# CS 440 - Homework 6

Matthew Liu

**1.**

Grade
Gain: 0.413

Freshman
6+ 0-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 3  | F     | N   | +        |
| 4  | F     | N   | +        |
| 11 | F     | N   | +        |
| 14 | F     | N   | +        |
| 17 | F     | N   | +        |
| 18 | F     | N   | +        |

0.0

→ +

Sophomore
2+ 3-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 1  | S     | N   | -        |
| 6  | S     | O   | -        |
| 8  | S     | N   | -        |
| 13 | S     | O   | +        |
| 20 | S     | O   | +        |

0.971

BMI
Gain: 0.4202

O
2+ 1-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 6  | S     | O   | -        |
| 13 | S     | O   | +        |
| 20 | S     | O   | +        |

0.918

By Majority

→ +

N
0+ 2-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 1  | S     | N   | -        |
| 8  | S     | N   | -        |

0.0

→ -

Juniors
2+ 7-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 2  | J     | O   | +        |
| 5  | J     | N   | -        |
| 7  | J     | N   | -        |
| 9  | J     | N   | -        |
| 10 | J     | N   | +        |
| 12 | J     | N   | -        |
| 15 | J     | O   | -        |
| 16 | J     | O   | -        |
| 19 | J     | O   | -        |

0.764

BMI
Gain: 0.002

O
1+ 3-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 2  | J     | O   | +        |
| 15 | J     | O   | -        |
| 16 | J     | O   | -        |
| 19 | J     | O   | -        |

0.811

By Majority

→ -

N
1+ 4-

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 5  | J     | N   | -        |
| 7  | J     | N   | -        |
| 9  | J     | N   | -        |
| 10 | J     | N   | +        |
| 12 | J     | N   | -        |

0.722

By Majority

→ -

(a) Decision Tree is attached above

(b) Classifications:

- Student 21: +
- Student 22: +
- Student 23: -
- Student 24: +
- Student 25: -
- Student 26: -
- Student 27: +
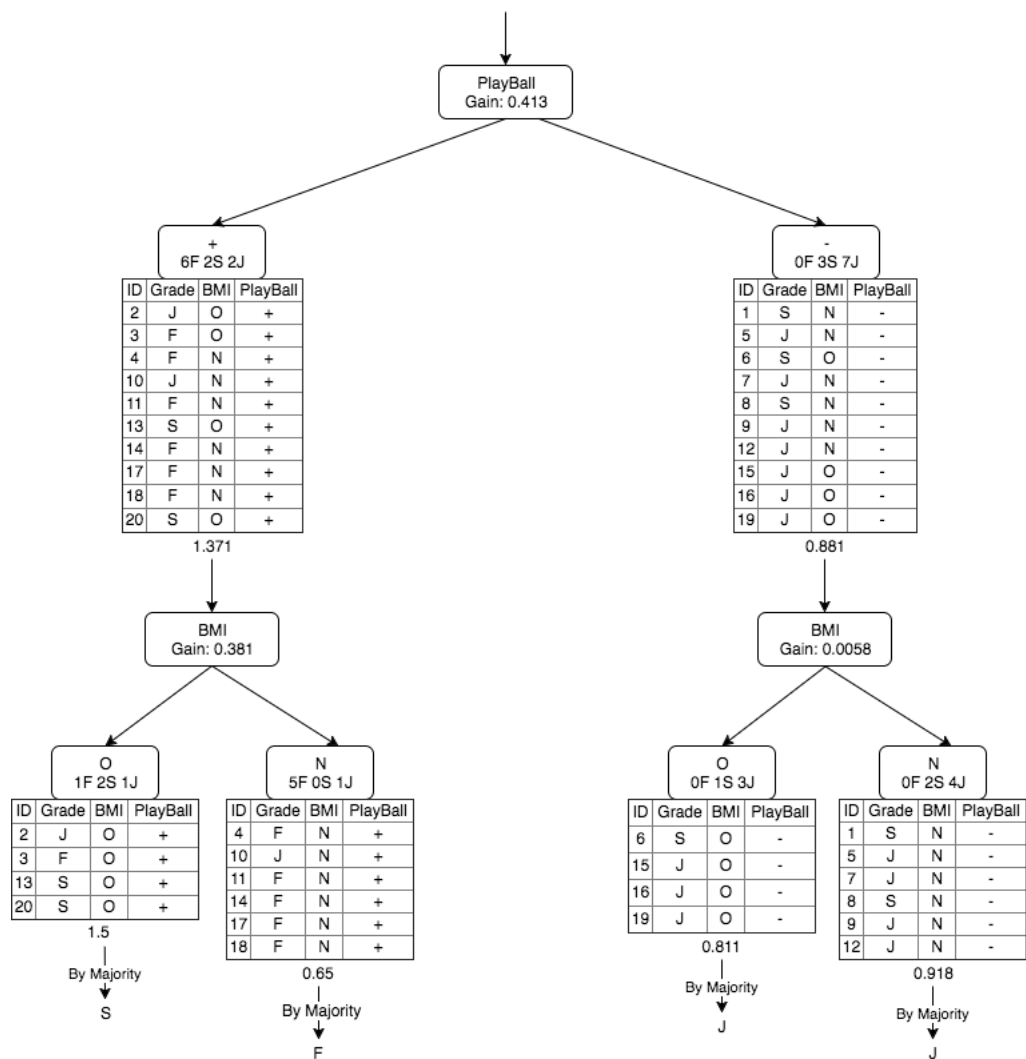- Student 28: +
- Student 29: -
- Student 30: +

