# Introduction to Big Data and Artificial Intelligence

Lecturer: Jinkyu Kim
Hanyang University, Ph.D. student

# Contents

# 1. Big Data Revolution: Trend

Figure 1: Factors leading to Big Data Revolution

| New Datasets | Methods of Analysis | Computing |
|---|---|---|
| Internet of Things | Machine Learning | Parallel/Distributed |
| Satellites, Phones | Deep Learning | Cheap Memory |
| Social Media, etc. | Artificial Intelligence | Cloud Computing |

**Big Data Revolution**

Source: J.P.Morgan Macro QDS.

☐ Exponential increase in amount of data available

☐ Increase in computing power and data storage capacity, at reduced cost

☐ Advancement in Machine Learning methods to analyze complex datasets

# 1. Big Data Revolution: Trend

□ Amount of data

- 90% of the data in the world today: created in the past two years

- Expected: 4.4 zettabytes (2015) => 44 zettabytes (2020)
  cf) 1 zettabyte = 1,000,000,000,000 Gb

- Further Acceleration:
  IoT (Internet of Things), smart phones, cost reduction in satellite techs

# 1. Big Data Revolution: Trend

□ Computing power & Storage "Cloud Computing" – Physical Ability


  - Remote, shared access to resources

    ex) A single web search on Google: coordinated across ~1000 computers
    ex) Software-as-a-service (SaaS), Platform-as-a-service (PaaS), etc.


  - Diminished the barriers to large-scale data processing and analytics

# 1. Big Data Revolution: Trend

□ Machine Learning methods – Software Ability

- Part of Statistics and Computer Science

- Enable analysis of large and unstructured datasets

- Classical Machine Learning, Deep Learning, Reinforcement Learning

- Ex) Amazon Echo, Google Home, Apple Siri, …

"Just 0.5% of the data produced is currently being analyzed [Regalado (2013)]"
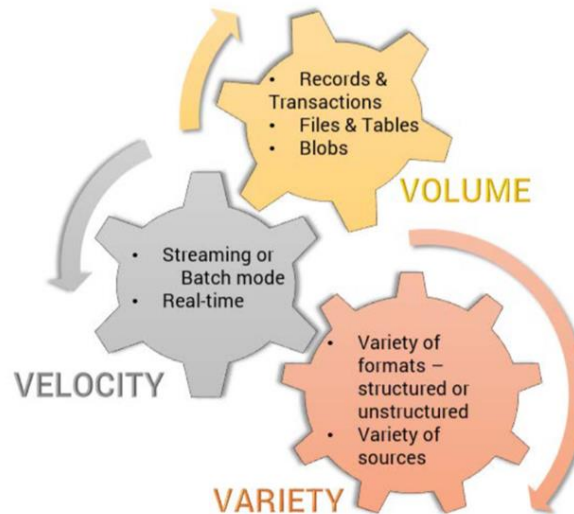- Compelling reason for learning about new datasets and ML toolkits.

# 2. Terminology Description

☐ "Big" Data

☐ Machine Learning

☐ Artificial Intelligence

# 2. Terminology Description

□ "Big" Data: What is Big?

Figure 2: Features of Big Data



Source: J.P.Morgan Macro QDS.

- Volume: the size of data is huge
- Velocity: the speed with which data is sent or received
- Variety: variety of formats (structured, semi-structured, unstructured)

# 2. Terminology Description

□ Machine Learning (ML)

  - Goal: enable computers to learn from their experience in certain tasks

  - Enable the machine to improve performance as their experience grows

  - Category: Supervised Learning, Unsupervised Learning, Deep and Reinforcement Learning
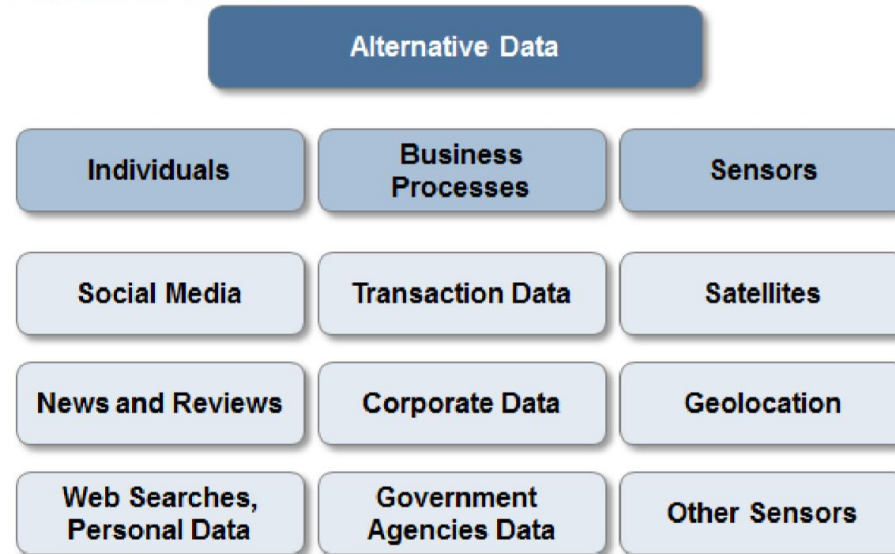
# 2. Terminology Description

☐ Artificial Intelligence (AI)

- Broader scheme of enabling machines with human-like intelligence

- At the first stage: hardcoding a large number of rules "Symbolic-AI"

- Machine Learning is another attempt to achieve AI

- Machine Learning & Deep Learning so far represent the most serious attempt at achieving AI (image&pattern recognition, language trans., etc.)

