

# Report for HW2

H24111057 姚博瀚

11/27/2023

## Section 1: 了解資料

### Part 1: 特徵 (Feature) 的資料屬性

首先我想先確認有無missing value (code 1)，再來分出哪些特徵為二元、連續還是離散的型態，code 2依照特徵的值是否只有兩種、整數、小數分為二元、離散、連續。接著我想看每個特徵的數值的分位數以及最大最小值，以describe()產生後再匯出成圖片 (code 3)。

Listing 1: Identify missing values

```
1 missing_values = df_train.isnull().sum()
2 df_missing_values = pd.DataFrame(missing_values)
3 df_missing_values.columns = ['Count']
4 print(df_missing_values)
5 -----
6                Count
7 profile pic         0
8 nums/length username 0
9 fullname words      0
10 nums/length fullname 0
11 name==username      0
12 description length   0
13 external URL         0
14 private              0
15 posts                0
16 followers            0
17 follows              0
18 fake                 0
```

Listing 2: Identify data type

```
1 import numpy as np
2 import pandas as pd
3
4 # Identify numeric features
5 numeric_features = X_train.select_dtypes(include=[np.number]).columns.tolist()
6
7 # Initialize lists for different types of features
8 binary_features = []
```

```

9 continuous_features = []
10 discrete_features = []
11
12 # Define a threshold to differentiate between binary and continuous/↔
    discrete
13 binary_threshold = 2
14
15 # Loop through each numeric feature
16 for feature in numeric_features:
17     unique_values = X_train[feature].nunique()
18     # Check if the feature is binary
19     if unique_values <= binary_threshold:
20         binary_features.append(feature)
21     else:
22         # Check if the feature has only integer values to classify as ↔
            discrete
23         if X_train[feature].dtype == np.int64:
24             discrete_features.append(feature)
25         else:
26             continuous_features.append(feature)
27
28 # Print the results
29 print("Numeric Features:")
30 if len(numeric_features) == len(X_train.columns):
31     print("All features are numeric")
32 else:
33     print(numeric_features)
34 print("\nBinary Features:")
35 print(binary_features)
36 print("\nContinuous Features:")
37 print(continuous_features)
38 print("\nDiscrete Features:")
39 print(discrete_features)
40 -----
41 Numeric Features:
42 All features are numeric
43
44 Binary Features:
45 ['profile pic', 'name==username', 'external URL', 'private']
46
47 Continuous Features:
48 ['nums/length username', 'nums/length fullname']
49
50 Discrete Features:
51 ['fullname words', 'description length', '#posts', '#followers', '#follows↔
    ']
```

Listing 3: Summary of features

```
1 import matplotlib.pyplot as plt
2
3 # X_train is a DataFrame
4 summary = X_train.describe()
5
6 # Format the values with 4 decimal places
7 formatted_values = summary.values.round(4)
8
9 # Create a figure and axis with specified DPI
10 fig, ax = plt.subplots(figsize=(16, 2), dpi=300)
11
12 # Plot the table with adjusted fontsize and expanded space for column ↵
    labels
13 ax.axis('off') # Turn off axis
14 table = ax.table(cellText=formatted_values,
15 rowLabels=summary.index,
16 colLabels=summary.columns,
17 loc='center',
18 cellLoc='center',
19 rowLoc='center',
20 colColours=['#f0f0f0'] * len(summary.columns)) # Set background color for↵
    column labels
21
22 table.auto_set_font_size(False)
23 table.set_fontsize(10) # Adjust the fontsize
24
25 # Adjust the space for column labels
26 table.auto_set_column_width([0] + list(range(len(summary.columns))))
27
28 # Save the figure as an image with adjusted DPI
29 fig.savefig('summary_features_table.png', dpi=300)
30
31 # Display the summary table
32 plt.show()
```

由Figure 1可得知訓練集（train.csv）共有576個樣本，並且post、followers、follows這三個特徵的數值皆偏大，後續可能進行標準化或是什麼轉換以降低此特徵的影響力過大。

