

An Overview of Machine Learning



Yuan-Fu Liao

National Taipei University of Technology

Outline & Content

- What is machine learning?
- Learning system model
- Training and testing
- Performance
- Algorithms
- Machine learning structure
- What are we seeking?
- Learning techniques
- Applications
- Conclusion

What is Machine Learning?

- A branch of **artificial intelligence**, concerned with the design and development of **algorithms** that allow computers to evolve behaviors based on empirical data.
- As intelligence requires knowledge, it is necessary for the computers to **acquire knowledge**.

From Human Learning to Machine Learning

Human learning: acquiring skill

with experience accumulated from observations



Machine learning: acquiring skill

with experience accumulated/computed from data

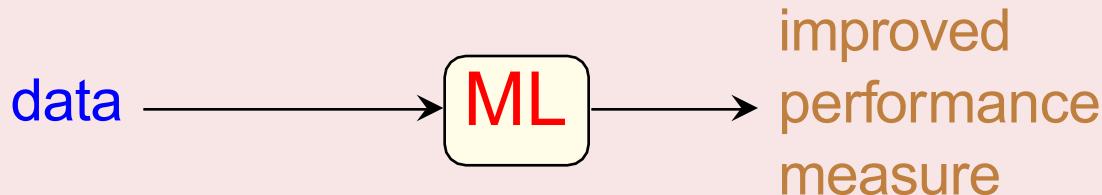


What is skill?

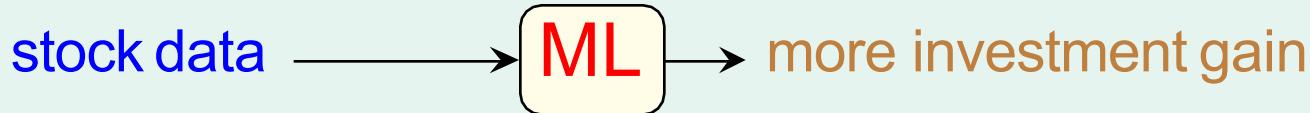
A More Concrete Definition

Skill \Leftrightarrow improve some performance measure (e.g. prediction accuracy)

machine learning: improving some performance measure with experience **computed** from **data**



An Application in Computational Finance



Why use machine learning?

Yet Another Application: Tree Recognition



- ‘define’ trees and hand-program: **difficult**
- learn from data (observations) and recognize: a **3-year-old can do so**
- ‘ML-based tree recognition system’ can be **easier to build** than hand-programmed system

ML: an **alternative route** to build complicated systems

The Machine Learning Route

ML: an **alternative route** to build complicated systems

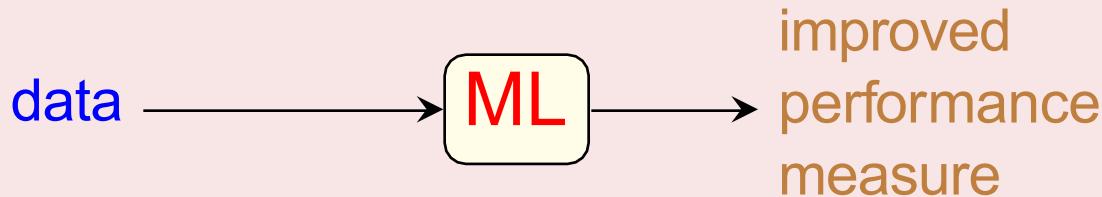
Some Use Scenarios

- when human cannot program the system manually
 - navigating on Mars
- when human cannot ‘define the solution’ easily
 - speech/visual recognition
- when needing rapid decisions that humans cannot do
 - high-frequency trading
- when needing to be user-oriented in a massive scale
 - consumer-targeted marketing

Give a **computer** a fish, you feed it for a day;
teach it how to fish, you feed it for a lifetime. :-)

Key Essence of Machine Learning

machine learning: improving some performance measure with experience **computed** from **data**



- ① exists some ‘underlying pattern’ to be learned
 - so ‘performance measure’ can be improved
- ② but **no** programmable (easy) **definition**
 - so ‘ML’ is needed
- ③ somehow there is **data** about the pattern
 - so ML has some ‘inputs’ to learn from

key essence: help decide whether to use ML

Daily Needs: Food, Clothing, Housing, Transportation



1 Food (Sadilek et al., 2013)

- **data**: Twitter data (words + location)
- **skill**: tell food poisoning likeliness of restaurant properly

2 Clothing (Abu-Mostafa, 2012)

- **data**: sales figures + client surveys
- **skill**: give good fashion recommendations to clients

3 Housing (Tsanas and Xifara, 2012)

- **data**: characteristics of buildings and their energy load
- **skill**: predict energy load of other buildings closely

4 Transportation (Stallkamp et al., 2012)

- **data**: some traffic sign images and meanings
- **skill**: recognize traffic signs accurately

ML is everywhere!

Education



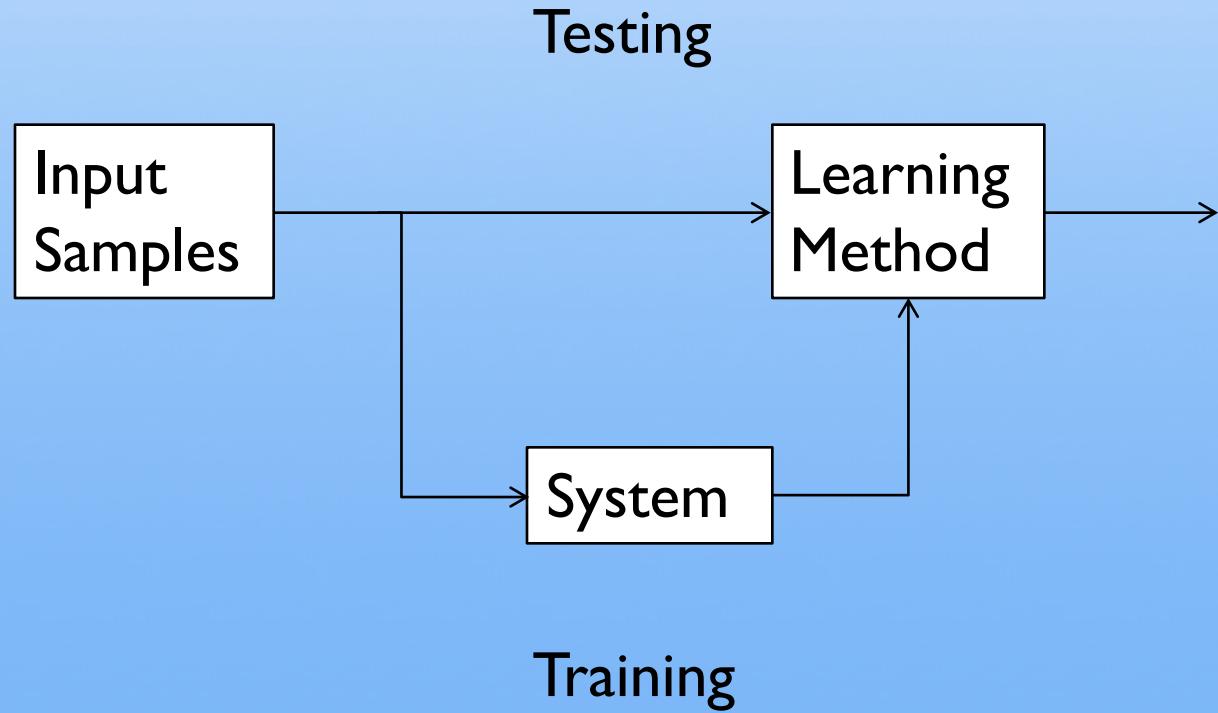
- **data**: students' records on quizzes on a Math tutoring system
- **skill**: predict whether a student can give a correct answer to another quiz question

A Possible ML Solution

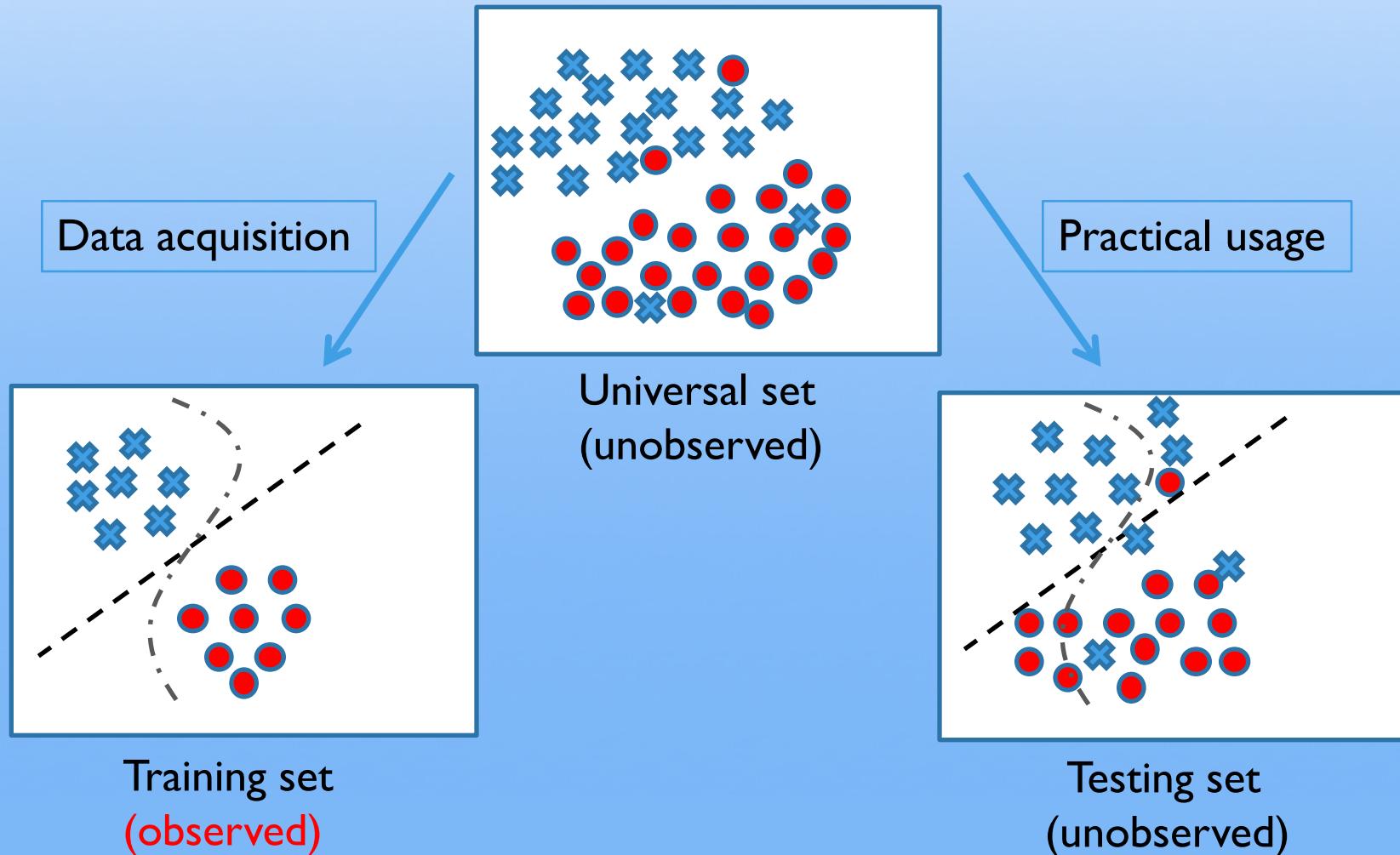
answer correctly \approx [recent **strength** of student > **difficulty** of question]

