

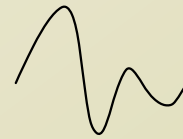
# Motivation and Basics

Il-Chul Moon  
Dept. of Industrial and Systems Engineering  
KAIST

[icmoon@kaist.ac.kr](mailto:icmoon@kaist.ac.kr)

# Weekly Objectives

- Motivate the study on
  - Machine learning, AI, Datamining....
  - Why? What?
  - Overview of the field
- Short questions and answers on a story
  - What consists of machine learning?
  - MLE
  - MAP
- Some basics
  - Probability
  - Distribution
  - And some rules...



# MOTIVATION

# Keywords

Substance가 중요

여러가지 분야에/서 다양한사람·방법 존재

→ 키워드 중요

- Many floating keywords
  - Data-mining, Knowledge discovery, Machine Learning, Artificial Intelligence...
- Comes from territory, perspectives, types of problems, researchers, etc
- We are going to focus on substance, not labeling.
- I am just going to call it “Machine Learning”
  - You can call it whatever you want

Statistics

AI in CS

Database in CS

Management

Industrial Engineering

.....

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING  
DEPARTMENT

Data 수집 (수집) ↑ : 어떠한 지식/정보/데이터는 여러곳 = 거대(거대) (ML)

# Abundance of Data

- Data are being collected everywhere

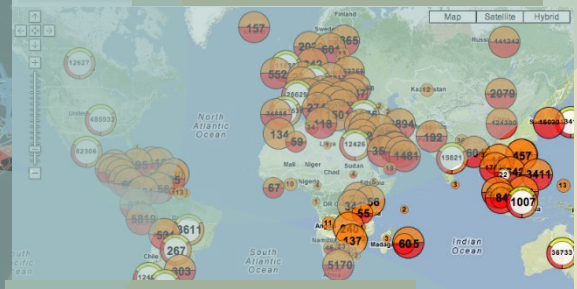
Image Data



Surveillance Data



Trajectory Data



Disease Outbreak Data

Geo-space



Vehicle

Time Series Data

1/17/2014

Text Data

조성민

Machine Logs

Social Networks

News Articles

Social Media

SNS

Blogs

amazon

Purchase+Review Data

10K Rep.



# Examples of Machine Learning Applications

- Machine Learning is everywhere...

Mail : Spam-filtering = document classification

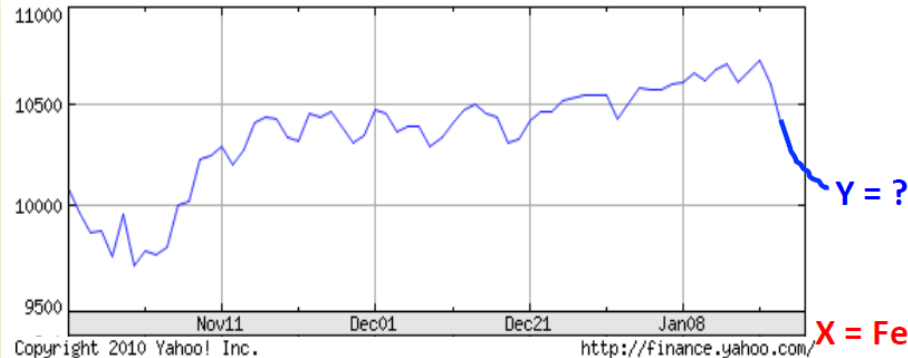
## Document Classification



Sports  
Science  
News

## Stock Market Prediction

DJ INDU AVERAGE (DOW JONES & CO)  
as of 22-Jan-2010



## Plate Num. Recognition



## SNS Recommendation

visit A Chartbanchachai and 7 others added you

Recommended Pages See All

클라라 Clara Lee  
Hong Seok Yang and 2 other friends like her.  
Like

SNL KOREA  
Yong Hoon Choi and 3 other friends like this.  
Like

TIME  
Woojoo Lee and 4 other friends like this.  
Like

People You May Know

Alexander H. Levis, University Professor of Electrical,  
Connect

Doo-Hwan Bae, --  
Connect

Paul Davidsson, Professor at Malmö University  
Connect

See more >

Ads You May Be Interested In

INSEAD The.CCP Plenum Deciphered  
Tsinghua-INSEAD EMBA Master Class on Mar. 7, Seoul.Discuss

