- **Hyperparameters are higher-level properties for a model**
  - Decides model's capacity (or complexity)
  - Not be learned during training
  - e.g., degree of regression model, weight decay

- **Parameters are properties of the training data**
  - Learned during training by a ML model
  - e.g., weights of regression model

*InfoLab* DGIST 대구경북과학기술원

# Setting Hyperparameters with Validation Set

■ **Setting hyperparameters in training step is not appropriate**

➤ Hyperparameters will be set to yield overfitting

(e.g., higher degree of regression model, $\lambda \to 0$)

■ **Test set will not be seen for training nor model choosing (hyperparameter setting)**

■ **So, we need validation set that the training algorithm does not observe**

1. Split validation set from training data
2. Train a model with training data (not including validation set)
3. Validate a model with validation set, update hyperparameters

# $k$-fold Cross Validation

Training set

Training folds

Validation fold

1st iteration $\qquad$ $E_1$

2nd iteration $\qquad$ $E_2$

3rd iteration $\qquad$ $E_3$

...

$k$-th iteration $\qquad$ $E_k$

$$E = \sum_{i}^{k} E_i$$

**InfoLab** DGIST 대구경북과학기술원

# Estimation on Statistics

## Point estimation

➤ Attempt to provide the singe best prediction of the true but unknown property of a model using sample data, $\{x^1, \dots, x^m\}$

## Function estimation

➤ Types of estimation that predict the relationship between input and target variables

## Point estimator

$$\widehat{\boldsymbol{\theta}}_{\boldsymbol{m}} = g(x^1, \dots, x^m)$$

$\widehat{\boldsymbol{\theta}}$: point estimator for the property of a model (e.g., expectation)

$m$: number of data elements

$\{x^1, \dots, x^m\}$: independent and identically distributed (i.i.d.) data points

$g(\cdot)$: any estimation function for the given data points

**InfoLab** DGIST 대구경북과학기술원

# Properties of an Estimator

- **Bias** measures the expected deviation from **the true value** of $\theta$

$$bias\left(\widehat{\boldsymbol{\theta}}_m\right) = \mathrm{E}[\widehat{\boldsymbol{\theta}}_m] - \boldsymbol{\theta}$$

- **Variance** measures the deviation from the expected estimator value that **any particular sampling of the data is likely to cause**

$$Var\left(\widehat{\boldsymbol{\theta}}\right)$$

*InfoLab* DGIST 대구경북과학기술원

# Graphical Illustration of Bias and Variance



Low Variance    High Variance

Low Bias

High Bias

High-bias model

High-variance model

Scott Fortmann-Roe, Understanding the Bias-Variance Tradeoff, 2012

InfoLab  DGIST 대구경북과학기술원

# Bias-Variance Trade-off with Capacity

InfoLab ᗪᏀᏙᎥᏕᎢ 대구경북과학기술원

# Bias & Variance on $MSE$

- **Let $h(\cdot; w)$ be a regression model defined by $w$**
  - $\theta = t(x)$: the true (but unseen) distribution
  - $\widehat{\theta} = E[h(x; w)]$: the estimation of $t(x)$

- **Then $MSE$ of $g$ is**

$$MSE = E\left[\left(h(x; w) - t(x)\right)^2\right] = E\left[\left(\widehat{\theta} - \theta\right)^2\right]$$

- **Add and subtract $E[\widehat{\theta}]$ on the internal term**

$$E\left[\left(\widehat{\theta} - E[\widehat{\theta}] + E[\widehat{\theta}] - \theta\right)^2\right]$$

$$E\left[\left(\widehat{\theta} - E[\widehat{\theta}] + E[\widehat{\theta}] - \theta\right)^2\right]$$

$$= E[\ \widehat{\theta}^2 - \widehat{\theta}E[\widehat{\theta}] + \widehat{\theta}E[\widehat{\theta}] - \widehat{\theta}\theta$$
$$-\widehat{\theta}E[\widehat{\theta}] + E[\widehat{\theta}]^2 - E[\widehat{\theta}]^2 + E[\widehat{\theta}]\theta$$
$$+\widehat{\theta}E[\widehat{\theta}] - E[\widehat{\theta}]^2 + E[\widehat{\theta}]^2 - E[\widehat{\theta}]\theta$$
$$-\widehat{\theta}\theta + E[\widehat{\theta}]\theta - E[\widehat{\theta}]\theta + \theta^2\ ]$$

$$= E\left[\left(\widehat{\theta} - E[\widehat{\theta}]\right)^2\right] + \left(E[\widehat{\theta}] - \theta\right)^2$$
$$+2\left(E[\widehat{\theta}E[\widehat{\theta}] - \widehat{\theta}\theta - E[\widehat{\theta}]^2 + E[\widehat{\theta}]\theta]\right)$$

$$= \underline{E\left[\left(\widehat{\theta} - E[\widehat{\theta}]\right)^2\right]} + \underline{\left(E[\widehat{\theta}] - \theta\right)^2}$$
$$\quad\quad Var(\widehat{\theta}) \quad\quad\quad\quad bias(\widehat{\theta})^2$$

**InfoLab** DGVIST 대구경북과학기술원

# Ways to Trade-off Bias & Variance

## Cross-validation

➢ Highly successful on many real-world tasks

## MSE of estimates

➢ MSE incorporates both bias and variance

$$MSE = E\left[\left(\widehat{\theta}_m - \theta\right)^2\right]$$
$$= bias\left(\widehat{\theta}_m\right)^2 + Var\left(\widehat{\theta}_m\right)$$

# Maximum Likelihood Estimation

- **Maximum likelihood estimation (MLE) is a method of estimating the parameters of a statistical model**

$$w_{ML} = \arg\max_{w} p_{model}(X; w)$$

$$= \arg\max_{w} \prod_{i}^{m} p_{model}(x^i; w)$$

$$\equiv \arg\max_{w} \sum_{i}^{m} \log(p_{model}(x^i; w))$$

- $p_{data}(x)$: true but unknown data-generating distribution
- $X = \{x^1, \dots x^m\}$: drawn independently from $p_{data}(x)$
- $p_{model}(x; w)$: a probability distribution estimating $p_{data}(x)$

**InfoLab** DGIST 대구경북과학기술원

# Linear Regression as Maximum Likelihood

- **Instead of single prediction $\hat{y}$, we consider of the model as producing a conditional distribution, $p(y|x)$**

- **Let $p(y|x) = \mathcal{N}(y \mid \hat{y}(x; w), \sigma^2)$**
  - $\hat{y}(x; w)$: the prediction of the mean of the Gaussian
  - $\sigma^2$: variance of $\mathcal{N}$ chosen by user

- **Log-likelihood:**

$$\sum_i^m \log(p_{model}(x^i; w)) = \sum_i^m \log(y^i | x^i; w)$$

$$= -m \log(\sigma) - \frac{m}{2} \log 2\pi - \sum_i^m \frac{\|\hat{y}^i - y^i\|^2}{2\sigma^2}$$

The same estimate of the parameter $w$

- **Cf. $MSE$:**

$$MSE_{\text{train}} = \frac{1}{m} \sum_i \|\hat{y}^i - y^i\|^2$$

InfoLab DGIST 대구경북과학기술원

# Contents on Part 2

- **Bayesian statistics**

- **Supervised learning algorithms**

- **Unsupervised learning algorithms**

- **Stochastic gradient descent**

- **Building a machine learning algorithm**

- **Challenges motivating deep learning**

**InfoLab**  DGIST 대구경북과학기술원

# Q&A

Thank you