- give ML **9 million records** from **3000 students**
- ML determines (**reverse-engineers**) **strength** and **difficulty** automatically

Learning system model

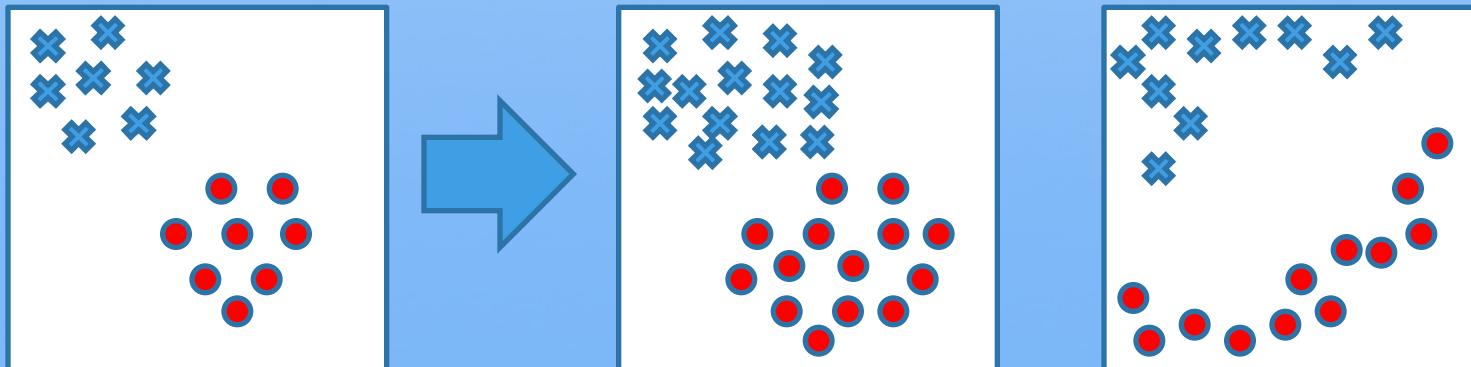


Training and testing



Training and testing

- Training is the process of making the system able to learn.
- No free lunch rule:
 - Training set and testing set come from the same distribution
 - Need to make some assumptions or bias



Performance

- There are several factors affecting the performance:
 - **Types of training** provided
 - The form and extent of any initial **background knowledge**
 - The **type of feedback** provided
 - The **learning algorithms** used
- Two important factors:
 - Modeling
 - Optimization

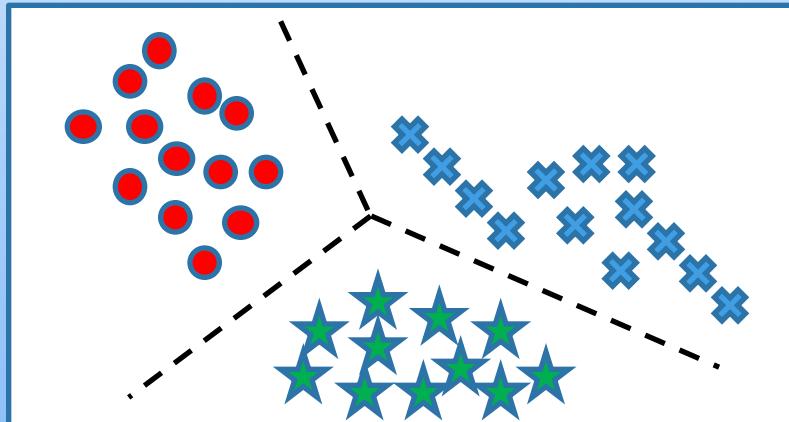
Algorithms

- The success of machine learning system also depends on the algorithms.
- The algorithms control the search to find and build the knowledge structures.
- The learning algorithms should extract useful information from training examples.

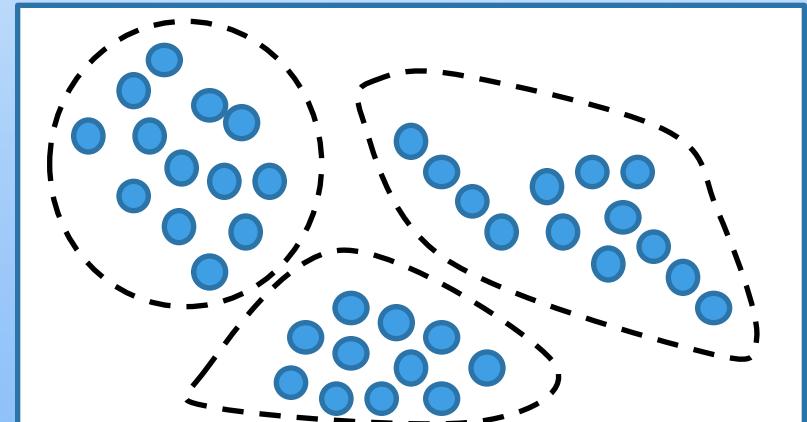
Algorithms

- **Supervised learning** ($\{x_n \in R^d, y_n \in R\}_{n=1}^N$)
 - Prediction
 - Classification (discrete labels), Regression (real values)
- **Unsupervised learning** ($\{x_n \in R^d\}_{n=1}^N$)
 - Clustering
 - Probability distribution estimation
 - Finding association (in features)
 - Dimension reduction
- **Semi-supervised learning**
- **Reinforcement learning**
 - Decision making (robot, chess machine)

