

Homework 2: Written Assignment

Handed Out: September 17, 2019

Due: October 7, 2019 11:55pm

Save your homework submission as *NETID-hw2-written.pdf*.

1 Naïve Bayes and Decision Tree (40 points)

Consider the following training set:

A	B	C	Y
0	1	1	0
1	1	1	0
0	0	0	0
1	1	0	1
0	1	0	1
1	0	1	1

Here A , B , and C are features; Y is label. Each row is a data object.

- (a) Make prediction for $(A = 1, B = 0, C = 0)$ using Naïve Bayes. Learn the Naïve Bayes classifier by estimating all necessary probabilities (including posteriori probabilities, prior probabilities, and likelihood probabilities).
- (b) Learn a Decision Tree from the above training set using the Information Gain criterion.

2 Kernel Function (20 points)

In class, we showed that the quadratic kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^2$ was equivalent to mapping each \mathbf{x} into a higher dimensional space where

$$\Phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$$

for the case where $\mathbf{x} = (x_1, x_2)$. Now consider the cubic kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^3$. What is the corresponding Φ function (again, for the special case where $\mathbf{x} = (x_1, x_2)$)?

3 Classification Evaluation (10 points)

(Choose one answer and prove it) With which of the following conditions, F1 score is equivalent to Accuracy in binary classification?

1. TP = TN; 2. TP = FP; 3. TP = FN; 4. FP = FN

TP: number of true positives; TN: true negatives; FP: false positives; FN: false negatives.

Homework 2: Programming Assignment

*Handed Out: September 17, 2019**Due: October 7, 2019 11:55pm*

Save your homework submission as *NETID-hw2-programming.zip*. The zip file has one pdf file *NETID-hw2-programming.pdf* and multiple code files.

Binary Classification with Categorical Features on ND Game Data using no Classification Packages (70 points)

Please use **Python** to solve the problems. You are NOT allowed to directly call any Classification functions (like the decision trees or naive bayes functions in Scikit).

Data files: Dataset-football-train.csv (**training set**), Dataset-football-test.csv (**test set**)

Data introduction: Given Notre Dame's football game data for the last two seasons (2015 and 2016), can we construct three classification models to predict game results on games in 2017? Can we evaluate the model performance? **See the table on the last page.** Each data object (or called instance) is a game. We have three attributes: (1) "Is Home/Away?", a 2-value attribute ("Home", "Away"), (2) "Is Opponent in AP Top 25 at Preseason?", a 2-value attribute ("In", "Out"), (3) "Media", a 5-value attribute ("1-NBC", "2-ESPN", "3-FOX", "4-ABC", "5-CBS"). The label "Win/Lose" is binary ("Win", "Lose").

- **Training set:** 24 games. Please use game ID 1-24 to construct classification models. (Background color: YELLOW)
- **Test set:** 12 games. Please use your models to predict labels of game ID 25-36 and evaluate the performance of the classification models. (Background color: BLUE)
- **Labels:** Suppose "Win" is the positive label and "Lose" is the negative label. Keep it in mind when you use Precision and Recall to evaluate the models.

For decision tree model construction: We stop splitting instances into child nodes when one of the criteria is satisfied:

- (1) All features have been used;
- (2) Information Gain or Gain Ratio will be zero with any feature that has not yet been used.

For decision tree model usage:

- (1) If the node is not pure, we predict with the majority: For example, if we have 5 positives and 1 negatives, we predict the testing case at this node to be a positive.
- (2) If the node has a balance (half/half labels), e.g., 2 positives and 2 negatives, we use the majority of the root node (the entire dataset) for prediction.

Q1: ID3 model using Information Gain (20 points)

Use ID3 to construct a decision tree based on the training set (24 games). Use the tree to predict labels of instances in the testing set (12 games) based on their attributes. Calculate Accuracy, Precision, Recall, and F1 score on the testing result.

Output: Write down (1) your decision tree (either hand-drawn or electronically drawn), (2) predicted labels of the 12 testing games, and (3) evaluation results in the pdf. Save your code as NETID-hw2-1.py.

Q2: C4.5 model using Gain Ratio (20 points)

Use C4.5 to construct a decision tree based on the training set (24 games). Use the tree to predict labels of instances in the testing set (12 games) based on their attributes. Calculate Accuracy, Precision, Recall, and F1 score on the testing result.