# 3. Classification of Alternative Data Sets

Figure 3: Classification of Big/Alternative Data sources

**Alternative Data**

| Individuals | Business Processes | Sensors |
|---|---|---|
| Social Media | Transaction Data | Satellites |
| News and Reviews | Corporate Data | Geolocation |
| Web Searches, Personal Data | Government Agencies Data | Other Sensors |

Source: J.P.Morgan Macro QDS.

☐ Data generated by individuals

☐ Data generated by Business Processes

☐ Data generated by sensors

# 3. Classification of Alternative Data Sets

☐ Data generated by individuals

- Mostly recorded through textual mediums
  (unstructured & distributed across multiple platforms)

☐ Possible sources

- Social media: Twitter, Facebook, LinkedIn, etc.

- Specialized sites: Business-reviewing website Yelp, E-commerce groups
  like Amazon, Mobile App analytics companies like App annie

- Web searches, and personalized data such as Google Search trends, etc.

# 3. Classification of Alternative Data Sets

☐ Data generated by Business Processes

   - produced by corporations and public entities (often highly structured)

☐ Examples

   - transaction records such as credit card data
   - banking records
   - supermarket scanner data
   - supply chain data

# 3. Classification of Alternative Data Sets

☐ Data generated by sensors

  - mechanically collected data through sensors embedded in various devices (typically unstructured & size is massive)

☐ Examples

  - satellite imaging
  - geolocation data
  - cameras fixed at a location (ex. CCTV)
  - weather and pollution data
  - IoT (Internet of Things)

# 4. Classification of Machine Learning Techniques

□ Necessity of New Methods

- Large and less structured data: hard to handle with spreadsheet program

- New methods of automatically analyzing unstructured data are needed

- ML yielded some spectacular results
  (image & pattern recognition, natural language processing, auto-driving)

"What is the application of ML in finance,
            and how do these methods differ from each other?"

# 4. Classification of Machine Learning Techniques

☐ First stage: Symbolic Artificial Intelligence – not a ML

  - We can instruct a computer to perform certain tasks based on rules.

  - Why don't we give machines a large number of rules so that it appears to be "intelligent"? (This is a so-called "Symbolic AI")

  - However, it will freeze the first time it encounters a situation that does not exactly match a set of pre-programmed rules.

  - This is not a Machine Learning. It's just making a lot of rules.

    => Not a big success

# 4. Classification of Machine Learning Techniques

☐ Machine Learning

- the computer is given an input (set of variables and datasets) and output (consequences of input variables)

- The machine then finds or "learn" a rule that link input and output

- ML has 3 important sub-categories, i.e., supervised learning, unsupervised learning, deep learning

# 4. Classification of Machine Learning Techniques

☐ Supervised Learning

- answer sheet
- trying to find a rule, an "equation" that we can use to predict a variable.

☐ Unsupervised Learning

- no answer sheet
- uncovering the structure of data

☐ Deep Learning

- similar to how people learn
- ex) A child will learn concept of face by identifying some simple features like eyes, nose, mouth, etc. From simple features one builds a more complex concepts such as face.

# 4. Classification of Machine Learning Techniques

☐ Classification of ML techniques

Figure 7: Classification of Machine Learning techniques



**Machine Learning / Artificial Intelligence**

| Supervised Learning | | Unsupervised Learning | | Deep Learning | | Other Approaches |
|---|---|---|---|---|---|---|
| Regression | Classification | Clustering | Factor Analysis | Time Series | Unstructured | Reinforcement Learning |
| Lasso, Ridge, Loess, KNN, Spline, XGBoost | Logistic, SVM, Random Forest, Hidden Markov | K-means, Birch, Ward Spectral Cluster | PCA, ICA, NMF | Multilayer Perceptron (MLP) Convolutional Neural Nets (CNN) Long Short-Term Memory (LSTM) Restricted Boltzmann Machine (RBM) | | Semi-Supervised / Active Learning |

Source: J.P.Morgan Macro QDS

# 5. Supervised Learning

☐ Supervised Learning

   - Algorithm is provided with historical data (both input and output variables) and is trying to find the relationship that has the best predictive power for out of the sample data

☐ Regression and Classification

   - Regressions: try to predict output variables based on a number of input variables (ex. predict asset returns)

   - Classification: attempt to group or classify output into categories (ex. decide buy/sell actions)

# 5. Supervised Learning-Regression

□ Linear Regression

  - Even a simple linear regression can be thought of as a Supervised Machine Learning method.

  - Sometimes, it is too simplistic to uncover the true relationship.

$$y = \beta_0 + \sum_{i=1}^{n} \beta_i x_i + \varepsilon.$$

OLS: Minimize Historical Sum of $\left(y - \left(\beta_0 + \sum_{i=1}^{n} \beta_i x_i\right)\right)^2$.

# 5. Supervised Learning-Regression

□ Penalized Regression

$$y = \beta_0 + \sum_{i=1}^{n} \beta_i x_i + \varepsilon.$$

OLS: Minimize Historical Sum of $\left(y - \left(\beta_0 + \sum_{i=1}^{n} \beta_i x_i\right)\right)^2.$

Lasso: Minimize Historical Sum of $\left(y - \left(\beta_0 + \sum_{i=1}^{n} \beta_i x_i\right)\right)^2 + \alpha \sum_{i=1}^{n} |\beta_i|,$

Ridge: Minimize Historical Sum of $\left(y - \left(\beta_0 + \sum_{i=1}^{n} \beta_i x_i\right)\right)^2 + \alpha \sum_{i=1}^{n} \beta_i^2,$

Elastic Net: Minimize Historical Sum of $\left(y - \left(\beta_0 + \sum_{i=1}^{n} \beta_i x_i\right)\right)^2 + \alpha_1 \sum_{i=1}^{n} |\beta_i| + \alpha_2 \sum_{i=1}^{n} |\beta_i|^2.$