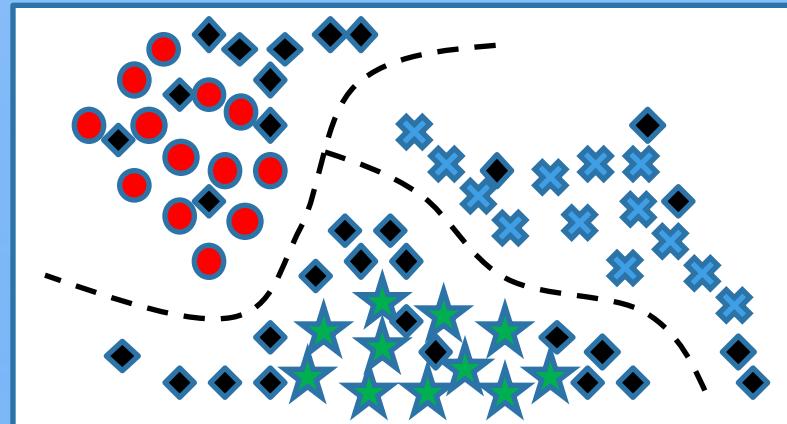
Algorithms



Supervised learning



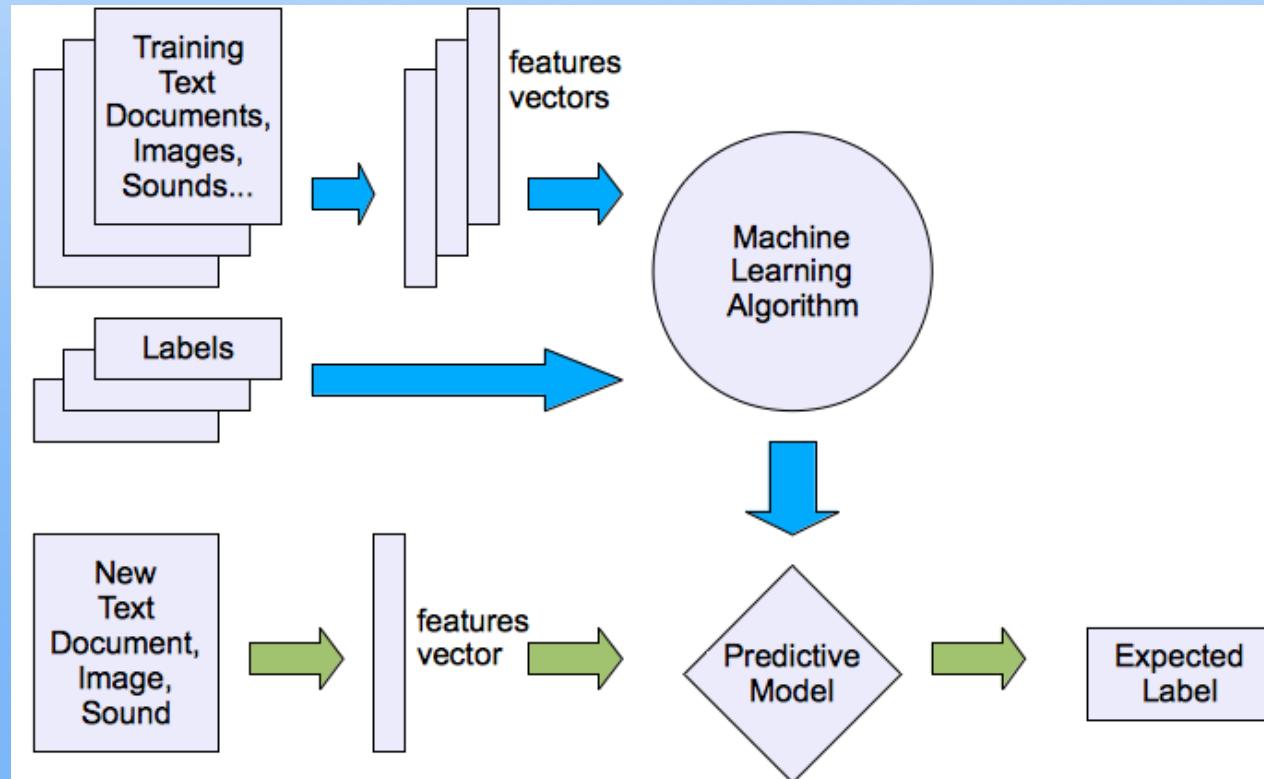
Unsupervised learning



Semi-supervised learning

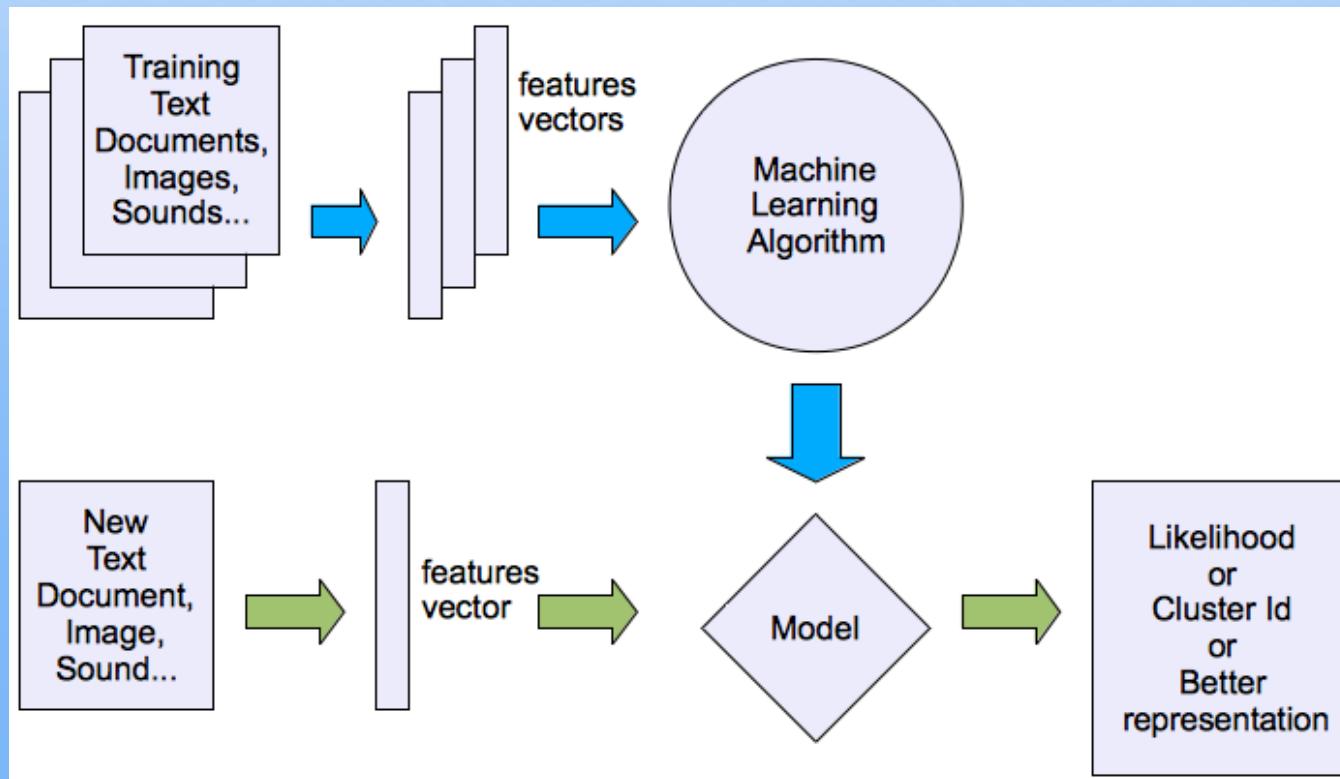
Machine learning structure

□ Supervised learning



Machine learning structure

❑ Unsupervised learning



What are we seeking?

- Supervised: Low E-out or maximize probabilistic terms

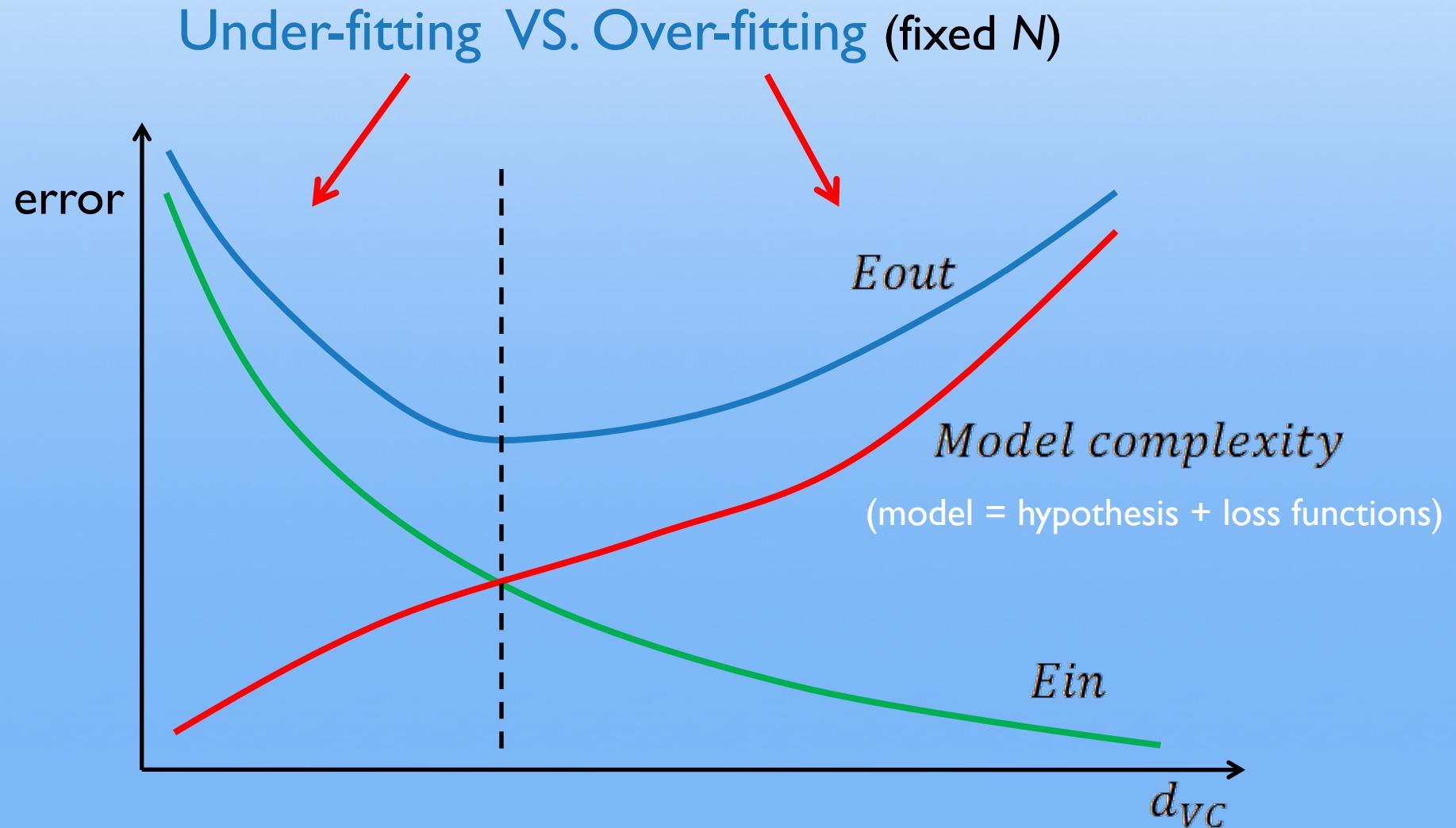
$$error = \frac{1}{N} \sum_{n=1}^N [y_n \neq g(x_n)]$$

E-in: for training set
E-out: for testing set

$$Eout(g) \leq Ein(g) \pm O\left(\sqrt{\frac{d_{VC}}{N} \ln N}\right)$$

- Unsupervised: Minimum quantization error, Minimum distance, MAP, MLE(maximum likelihood estimation)

What are we seeking?

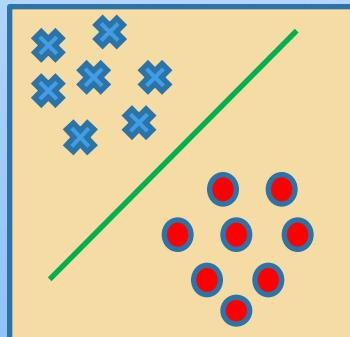


Learning techniques

- Supervised learning categories and techniques
 - **Linear classifier** (numerical functions)
 - **Parametric** (Probabilistic functions)
 - Naïve Bayes, Gaussian discriminant analysis (GDA), Hidden Markov models (HMM), Probabilistic graphical models
 - **Non-parametric** (Instance-based functions)
 - K-nearest neighbors, Kernel regression, Kernel density estimation, Local regression
 - **Non-metric** (Symbolic functions)
 - Classification and regression tree (CART), decision tree
 - **Aggregation**
 - Bagging (bootstrap + aggregation), Adaboost, Random forest

Learning techniques

- Linear classifier



$$g(x_n) = \text{sign}(w^T x_n)$$

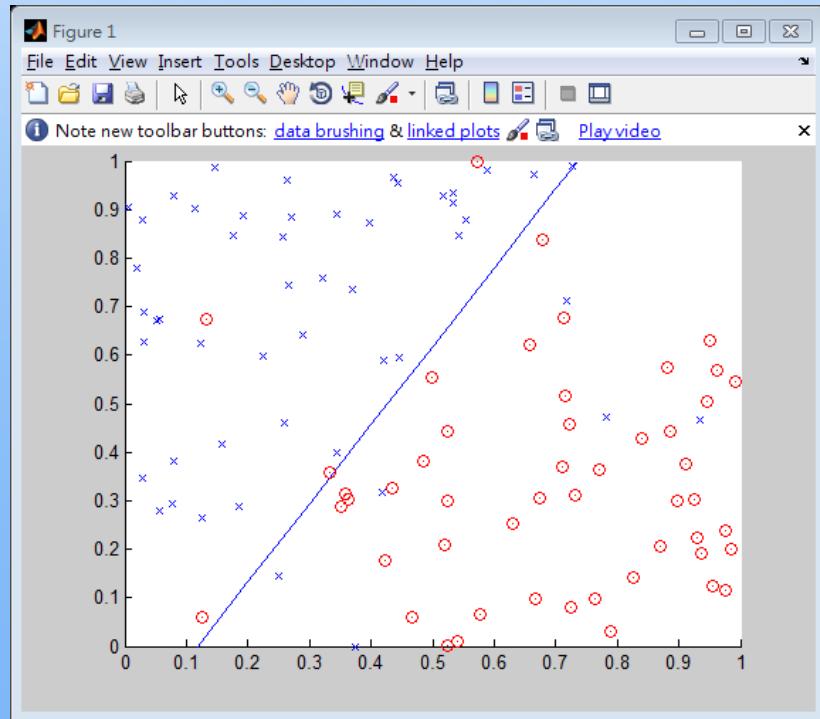
, where w is an d -dim vector (learned)

- Techniques:

- Perceptron
- Logistic regression
- Support vector machine (SVM)
- Ada-line
- Multi-layer perceptron (MLP)