## Helicopter Control



# Spam Filtering and more



Table 2 Detailed evaluation results of SVMs with each representation scheme and varying training-set sizes. Macro-averaged MAE scores are provided with p-values, indicating the statistical significances of performance improvement over that of BF (using basic features alone). Numbers in bold font indicate the best method for each fixed training-set size. One star indicates the p-values in (0.01, 0.05]; two stars indicate the p-values equal or less than 1%.

	BF		BF+NC		BF+SI		BF+SIP		BF+SI+NC		BF+SI+NC+SIP	
# of tr	MAE	p-value	MAE	p-value	MAE	p-value	MAE	p-value	MAE	p-value	MAE	p-value
10	0.9666		0.9063	* 0.0382	0.8837	* 0.0106	0.8968	* 0.0311	0.9112	* 0.0211	<b>0.8827</b>	** 0.0087
20	0.9720		0.8969	0.0506	<b>0.8596</b>	* 0.0315	0.9095	* 0.0435	0.9071	0.0558	0.8659	* 0.0235

SVM?

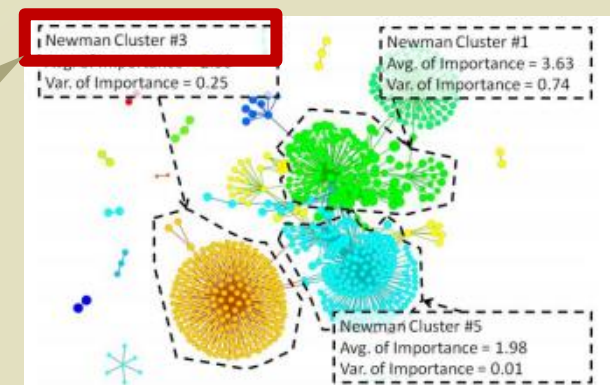
- Spam filter
- More?
  - Importance vs. Urgency
- How to predict an important email?
  - Social networks
  - Contents : *Text*
- Shinjae Yoo, Yiming Yang, Frank Lin, and Il-Chul Moon, Mining Social Networks for Personalized Email Prioritization, ACM SIGKDD Conference, Paris, France, Jun, 28, 2009

Features

## 5.3 Features

The basic features are the tokens in the sections of *from*, *to*, *cc*, *title*, and *body text* in email messages. Let us use a  $v$ -dimensional vector to represent these features for each email message where  $v$  is the vocabulary size. We call it the *basic feature* (BF) sub-vector.

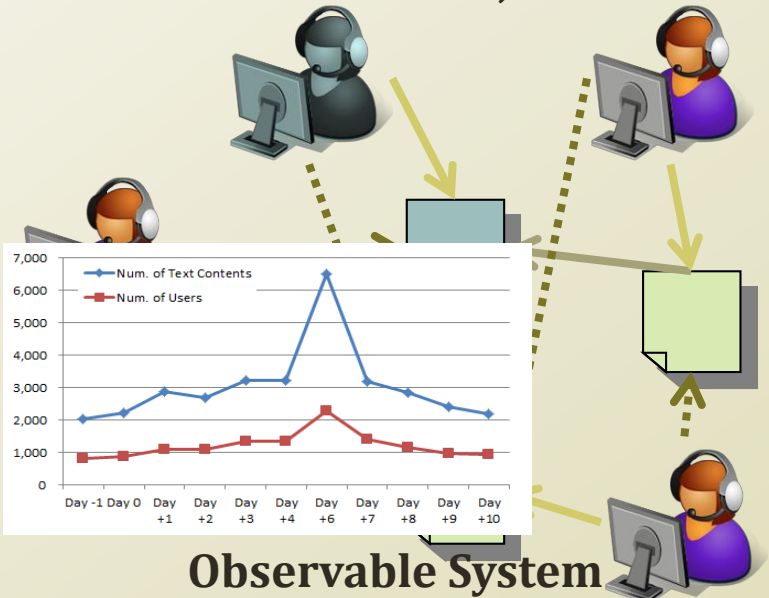
Clusters?  
Is this a machine learning technique?



