Course Overview and Introduction Machine Learning

Hamid R Rabiee – Zahra Dehghanian Spring 2025



Prerequisites:

- Programming skills
 - Display="block-color: block-space; block-spa
- Probability and statistics
- Basic linear algebra



Marking Scheme

Midterm Exam:	20%	
Final Exam:	30%	
Homeworks (Includes programming assignments):		40%
Quizzes:	10%	
Class Activity & Exam Bonus:	7.5%	
(1.5 extra points)		



Assignments

6 Problem sets

contain both theoretical and programming assignments





Invited Talks from Machine Learning & Al Industry
 The Schedule and the Speaker will be announced





- Deattern Recognition and Machine Learning, C. Bishop, Springer, 2006.
- Machine Learning, T. Mitchell, MIT Press, 1998.
- Other books:
 - The elements of statistical learning, T. Hastie, R. Tibshirani, J. Friedman, Second Edition, 2008.
 - Machine Learning: A Probabilistic Perspective, K. Murphy, MIT Press, 2012.
 - Richard Sutton and Andrew Barto, Reinforcement Learning: An introduction. MIT Press, Second edition, 2017.
- Recommended Video Lectures:
- Andrew Ng:

https://www.youtube.com/watch?v=jGwO_UgTS7I



A Definition of ML

□ Tom Mitchell (1998):

- A computer program is said to learn a task from <u>experience</u> if its performance improves with experience
- Using the observed data to make better decisions
 - Generalizing from the observed data



ML has a wide reach

- Wide applicability
- Very-large-scale complex systems
- Internet (billions of nodes), sensor network (new multi-modal sensing devices), genetics (human genome)
- Huge multi-dimensional data sets
- 20,000 genes x 10,000 drugs x 100 species x ...
- Improved machine learning algorithms
- Improved data capture (Terabytes, Petabytes of data), networking, faster computers
- The New York Times is regularly talks about machine learning!



ML Definition: Example

Consider an email program that learns how to filter spam according to emails you do or do not mark as spam.

- Task: Classifying emails as spam or not spam.
- Experience: Watching you label emails as spam or not spam.
- Performance: The number (or fraction) of emails correctly classified as spam/not spam.



The essence of machine learning

- A pattern exist
- We do not know it mathematically
- We have data on it



Example: Home Price

Housing price prediction





Sharif University of Technology

Example: Home Price

Predicting house price from 3 attributes

	Age (year)	Region	
100	2	5	500
80	25	3	250



Sharif University of Technology

Example: Home Price





Sharif University of Technology

Example: Bank loan

Applicant form as the input:

- salary
- 🛛 age
- gender
- current debt
- ...
- Output: approving or denying the request



Experience (E) in ML

- Basic premise of learning:
 - "Using a set of observations to uncover an underlying process"
- We have different types of (getting) observations in different types or paradigms of ML methods



Early Paradigms of ML

- <u>Supervised learning</u> (regression, classification)
 - predicting a target variable for which we get to see examples.
- Unsupervised learning
 - revealing structure in the observed data



Data in Supervised Learning

- \blacktriangleright Data are usually considered as vectors in a d dimensional space
 - Now, we make this assumption for illustrative purpose
 - We will see it is not necessary

Columns: Features/attributes/dimensions
Rows:
Data/points/instances/examples/sample
s
Y column:
Target/outcome/response/label



• • •

Sharif University of Technology

Sample

Sample 2

. . .

Sample n-1

Sample n

Supervised Learning: Regression vs. Classification

- Supervised Learning
 - **Regression**: predict a <u>continuous</u> target variable
 - ▶ E.g., *y* ∈ [0,1]
 - Classification: predict a <u>discrete</u> (unordered) target variable
 E.g., y ∈ {1,2, ..., C}



Regression: Example

Housing price prediction





Sharif University of Technology

Classification: Example

 Classification of tumors to Benign/Malignant according to attributes





Sharif University of Technology

Classification: Example

Classification of tumors to Benign/Malignant according to attributes





Sharif University of Technology

Training data: Example



Training data





Sharif University of Technology

Handwritten Digit Recognition Example

Data: labeled samples



Handwritten Digit Recognition Example





Sharif University of Technology

Components of (Supervised) Learning

- ▶ Unknown target function: $f: \mathcal{X} \to \mathcal{Y}$
 - Input space: \mathcal{X}
 - Output space: *Y*
- Training data: $(x_1, y_1), (x_2, y_2), ..., (x_N, y_N)$
- Pick a formula $g\colon \mathcal{X} \to \mathcal{Y}$ that approximates the target function f
 - \blacktriangleright selected from a set of hypotheses ${\cal H}$