## Part 2: 目標（Target）的資料屬性

由以下code可進一步得知訓練集的576個樣本之中各有288個假帳號與288個真帳號。

Listing 4: Check Target's value

```
1 # y_train only contains fake columns extracted from train.csv
```

```

2 y_train.value_counts()
3 -----
4 0      288
5 1      288
6 Name: fake, dtype: int64

```

---

	profile pic	nums/length username	fullname words	nums/length fullname	name==username
count	576.0	576.0	576.0	576.0	576.0
mean	0.7014	0.1638	1.4601	0.0361	0.0347
std	0.458	0.2141	1.0526	0.1251	0.1832
min	0.0	0.0	0.0	0.0	0.0
25%	0.0	0.0	1.0	0.0	0.0
50%	1.0	0.0	1.0	0.0	0.0
75%	1.0	0.31	2.0	0.0	0.0
max	1.0	0.92	12.0	1.0	1.0

	description length	external URL	private	#posts	#followers	#follows
count	576.0	576.0	576.0	576.0	576.0	576.0
mean	22.6233	0.1163	0.3819	107.4896	85307.2361	508.3819
std	37.703	0.3209	0.4863	402.0344	910148.4577	917.9812
min	0.0	0.0	0.0	0.0	0.0	0.0
25%	0.0	0.0	0.0	0.0	39.0	57.5
50%	0.0	0.0	0.0	9.0	150.5	229.5
75%	34.0	0.0	1.0	81.5	716.0	589.5
max	150.0	1.0	1.0	7389.0	15338538.0	7500.0

Figure 1: Summary table of features

## Section 2: 前處理

由code 1 已知所有的特徵以及目標皆無遺失值，並且已知特徵皆是數字型態（code 2）。由於我後續建立的Decision Tree在decision code要分割樣本的時候可以直接處理二元、離散、連續性的特徵（皆以兩個樣本的數值中點當分割依據），因此特徵不做任何的轉換皆已適合模型的輸入。

但為了避免post、followers、follows這三個特徵的影響力過大，我對這三個特徵的數值進行0.3次方的轉換，目的是把數值的尺度縮小。原本是想進行0.5次方的轉換（也就是把特徵的數值開根號），但做0.3次方的轉換讓followers的最大值比較接近其餘沒做轉換的特徵的最大值（約150）。如Figure 2所示，左邊為原來的資料的分布圖，右邊則為做完0.3次方轉換後的分布圖。Figure 3則為做完轉換後的summary table。後續模型的輸入皆為直接用轉換過後的資料（X\_train\_root）。

## Section 3: 建立模型

依照講義的作法，使用Entropy來計算node的impurity，再依照Information gain 來決定要用哪個特徵以及那個特徵的哪個threshold來分割目前的node。簡而言之，當使用fit fuction時，會呼叫\_grow\_tree(呼叫\_best\_split來找出最好的分割feature與threshold)來遞迴地生成下面的child node。註1：左child node是小於threshold；右child node是大於threshold。註2：二元

