

Bookkeeping

- Today: Variable Elimination (quick redux)
- Today: ML 1
 - What is machine learning?
 - Classification
 - Intro to decision trees?

Next class

- In-class midterm review
 - What kinds of questions?
 - What material will be covered?







Variable Elimination

- Variable elimination: carry out summations right-to-left, storing intermediate results (factors) to avoid recomputation
- A factor is a function from some set of variables into a specific value
 - A factor can be (and often is) a CPT
 - A set of values (because it's a function)
- Variable elimination works by eliminating all variables in turn until there is a factor with only the query variable
- To eliminate a variable:
 - Join all factors containing that variable (like DB)
 - Sum out the influence of the variable on new factor

Factors

- A factor is a function from a set of random variables to a number
- Formally, let f denote a factor and let X1,..,Xj denote the variables in the factor
- A factor can represent a joint probability or a conditional probability
 - Ex: a factor with two variables X1 and X2 can represent $P(X1 \wedge X2)$ the joint probability
 - Or, it can represent P(X1|X2) this is a conditional probability
- For variable elimination define a factor for every variable/node in the Bayesian network
- The initial factor for each variable/node captures the conditional probability distribution for that variable/node



Variable Elimination Algorithm

function ELIMINATION-ASK(X, e, bn) returns a distribution over X inputs: X, the query variable e, evidence specified as an event bn, a belief network specifying joint distribution $P(X_1, ..., X_n)$ factors \leftarrow []; vars \leftarrow REVERSE(VARS[bn]) for each var in vars do factors \leftarrow [MAKE-FACTOR(var, e)|factors] if var is a hidden variable then factors \leftarrow SUM-OUT(var, factors) return NORMALIZE(POINTWISE-PRODUCT(factors))





Why "Learn" ?

- Machine learning is programming computers to optimize a performance criterion using example data or past experience.
- There is no need to "learn" to calculate payroll
- Learning is used when:
 - Human expertise does not exist (navigating on Mars)
 - Humans are unable to explain their expertise (speech recognition)
 - Solution changes in time (routing on a computer network)
 - Solution needs to be adapted to particular cases (user biometrics)



What is Machine Learning?

- Optimize a performance criterion using example data or past experience
- Role of Statistics: Inference from a sample
- Role of Computer science: Efficient algorithms to:
- 1. Solve the optimization problem
- 2. Represent and evaluate the model for inference



Associations

- Basket analysis:
- *P*(*Y* / *X*) probability that somebody who buys *X* also buys *Y* where *X* and *Y* are products/services.
- Example: P (chips | beer) = 0.7





Classification: Applications

- AKA Pattern recognition
- Face recognition: Pose, lighting, occlusion (glasses, beard), make-up, hair style
- Character recognition: Different handwriting styles.
- Speech recognition: Temporal dependency.
 - Use of a dictionary or the syntax of the language.
 - Sensor fusion: Combine multiple modalities; eg, visual (lip image) and acoustic for speech
- Medical diagnosis: From symptoms to illnesses
- ...



www.thestartupfounder.com/best-and-most-accurate-factal-recognition-engines-in-2022 Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)





Supervised Learning: Uses

- **Prediction of future cases:** Use the rule to predict the output for future inputs
- Knowledge extraction: The rule is easy to understand
- **Compression:** The rule is simpler than the data it explains
- **Outlier detection:** Exceptions that are not covered by the rule, e.g., fraud

[Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)





So... What is Learning?

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

Why Learn?

- Discover previously-unknown new things or structure
 - Data mining, scientific discovery
- Fill in skeletal or incomplete domain knowledge
- Build agents that can adapt to users or other agents
- Understand and improve efficiency of human learning
 - Use to improve methods for teaching and tutoring people (e.g., better computer-aided instruction)

Machine Leaning Successes

- Sentiment analysis
- Spam detection
- Machine translation
- Spoken language understanding
- Named entity detection
- Self driving cars

- Motion recognition
- Identifying places in digital images
- Recommender systems (Netflix, Amazon)
- Credit card fraud detection
- Etc. (lots of etc.)