- the optimizer will set unnecessary and very large betas to zero

# 5. Supervised Learning-Regression

☐ K-Nearest Neighbor (KNN)

- Given a new datapoint, we search through historical data and identify a number of K of similar instances

- Once we have located K nearest neighbors, we can average the output variable y for this subset and use that as our prediction

☐ LOESS

- Using data for K-nearest neighbors, for each new point we fit a linear regression based on the K nearest neighbors, and predict the output using those coefficients

=> Simple way to extrapolate analysis on past similar events

# 5. Supervised Learning-Classification

□ Logistic Regression

- produces output as a binary decision such as "buy" or "sell"

- forecast the probability of an event given historical sample data

# 5. Supervised Learning-Classification

□ Support Vector Machine (SVM)

- the most popular classification algorithm

- enhance with new derived features / axis-transform / higher dimension

Figure 58: Simplified illustration of SVM classification



Points 'P' are support vectors

Source: J.P.Morgan Macro QDS

# 5. Supervised Learning-Classification

□ Decision Tree

- Essentially, flow charts used commonly in business management

- To fit a decision tree, the algorithm usually looks for the best variable and the best splitting value among all possibilities, so that a particular loss function is minimized

Figure 63: Example of decision tree to classify whether future stock returns are good or poor. In each node, we split the data into two subsets depending on the value of the splitting variable

```
                        All stocks

           PB < 1.5                    PB > 1.5

      Vol < 0.8    Vol > 0.8    Momentum < 1    Momentum > 1

   Strong returns  Weak returns  Weak returns   Strong returns
```

# 5. Supervised Learning-Classification

□ Random Forests

   - A single decision tree is capture the structure of the in-sample data
   very well, but it tends to work poorly out-of-sample (overfitting)

   - Bagging: draw many random subsets of the original data

   - Fit a decision tree on each subset

   - Average the predictions of all decision trees

=> Reducing the variances of the prediction

# 6. Unsupervised Learning

☐ Unsupervised Learning

  - Identify relationships between variables and their common drivers

  - In unsupervised learning, the machine is simply given the entire data sets and it does not have a notion of what are independent and what are the dependent variables

☐ Clustering and Factor analysis

  - Clustering: splitting a dataset into smaller groups based on similarity

  - Factor Analysis: identify the main drivers of the data

# 6. Unsupervised Learning-Clustering

□ Clustering

- In the supervised learning, the algorithm tries to minimize the sum of squared errors, i.e., (predicted value – actual value)^2.

- Now, we have no answer sheet. We don't know what actual values are.

- In this situation, how can we compute errors?

# 6. Unsupervised Learning-Clustering

☐ K-means clustering

  - simplest clustering algorithm

  - algorithm starts by initially marking random points as centroids

  - it iteratively does a two-step calculations

    1. maps points to closest centroids
    2. redefines the centroid as the mean of the points mapped to it



(a)         (b)         (c)

(d)         (e)         (f)

# 6. Unsupervised Learning-Clustering



(a)     (b)     (c)

(d)     (e)     (f)

# 6. Unsupervised Learning-Factor Analysis

☐ Factor Analysis

- Goal: identify the main drivers of the data or identify the best representation of data


☐ Principal Component Analysis

- a very well-known method of factor analysis

- originally a statistical tool for dimensionality reduction

# 6. Unsupervised Learning-Factor Analysis

□ Principal Component Analysis - Intuition

 - Assume we want to categorize observations based on similarities

 - If we have only one characteristic, that characteristic is the only reference so that we can compare each group based on that reference

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |

# 6. Unsupervised Learning-Factor Analysis

□ Principal Component Analysis - Intuition

- If we have two characteristics, we can compare it with two axes.

- KEY IDEA: also, we can compare it with the single reference line "$l$".

- This single line is roughly the concept of "Principal Component" (PC)

- We can compare two characteristics with just one PC. But, what if there are $n$ characteristics?



| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

# 6. Unsupervised Learning-Factor Analysis

☐ Principal Component Analysis - Intuition

- PCA is the technique that extract the most valuable characteristics or the combination of characteristics for clustering the data. (high variation)

- That is, PC can be thought of as the most important part of the variation that explains the variation of characteristics of entire samples

- Once we finish the extraction of PC, we can compare observations based on the extracted PC.

- For example, PCA might tell us that Gene3 is responsible for separating samples along the x-axis.

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |
| Gene 3 | 12 | 9 | 10 | 2.5 | 1.3 | 2 |
| Gene 4 | 5 | 7 | 6 | 2 | 4 | 7 |

PC 2 (4%)

PC 1 (91%)

# 7. Deep Learning

☐ Deep Learning

- analyzes data in multiple layers of learning

- simpler concepts => combining => complex and abstract

☐ Deep Learning vs. Classical Machine Learning

- Classical machine learning perform well on numerical prediction and analysis tasks

- But not at "intuitive" things (ex. recognizing images, analyzing texts)

- Goal of Deep Learning: perform tasks that are difficult for people to define, but easy to perform

# 7. Deep Learning

□ How it works

- Neuron: a simple calculator that computes **first** the weighted average of inputs given to it and **then** output the result if the value exceeds a pre-determined threshold

Figure 85: Example of a Neural Network with 2 hidden layers (4 neurons each), 3 input variables and 1 output variable



Input Layer

Hidden Layer 1

Hidden Layer 2

Output Layer

Source: J.P.Morgan Macro QDS

# 7. Deep Learning

□ Simple illustration : Hand-written number recognition



Each Neuron

$$\max\left[\sigma\left(w_1 a_1 + w_2 a_2 + \cdots + w_n a_n - 10\right), 0\right] \Leftrightarrow a^{(1)} = \sigma\left(W a^{(0)} + b\right)$$