特徵的threshold只有一個，為1；離散與連續性特徵的thresholds是每兩個unique value的平均。

Listing 5: Decision Tree Class

```
1
2 import numpy as np
3 import pandas as pd
4
5 # Decision Tree Node class
6 class Node:
7     def __init__(self, entropy, num_samples, num_samples_per_class, ↵
        predicted_class):
8         self.entropy = entropy
9         self.num_samples = num_samples
10        self.num_samples_per_class = num_samples_per_class
11        self.predicted_class = predicted_class
12        self.feature_index = 0
13        self.threshold = 0
14        self.children = {}
15
16 # Decision Tree class
17 class DecisionTree:
18     def __init__(self, max_depth = 30):
19         self.max_depth = max_depth
20
21     def fit(self, X, y):
22         self.n_classes_ = len(set(y))
23         self.n_features_ = X.shape[1]
24         self.tree_ = self._grow_tree(X, y)
25
26     def _entropy(self, y):
27         m = len(y)
28         class_probs = [np.sum(y == c) / m for c in range(self.n_classes_)]
29         return -np.sum(p * np.log2(p) if p > 0 else 0 for p in class_probs)↵
        )
30
31     def _information_gain(self, y, y_left, y_right):
32         ent_parent = self._entropy(y)
33         ent_left = self._entropy(y_left)
34         ent_right = self._entropy(y_right)
35         weight_left = len(y_left) / len(y)
36         weight_right = len(y_right) / len(y)
37         return ent_parent - (weight_left * ent_left + weight_right * ↵
            ent_right)
38
39     def _best_split(self, X, y):
```

```

40     m, n = X.shape
41     if m <= 1:
42         return None, None
43     best_info_gain = 0
44     best_idx, best_thr = None, None
45
46     for idx in range(self.n_features_):
47         feature_values = X.iloc[:, idx]
48         unique_values = feature_values.unique()
49
50         # If the feature is binary, consider only one threshold
51         if len(unique_values) == 2:
52             thresholds = [unique_values[0]]
53         else:
54             # For continuous or discrete features, consider midpoints ↵
55             # between unique values
56             thresholds = [(unique_values[i] + unique_values[i + 1]) / ↵
57                           2 for i in range(len(unique_values) - 1)]
58
59         for thr in thresholds:
60             # Split the dataset based on the current threshold
61             y_left = y[feature_values < thr]
62             y_right = y[feature_values >= thr]
63
64             # Calculate information gain
65             info_gain = self._information_gain(y, y_left, y_right)
66
67             # Update best split if current information gain is higher
68             if info_gain > best_info_gain:
69                 best_info_gain = info_gain
70                 best_idx = idx
71                 best_thr = thr
72
73     return best_idx, best_thr
74
75 def _grow_tree(self, X, y, depth=0):
76     num_samples_per_class = [np.sum(y == i) for i in range(self.↵
77         n_classes_)]
78     predicted_class = np.argmax(num_samples_per_class)
79     node = Node(
80         entropy=self._entropy(y),
81         num_samples=len(y),
82         num_samples_per_class=num_samples_per_class,
83         predicted_class=predicted_class
84     )
85
86     if depth < self.max_depth:

```

```

84         idx, thr = self._best_split(X, y)
85         if idx is not None:
86             feature_values = X.iloc[:, idx]
87             X_left, y_left = X[feature_values < thr], y[feature_values < thr]
88             X_right, y_right = X[feature_values >= thr], y[feature_values >= thr]
89             node.feature_index = idx
90             node.threshold = thr
91             node.children['left'] = self._grow_tree(X_left, y_left, depth + 1)
92             node.children['right'] = self._grow_tree(X_right, y_right, depth + 1)
93         return node
94
95     def predict(self, X):
96         return [self._predict(inputs) for inputs in X.values]
97
98     def _predict(self, inputs):
99         node = self.tree_
100         while 'left' in node.children:
101             feature_value = inputs[node.feature_index]
102
103             # Check if the feature value is numeric
104             if isinstance(feature_value, (int, float)):
105                 if feature_value < node.threshold:
106                     node = node.children['left']
107                 else:
108                     node = node.children['right']
109             else:
110                 # Handle non-numeric feature values
111                 raise ValueError("Unsupported feature type in decision tree prediction.")
112
113         return node.predicted_class

```

---

## Section 4: 優化

為了使整個模型不過於複雜以及為了避免overfitting，繼承了原本的Decision Tree class來修改\_best\_split、\_grow\_tree以增加以下四個conditions於PrePrunedDecisionTree class：

1. min\_samples\_split(default=4)：當要分割當前node所需的最小樣本數。
2. min\_samples\_leaf(default=2)：要成為leaf node(terminal node)所需的最小樣本數。
3. min\_info\_gain(default=0.05)：每次分隔node的最小所需information gain。