Questions

- What's supervised learning?
 - What's classification? What's regression?
 - What's a hypothesis? What's a hypothesis space?
 - What are the training set and test set?
 - What is Ockham's razor?
- What's unsupervised learning?

	Suponvisod Loarning	
-	Supervised Learning	onsupervised Learning
Discrete	classification or categorization	clustering
Continuous	regression	dimensionality reduction





The Machine Learning Framework

• Apply a prediction function to a feature representation of the data to get the desired output:

$$f(round, stem) = "apple"$$
$$f(round, \neg stem) = "tomato"$$
$$f(\neg round, \neg stem) = "cow"$$

Slide credit: Svetlana Lazebnik







Many classifiers to choose from SVM **Randomized Forests** • • Neural networks **Boosted Decision Trees** • • Naïve Bayes K-nearest neighbor • • **Bayesian network RBMs** • • Logistic regression Etc. • • Which is the best one? Slide credit: Derek Hoiem 42

Major Paradigms of ML (1)

- **Rote learning**: 1-1 mapping from inputs to stored representation, learning by memorization, association-based storage & retrieval
- Induction: Use specific examples to reach general conclusions
- **Clustering**: Unsupervised discovery of natural groups in data

Major Paradigms of ML (2)

- Analogy: Find correspondences between different representations
- Discovery: Unsupervised, specific goal not given
- Genetic algorithms: Evolutionary search techniques, based on an analogy to survival of the fittest
- **Reinforcement:** Feedback (positive or negative reward) given at the end of a sequence of steps

<text><list-item><list-item><list-item><list-item><list-item><list-item>







Unsupervised Learning

Given only *unlabeled* data as input, learn some sort of structure, e.g.:

- Cluster your Facebook friends based on similarity of posts and friends
- Find sets of words whose meanings are related (e.g., doctor, hospital)
- Induce N topics and the words that are common in documents that are about that topic







Inductive Learning Framework

- Raw input data from sensors preprocessed to obtain feature vector,
 X, of relevant features for classifying examples
- Each X is a list of (attribute, value) pairs
- *n* attributes (a.k.a. features): fixed, positive, and finite
- Features have fixed, finite number # of possible values
 - Or continuous within some well-defined space, e.g., "age"
- Each example is a point in an *n*-dimensional feature space
 - X = [Name:Sue, EyeColor:Brown, Age:Young, Gender:Female]
 - X = [Cheese:f, Sauce:t, Bread:t]
 - X = [Texture:Fuzzy, Ears:Pointy, Purrs:Yes, Legs:4]



Inductive Learning as Search

- Instance space, I, is the set of all possible examples
 - Defines the language for the training and test instances
 - Usually each instance $i \in I$ is a feature vector
 - Features are also sometimes called attributes or variables
- Class variable C gives an instance's class (to be predicted)
- Model space M defines the possible classifiers
 - $M: I \rightarrow C, M = \{m_1, \dots, m_n\}$ (possibly infinite)
 - Model space is **sometimes** defined using same features as instance space











Model Spaces (1)

- Decision trees
 - Partition the instance space I into axis-parallel regions
 - Labeled with class value
- Nearest-neighbor classifiers
 - Partition the instance space I into regions defined by centroid instances (or cluster of *k* instances)
- Bayesian networks
 - Probabilistic dependencies of class on attributes
 - Naïve Bayes: special case of BNs where class ightarrow each attribute

Model Spaces (2)

- Neural networks
 - Nonlinear feed-forward functions of attribute values
 - Can be "deep"
 - Much learning today falls under neural approaches
- Support vector machines
 - Find a separating plane in a high-dimensional feature space
- Associative rules (feature values → class)
- First-order logical rules





Summary: Machine Learning 1

- Core idea: given (possibly labeled) training data, learn a model of how the world works that lets you make predictions about new observations at test time
- Supervised vs. unsupervised, continuous vs. discrete
- Supervised learning over discrete data = classification
- Decision trees are one approach for discrete data







Decision Tree Induction

• The Big Idea: build a tree of **decisions**, each of which splits training data into smaller groups

• Very common machine learning technique!

