

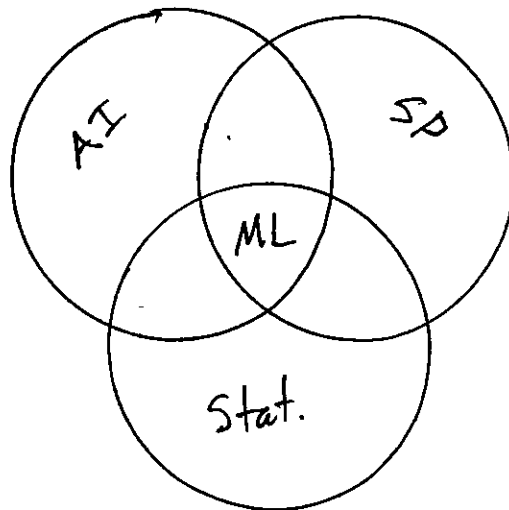
STATISTICAL MACHINE LEARNING

Machine Learning

Machine learning is a multi-disciplinary field concerned with the study of algorithms that learn from examples.

ML theory and methodology emerged historically out of three areas: artificial intelligence, signal processing, and statistics.

By now, the best practices of these areas have spread to the others, and ML is its own well-defined field.



○ Learning from examples \Leftrightarrow learning from data

Notation for data:

X random variable

x realization of X

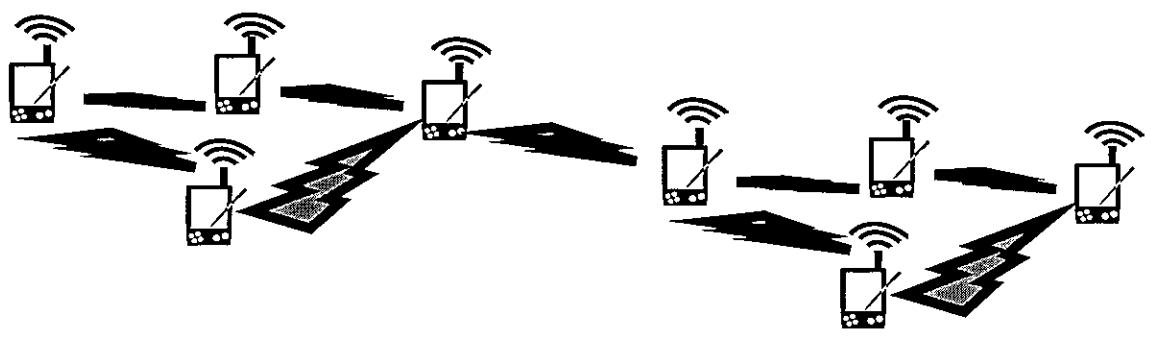
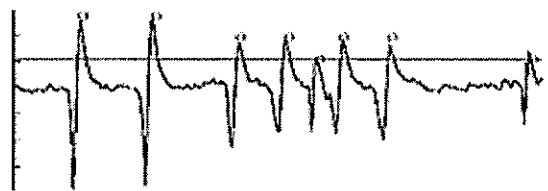
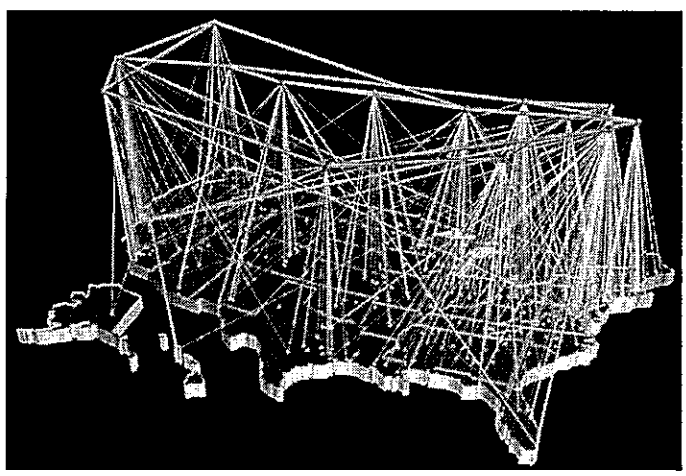
Typically, $X \in \mathbb{R}^d$.

X represents a measurement or observation of some natural or man-made phenomenon, and may be called a . The coordinates of X are called .

- pattern \nearrow
- signal
- feature vector
- input
- instance
- independent variable

- features \nearrow
- attributes
- predictors
- covariates

0	1	2	3	4
5	6	7	8	9



Statistical Machine Learning

In most applications, there is some uncertainty or randomness inherent in the data.