4. `max_features(default=None)`：每次分隔時考慮的特徵數量，會隨機抽出所設的特徵數量。`None`為考慮所有的特徵。

註：`max_depth`為原本就有的超參數。

Listing 6: Pre-Pruned Decision Tree Class

```
1 import numpy as np
2 import pandas as pd
3 import random
4
5 # Pre-Pruned Decision Tree class
6 class PrePrunedDecisionTree(DecisionTree):
7     def __init__(self, max_depth=30, min_samples_split=4, min_samples_leaf=
8         =2, min_info_gain=0.05, max_features=None):
9         super().__init__(max_depth=max_depth)
10        self.min_samples_split = min_samples_split
11        self.min_samples_leaf = min_samples_leaf
12        self.min_info_gain = min_info_gain
13        self.max_features = max_features
14
15    def _best_split(self, X, y):
16        m, n = X.shape
17        if m <= 1:
18            return None, None
19
20        best_info_gain = 0
21        best_idx, best_thr = None, None
22
23        # Select a random subset of features if max_features is specified
24        if self.max_features is not None:
25            all_features = range(self.n_features_)
26            selected_features = random.sample(all_features, min(self.
27                max_features, self.n_features_))
28        else:
29            selected_features = range(self.n_features_)
30
31        for idx in selected_features:
32            feature_values = X.iloc[:, idx]
33            unique_values = feature_values.unique()
34
35            # If the feature is binary, consider only one threshold
36            if len(unique_values) == 2:
37                thresholds = [unique_values[1]]
38            else:
39                # For continuous or discrete features, consider midpoints
40                # between unique values
41                thresholds = [(unique_values[i] + unique_values[i + 1]) /
```



```

        2 for i in range(len(unique_values) - 1)]
39
40     for thr in thresholds:
41         # Split the dataset based on the current threshold
42         y_left = y[feature_values < thr]
43         y_right = y[feature_values >= thr]
44
45         # Apply pre-pruning conditions
46         if (len(y_left) < self.min_samples_leaf or len(y_right) < ←
            self.min_samples_leaf):
47             continue
48
49         # Calculate information gain
50         info_gain = self._information_gain(y, y_left, y_right)
51
52         # Update best split if current information gain is higher
53         if info_gain > best_info_gain:
54             best_info_gain = info_gain
55             best_idx = idx
56             best_thr = thr
57
58     # Apply pre-pruning conditions
59     if best_info_gain < self.min_info_gain:
60         return None, None
61
62     return best_idx, best_thr
63
64
65 def _grow_tree(self, X, y, depth=0):
66     num_samples_per_class = [np.sum(y == i) for i in range(self.←
        n_classes_)]
67     predicted_class = np.argmax(num_samples_per_class)
68     node = Node(
69         entropy=self._entropy(y),
70         num_samples=len(y),
71         num_samples_per_class=num_samples_per_class,
72         predicted_class=predicted_class
73     )
74
75     if depth < self.max_depth:
76         idx, thr = self._best_split(X, y)
77         if idx is not None:
78             feature_values = X.iloc[:, idx]
79             X_left, y_left = X[feature_values < thr], y[feature_values←
                < thr]
80             X_right, y_right = X[feature_values >= thr], y[←
                feature_values >= thr]

```

```

81
82         # Apply pre-pruning conditions
83         if (
84             len(y_left) >= self.min_samples_split
85             and len(y_right) >= self.min_samples_split
86         ):
87             node.feature_index = idx
88             node.threshold = thr
89             node.children['left'] = self._grow_tree(X_left, y_left,
90             , depth + 1)
91             node.children['right'] = self._grow_tree(X_right, y_right,
92             depth + 1)
93
94     return node

