

# 1.Implementation

a) I'll try to describe at a high level the step I took in order to get to a decision tree:

- A. Using the class ***DataSetLoader***, I am loading the training data as an array having the following form:

```
[
    ['Wil van der Aalst', '-'],
    ['Scott Aaronson', '+'],
    ...
]
```

- B. I use then the ***DataSetFeaturesEnricher*** to create the features for each row. At the end I get an array having the following form:

```
[
    [False, False, False, False, True, False, '+'],
    [False, True, False, True, False, True, '-'],
    ...
]
```

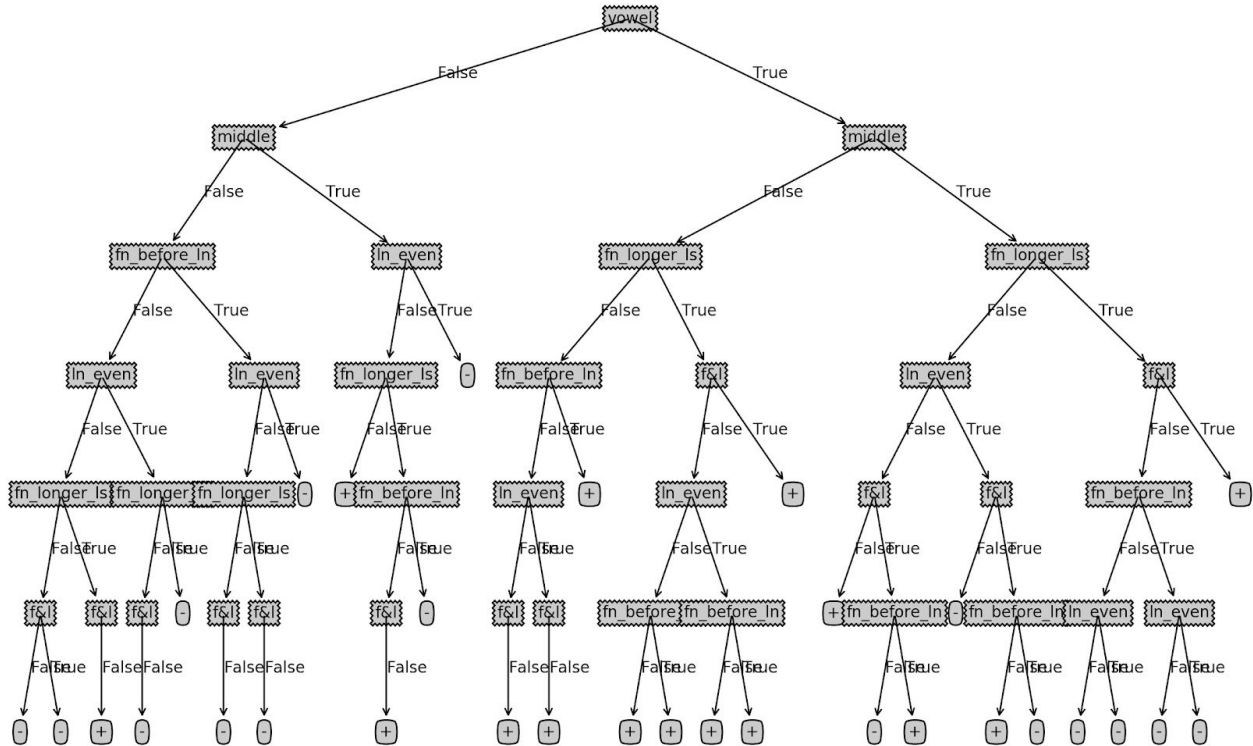
The first 6 columns are the features values, while the last is the label for that entry.