- At each split, an attribute of the training data a feature is chosen to divide data into classes
- Goal: each leaf group in the tree consists entirely of one class
- Learning: creating that tree



	Class label	Features				
Object	Yellow?	R	G	В	Fuzzy?	
Duckie1	N	0	255	0	N	
Fish	Y	240	240	0	Y	
Tomato	N	250	0	0	N	
Banana	Y	255	230	0	N	
Duckie2	Y	250	255	0	N	
Frog	N	0	120	0	Y	















Preference Bias: Ockham's Razor

- A.k.a. Occam's Razor, Law of Economy, or Law of Parsimony
- Stated by William of Ockham (1285-1347/49):
 - "Non sunt multiplicanda entia praeter necessitatem"
 - "Entities are not to be multiplied beyond necessity"
- The simplest consistent explanation is the best.
- Smallest decision tree that correctly classifies all training examples
- Finding the provably smallest decision tree is NP-hard!
- So, instead of constructing the absolute smallest tree consistent with the training examples, construct one that is "pretty small"

76

R&N's Restaurant Domain

- Model the decision a patron makes when deciding whether to wait for a table or leave the restaurant
 - Two classes (outcomes): wait, leave
 - Ten attributes:
 - Alternative available? ∃ Bar? Is it Friday? Hungry? How full is restaurant? How expensive? Is it raining? Do we have a reservation? What type of restaurant is it? What's purported waiting time?
- Training set of 12 examples
- ~ 7000 possible cases

Datum	Attributes										Outcome (Label)
	altern- atives	bar	Friday	hungry	people	\$	rain	reser- vation	type	wait time	Wait?
X1	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	Yes
X ₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	No
X ₃	No	Yes	No	No	Some	\$	No	No	Burger	0-10	Yes
X4	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	Yes
X ₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	No
X ₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	Yes
X ₇	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	No
X ₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	Yes
X ₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	No
X ₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	0-30	No
X ₁₁	No	No	No	No	None	\$	No	No	Thai	0-10	No
X ₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	Yes



lssues

- It's like 20 questions:
- We can generate many decision trees depending on what attributes we ask about and in what order
- How do we decide?
- What makes one decision tree better than another: number of nodes? number of leaves? maximum depth?

80

ID3/C4.5

- A greedy algorithm for decision tree construction
 - Ross Quinlan, 1987
- Construct decision tree top-down by recursively selecting the "best attribute" to use at current node
 - Select attribute for current node
 - Generate child nodes (one for each possible value of attribute)
 - Partition training data using attribute values
 - Assign subsets of examples to the appropriate child node
 - Repeat for each child node until all examples associated with a node are either all positive or all negative





		Outlook	Тетр	Humidity	Windy	Play golf?
1		Rainy	Hot	High	False	No
2		Rainy	Hot	High	True	No
3		Overcast	Hot	High	False	Yes
4		Sunny	Mild	High	False	Yes
5		Sunny	Cool	Normal	False	Yes
6		Sunny	Cool	Normal	True	No
7		Overcast	Cool	Normal	True	Yes
8		Rainy	Mild	High	False	No
9		Rainy	Cool	Normal	False	Yes
10)	Sunny	Mild	Normal	False	Yes
11	L	Rainy	Mild	Normal	True	Yes
12	2	Overcast	Mild	High	True	Yes
13	3	Overcast	Hot	Normal	False	Yes
14	1	Sunny	Mild	High	True	No



Choosing the Best Attribute

- Key problem: what attribute to split on?
- Some possibilities are:
 - Random: Select any attribute at random
 - Least-Values: Choose attribute with smallest number of values
 - Most-Values: Choose attribute with largest number of values
 - Max-Gain: Choose attribute that has the largest expected information gain the attribute that will result in the smallest expected size of the subtrees rooted at its children
- ID3 uses Max-Gain to select the best attribute

86



Restaurant Example

- What do these approaches split restaurants on, given the data in the table?
 - Random: Patrons or Type
 - Least-values: Patrons
 French
 - Most-values: Type
 - Max-gain: ???

French		Y	Ν
Italian		Y	Ν
Thai	Ν	Y	Ν
Burger	Ν	Y	Y
	Empty	Some	Full