```

---

Listing 7: Accuracy before pruning

```

1  import numpy as np
2  import pandas as pd
3
4  df_train = pd.read_csv('train.csv')
5  X_train = df_train.drop('fake', axis=1)
6  y_train = df_train['fake']
7
8  df_test = pd.read_csv('test.csv')
9  X_test = df_test.drop('fake', axis=1)
10 y_test = df_test['fake']
11
12 X_train_root = X_train.copy()
13 X_train_root['#followers'] = X_train_root['#followers']**0.3
14 X_train_root['#follows'] = X_train_root['#follows']**0.3
15 X_train_root['#posts'] = X_train_root['#posts']**0.3
16
17 X_test_root = X_test.copy()
18 X_test_root['#followers'] = X_test_root['#followers']**0.3
19 X_test_root['#follows'] = X_test_root['#follows']**0.3
20 X_test_root['#posts'] = X_test_root['#posts']**0.3
21
22 model_root = DecisionTree(max_depth=7)
23 model_root.fit(X_train_root, y_train)
24
25 y_pred_root = model_root.predict(X_test_root)
26 accuracy_root = np.mean(y_pred_root == y_test.values.flatten())
27 print(f"Accuracy before prepruning: {accuracy_root:.5f}")
28 -----
29 Accuracy before prepruning: 0.88333

```

---

```

1
2 import itertools
3
4 # Define the parameter grid to search
5 param_grid = {
6     'max_depth': [5, 7, 9],
7     'min_samples_split': [3, 4, 5],
8     'min_samples_leaf': [2, 3, 4],
9     'min_info_gain': [0.01, 0.04, 0.08],
10    'max_features': [4, 6, 8]
11 }
12
13 # Create the decision tree model
14 base_decision_tree = PrePrunedDecisionTree()
15
16 best_accuracy = 0
17 best_params = None
18
19 # Perform grid search
20 for params in itertools.product(*param_grid.values()):
21     param_dict = dict(zip(param_grid.keys(), params))
22
23     # Create the decision tree model with current parameters
24     decision_tree = PrePrunedDecisionTree(**param_dict)
25
26     # Fit the model on the training data
27     decision_tree.fit(X_train_root, y_train)
28
29     # Make predictions on the test set
30     y_pred = decision_tree.predict(X_test_root)
31
32     # Evaluate accuracy
33     #accuracy = accuracy_score(y_test, y_pred)
34     accuracy = np.mean(y_pred == y_test)
35
36
37     # Update best parameters if the current model is better
38     if accuracy > best_accuracy:
39         best_accuracy = accuracy
40         best_params = param_dict
41
42 # Print the best hyperparameters and accuracy
43 print("Best Hyperparameters:", best_params)
44 print("Test Accuracy after pruning:", best_accuracy)
45 -----
46 Best Hyperparameters: {'max_depth': 7, 'min_samples_split': 3, '↵

```

```
min_samples_leaf': 2, 'min_info_gain': 0.01, 'max_features': 8}
```

47 Test Accuracy after pruning: 0.95

---

## Section 5: 解釋

Figure 4 為最終的決策樹結構，依此圖來看，follows、followers、posts為用來做決策的主要特徵，如果要再做優化可能仍須再把這三個的數值尺度降低，或者是max\_features要設小一點。

Listing 9: output structure

---

```
1  #from sklearn.tree import export_graphviz
2  import graphviz
3
4  # Helper function to convert ID3 tree to a format compatible with ↵
   export_graphviz
5  def tree_to_graphviz(tree, feature_names, class_names, dpi=300, size=None)↵
   :
6      def recurse(node, graph):
7          # Combine nodes for decision split
8          if 'left' in node.children:
9              #decision_label = f"If feature {node.feature_index} < {node.↵
               threshold}"
10             #decision_label = f"If feature {node.feature_index} < {node.↵
               threshold}\nClass: {node.predicted_class}\nEntropy = {node↵
               .entropy}\nSamples = {node.num_samples}\nClass ↵
               Distribution: {node.num_samples_per_class}"
11             decision_label = f"If {feature_names[node.feature_index]} < {↵
               node.threshold}\nClass: {node.predicted_class}\nEntropy = ↵
               {node.entropy}\nSamples = {node.num_samples}\nClass ↵
               Distribution: {node.num_samples_per_class}"
12
13             graph.node(str(id(node)), label=decision_label, shape='box', ↵
               fillcolor='#78bceb')
14             left_child = node.children['left']
15             right_child = node.children['right']
16
17             recurse(left_child, graph)
18             recurse(right_child, graph)
19
20             graph.edge(str(id(node)), str(id(left_child)), label='True')
21             graph.edge(str(id(node)), str(id(right_child)), label='False')
22         else:
23             # Leaf node
```

```

24         leaf_label = f"Class: {node.predicted_class}\nEntropy = {node.↵
           entropy}\nSamples = {node.num_samples}\nClass Distribution↵
           : {node.num_samples_per_class}"
25         graph.node(str(id(node)), label=leaf_label, shape='box', ↵
           fillcolor='#78bceb')
26
27     graph = graphviz.Digraph(format='png')
28     root_label = f"Root\nEntropy = {tree.entropy}\nSamples = {tree.↵
           num_samples}\nClass Distribution: {tree.num_samples_per_class}"
29     graph.node(str(id(tree)), label=root_label, shape='box', fillcolor='↵
           #78bceb')
30     recurse(tree, graph)
31
32     # Set DPI and size attributes
33     if dpi is not None:
34         graph.attr(dpi=str(dpi))
35     if size is not None:
36         graph.attr(size=size)
37
38     return graph
39
40 -----
41 feature_names = list(X_train.columns)
42 class_names = list(map(str, set(y_train)))
43 # Set the desired DPI
44 dpi = 300
45 # Set the desired size in inches (width, height)
46 size = None
47 graph = tree_to_graphviz(prepruned_model_root.tree_, feature_names, ↵
           class_names, dpi=dpi, size=size)
48 graph.render("prepruned_model_root")