- C. The array from point B is sent to the ***DecisionTree*** class where using the ID3 algorithm I return the tree having the form:

```
{
  'vowel': {
    False: {
      'middle': {
        False: '-',
        True: {... }
      }
    },
    True: {
      'middle': {
        False: '+',
        True: {
          'ln_even': {
            False: {... },
            True: {... }
          }
        }
      }
    }
  }
}
```

}

D. I then render my tree using the **mathplotlib** module. I get the following image after I train a decision tree using the 6 features listed in the **feature\_description.txt** file:

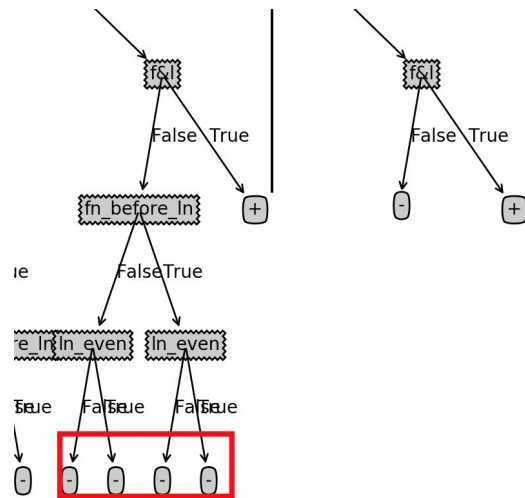


### Legend:

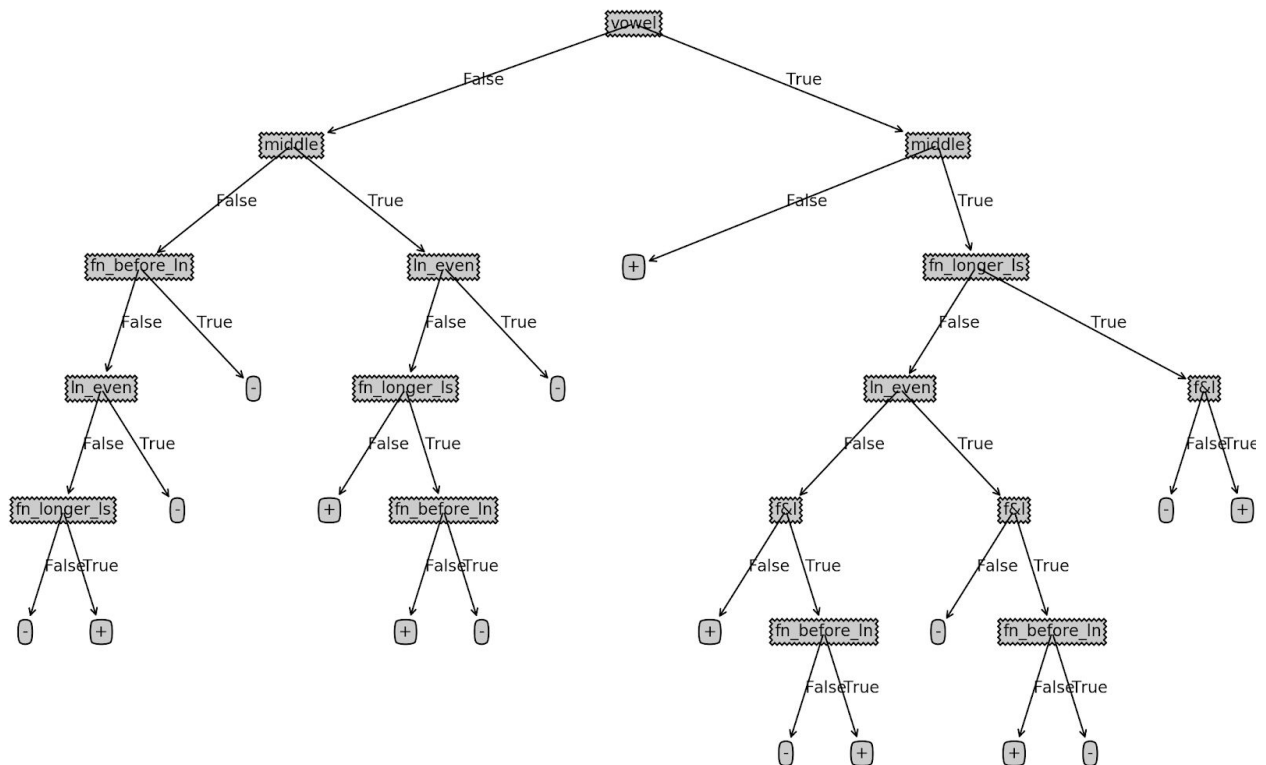
1. vowel - Is the second letter of their first name a vowel (a,e,i,o,u)?
2. middle - Do they have a middle name?
3. fn\_even - Is the number of letters in their last name even?
4. fn\_longer\_Is - Is their first name longer than their last name?
5. fn\_before\_In - Does their first name come alphabetically before their last name?
6. f&l - Does their first name start and end with the same letter?