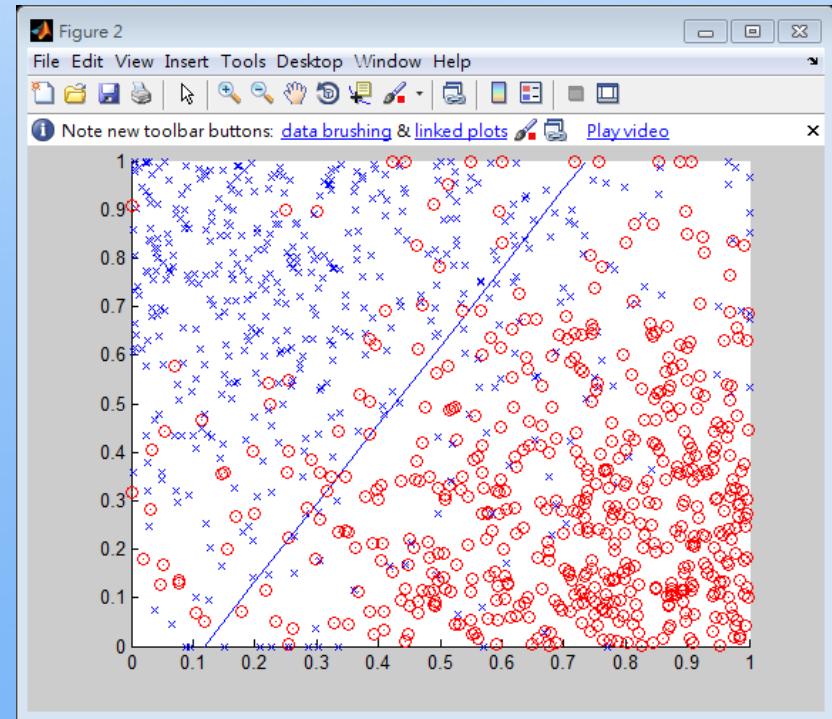
Learning techniques

Using perceptron learning algorithm(PLA)



Training

Error rate: 0.10

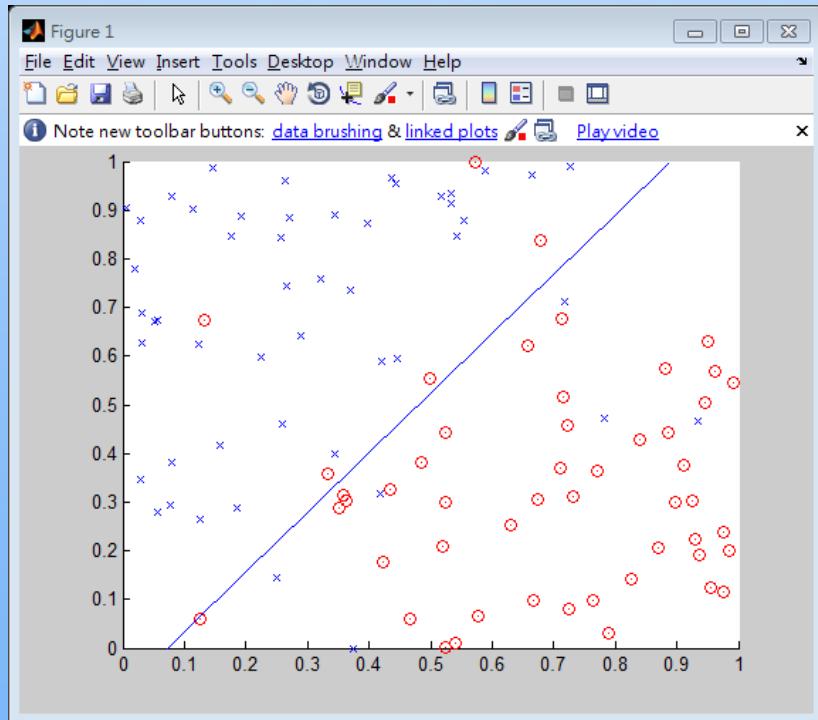


Testing

Error rate: 0.156

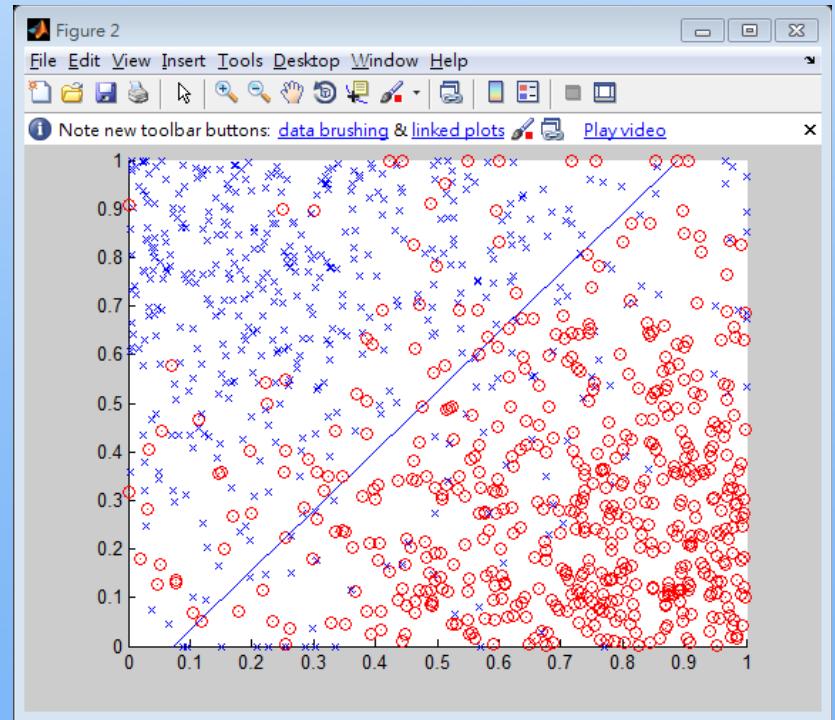
Learning techniques

Using logistic regression



Training

Error rate: 0.11

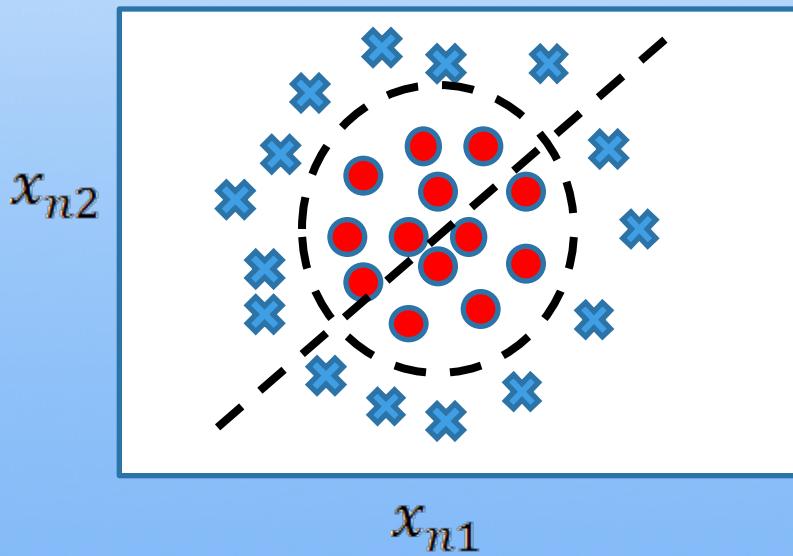


Testing

Error rate: 0.145

Learning techniques

- Non-linear case



$$x_n = [x_{n1}, x_{n2}]$$



$$x_n = [x_{n1}, x_{n2}, x_{n1} * x_{n2}, x_{n1}^2, x_{n2}^2]$$
$$g(x_n) = \text{sign}(w^T x_n)$$

- Support vector machine (SVM):
 - Linear to nonlinear: **Feature transform** and **kernel function**

Learning techniques

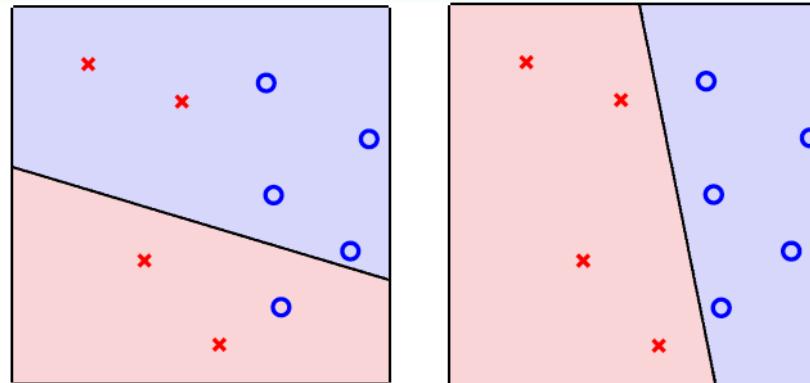
- Unsupervised learning categories and techniques
 - **Clustering**
 - K-means clustering
 - Spectral clustering
 - **Density Estimation**
 - Gaussian mixture model (GMM)
 - Graphical models
 - **Dimensionality reduction**
 - Principal component analysis (PCA)
 - Factor analysis

Example



Perceptrons in \mathbb{R}^2

$$h(\mathbf{x}) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$



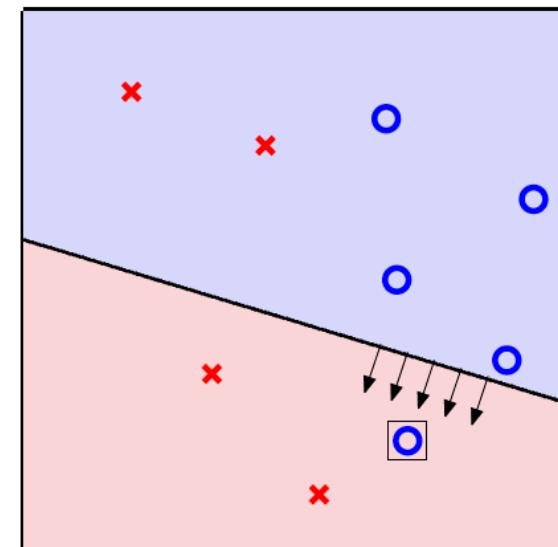
- customer features \mathbf{x} : points on the plane (or points in \mathbb{R}^d)
- labels y : $\circ (+1)$, $\times (-1)$
- hypothesis h : **lines** (or hyperplanes in \mathbb{R}^d)
—**positive** on one side of a line, **negative** on the other side
- different line classifies customers differently

perceptrons \Leftrightarrow **linear (binary) classifiers**

Select g from \mathcal{H}

\mathcal{H} = all possible perceptrons, $g = ?$

- want: $g \approx f$ (hard when f unknown)
- almost necessary: $g \approx f$ on \mathcal{D} , ideally
 $g(\mathbf{x}_n) = f(\mathbf{x}_n) = y_n$
- difficult: \mathcal{H} is of **infinite** size
- idea: start from some g_0 , and ‘correct’ its mistakes on \mathcal{D}



will represent g_0 by its weight vector \mathbf{w}_0

Perceptron Learning Algorithm

start from some \mathbf{w}_0 (say, $\mathbf{0}$), and ‘correct’ its mistakes on \mathcal{D}

For $t = 0, 1, \dots$

- ➊ find a mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$

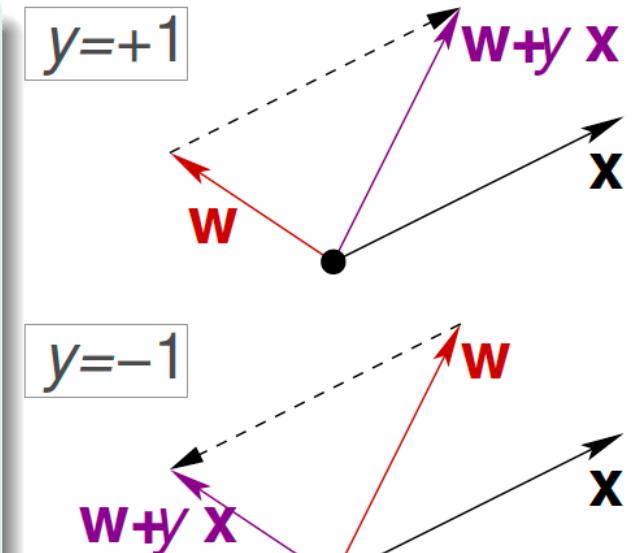
$$\text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_{n(t)}$$

- ➋ (try to) correct the mistake by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$$

... until no more mistakes

return last \mathbf{w} (called \mathbf{w}_{PLA}) as g



That's it!