↑ Sigmoid        ↑ bias

$$\min \| x_i - y_k \|$$

- Learning = Finding the appropriate weights and biases, i.e., $W$ & $b$
- MLP (Multilayer Perceptron)

# 7. Deep Learning

□ CNN (Convolutional Neural Network): Extension1

- Let's use convolutional filter as a layer to detect the edge of images well

# 7. Deep Learning

□ CNN (Convolutional Neural Network): Extension1

- One simple convolution layer
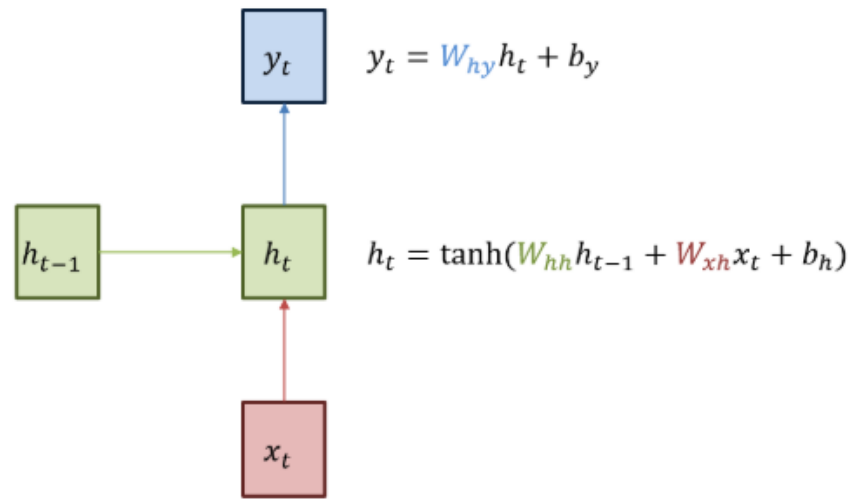
# 7. Deep Learning

□ Auto-encoder: Extension2

- Why not reverse-engineer the extracted feature and reconstruct the original input and compare it with the "real" original one?

Input layer        Hidden layers        Output layer

$x_1$   $a_1$   $a_1$   $a_1$   $\hat{x}_1$

$x_2$   $a_2$   $a_2$   $a_2$   $\hat{x}_2$

$x_3$   $a_3$   $a_3$   $a_3$   $\hat{x}_3$

$x_4$   $a_4$   $a_4$   $\hat{x}_4$

$x_5$   $\hat{x}_5$

$x_6$   $\hat{x}_6$

encoder        decoder

# 7. Deep Learning

□ RNN (Recurrent Neural Network): Extension3

- Why not incorporate historical contexts to the learning method?
  ex) predicting consecutive letter of word? h,e,l,l,o

- Vanishing gradient problem: cannot deal with long term memory
  (structure of the product, characteristics of tanh)

$$y_t = W_{hy}h_t + b_y$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

# 7. Deep Learning

□ LSTM (Long-Short Term Memory): Extension4

- Add special boxes to the RNN to adjust vanishing gradient problem
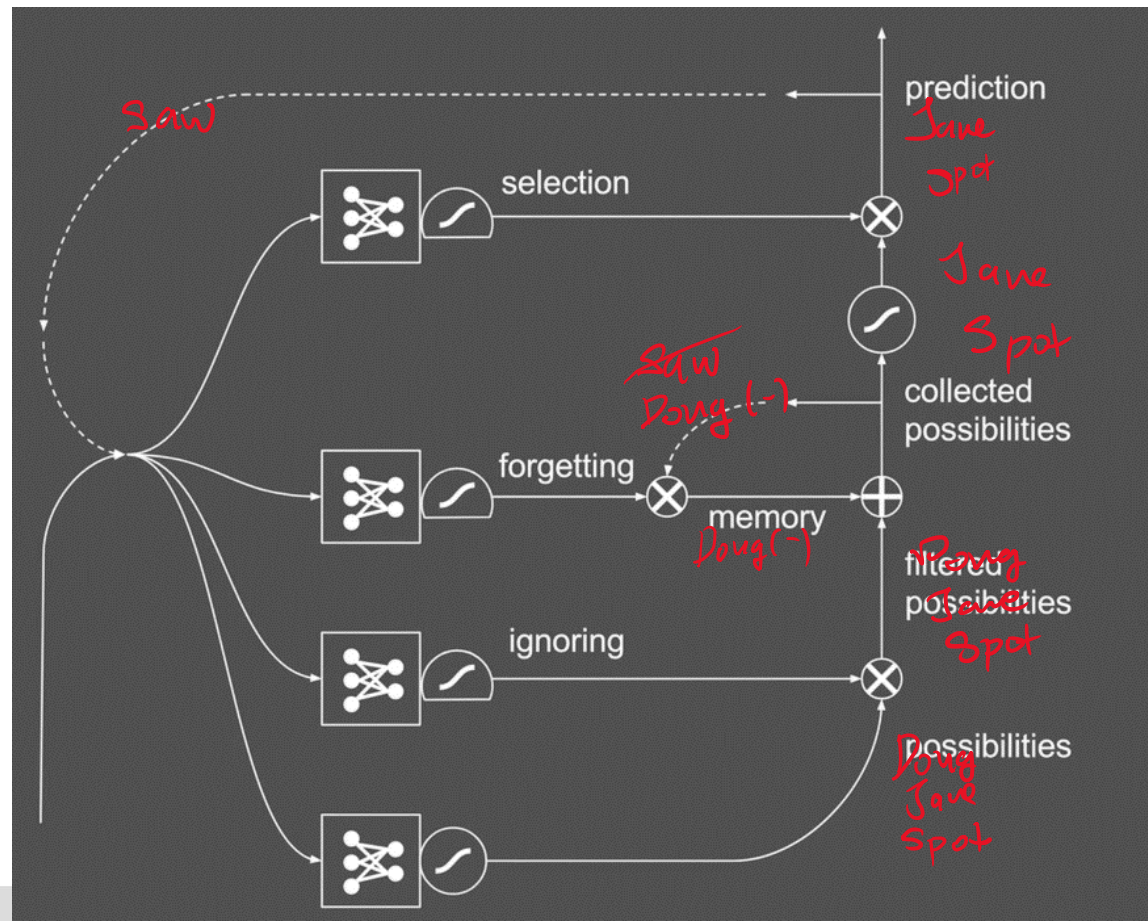
# 7. Deep Learning