Components of (Supervised) Learning

We have some example pairs of (input, output) called training samples

•
$$(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})$$

- We want to select a function from the input space to the output space
 - $\flat \ f \colon \mathcal{X} \to \mathcal{Y}$
- We choose a set of hypotheses (candidate formulas)
 - e.g., linear functions
- We use a learning algorithm to select a function from hypothesis set that approximates the target function



Components of (Supervised) Learning



Source: Yaser Abou Mostafa, Learning from Data



Sharif University of Technology

(Supervised) Learning problem

Selecting a hypothesis space

- Hypothesis space: a set of mappings from feature vector to target
- **Learning:** find mapping \hat{f} (from hypothesis set) based on the training data
 - Which notion of error should we use? (loss functions)
 - \blacktriangleright Optimization of loss function to find mapping \hat{f}
- **Evaluation**: we measure how well \hat{f} generalizes to unseen examples (generalization)



Solution Components

Learning model composed of:

- Hypothesis set
- Learning algorithm

Perceptron example

Handwritten Digit Recognition Example

Data: labeled samples





Example: Input representation

'raw' input
$$\mathbf{x}=(x_0,\!x_1,x_2,\cdots,x_{256})$$

linear model: $(w_0, w_1, w_2, \cdots, w_{256})$

Features: Extract useful information, e.g.,

intensity and symmetry $\mathbf{x} = (x_0, x_1, x_2)$

linear model: (w_0, w_1, w_2)

Data Source: Yaser Abou Mostafa, Learning from







Example: Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$
 x_1 : intensity x_2 : symmetry



Data Source: Yaser Abou Mostafa, Learning from

Course Overview and Introduction



Perceptron classifier

- Input $\boldsymbol{x} = [x_1, \dots, x_d]$
- Classifier:
 - If $\sum_{i=1}^{d} w_i x_i$ > threshold then output 1
 - ▶ else output −1
- The linear formula $g \in \mathcal{H}$ can be written: $g(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{d} \mathbf{w}_{i} x_{i} - \operatorname{threshold}\right)$





Perceptron classifier

Input
$$\boldsymbol{x} = [x_1, \dots, x_d]$$

- Classifier:
 - If $\sum_{i=1}^{d} w_i x_i$ > threshold then output 1
 - ▶ else output −1
- The linear formula $g \in \mathcal{H}$ can be written:

$$g(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{d} \mathbf{w}_{i} x_{i} + \mathbf{w}_{0}\right)$$





Perceptron classifier

Input
$$x = [x_1, ..., x_d]$$

- Classifier:
 - If $\sum_{i=1}^{d} w_i x_i$ > threshold then output 1
 - ▶ else output −1
- The linear formula $g \in \mathcal{H}$ can be written: $g(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{d} \mathbf{w}_{i} x_{i} + \mathbf{w}_{0}\right)$



If we add a coordinate $x_0 = 1$ to the input:

$$g(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=0}^{d} \mathbf{w}_{i} x_{i}\right)$$

Vector form





Sharif University of Technology

Perceptron learning algorithm: linearly separable data

Give the training data
$$(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})$$

Misclassified data
$$(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$$
:
 $\operatorname{sign}(\mathbf{w}^T \mathbf{x}^{(n)}) \neq y^{(n)}$

Repeat

Pick a misclassified data $(x^{(n)}, y^{(n)})$ from training data and update w:

$$\boldsymbol{w} = \boldsymbol{w} + \boldsymbol{y}^{(n)} \boldsymbol{x}^{(n)}$$

Until all training data points are correctly classified by g



Perceptron learning algorithm: Example of weight update



Example: linear classifier

$$\mathbf{x} = (x_0, x_1, x_2)$$
 x_1 : intensity x_2 : symmetry



Course Overview and Introduction



(Supervised) Learning problem

Selecting a hypothesis space

- Hypothesis space: a set of mappings from feature vector to target
- Learning (estimation): optimization of a cost function
 - Based on the training set $D = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$ and a cost function we find (an estimate) $f \in F$ of the target function
- **Evaluation**: we measure how well \hat{f} generalizes to unseen examples



Generalization

- We don't intend to memorize data but want to distinguish the pattern.
- A core objective of learning is to generalize from the experience.
 - Generalization: ability of a learning algorithm to perform accurately on new, unseen examples after having experienced.