E. In the previous picture you can notice that some features leaves are returning the same label on both branches (True/False). I created a class called **TreePruner** where I try to fix this.

For example the 4 leaves from below are becoming only one (having label '-') under the 'f&l' label.



After the pruning operation I get to this final tree:



F. After I get this tree, I send it to a **Classifier** class that will return labels for any new entry. Each new entry should be an array of 6 features.

b) Suggest at least 4 other features you could have extracted from this dataset.

1. The number of vowels in the full name is odd;
2. Last letter of the last name is vowel;
3. Does the full name contains more than 2 names;
4. Does the full name contains less than 15 characters.

c) Report the error of your decision tree on the **Dataset/training.data** file.

The error rate that I got on the training data was: **7.19%**

d) Report the error of your decision tree on the **Dataset/test.data** file.

The error rate that I got on the test data was: **8.11%**

e) Report the maximum depth of your decision tree.

The maximum depth of my decision tree is 6 (equal with the number of features used).

## 2. Limiting depth

a) Using the steps described in the homework description I used 10 features (6 feature listed in the **feature\_descriptions.txt** file and the 4 features that I listed inside the **Implementation** section (b)) to train the tree using the k-fold cross validation method. I got to the following values:

Depth	Error rate when <b>training00.data</b> was used as test set	Error rate when <b>training01.data</b> was used as test set	Error rate when <b>training02.data</b> was used as test set	Error rate when <b>training03.data</b> was used as test set	Average error rate	Standard deviation
1	13.51%	18.92%	18.92%	18.02%	17.34%	2.59
2	14.41%	11.71%	10.81%	10.81%	11.94%	1.7
3	14.41%	9.91%	9.01%	10.81%	11.04%	2.37
4	14.41%	8.11%	7.21%	6.31%	9.01%	3.67
<b>5</b>	<b>9.01%</b>	<b>6.31%</b>	<b>8.11%</b>	<b>8.11%</b>	<b>7.88%</b>	<b>1.13</b>
10	10.81%	10.81%	13.51%	7.21%	10.59%	2.59

We can see that a tree having depth 5 has the lowest error rate (7.88%) while also has the lowest standard deviation. That means that the values used for calculating the average are not spread.

b) Using the depth with the greatest cross-validation accuracy from your experiments: train your decision tree on the Dataset/training.data file. Report the accuracy of your decision tree on the Dataset/test.data file.

I tested a decision tree using all the training data, but limiting the depth of the tree to 5. The error rate that I got on the testing data using the tree with depth 5 was: **5.41%**

c) Discuss the performance of the depth limited tree as compared to the full decision tree. Do you think limiting depth is a good idea? Why?

When we limited the depth to 5 (the initial tree had the depth 6), the error rate on the test data went from 8.11% down to 5.41%. I think the initial tree was overfitted on the training data. By cutting the depth with one level, we removed a feature that was not bringing a lot of ordering in the data (that feature had the lowest entropy gain) but was making at the same time the tree overfitted. Not using all the features, make the tree more generic.

### 3. Extra validation

For validation my tree, I used also the decision tree classifier from the SciKit project (<http://scikit-learn.org/stable/modules/tree.html>). The SciKit classifier uses a optimized version of CART. The tree that I got is similar with the one created by my script:



