

Dreamers. Thinkers. Doers.



CIVL 7012/8012

Introduction to Machine Learning and basic algorithms

www.memphis.edu

What is Machine Learning?

- Field of study that enables computers to learn by itself without being programmed.
- "Learning" implies improving on certain task T, with respect to performance metric P based on training experience E

Some examples

- Identifying spams in email
 - T: Classify emails as legitimate or spam
 - P: Percentage of accurate classifications
 - E: Email database
- Detecting damages in structures
 T: Detect the location and extent of damages in structure
 P: Percentages of accurate detections
 E: Visual and sensory data

Some more examples

- Classifying soil using photographs
 T: Classifying soil class and type
 P: Percentages of accurate classifications
 E: Visual data such as photographs
 - E: Visual data such as photographs
- Identifying spoken words (e.g. Apple Siri, Amazon Alexa, etc.)
 - T: Classify and identify words spoken by someone
 - P: Percentage of accurate identifications
 - E: Database of words spoken by people with different accents



Dreamers. Thinkers. Doers.

How is ML different from traditional modeling?



ML Terminology

- Generalization: The ability of ML algorithm to fit new data.
- Features, class and examples:

See figure below for a sample data with dependent and independent variables from HW 4 Q3.

		Fe	eatures (li	ndepe	endent va	ariables)				
Class label or										
target class (Dependent	Speed	Route	Seatbelt	Age	Income	Marital status				
variable)	Low	1	1	1	2	Married				
Examples or samples (Observations)	Medium	2	0	3	3	Married	5			
	High	1	0	1	1	Unmarried				
	High	1	1	2	2	Married				
							5			



Types of learning

MEMPHIS

- 1. Supervised learning
- Used when the data has desired outputs
- E.g. regression, decision trees
- 2. Unsupervised learning
- Used when the data does not have desired outputs
- E.g. Clustering algorithm such as K-means
- 3. Reinforcement learning
- Commonly used in robotics
- Learning based on rewards gained from actions







Types of learning

- 1. Supervised learning
- Train a model using training data
- Predict on a new unseen data
- Training data includes desired outputs (Remember regression?)



7



Supervised learning

- For a continuous data prediction-regression
- For categorical prediction-classification (e.g. binary logit)



Types of learning

- 2. Unsupervised learning
- Training data does not include desired outputs
- Creates 'clusters' of similar data

Example.

Figure below shows clustering of unlabeled data into three groups based on their similarity with features X, and Y



Types of learning

- 3. Reinforcement learning
- Considerations in reinforcement learning
 - There are finite number of states with finite number of actions
 - The *agent* improves its's performance based on *reward* from *actions* or interaction with the *environment*.



Example 1. A chess engine (*agent*) training on chess moves will base its moves on the state of the chess board (*environment*) with the reward positive if win, negative if lose.



Reinforcement learning

Example 2

- Consider a self driving car stopped at a red-light signal. The
- lights turn green and the car starts moving. Here the agent is
- the car and it is moving from one state to another. This
- actions produces some reward.

On the contrary if the car does not start and remains at the intersection despite green signal, it receives a traffic ticket. This is negative reward.



Some common learners

- Naïve Bayes (Supervised)
- Decision trees (Supervised)
- Perceptron and Artificial Neural Networks (Supervised/unsupervised)
- K-means clustering (Unsupervised)



Naïve Bayes

 Naïve Bayes is based on Bayes theorem for conditional probability:

$$P(y_i|X) = \frac{P(X|y_i)P(y_i)}{P(X)}$$

 As we are calculating probabilities for all classes the above equation can be written as

 $P(y_i|X) \propto P(X|y_i)P(y_i)$



Naive Bayes

- Even with this, it becomes difficult to learn $P(X|y_i)$
- Example
- Consider X has only 2 features: gender(male, female) and marital status (married and unmarried), and there are two possible outcomes for target class wage(low and high).
- To apply Bayes theorem we will need to learn
 2ⁿ conditional probabilities.



Naive Bayes

- With large number of features Bayes theorem becomes difficult to apply
- To learn $P(X|y_i)$ where X has n features, we need to learn the conditional probabilities $P(x_1|y_i), P(x_2|y_i) \dots P(x_n|y_i)$ which is much easier
- We make a naïve assumption to simplify problems i.e. features are conditionally independent given class.



Naïve Bayes

$$P(x_{1}, x_{2}|y) = P(x_{1}|y)P(x_{2}|y)$$

OR
$$P(x_{1}, x_{2} \dots x_{n}|y) = \prod_{i} P(x_{i}|y)$$

• Classification is then based on class selection with largest probability $y^* = argmax_y P(y)^* \prod_i P(x_i|y)$



Naïve Bayes

Example. For the given features W,X and Y and Class label, predict the class label for W=F, X=T, Y=F.

W	X	Y	Class	
Т	т	т	Т	
Т	F	т	F	
Т	F	F	F	
F	т	т	F	
F	F	F	Т	
F	т	F	???	