3 3 3 3 3

In statistical machine learning, we will

- view a pattern X as a random variable
- use the tools of probability and statistics to provide a mathematical framework for
 - posing machine learning problems
 - formulating solutions to those problems.
 - evaluation

The following terms are often used interchangeably:

- Machine learning
- Statistical machine learning
- Statistical learning
- Pattern recognition
- Multivariate data analysis

Types of Learning Problems

We will consider various paradigms for statistical learning.

Supervised Learning

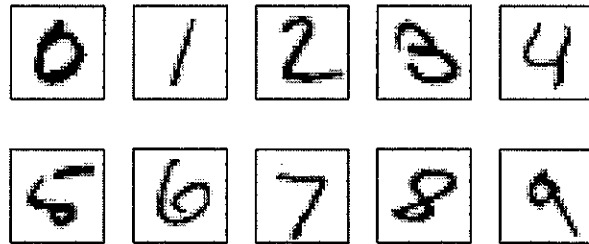
In addition to patterns

$$X_1, \dots, X_n$$

we also have access to variables

$$Y_1, \dots, Y_n.$$

Examples of patterns



Training data (suppose correct labels are provided)

7 2 1 0 4 1 4 9 5 9
0 6 9 0 1 5 9 7 3 4
9 6 6 5 4 0 7 4 0 1
3 1 3 4 7 2 7 1 2 1
1 7 4 2 3 5 1 2 4 4
6 3 5 5 6 0 4 1 9 5
7 8 9 3 7 4 6 4 3 0
7 0 2 9 1 7 3 2 8 7
7 6 2 7 8 4 7 3 6 1
3 6 9 3 1 4 1 7 6 9

Goal: predict label of a future pattern

We may think of the pair (X, Y) as obeying a (possibly noisy) input-output relationship.

The goal of supervised learning is usually to generalize the input-output relationship, facilitating the prediction of the output associated with previously unseen inputs X .

The data

$$(X_1, Y_1), \dots, (X_n, Y_n)$$

are called training data.

The Y variable may be called a

- response
- output
- label
- dependent variable

The primary supervised learning problems are

- classification : $Y \in \{1, \dots, M\}$
- regression : $Y \in \mathbb{R}$

Unsupervised Learning

The patterns

$$X_1, \dots, X_n$$

are not accompanied by output variables.

The goal of unsupervised learning is typically not related to future observations. Instead, one seeks to understand structure in the data sample itself, or to infer some characteristic of the underlying probability distribution.

The primary unsupervised learning problems are

- clustering
- density estimation
- dimensionality reduction

↑ can also be supervised

Reinforcement Learning

The patterns X_1, X_2, \dots are observed sequentially. After each X_i is observed, the learner must take an action. After each action, the learner receives a reward from the environment. The goal of the learner is to determine a policy (for selecting actions based on observations) to maximize long-term reward.

RL is important in robotics (navigation, path planning), economics, and other areas.

Other Learning Problems

- Semi-supervised learning
- Active learning
- Online learning
- Novelty detection
- Ranking
- Transfer learning
- Multi-task learning

Types of Learning Methods

Distributional assumptions

- Generative: full probabilistic model
- Discriminative: partial or no probabilistic model; typically models only the desired function or set (e.g. the decision boundary in classification)

Computational Form

- Linear
- Nonlinear
 - polynomial
 - partition-based
 - kernel based

Complexity

- parametric: # of model parameters is independent of sample size
- nonparametric: # of model parameters grows with sample size

Prerequisites

Probability

- joint distributions, multivariate densities and mass functions, expectation, independence, conditional distributions, Bayes rule, the multivariate normal (Gaussian) distribution

Ex] In classification, we view (X, Y) as a jointly distributed pair. We will show that the optimal classifier, in terms of probability of error, is

$$x \mapsto \arg \max_y P(Y=y | X=x)$$

Linear Algebra

Linear methods are widely used in all learning problems. They are typically simpler to understand, and also form the basis of more sophisticated nonlinear methods

- rank, nullspace, linear independence, span, dimension, inner product, orthogonality, positive (semi-) definite matrices, eigenvalue decomposition, projection

Parameter Estimation

Suppose

$$X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta)$$

← independent and identically distributed

where θ is an unknown parameter to be estimated.

The maximum likelihood estimate (MLE) of θ is

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} f(x_1, \dots, x_n; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n f(x_i; \theta) \end{aligned}$$

If θ is viewed as a random variable, we write

$X_i \sim f(x | \theta)$, and base our estimates on

the posterior distribution $f(\theta | x_1, \dots, x_n)$

$$\propto f(x_1, \dots, x_n | \theta) \cdot f(\theta)$$

← Bayes' rule ← prior

Ex] MAP estimate: $\hat{\theta} = \arg \max_{\theta} f(\theta | x_1, \dots, x_n)$