# Opinion Mining and more

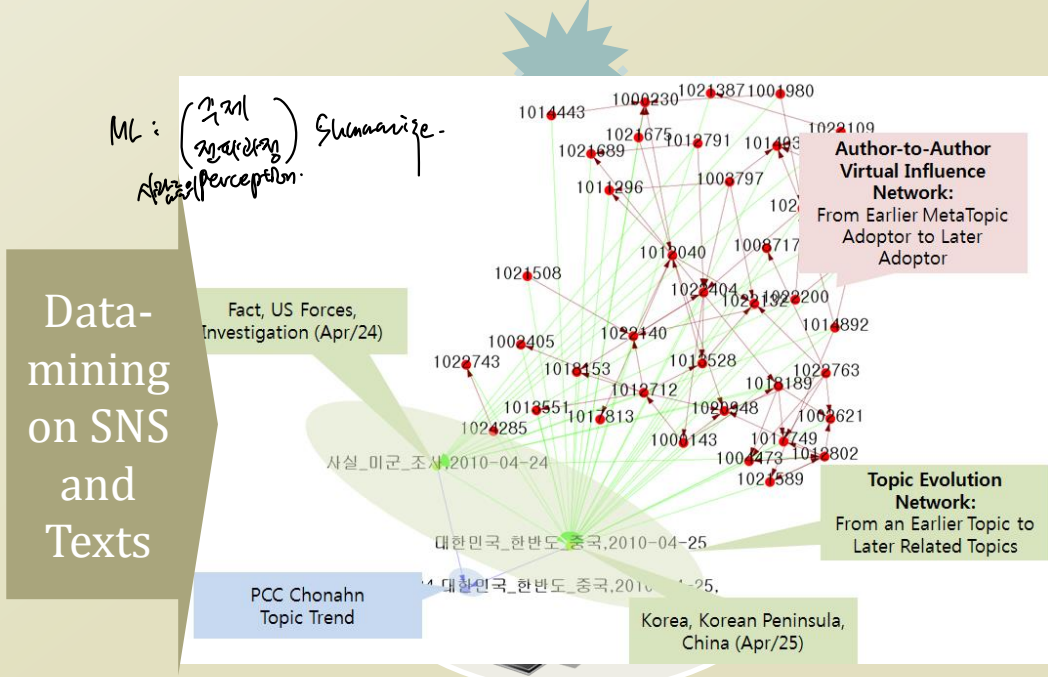


PCC Cheonan  
Sank on Mar 26, 2010



Observable System

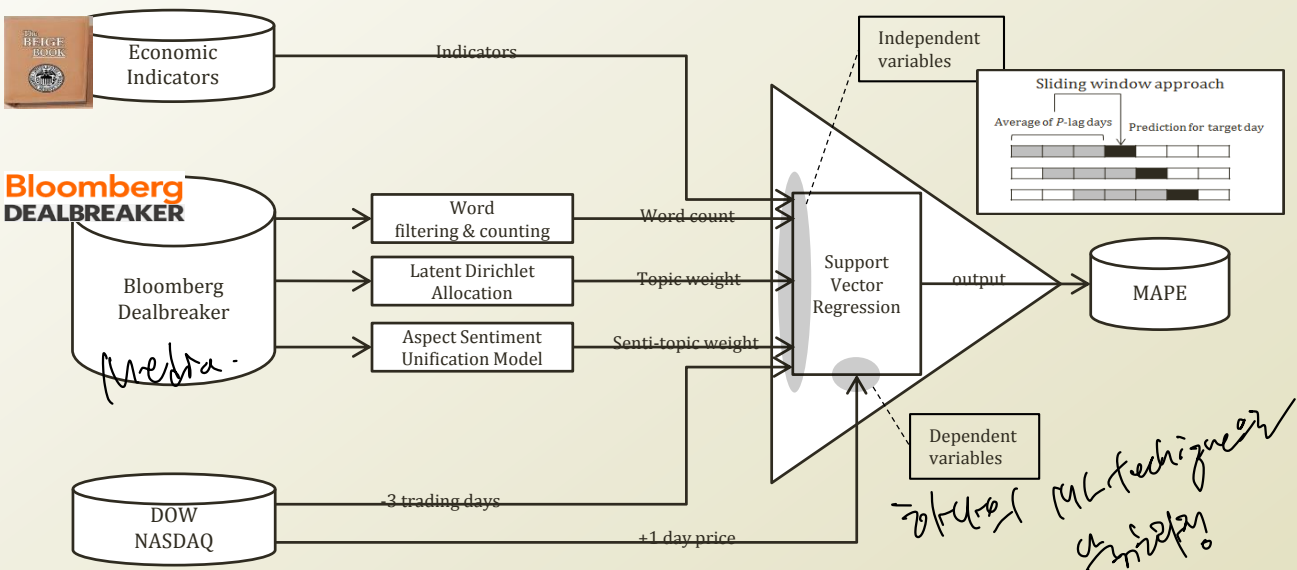
- Finding out consensus of the population
- Mining population's perception of the event