MEMPHIS Naive Bayes

THE UNIVERSITY OF

- We calculate conditional probabilities for the outcomes and choose the one that is most likely:
- P(W = F, X = T, Y = F | Class = T) = P(W = F | Class = T) * P(X = T | Class = T) * P(Y = F | Class = T) =0.5 * 0.5 * 0.5 = 0.125
- P(Class = T | W = F, X = T, Y = F) = 0.125 * P(Class = T)= 0.125 * 0.4 = 0.05 Similarly;
- P = (Class = F|W = F, X = T, Y = F) = P(W = F|Class = F) * P(X = T|Class = F) * P(Y = F|Class = F) *P(Class = F) = 0.022
- The predicted class is therefore T

Decision trees

- Uses divide and conquer approach to split data to smaller subsets
- More intuitively it is like a series of IF-THEN statements Example. Consider decision tree for playing badminton based on conditions outside



www.memphis.edu

Decision trees

- The root node represents the entire dataset.
- The algorithm chooses a feature that is most predictive of the target class and then partitions the examples into groups based on values of the feature. Here, "sunny", "cloudy", and "rainy".
- The algorithm continues this divide and conquer approach until there are no feature remaining or the size limit for the tree is reached.

Decision trees

- How to decide feature for split?
- Different methods are used to decide the feature for branching.
- Some algorithms used to determine feature upon which to branch are Chi-square, entropy, and Gini index
- For example using Chi-square;
 - Calculate Chi-square for all the potential features and their target class expected on branching

$$Chi - square = \sqrt{\frac{(Target_{Actual} - Target_{Expected})^2}{Target_{Expected}}}$$

Branch along the feature with the highest Chi-square

Dreamers. Thinkers. Doers.

THE UNIVERSITY OF **MEMPHIS**

Decision trees

Example of split (Taken from a Medium <u>post</u> by Rishabh Jain.) We want to segregate the students based on target class (play cricket or not). There are two possible feature on which to split; 1) gender, 2) class. Select feature to be branched (split) upon.



MEMPHIS

Decision trees

Calculate Chi-square for each of the features: Gender

Node Play Cricket	Not Play Cricket	Total	Expected Play Cricket	Expected Not Play Cricket	Deviation Play Cricket	Deviation Not Play Cricket	Chi-Square		
							Play	Not Play	
				•	•		•	Cricket	Cricket
Female	2	8	10	5	5	-3	3	1.34	1.34
Male	13	7	20	10	10	3	-3	0.95	0.95
			_	-	-	-	Total Chi-Square	4.	58

Class

		Not Play		Expected	Expected Not	Deviation	Deviation Not	Chi-Square	
Node	Play Cricket	Cricket	Total	Play Cricket	Play Cricket	Play Cricket	Play Cricket	Play	Not Play
					,	,		Cricket	Cricket
IX	6	8	14	7	7	-1	1	0.38	0.38
X	9	7	16	8	8	1	-1	0.35	0.35
							Total Chi-Square	1.	46

Chi-square for gender is more so Gender split is more significant.



Perceptron

- Perceptron is a binary linear classifier that can separate only linearly separable data(could be a line, plane, etc.).
- Visualize perceptron decision surface as;
 - Line if for two features
 - Plane for three features
 - N-dimensional hyperplane for n features
- Given the equation for the decision surface, classification can is done (think of regression using categorical outcomes)



Dreamers. Thinkers. Doers.

Perceptron

• Given inputs x_1, x_2, \dots, x_n , with weights $w_0, w_1, w_2, \dots, w_n$ perceptron produces output based on a function called activation function.





Perceptron

Example. Find a perceptron to classify the following data with features X1, X2 and target class Y

X1	X2	Y	
0	0	-1	
0	1	1	
1	0	1	
1	1	1	

Perceptron

• The figure shows three possible perceptrons which can be represented as $w_0 + w_1 * X1 + w_2 * X2$.

Here w_0, w_1 and w_2 are weights that need to be estimated that can classify the points



makin od

27

Perceptron

EMPHIS

THE UNIVERSITY OF

- Note that infinite possible weight settings are possible.
- Let us look at one possible set of weights weight settings for the problem.
- We use the weights; $w_0 = -0.5, w_1 = 1 \text{ and } w_2 = 1$.

The perceptron: -0.5 + X1 + X2

X1	X2	Y (True)	Y (Predicted)	
0	0	-1	-0.5 (-ve so output -1)	
0	1	1	0.5 (+ve so output 1)	X
1	0	1	0.5 (+ve so output 1)	
1	1	1	1.5 (+ve so output 1)	



Perceptron learning

• In the previous example; simple hit and trial with weights could help us determine the perceptron.

But What if we have 1000 features?

- Gradient descend is used to estimate the weights for the features
- Given a set of n training examples (x₁, x₂, x₃...x_n) we compute corresponding weights (w₀, w₁, w₂,w_n) such that perceptron classifies the training examples correctly
- We minimize the squared difference between expected and observed output slowly by moving along the function and updating weights as we move along



Dreamers. Thinkers. Doers.