```

---

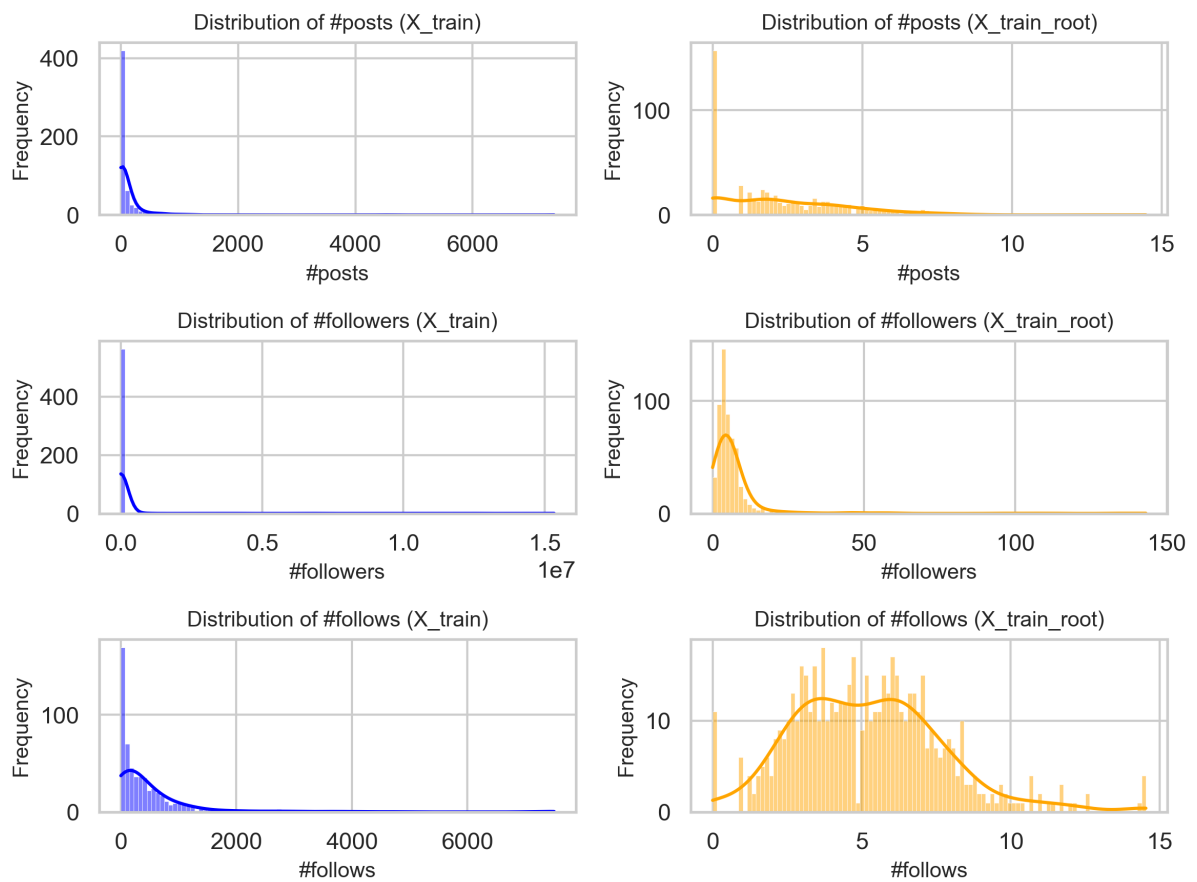


Figure 2: Distribution comparison(post, followers, follows are transformed to the power of 0.3)

	profile pic	nums/length username	fullname words	nums/length fullname	name==username
count	576.0	576.0	576.0	576.0	576.0
mean	0.7014	0.1638	1.4601	0.0361	0.0347
std	0.458	0.2141	1.0526	0.1251	0.1832
min	0.0	0.0	0.0	0.0	0.0
25%	0.0	0.0	1.0	0.0	0.0
50%	1.0	0.0	1.0	0.0	0.0
75%	1.0	0.31	2.0	0.0	0.0
max	1.0	0.92	12.0	1.0	1.0

	description length	external URL	private	#posts	#followers	#follows