  - Mining key opinion buried in a data chunk
  - Estimating future polarity of the population
  - Strategy to maintain the unity of the population



Implicit System

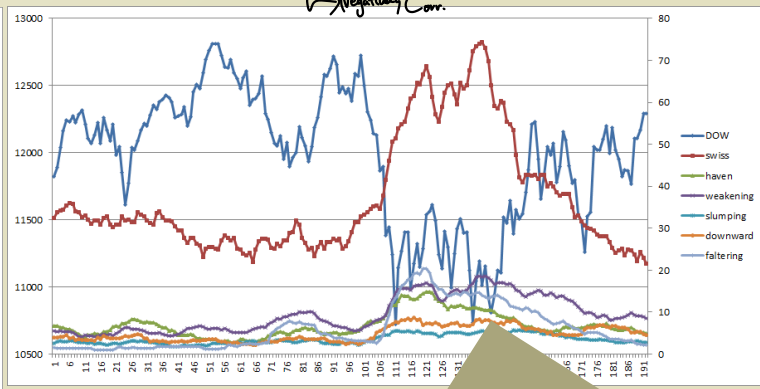
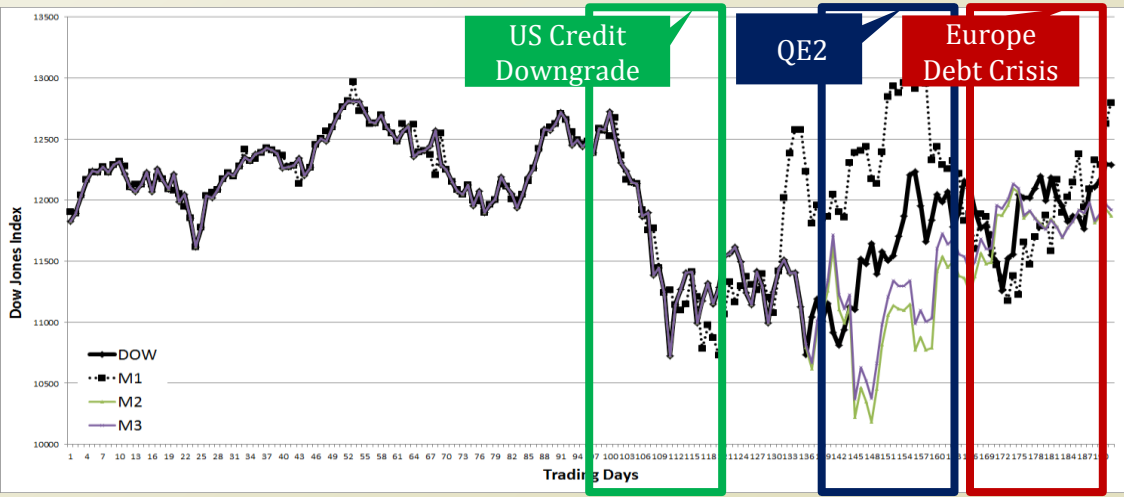