Gradient descent (Taken from "Machine learning" by Tom Mitchell (1997))

Gradient Descent

To understand, consider simpler *linear unit*, where

$$o = w_0 + w_1 x_1 + \dots + w_n x_n$$
 Perceptron

Let's learn w_i 's that minimize the squared error

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Where D is set of training examples

Sum of $(true \ label - \ observed \ label)^2$ for training example, d

Dreamers. Thinkers. Doers.

Gradient descent (Taken from "Machine learning" by Tom Mitchell (1997))

 Direction of steepest gradient on the surface is found by finding the derivative of E w.r.t. the components of vector of weights

$$abla E[\vec{w}] \equiv \left[rac{\partial E}{\partial w_0}, rac{\partial E}{\partial w_1}, \cdots rac{\partial E}{\partial w_n}
ight]$$

• This can be iterated over each weights by differentiation

$$\frac{\partial E}{\partial w_{i}} = \frac{\partial}{\partial w_{i}} \frac{1}{2} \sum_{d} (t_{d} - o_{d})^{2}$$

$$= \frac{1}{2} \sum_{d} \frac{\partial}{\partial w_{i}} (t_{d} - o_{d})^{2}$$
For training example d:

$$t_{d} - true \ label$$

$$o_{d} - observed \ or \ predicted \ label$$

$$x_{i,d} - feature \ input \ i$$

$$\frac{\partial E}{\partial w_{i}} = \sum_{d} (t_{d} - o_{d}) \frac{\partial}{\partial w_{i}} (t_{d} - \vec{w} \cdot \vec{x_{d}})$$

$$\frac{\partial E}{\partial w_{i}} = \sum_{d} (t_{d} - o_{d}) (-x_{i,d})$$
3



Dreamers. Thinkers. Doers.



www.memphis.edu

Artificial Neural Networks

- ANN is a multilayered perceptron
- There are three distinct layers; the input, hidden and output.
- Each input node represents a feature and output node represents output
- ANN are called black box because the operations in the hidden layers cannot be observed.
- Variable output produced from the layers depends on threshold functions as with perceptrons
- More complex backward and forward propagation methods are used in estimating weights by distributing errors to the nodes





Dreamers. Thinkers. Doers.

K-means clustering

EMPHI

- Clustering clusters data in groups by grouping similar instances
- Useful in detecting patterns e.g. in
 - consumer's spending pattern
 - group emails based on keywords
- Most useful when label is not known, or we don't know what we are looking for



K-means clustering

- We group similar instances together, but what does "similar" mean?
- Similarity can be based on some distance measure
- For 2D data points, grouping can be based on Euclidean distance between the points.
- Squared Euclidian distance for points $x = (x_1, x_{2...}, x_n)$ and $y = (y_1, y_{2...}, y_n)$: $d(y, x) = \sum_{i=1}^{n} (y_i - x_i)^2$

Dreamers. Thinkers. Doers.

How does K-means clustering work?

- Steps in K-means clustering
 - 1. Pick K random points
 - 2. Assign data instances to the closest cluster center
 - 3. Update cluster center to the average of assigned points
 - 4. Stop when all points are assigned else go to Step 2

Dreamers. Thinkers. Doers.

K-means clustering



Dreamers. Thinkers. Doers.

K-means clustering





Dreamers. Thinkers. Doers.

Evaluation metrics for ML classifiers

- Accuracy alone is not a good measure of the accuracy of an algorithm.
- Consider an algorithm for detecting email spams. Out of 10,000 emails we know there are 100 spams.
 Case 1: Any email the algorithm sees is classified as Not-spam.
 - Accuracy for spam emails only = 0/100

Case 2: Any email the algorithm sees is classified as spam

- Accuracy for spam = 100/100

Is the second algorithm better?

 Cases where an outcome is less frequent is a common problem warranting use of accuracy based on rate of false positive and false negatives predictions



Dreamers. Thinkers. Doers.

Precision, recall and accuracy

EMPHIS

- Two standard measures to evaluate quality of prediction.
 - Precision: Accuracy in terms of classifying relevant instances

 $Precision = \frac{True \ positive}{Actual \ results} = \frac{True \ positive}{True \ positive + False \ positive}$

- Recall : Accuracy in classifying all relevant instances correctly

 $Recall = \frac{True \ positive}{Predicted \ results} = \frac{True \ positive}{True \ positive + False \ negative}$

- Also, $Accuracy = \frac{True \ positive + True \ negative}{Total \ instances}$



MEMPHIS

F1 score

- Combines precision and recall
- Calculated as harmonic mean of precision and recall
- **F1 score,** $F1 = \frac{2*Precision*Recall}{Precision+Recall}$



Example 1

- Consider there are 100 known spam emails in a dataset with
- 10,000 emails. Find precision, recall assuming all the emails
- the algorithm sees are classified as Not-spam.
 - True positive=0, False positive=0, false negative=100
 - Precision=undefined
 - Recall=0
 - F1 score =undefined;
 - Conclusion not a good algorithm

Example 2

- Consider there are 500 known spam emails in a dataset with
- 1000 emails. Find precision, recall assuming all the emails the algorithm sees are classified as spams.
 - True positive=500, False positive=500, false negative=0
 - Precision=0.5
 - Recall=1
 - F1 score =0.667;
 - Conclusion relatively good algorithm