count	576.0	576.0	576.0	576.0	576.0	576.0
mean	22.6233	0.1163	0.3819	2.3739	7.1458	5.224
std	37.703	0.3209	0.4863	2.1771	12.6298	2.5185
min	0.0	0.0	0.0	0.0	0.0	0.0
25%	0.0	0.0	0.0	0.0	3.0014	3.372
50%	0.0	0.0	0.0	1.9332	4.5005	5.1078
75%	34.0	0.0	1.0	3.7441	7.1858	6.7787
max	150.0	1.0	1.0	14.4736	143.1315	14.5385

Figure 3: Summary table of features (post, followers, follows are transformed to the power of 0.3)

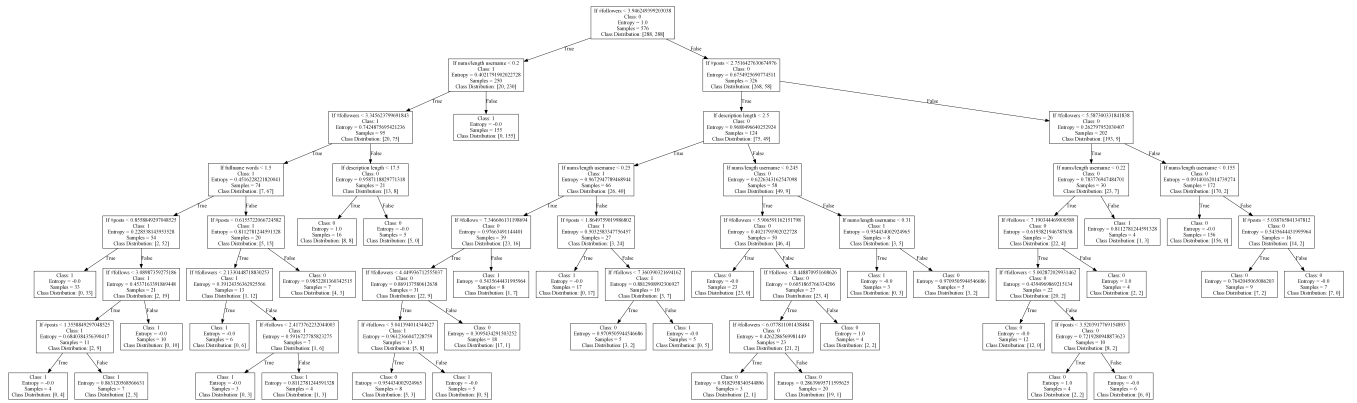


Figure 4: Tree Structure