

# Lecture 5: Supervised Learning

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning  
August–December 2020

# Supervised learning

- A set of items
  - Each item is characterized by attributes  $(a_1, a_2, \dots, a_k)$
  - Each item is assigned a class or category  $c$
- Given a set of examples, predict  $c$  for a new item with attributes  $(a'_1, a'_2, \dots, a'_k)$
- Examples provided are called **training data**
- Aim is to **learn** a mathematical model that **generalizes** the training data
  - Model built from training data should extend to previously unseen inputs
- **Classification** problem
  - Usually assumed to binary — two classes

## Example: Loan application data set

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	<b>No</b>
2	young	false	false	good	<b>No</b>
3	young	true	false	good	<b>Yes</b>
4	young	true	true	fair	<b>Yes</b>
5	young	false	false	fair	<b>No</b>
6	middle	false	false	fair	<b>No</b>
7	middle	false	false	good	<b>No</b>
8	middle	true	true	good	<b>Yes</b>
9	middle	false	true	excellent	<b>Yes</b>
10	middle	false	true	excellent	<b>Yes</b>
11	old	false	true	excellent	<b>Yes</b>
12	old	false	true	good	<b>Yes</b>
13	old	true	false	good	<b>Yes</b>
14	old	true	false	excellent	<b>Yes</b>
15	old	false	false	fair	<b>No</b>

# Basic assumptions

## Fundamental assumption of machine learning

- Distribution of training examples is identical to distribution of unseen data

## What does it mean to learn from the data?

- Build a model that does better than random guessing
  - In the loan data set, always saying **Yes** would be correct about 9/15 of the time
- Performance should ideally improve with more training data

## How do we evaluate the performance of a model?

- Model is optimized for the training data. How well does it work for unseen data?
- Don't know the correct answers in advance to compare — different from normal software verification

# The road ahead

## Many different models

- Decision trees
- Probabilistic models — naïve Bayes classifiers
- Models based on geometric separators
  - Support vector machines (SVM)
  - Neural networks

## Important issues related to supervised learning

- Evaluating models
- Ensuring that models generalize well to unseen data
  - A theoretical framework to provide some guarantees
- Strategies to deal with the training data bottleneck