Machine Learning: Introduction and Unsupervised Learning

Chapter 18.1 – 18.2 and "Introduction to Statistical Machine Learning"

What is Learning?

- "Learning is making useful changes in our minds" – Marvin Minsky
- "Learning is constructing or modifying representations of what is being experienced" – Ryszard Michalski
- "Learning denotes changes in a system that ... enable a system to do the same task more efficiently the next time" – Herbert Simon

Why do Machine Learning?

- Solve classification problems
- Learn models of data ("data fitting")
- Understand and improve efficiency of human learning (e.g., Computer-Aided Instruction (CAI))
- Discover new things or structures that are unknown to humans ("data mining")
- Fill in skeletal or incomplete specifications about a domain

Major Paradigms of Machine Learning

- Rote Learning
- Induction
- Clustering
- Analogy
- Discovery
- Genetic Algorithms
- Reinforcement

Inductive Learning

- Generalize from a given set of (training) examples so that accurate predictions can be made about future examples
- Learn unknown function: f(x) = y
 - -x: an input example (aka instance)
 - -y: the desired output
 - -h (hypothesis) function is learned that approximates f

Representing "Things" in Machine Learning

- An example or instance, x, represents a specific object ("thing")
- **x** often represented by a *D*-dimensional **feature vector** $\mathbf{x} = (x_1, \dots, x_D) \in \Re^D$
- Each dimension is called a feature or attribute
- Continuous or discrete
- **x** is a point in the **D-dimensional feature space**
- Abstraction of object. Ignores any other aspects (e.g., two people having the same weight and height may be considered identical)

Feature Vector Representation

- · Preprocess raw sensor data
 - extract a feature (attribute) vector, *x*, that describes all attributes relevant for an object
- Each *x* is a list of (*attribute*, *value*) pairs
 - $\boldsymbol{x} = [(Rank, queen), (Suit, hearts), (Size, big)]$
 - number of attributes is fixed: Rank, Suit, Size
 - number of possible values for each attribute is fixed (if discrete)
 - Rank: 2, ..., 10, jack, queen, king, ace Suit: diamonds, hearts, clubs, spades
 - Size: big, small

Feature Vector Representation

Each example can be interpreted as a point in a *D*-dimensional feature space, where *D* is the number of features/attributes



Feature Vector Representation Example

- Text document
 - Vocabulary of size D (~100,000): aardvark, ..., zulu
- "bag of words": counts of each vocabulary entry
 - − To marry my true love → (3531:1 13788:1 19676:1)
 - I wish that I find my soulmate this year → (3819:1 13448:1 19450:1 20514:1)
- Often remove stopwords: the, of, at, in, ...
- Special "out-of-vocabulary" (OOV) entry catches all unknown words

More Feature Representations

- Image
 - Color histogram
- Software
 - Execution profile: the number of times each line is executed
- Bank account
 - Credit rating, balance, #deposits in last day, week, month, year, #withdrawals, ...
- You and me
 - Medical test1, test2, test3, ...

Training Sample

- A training sample is a collection of examples (aka instances), x_1, \ldots, x_n , which is the input to the learning process
- $\mathbf{x}_i = (x_{i1}, ..., x_{iD})$
- Assume these instances are sampled independently from an unknown (population) distribution, P(x)
- We denote this by x_i ^{i.i.d.}P(x), where i.i.d. stands for independent and identically distributed

Training Sample

- A training sample is the "experience" given to a learning algorithm
- What the algorithm can learn from it varies
- Two basic learning paradigms:
 - -unsupervised learning
 - -supervised learning

Inductive Learning

- Supervised vs. Unsupervised Learning
 - supervised: "teacher" gives a set of (x, y) pairs
 - unsupervised: only the *x* 's are given
- In either case, the goal is to estimate *f* so that it generalizes well to correctly deal with "future examples" in computing *f*(*x*) = *y*

Unsupervised Learning

- Training sample is $\mathbf{x}_1, \ldots, \mathbf{x}_n$, that's it
- No teacher providing supervision as to how individual examples should be handled
- · Common tasks:
 - Clustering: separate the *n* examples into groups
 - Novelty detection: find examples that are very different from the rest
 - Dimensionality reduction: represent each example with a lower dimensional feature vector while maintaining key characteristics of the training samples







Detecting Events on Twitter

- Use real-time text and images from tweets to discover new social events
- Clusters defined by similar words and word cooccurences, plus image features

Digital Photo Collections You have 1000's of digital photos stored in various folders ... Organize them better by grouping into clusters Simplest idea: use image creation time (EXIF tag) More complicated: extract image features

Three Frequently Used Clustering Methods

- Hierarchical clustering
 - Build a binary tree over the dataset
- K-means clustering
 - Specify the desired number of clusters and use an iterative algorithm to find them
- Mean Shift clustering











Hierarchical Agglomerative Clustering Algorithm

Input: a training sample $\{\mathbf{x}_i\}_{i=1}^n$; a distance function d().

Initially, place each instance in its own cluster (called a singleton cluster).
 while (number of clusters > 1) do:

- 3. Find the closest cluster pair A, B, i.e., they minimize d(A, B).
- 4. Merge A, B to form a new cluster.

Output: a binary tree showing how clusters are gradually merged from singletons to a root cluster, which contains the whole training sample.