□ LSTM Example: Let's make a children's book

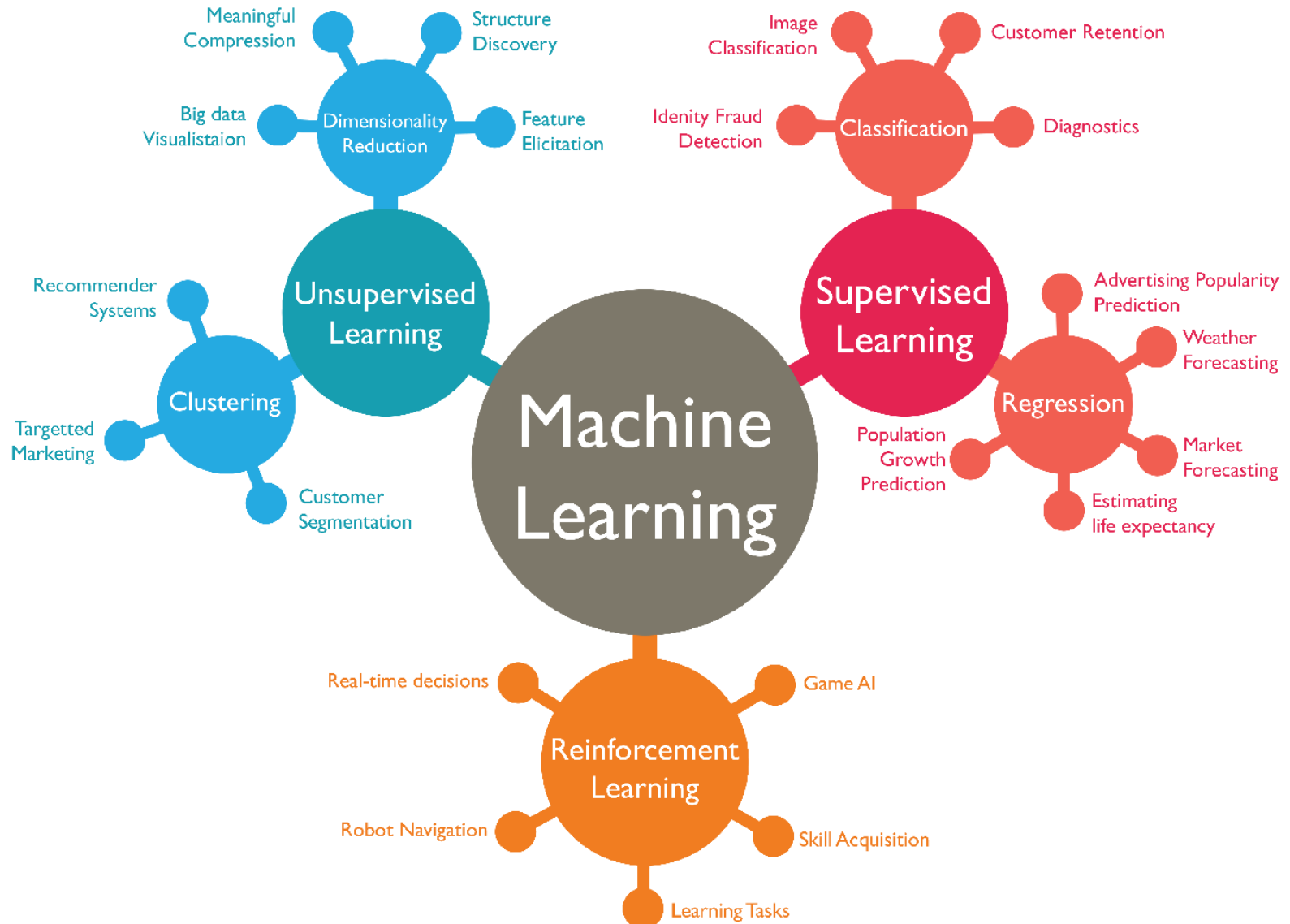- component: Jane, Spot, Doug, saw, .

- Goal:
  Jane saw Spot.
  Spot saw Doug.
  Doug saw Spot. ...