Paradigms of ML

- <u>Supervised learning</u> (regression, classification)
 - predicting a target variable for which we get to see examples.
- Unsupervised learning
 - revealing structure in the observed data



Supervised Learning vs. Unsupervised Learning

- Supervised learning
 - Given:Training set
 - ▶ labeled set of N input-output pairs $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{N}$
 - Goal: learning a mapping from **x** to y
- Unsupervised learning
 - Given:Training set
 - $\blacktriangleright \left\{ \boldsymbol{x}^{(i)} \right\}_{i=1}^{N}$
 - Goal: find groups or structures in the data
 - Discover the intrinsic structure in the data



Supervised Learning: Samples



Unsupervised Learning: Samples

Wants to use data to improve their knowledge on a task





Sample Data in Unsupervised Learning

Unsupervised Learning:

Columns: *Features/attributes/dimensions*

Rows: Data/points/instances/examples /samples

•••		
		Sample 1
		Sample 2
		Sample n-1
		Sample n



Unsupervised learning

- Clustering: partitioning of data into groups of similar data points.
- Dimensionality reduction: data representation using a smaller number of dimensions while preserving (perhaps approximately) some properties of the data.
- Density estimation & generative models



Some clustering purposes

- Preprocessing stage to index, compress, or summarize the data
- As a tool to understand the hidden structure in data or to group them
 - To gain knowledge (insight into the structure of the data) or
 - To group the data when no label is available



Clustering: Example Applications

- Clustering docs based on their similarities
 Grouping new stories in the Google news site
- Market segmentation: group customers into different market segments given a database of customer data.
- Community detection in social networks



Clustering of docs







Course Overview and Introduction

Dimensionality reduction: Example

How to map the high dimensional data into a lower dimensional space in which the distance is more meaningful.



Sharif University of Technology

Generative models: Example





Sharif University of Technology

Paradigms of ML

- <u>Supervised learning</u> (regression, classification)
 - predicting a target variable for which we get to see examples.
- Unsupervised learning
 - revealing structure in the observed data
- Other paradigms: semi-supervised learning, weakly supervised, self-supervised learning, active learning, etc.





Task (i.e. what is the type of knowledge that we seek from data)

Data

Algorithm



Three axes of ML

- Task (i.e. what is the type of knowledge that we seek from data)
 - Prediction (i.e. classification or regression)
 - Control
 - Description
- 🛛 Data
 - Fully observed / Partially observed
 - Passively / Actively collecting data
 - Online / Offline
- □ Algorithm
 - Parametric / Nonparametric models
 - ····



Parametric models

- We consider a parametric boundary (e.g., hyper-plane, hyperbola, ...) and learn its parameters form data
 - The set of parameters does not grow with increasing the data







Nonparametric models

- We must store data and for each prediction, we need to process training data
- More data means a more complex model
 - Models that grow with the data



Nonparametric models

- k-NN classifier
 - Label for x predicted by majority voting among its k-NN.



Find k nearest training data to the new input and predict its label from the labels of its k nearest neighbors

The number of points to search scales with the training data

Course Overview and Introduction



ML in Computer Science

Why ML applications are growing?

- Improved machine learning algorithms
- Availability of data (Increased data capture, networking, etc)
- Software too complex to write by hand
 - Demand for complex systems (on high-dimensional, multi-modal, or heterogeneous data)
 - Demand for self-customization to user or environment



Relation to other fields

- **Statistics:** the goal is the understanding of the data at hand
- **Artificial Intelligence:** the goal is to build an intelligent agent
- **Data Mining:** the goal is to extract patterns from large-scale data
- Data Science: the science encompassing collection, analysis, and interpretation of data
- The goal of machine learning is the underlying mechanisms and algorithms that allow improving our knowledge with more data



Some Learning Application Areas

- Computer Vision (Photo tagging, face recognition, ...)
- Natural language processing (e.g., machine translation)
- Robotics
- Speech recognition
- Autonomous vehicles
- Social network analysis
- Web search engines
- Medical outcomes analysis
- Market prediction (e.g., stock/house prices)
- Computational biology (e.g., annotation of biological sequences)
- Self-customizing programs (recommender systems)



Topics of this course

- Regression & generalization
- Evaluation & Model Selection
- □ Classification
 - Linear classifier
 - Probabilistic classifiers
 - SVM & kernel
 - Decision tree
- Neural Networks
- Learning Theory
- Non-parametric methods
- Ensemble learning
- Dimensionality reduction
- Clustering