—A fault confessed is half redressed. :-)

Practical Implementation of PLA

start from some \mathbf{w}_0 (say, $\mathbf{0}$), and ‘correct’ its mistakes on \mathcal{D}

Cyclic PLA

For $t = 0, 1, \dots$

- 1 find **the next** mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$

$$\text{sign} \left(\mathbf{w}_t^T \mathbf{x}_{n(t)} \right) \neq y_{n(t)}$$

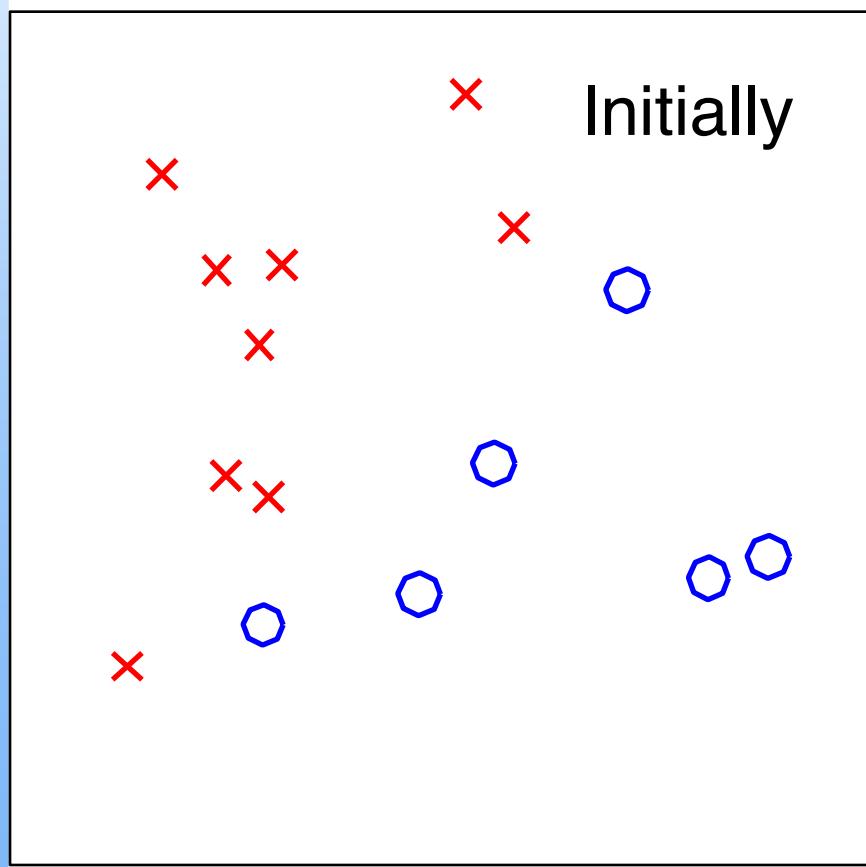
- 2 correct the mistake by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$$

... until **a full cycle of not encountering mistakes**

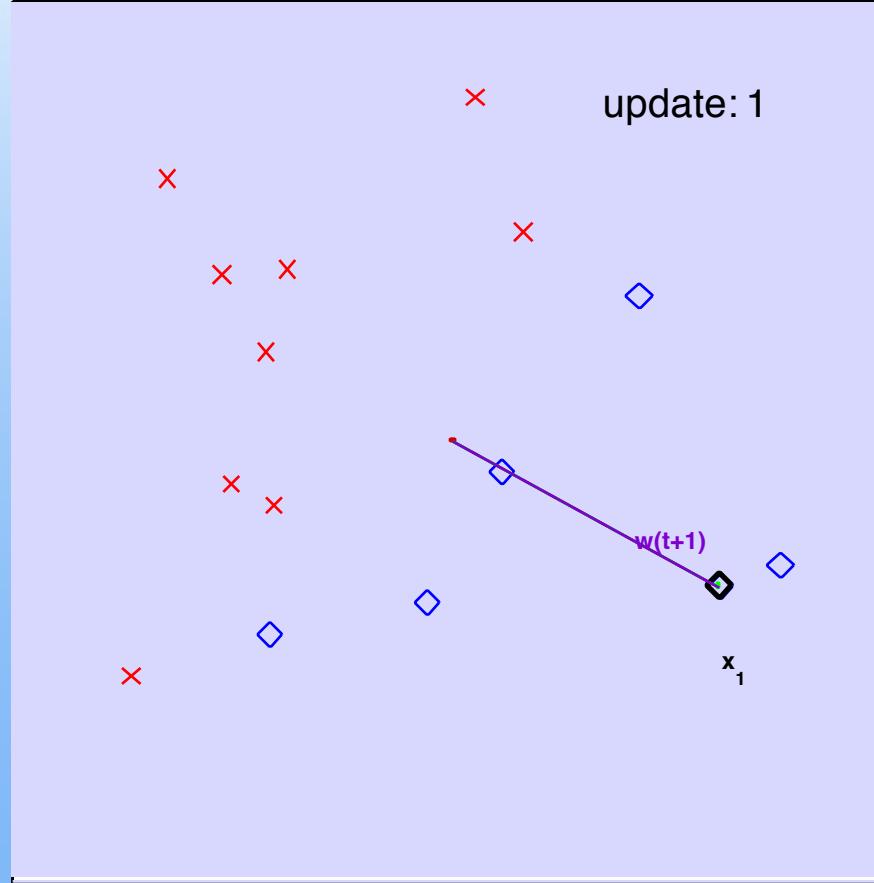
next can follow naïve cycle $(1, \dots, N)$
or precomputed random cycle

Seeing is Believing



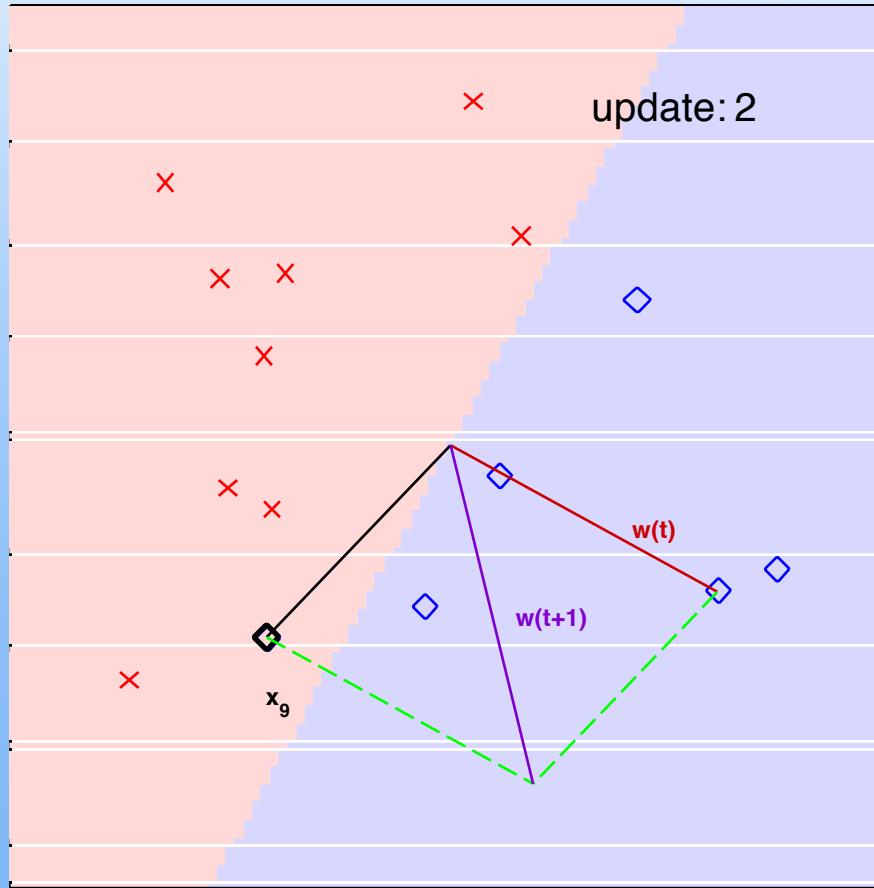
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



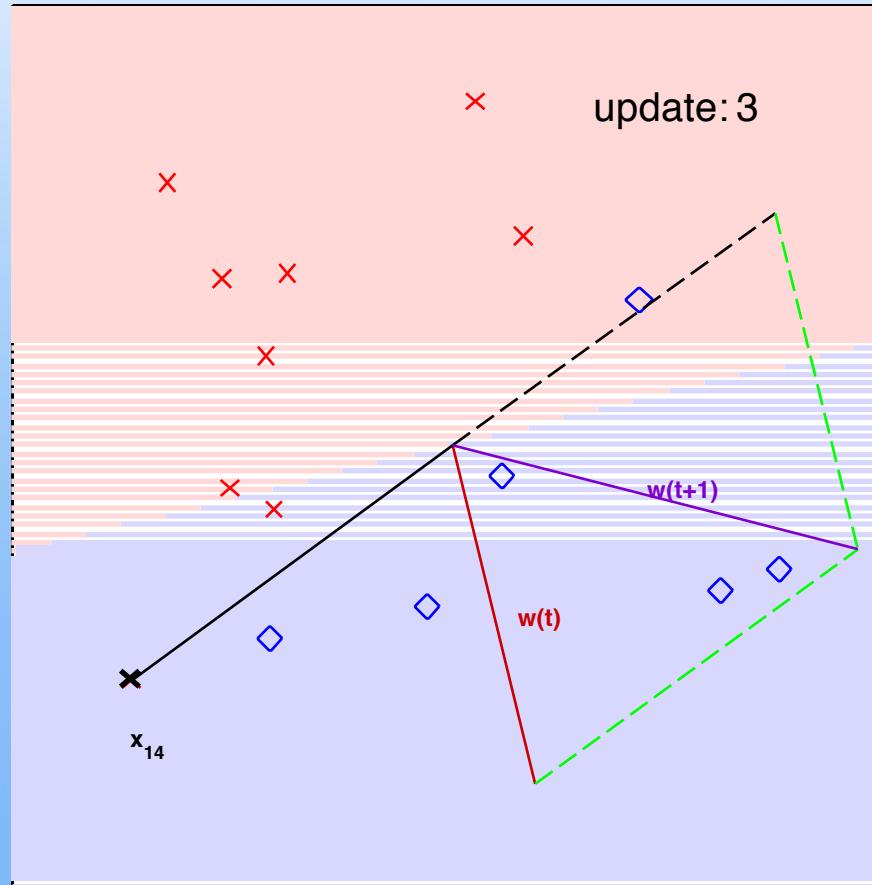
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