- Possible mistakes:
  Jane saw Jane.
  Doug saw Spot saw ...
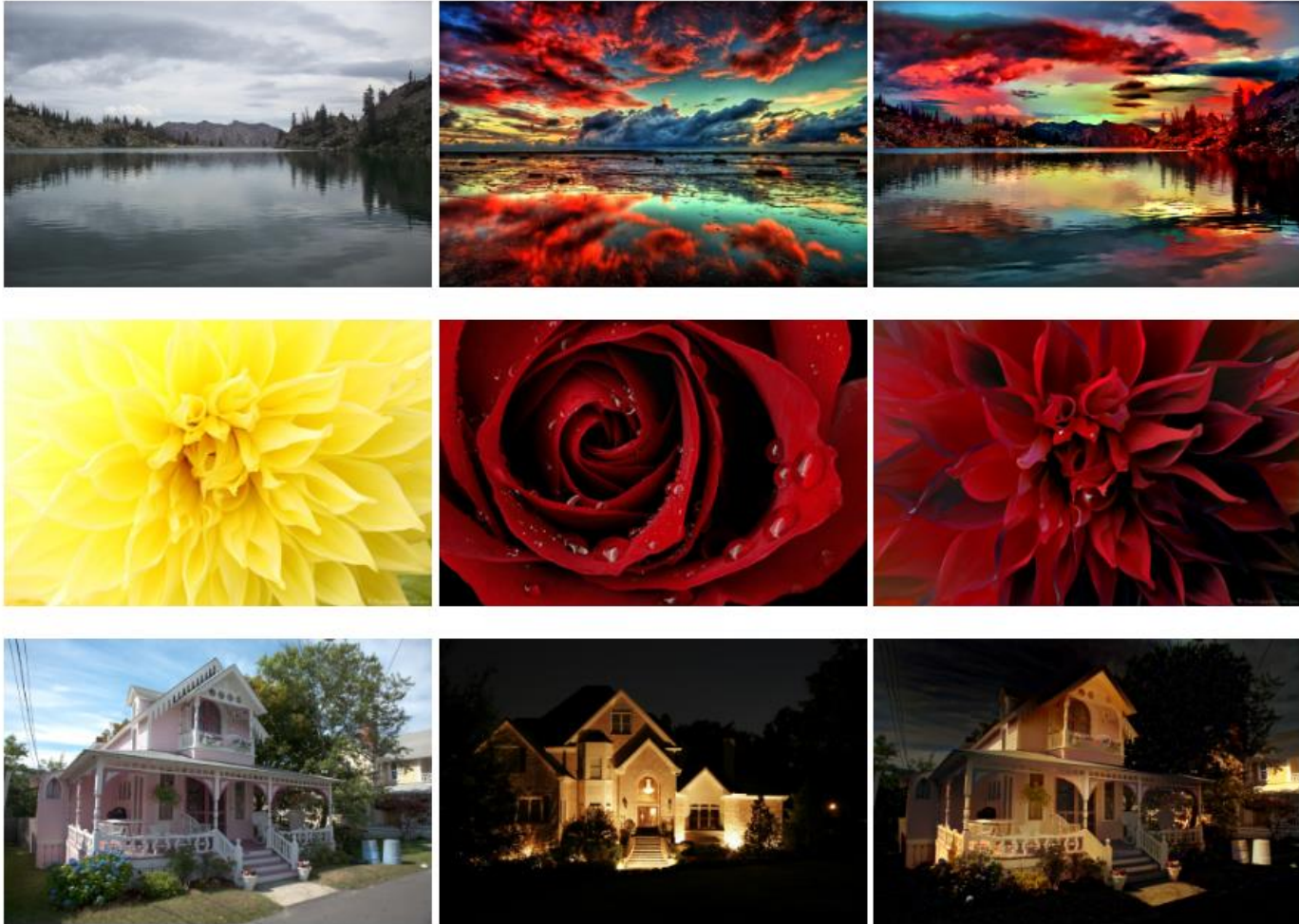  Doug.

- Current stage:
  Jane saw Spot.
  Doug ...

# 7. Deep Learning

□ LSTM Example: Let's make a children's book

- component: Jane, Spot, Doug, saw, .

- Goal:
  Jane saw Spot.
  Spot saw Doug.
  Doug saw Spot. ...

- Possible mistakes:
  Jane saw Jane.
  Doug saw Spot saw ...
  Doug.

- Current stage:
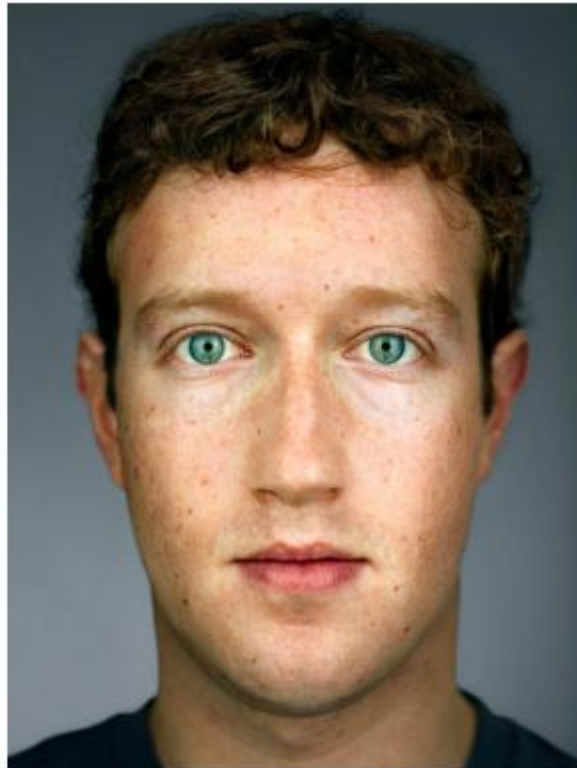  Jane saw Spot.
  Doug saw ...