(c) Classifications:

- Student 11: +
- Student 12: -
- Student 13: +
- Student 14: +
- Student 15: -
- Student 16: -
- Student 17: +
- Student 18: +
- Student 19: -
- Student 20: +

(d) The results from **(b)** would be the better estimate as **(c)** was part of the training set for the classifier

(e) The classifier from Scenario 2 would be significantly more accurate then the classifier from Scenario 1, as the Scenario 1 does not have enough training data to produce meaningful results

(f) Decision Tree is attached below

PlayBall
Gain: 0.413

+
6F 2S 2J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 2  | J     | O   | +        |
| 3  | F     | O   | +        |
| 4  | F     | N   | +        |
| 10 | J     | N   | +        |
| 11 | F     | N   | +        |
| 13 | S     | O   | +        |
| 14 | F     | N   | +        |
| 17 | F     | N   | +        |
| 18 | F     | N   | +        |
| 20 | S     | O   | +        |

1.371

-
0F 3S 7J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 1  | S     | N   | -        |
| 5  | J     | N   | -        |
| 6  | S     | O   | -        |
| 7  | J     | N   | -        |
| 8  | S     | N   | -        |
| 9  | J     | N   | -        |
| 12 | J     | N   | -        |
| 15 | J     | O   | -        |
| 16 | J     | O   | -        |
| 19 | J     | O   | -        |

0.881

BMI
Gain: 0.381

O
1F 2S 1J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 2  | J     | O   | +        |
| 3  | F     | O   | +        |
| 13 | S     | O   | +        |
| 20 | S     | O   | +        |

1.5

By Majority

S

N
5F 0S 1J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 4  | F     | N   | +        |
| 10 | J     | N   | +        |
| 11 | F     | N   | +        |
| 14 | F     | N   | +        |
| 17 | F     | N   | +        |
| 18 | F     | N   | +        |

0.65

By Majority

F

BMI
Gain: 0.0058

O
0F 1S 3J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 6  | S     | O   | -        |
| 15 | J     | O   | -        |
| 16 | J     | O   | -        |
| 19 | J     | O   | -        |

0.811

By Majority

J

N
0F 2S 4J

| ID | Grade | BMI | PlayBall |
|----|-------|-----|----------|
| 1  | S     | N   | -        |
| 5  | J     | N   | -        |
| 7  | J     | N   | -        |
| 8  | S     | N   | -        |
| 9  | J     | N   | -        |
| 12 | J     | N   | -        |

0.918

By Majority

J

**2.**

(a) The following are three possible approaches

- Limit the depth: limit how deep the decision tree can go
- Limit the minimum number of example used to select a split: Require a minimum amount of data to continue splitting the decision tree. This is similar to limiting the depth, but does so differently
- Cross Validation: Partition data into training set and testing set, and rotate data for both. sizes of both sets can vary.

(b) Cross Validation

(c) I would break up the data into sets of 5 (into a total of N sets), and rotate which N-1 sets will be used for training use the remaining set as a evaluation/testing set. This is better than limiting the depth because that forcefully prunes the tree, and in a tree like this one where it can't really go all that deep, all it does is create a poor classifier. Limiting the minimum number of examples used to select a split would work if there was a plethora of data, however this dataset does not have a wealth of data, hence this method would again limit the effectiveness of our classifier

(d) As this algorithm rotates through and considers all possible combinations of data subsets for training and evaluation, it will always produce the shortest decision tree because of it. Of course the cost is that depending on the size of the dataset, this algorithm can take a while to complete.