Example: Histogram-Based Image Segmentation Goal: Segment the image into K regions Reduce the number of colors to K and map each pixel to the closest color

Example: Histogram-Based Image Segmentation

- Goal: Segment the image into K regions
 - Reduce the number of colors to K and map each pixel to the closest color



Hierarchical Clustering

- How do you measure the closeness between two clusters? At least three ways:
 - Single-linkage: the shortest distance from any member of one cluster to any member of the other cluster
 - Complete-linkage: the greatest distance from any member of one cluster to any member of the other cluster
 - Average-linkage: you guessed it

Hierarchical Clustering

• How do you measure the closeness between two clusters?

Hierarchical Clustering

- The binary tree you get is often called a **dendrogram**, or taxonomy, or a hierarchy of data points
- The tree can be cut at various levels to produce different numbers of clusters: if you want *k* clusters, just cut the (*k*-1) longest links
- Sometimes the hierarchy itself is more interesting than the clusters















































K-Means Algorithm

- Input: *x*₁, ..., *x*_n, *k*
- **Step 1**: select *k* cluster centers *c*₁ ... *c*_k
- **Step 2**: for each point *x*, determine its cluster: find the closest center in Euclidean space
- Step 3: update all cluster centers as the centroids
 - $c_i = \sum_{\{x \text{ in cluster } i\}} x / \text{SizeOf(cluster } i)$
- Repeat steps 2 and 3 until cluster centers no longer change

K-Means Questions

- Will it always terminate?
 - Yes (finite number of ways of partitioning a finite number of points into k groups)
- Is it guaranteed to find an "optimal" clustering?
 - No, but each iteration will reduce the error/distortion of the clustering

Example: Image Segmentation







Input image

Clusters on intensity

Clusters on color

Say k=3 and you are given the following points:

Non-Optimal Clustering Given a poor choice of the initial cluster centers, the following result is possible:

Picking the Number of Clusters

- Difficult problem
- Heuristic approaches depend on number of points and number of dimensions

Picking Starting Cluster Centers

Which local optimum *k*-Means goes to is determined solely by the starting cluster centers

- Idea 1: Run k-Means multiple times with different starting, random cluster centers (hill climbing with random restarts)
- **Idea 2**: Pick a random point x_1 from dataset
 - 1. Find the point x_2 farthest from x_1 in the dataset
 - 2. Find x_3 farthest from the closer of x_1, x_2
 - 3. ... Pick *k* points like this, and use them as the starting cluster centers for the *k* clusters

Uses of K-Means

- Often used as an exploratory data analysis tool
- In one-dimension, a good way to quantize realvalued variables into *k* non-uniform buckets
- Used on acoustic data in speech recognition to convert waveforms into one of k categories (known as Vector Quantization)
- Also used for choosing color palettes on graphical display devices

Copyright © 2001, 2004 Andrew W. Moore

Mean Shift Clustering

- 1. Choose a search window size
- 2. Choose the initial location of the search window
- 3. Compute the mean location (centroid of the data) in the search window
- 4. Center the search window at the mean location computed in Step 3
- 5. Repeat Steps 3 and 4 until convergence

The mean shift algorithm seeks the *mode*, i.e., point of highest density of a data distribution:



































Concept Learning

- Determine from a given a set of examples if a given example is or isn't an instance of the concept/class/category
 - If it is, call it a **positive** example
 - If not, called it a **negative** example



Supervised Concept Learning by Induction

- Given a training set of positive and negative examples of a concept:
 - $-\{(\boldsymbol{x}_{1},\,y_{1}),\,(\boldsymbol{x}_{2},\,y_{2}),\,...,\,(\boldsymbol{x}_{n},\,y_{n})\}$
 - each y_i is either + or -
- Construct a description that accurately classifies whether future examples are positive or negative:
 - $-h(\boldsymbol{x}_{n+1}) = y_{n+1}$
 - where y_{n+1} is the + or prediction

Inductive Learning by Nearest-Neighbor Classification

- A simple approach:
 - save each training example as a point in Feature Space
 - classify a new example by giving it the same classification as its *nearest neighbor* in Feature Space

k-Nearest-Neighbor (k-NN)

Input: Training data $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)$; distance function d(); number of neighbors k; test instance \mathbf{x}^*

1. Find the k training instances $\mathbf{x}_{i_1}, \ldots, \mathbf{x}_{i_k}$ closest to \mathbf{x}^* under distance d(). 2. Output y^* as the majority class of y_{i_1}, \ldots, y_{i_k} . Break ties randomly.







- What if we want regression?
 - Instead of majority vote, take average of neighbors' y
- How to pick k?
 - Split data into training and tuning sets
 - Classify tuning set with different k values
 - Pick k that produces least tuning-set error

Inductive Bias

- Inductive learning is an inherently conjectural process. Why?
 - any knowledge created by generalization from specific facts cannot be proven true
 - it can only be proven false
- Hence, inductive inference is "falsity preserving," not "truth preserving"

Inductive Bias

- · Biases commonly used in machine learning
 - Restricted Hypothesis Space Bias: allow only certain types of *h* 's, not arbitrary ones

- Preference Bias:

define a metric for comparing h 's so as to determine whether one is better than another

Inductive Bias

- Learning can be viewed as searching the Hypothesis Space *H* of possible *h* functions
- Inductive Bias
 - is used when one h is chosen over another
 - is needed to generalize beyond the specific training examples
- · Completely unbiased inductive algorithm
 - only memorizes training examples
 - can't predict anything about unseen examples

Supervised Learning Methods

- k-nearest-neighbor (k-NN)
- Decision trees
- Neural networks (NN)
- Support vector machines (SVM)
- etc.