# 8. Application

□ deep-photo-styletransfer: input, style, output

# 8. Application

□ deep-photo-styletransfer: input, style, output

  - how to map style to the original contents?
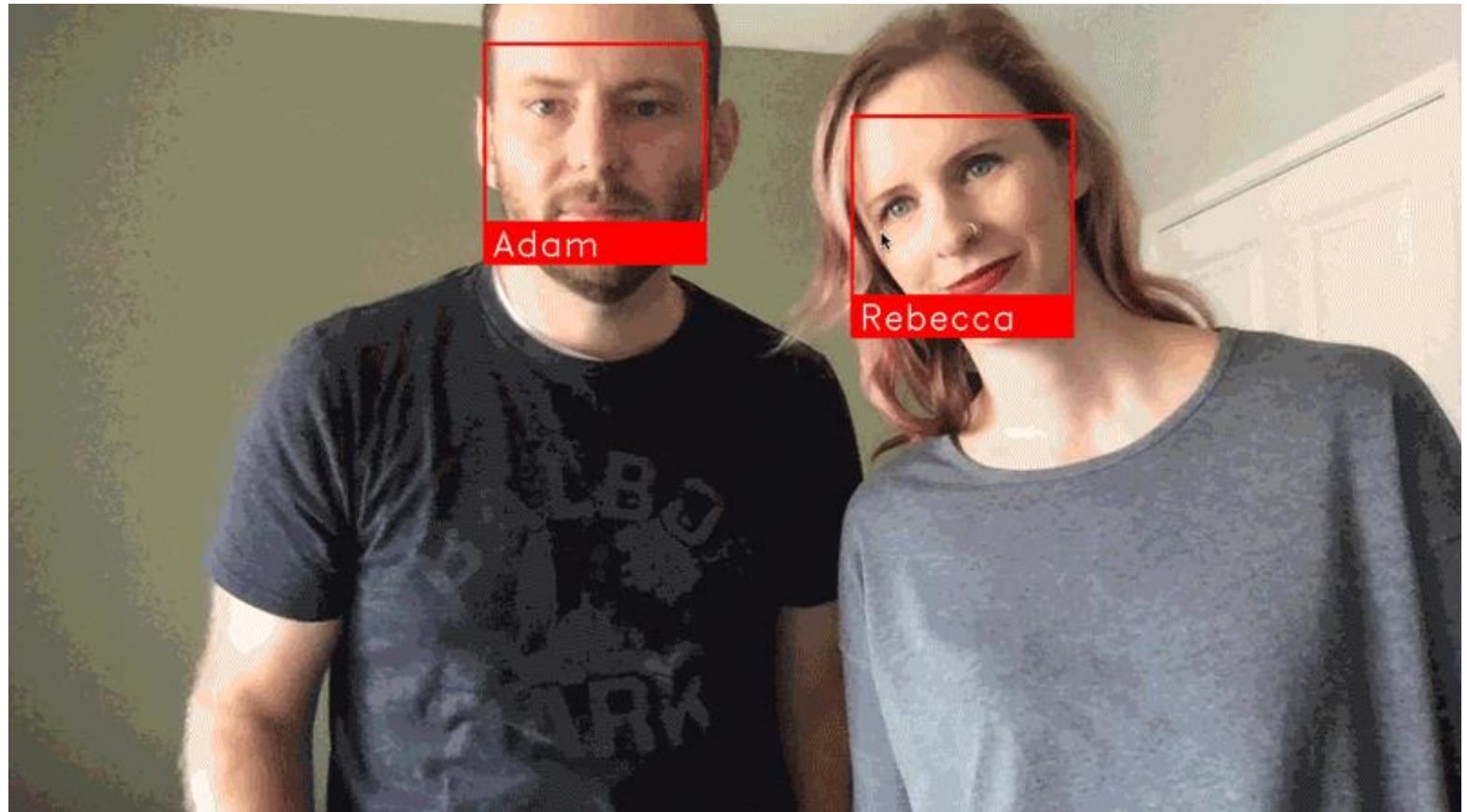  - extension of CNN

# 8. Application

☐ deep-photo-styletransfer extension



Freiburger Münsterplatz

# 8. Application

☐ Face recognition

# 8. Application

□ Face recognition

- Overcome this!



(a)  (b)

(c)  (d)

# 8. Application

□ Face recognition

- Face detection: find the position of the faces

- Face alignment: scale and crop face images to compare with references

- Face representation: pixel values of face images->feature vector

- Face matching: compare the similarity score
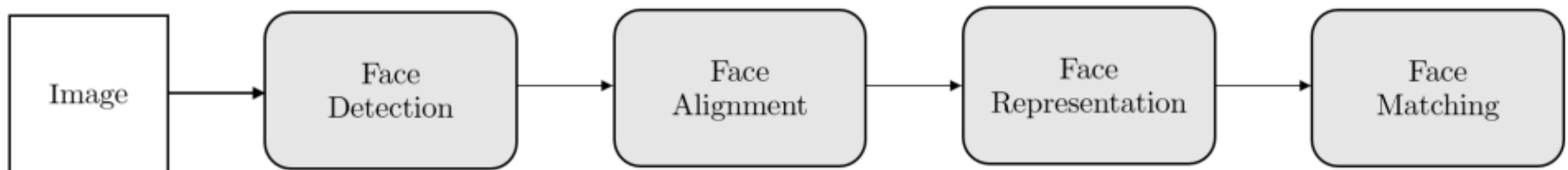
- This time again, CNN



Fig. 2: Face recognition building blocks.