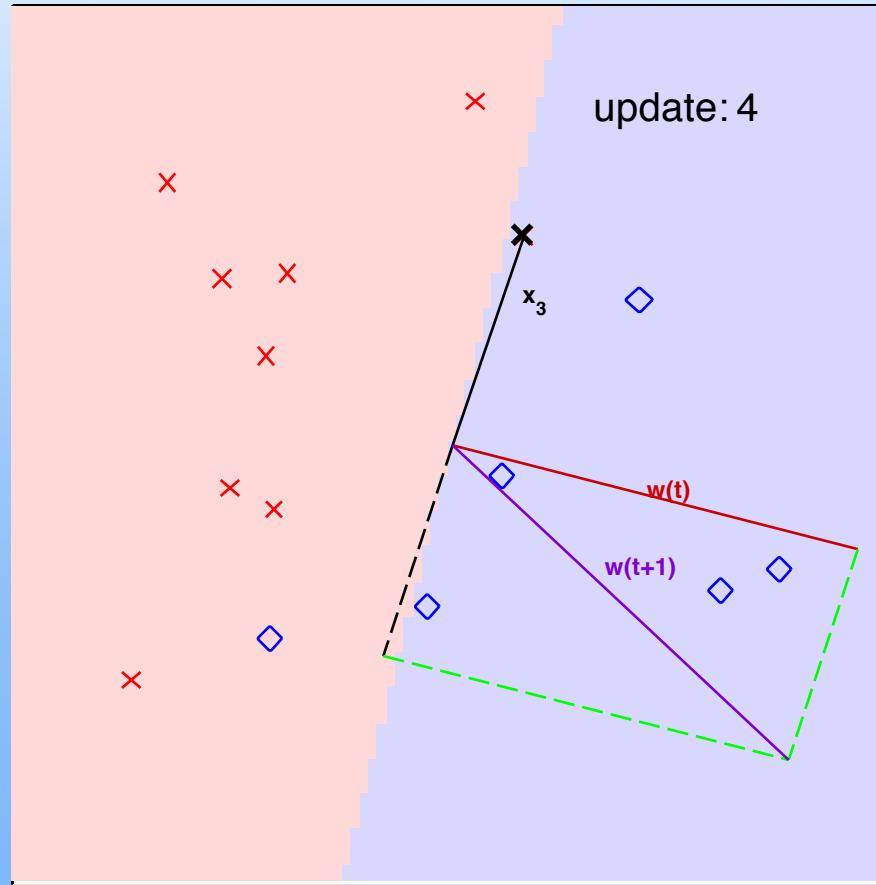
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



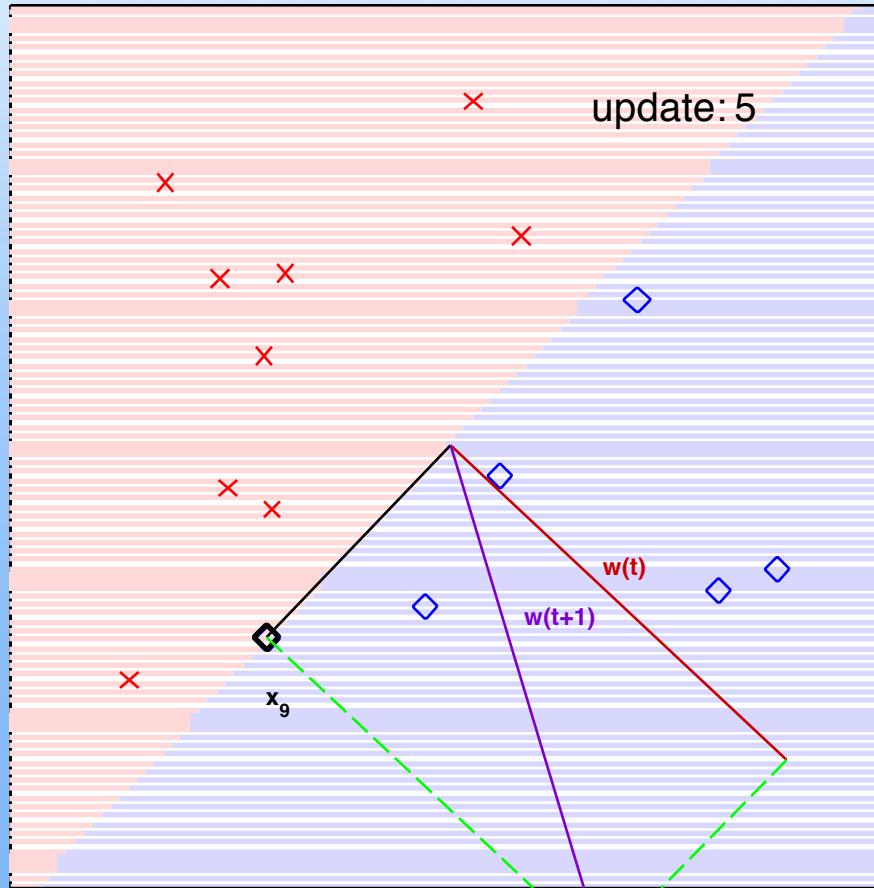
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



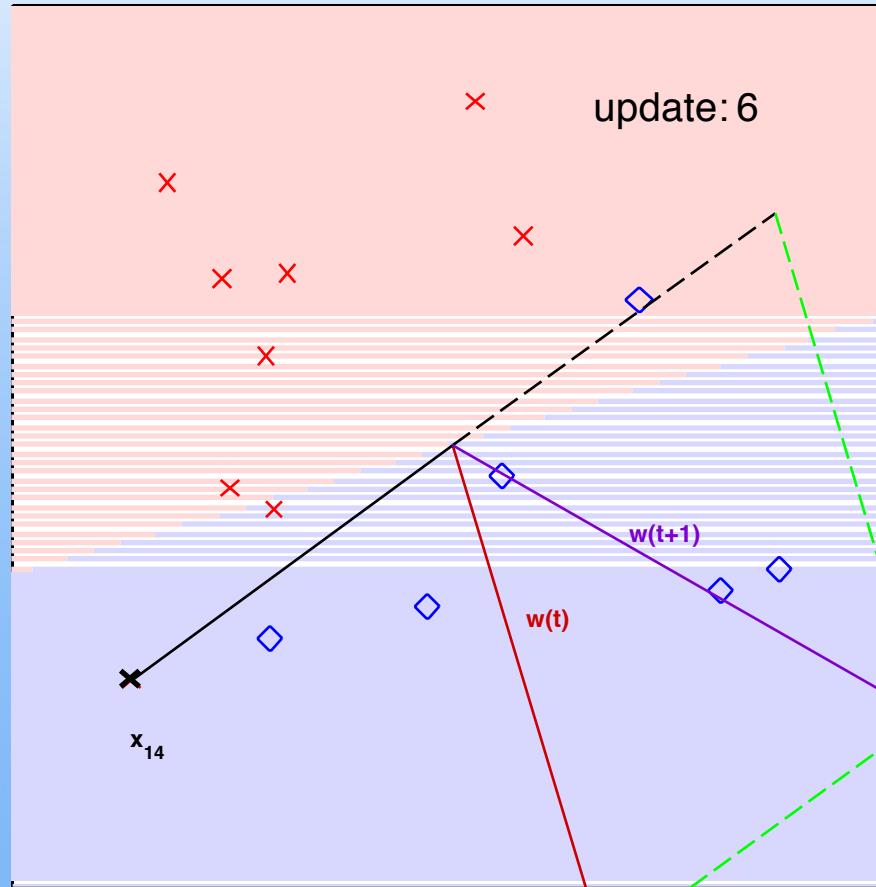
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



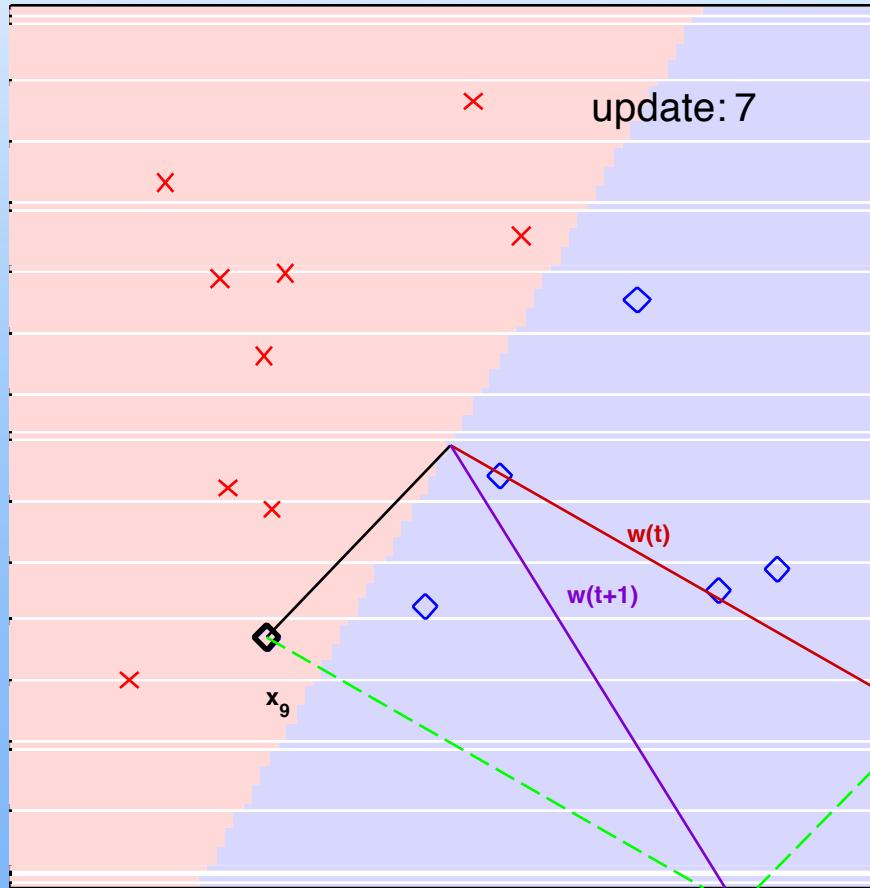
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