# Stock Market Prediction and more



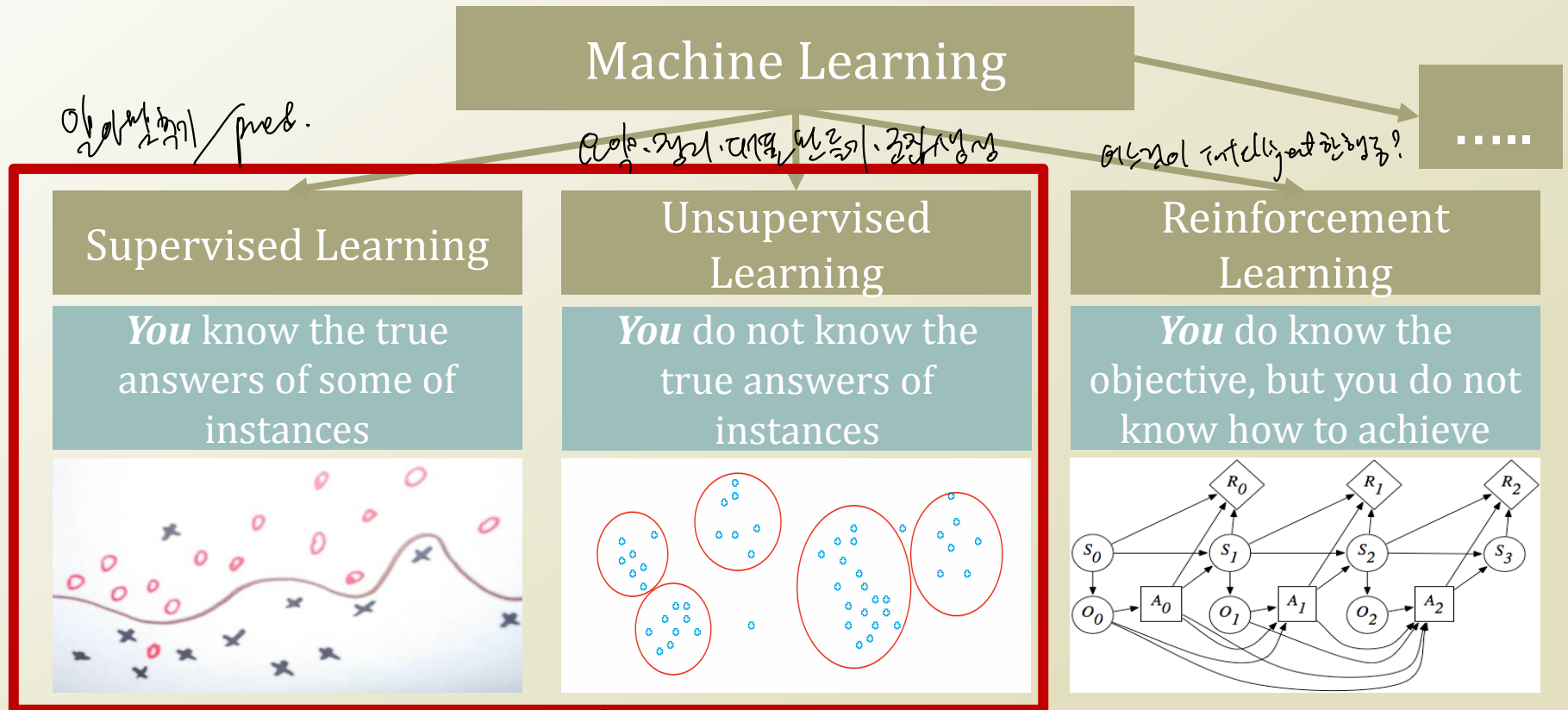
High Coefficients on Prediction

TopicWeight 26	TopicWeight 1	TopicWeight 7
-0.609	0.520	0.508
notes	obama	jun
moodys	republican	pence
swaps	republicans	na
treasuries	congress	swiss
versus	senate	chg
ratings	bill	francs
auCTION	barack	spa
default	lawmakers	fullyear
strategist	administration	nv
franc	democrats	dividend
twoyear	taxes	firstquarter
samp	white	ks
currencies	workers	paris
yen	democrat	reporting
swiss	obamas	tech



Heavy negative correlation between "swiss" and DJIA

# Types of Machine Learning



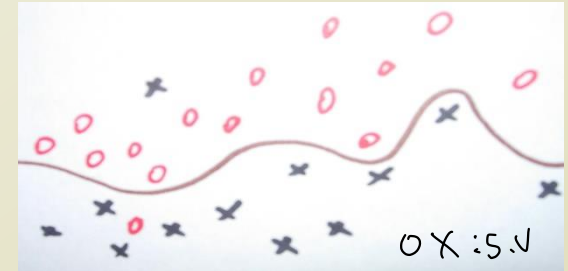
- **You** can
  - Machine learning
  - Dataset provider
  - Machine learning users
  - etc

알려진 답이 있는 문제

- Various classifications by different professors
  - Purpose, data types, etc
- Other learning classifications also exist

# Supervised Learning

*You* know the true answers of some of instances



Supervision이란 ML.  
 알려주는 data (tag 있는 - guide 있는)  
 찾는 것만 2차원 (안녕하세요 2차원)  
 ex) Spam or not.

- **You know the true value, and you can provide examples of the true value.**

- Cases, such as

- Spam filtering
- Automatic grading
- Automatic categorization (ex) 아사지 상조영지.

- Classification or Regression of

- Hit or Miss: Something has **either disease or not**.
- Ranking: Someone received **either A+, B, C, or F**.
- Types: An article is **either positive or negative**.
- Value prediction: The price of this artifact is **X**. : Conf. value -

- Methodologies

- Classification: estimating a discrete dependent value from observations
- Regression: estimating a (continuous) dependent value from observations

ML  
 분야

ex)

# Unsupervised Learning

데이터의 숨겨진 구조. supervision X

- **You don't know the true value, and you cannot provide examples of the true value.**
- Cases, such as
  - Discovering clusters
  - Discovering latent factors
  - Discovering graph structures
- Clustering or filtering or completing of
  - Finding **the representative topic words from text data**
  - Finding **the latent image from facial data**
  - Completing the incomplete **matrix of product-review scores**
  - Filtering the **noise from the trajectory data**
- Methodologies
  - Clustering: estimating sets and affiliations of instances to the sets
  - Filtering: estimating underlying and fundamental signals from the mixture of signals and noises

군집

잠재요인 / factor.

군집화-라벨이 주어지지 않음.  
매일 찾기 ↑



↓  
latent image  
이러한 이미지를 생성하는  
→ 어떻게 할까?  
그 답이 바로.

## Unsupervised Learning

*You* do not know the true answers of instances