Output: Write down (1) your decision tree (either hand-drawn or electronically drawn), (2) predicted labels of the 12 testing games, and (3) evaluation results in the pdf. Save your code as NETID-hw2-2.py.

Q3: Naive Bayes model without Zero Correction (30 points)

Use Naive Bayes to predict labels of instances in the test set (12 games) based on the training set (24 games). Calculate Accuracy, Precision, Recall, and F1 score on the testing result.

Output: Write down (1) predicted labels of the 12 testing games and (2) evaluation results in the pdf. Save your code as NETID-hw2-3.py.

Multi-Class Classification with Numerical Features on Film Data using CART (30 points)

You are allowed to use *any* programming language (Python recommended; R, C++, Java, etc.), however, the solutions will be in **Python**. You are allowed to use *any* public package (including Numpy and Scikit-learn) and any other kind of tools (Excel).

Please self learn the Section 1.10.1 on Web Page:

<http://scikit-learn.org/stable/modules/tree.html>.

Construct the CART model following the instruction and apply on the film data we used in Homework 1. All the 150 films will be used as training data points. You are not asked to do model usage or model evaluation.

Output: Write down the CART model (the decision tree) that was learned from the film

data and in the pdf. We strongly recommend you export your tree by in Graphviz format. For more detail of this package, here is the link: <http://www.graphviz.org>. Save your code as NETID-hw2-4.py.

ID	Date	Opponent	Is Home or Away?	Is Opponent in AP Top 25 at Preseason?	Media	Label: Win/Lose
1	9/5/15	Texas	Home	Out	1-NBC	Win
2	9/12/15	Virginia	Away	Out	4-ABC	Win
3	9/19/15	Georgia Tech	Home	In	1-NBC	Win
4	9/26/15	UMass	Home	Out	1-NBC	Win
5	10/3/15	Clemson	Away	In	4-ABC	Lose
6	10/10/15	Navy	Home	Out	1-NBC	Win
7	10/17/15	USC	Home	In	1-NBC	Win
8	10/31/15	Temple	Away	Out	4-ABC	Win
9	11/7/15	PITT	Away	Out	4-ABC	Win
10	11/14/15	Wake Forest	Home	Out	1-NBC	Win
11	11/21/15	Boston College	Away	Out	1-NBC	Win
12	11/28/15	Stanford	Away	In	3-FOX	Lose
13	9/4/16	Texas	Away	Out	4-ABC	Lose
14	9/10/16	Nevada	Home	Out	1-NBC	Win
15	9/17/16	Michigan State	Home	Out	1-NBC	Lose
16	9/24/16	Duke	Home	Out	1-NBC	Lose
17	10/1/16	Syracuse	Home	Out	2-ESPN	Win
18	10/8/16	North Carolina State	Away	Out	4-ABC	Lose
19	10/15/16	Stanford	Home	In	1-NBC	Lose
20	10/29/16	Miami Florida	Home	Out	1-NBC	Win
21	11/5/16	Navy	Home	Out	5-CBS	Lose
22	11/12/16	Army	Home	Out	1-NBC	Win
23	11/19/16	Virginia Tech	Home	In	1-NBC	Lose
24	11/26/16	USC	Away	In	4-ABC	Lose
25	9/2/17	Temple	Home	Out	1-NBC	Win
26	9/9/17	Georgia	Home	In	1-NBC	Lose
27	9/16/17	Boston College	Away	Out	2-ESPN	Win
28	9/23/17	Michigan State	Away	Out	3-FOX	Win
29	9/30/17	Miami Ohio	Home	Out	1-NBC	Win
30	10/7/17	North Carolina	Away	Out	4-ABC	Win
31	10/21/17	USC	Home	In	1-NBC	Win
32	10/28/17	North Carolina State	Home	Out	1-NBC	Win
33	11/4/17	Wake Forest	Home	Out	1-NBC	Win
34	11/11/17	Miami Florida	Away	In	4-ABC	Lose
35	11/18/17	Navy	Home	Out	1-NBC	Win
36	11/25/17	Stanford	Away	In	4-ABC	Lose