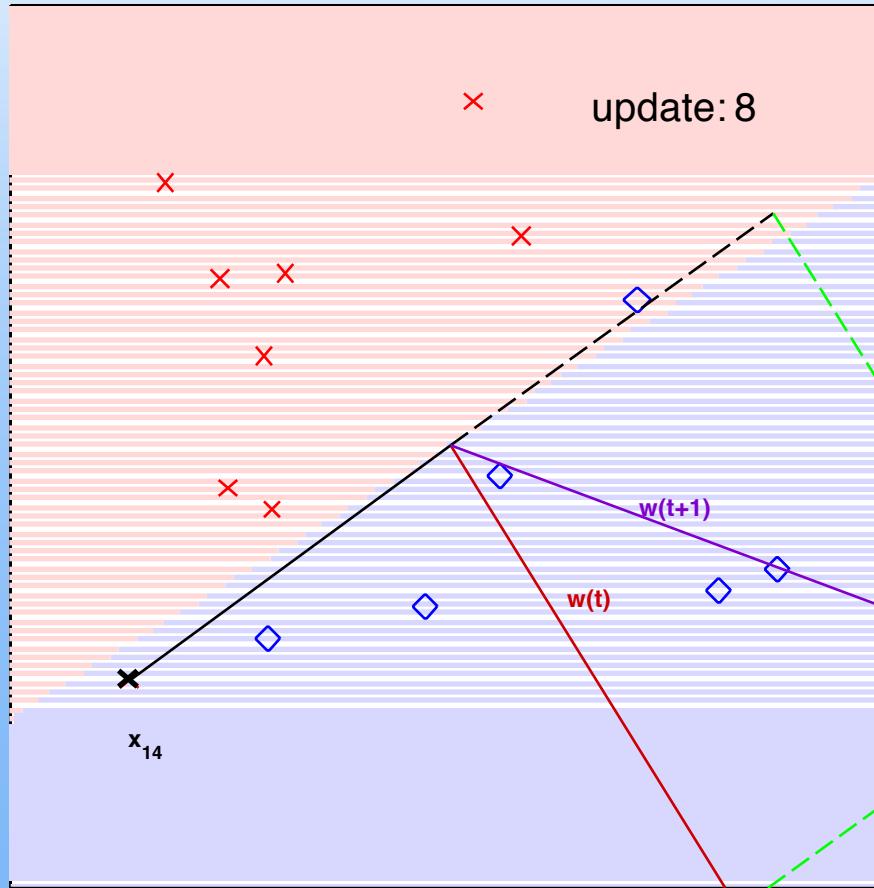
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



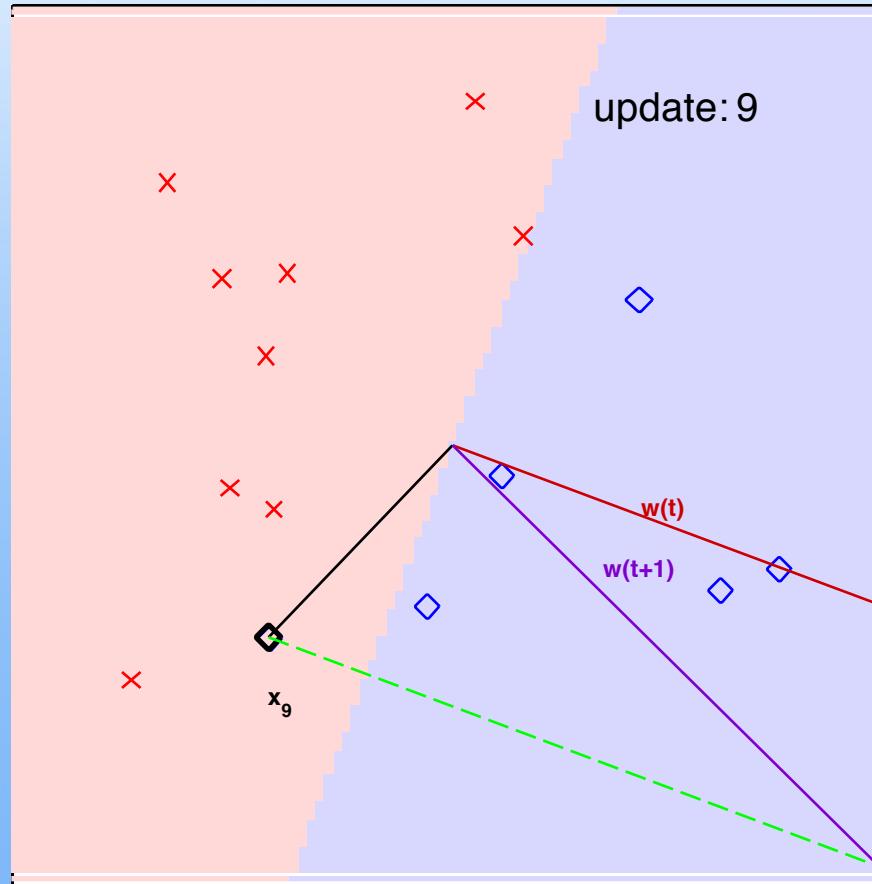
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



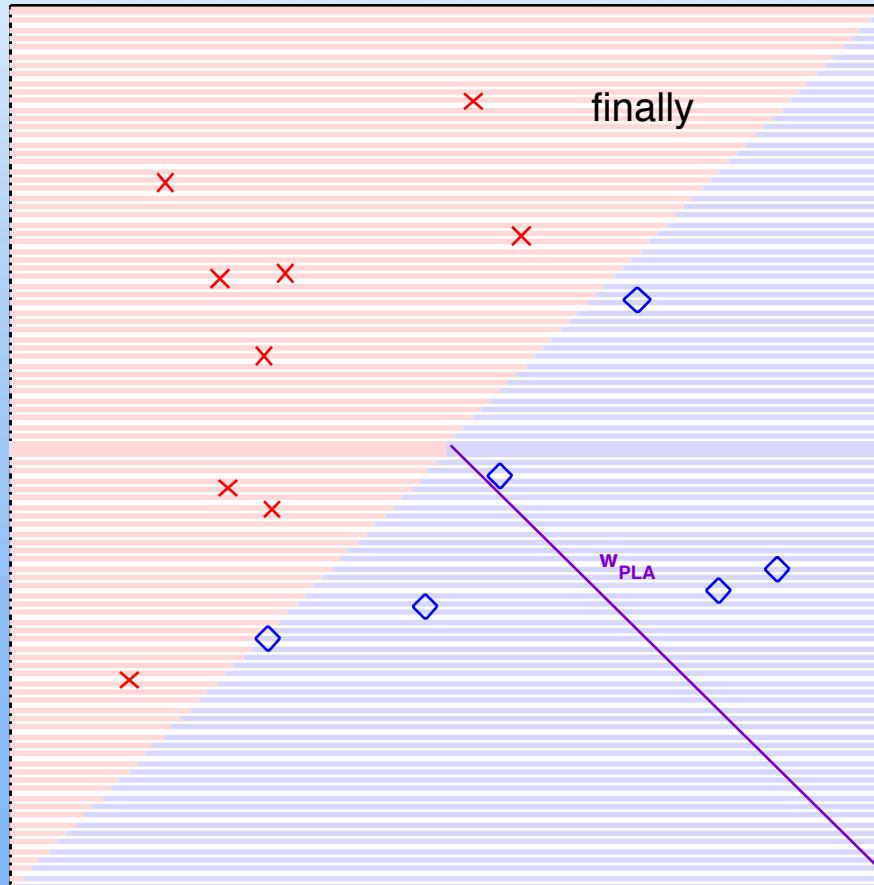
worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing



worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

Seeing is Believing

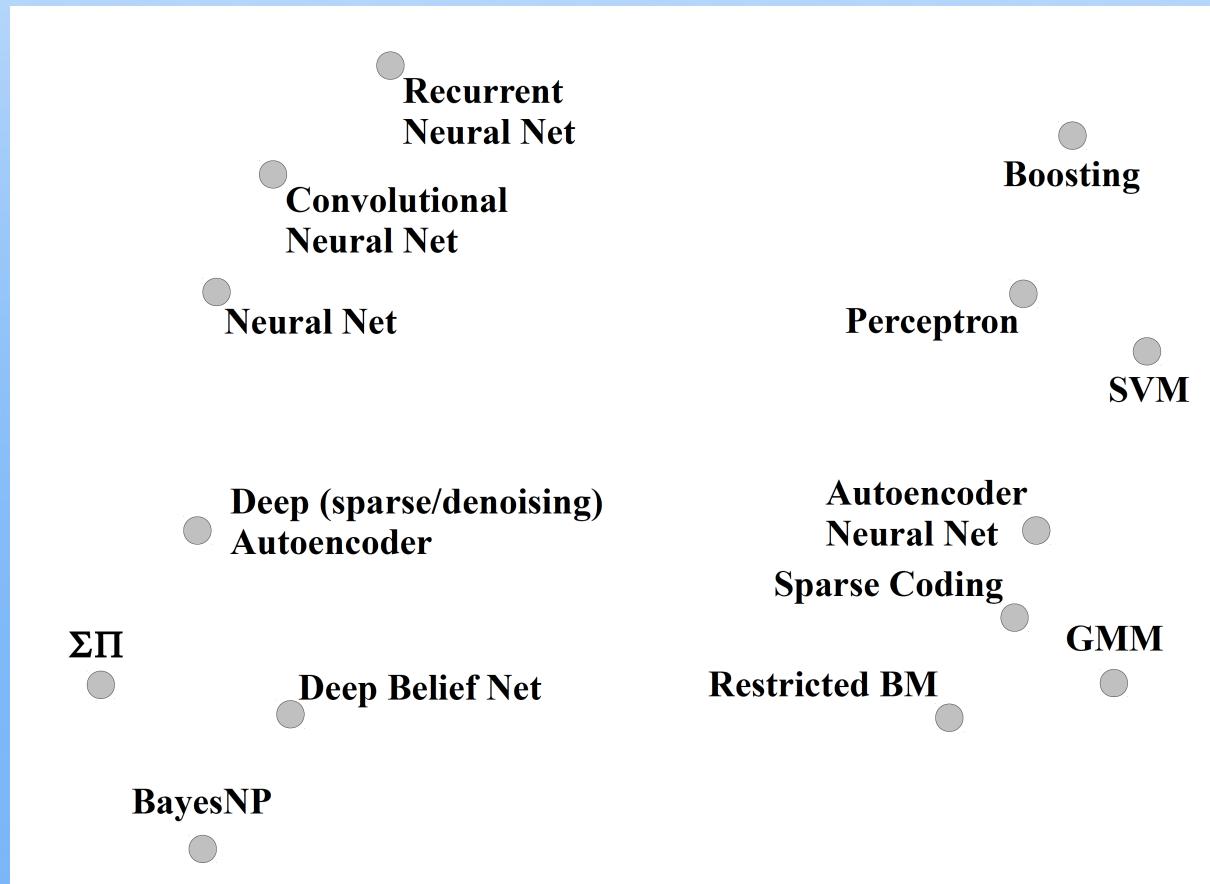


worked like a charm with < 20 lines!!
(note: made $x_i \gg x_0 = 1$ for visual purpose)