# 8. Application

□ Magenta: Art and Music composer

- RNN

# 8. Application

□ Google: Neutral Networks and 'Machine that Dream'



Selection of images created by Google's neural network

- These days, it's probably easier to list areas that Google isn't working on with machine learning.
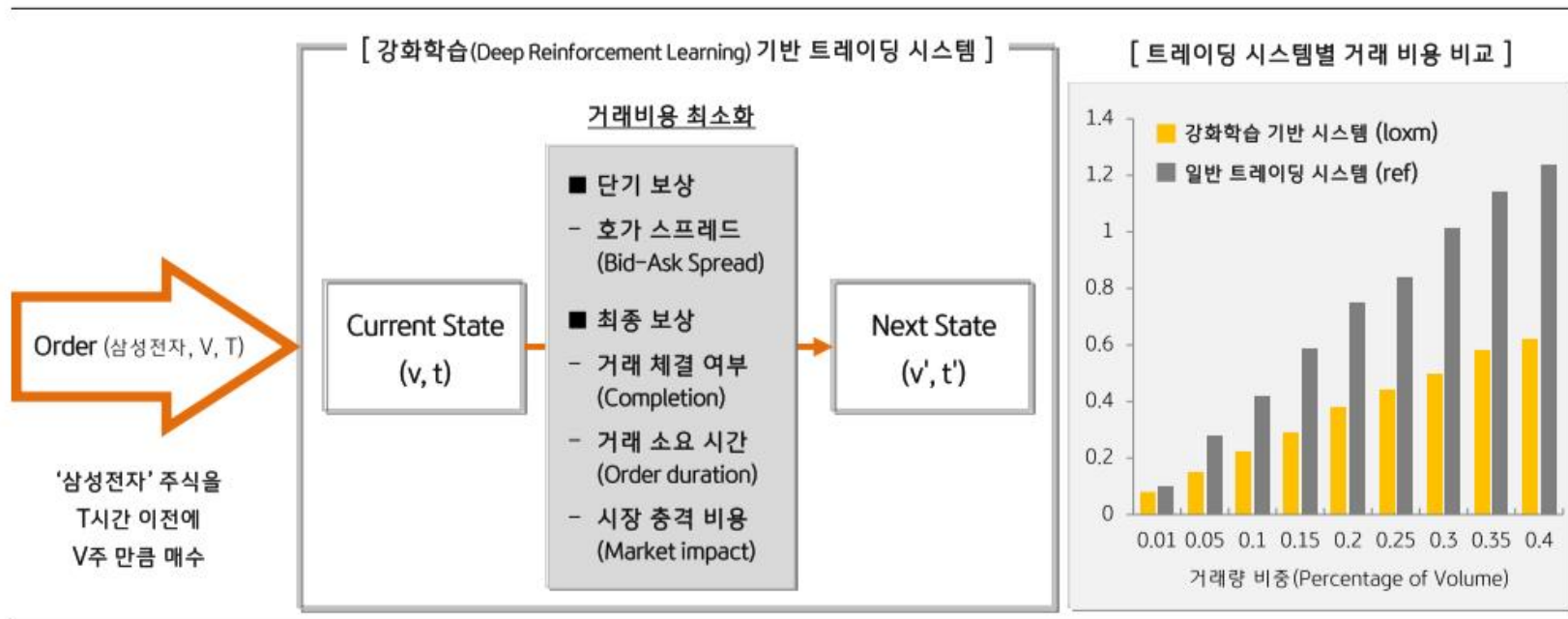
# 8. Application

□ Deep Voice Conversion: Neural Voice Cloning

# 8. Application

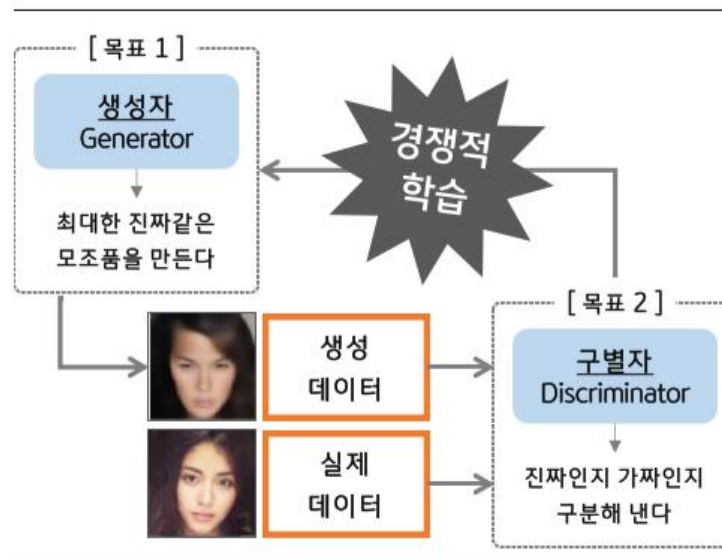☐ LOXM: JP Morgan Trading system

- minimize the cost function with the inputs as bid-ask spread, the success of transactions, duration, market impact costs

- great performance compared to the benchmark system



자료: 'Active Learning in Trading Algorithms(JPMorgan, 2016)' 및 KB경영연구소

# 8. Application

□ Fraud Detection

- Making a competition between two model: Generator and Discriminator

- Generator generates fraud-like data
- Discriminator discriminates real one with the generated one

- Good choice when not enough data



자료: Tech M 및 KB경영연구소

# 8. Application

□ Credit Rating System

  - need to know "why"

  - make a small perturbation to the specific variable so that figure out
    which factor has the most influence for categorization

  - FICO: Credit rating company
    "Reasons Reporter"

  - Abnormality detection and
    Credit Rating



원본 이미지
P(개구리) = 0.54

| 그림 요소 변형 | | | |
|---|---|---|---|
| 확률 | 0.85 | 0.00001 | 0.52 |

회색으로 가려진 부분은 해당 이미지가 입력값에서 삭제된 것
얼굴과 눈 부분이 있을 때 개구리로 분류할 확률이 증가(85%)하고
없을 때 거의 0으로 감소, '얼굴과 눈'이 개구리 분류의 핵심 요인

자료: Marco Tulio Ribeiro et al.(O'REILLY, 2016)

# 00. Epilogue

THANK YOU

FOR LISTENING