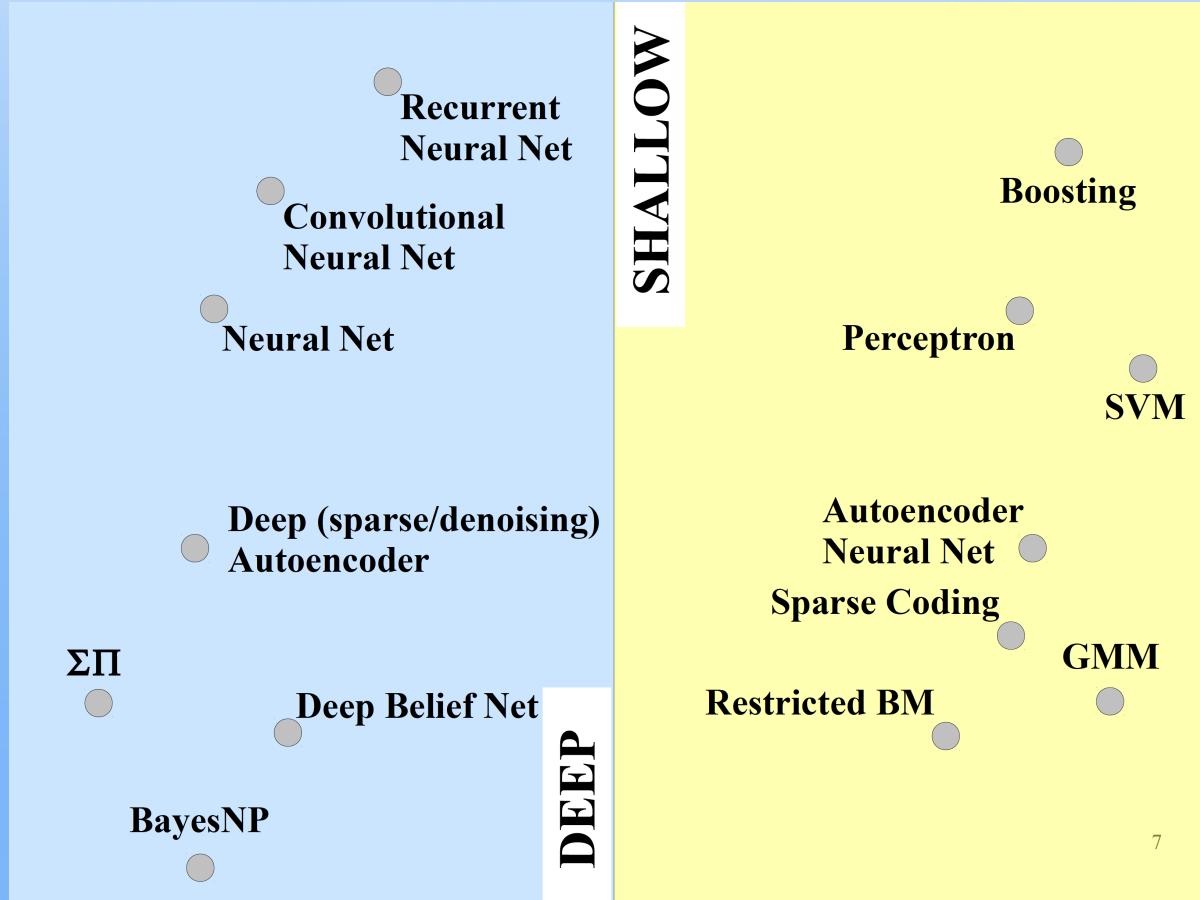
Category



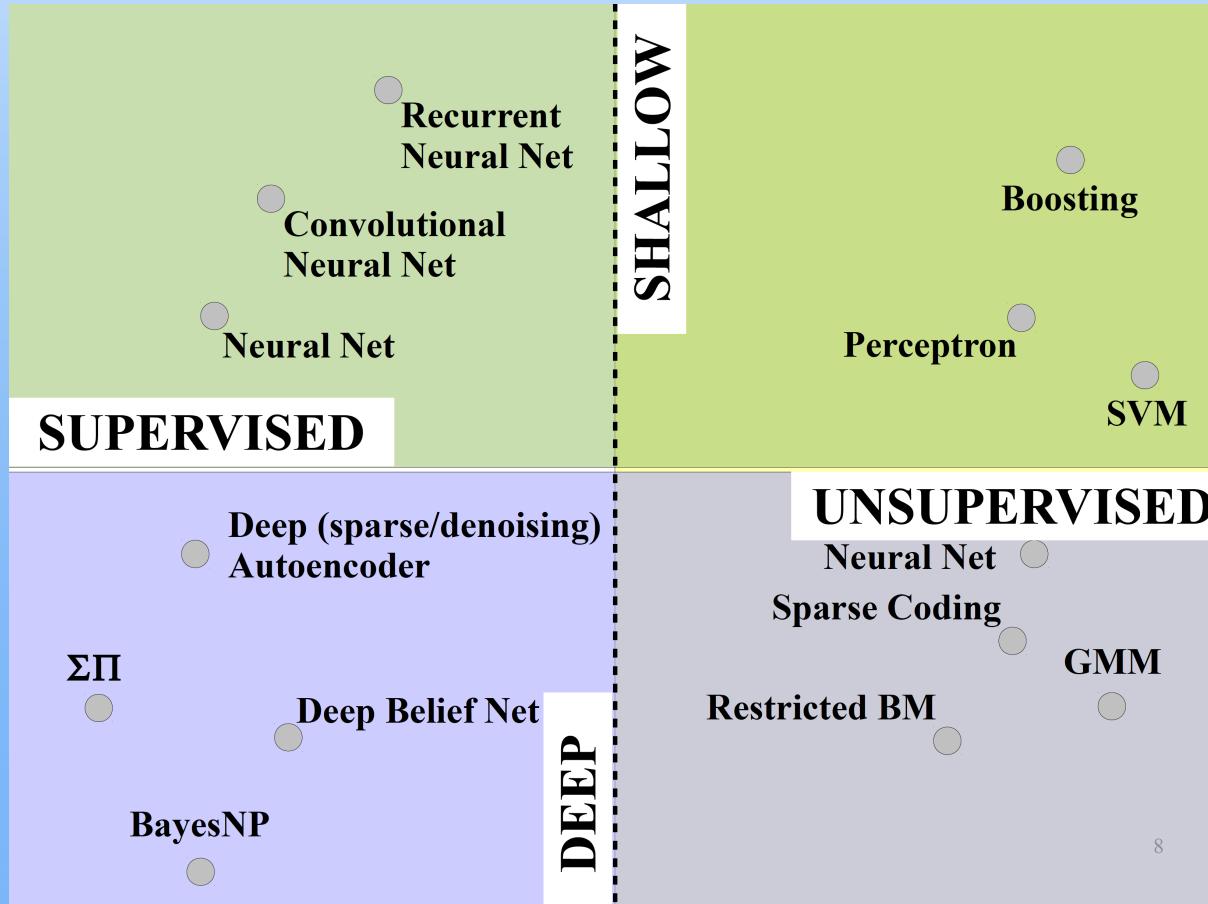
Machine Learning Algorithms



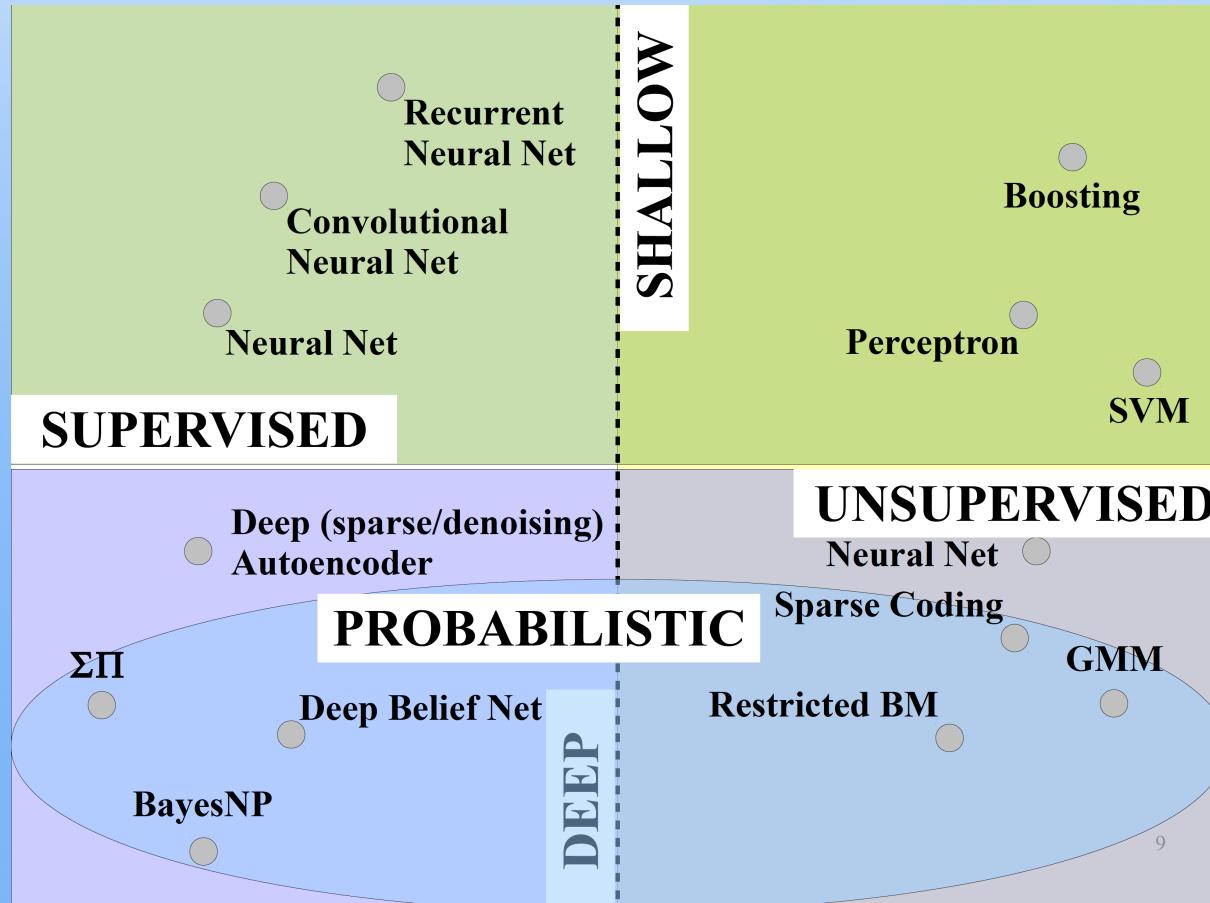
Machine Learning Algorithms



Machine Learning Algorithms



Machine Learning Algorithms



Applications

- Face detection
- Object detection and recognition
- Image segmentation
- Multimedia event detection
- Economical and commercial usage

Conclusion

We have a simple overview of some techniques and algorithms in machine learning. Furthermore, there are more and more techniques apply machine learning as a solution. In the future, machine learning will play an important role in our daily life.