



---

VASILIS ODYSSEOS

# DECISION TREES

ID3 Algorithm and Post-Pruning

---

# Datasets Used

- Hepatitis dataset
  - 155 instances
  - 19 attributes
  - Classification of patient's survival
- Horse Colic Dataset
  - 368 instances
  - 27 attributes
  - Classification of whether lesion was surgical or not
- Mushroom Dataset
  - 8124 instances
  - 22 attributes
  - Classification of whether mushroom is poisonous or edible

+others

```
def open_file(filename):
    df = pd.read_csv(filename)
    cont_attr = get_continuous(df)
    continuous(df, cont_attr, 4)
    df = df.sample(frac=1).reset_index(drop=True)
    training = df.sample(frac=0.8)
    rest = df.drop(training.index).reset_index(drop=True)
    training = training.reset_index(drop=True)
    testing = rest.sample(frac=0.5)
    validation = rest.drop(testing.index)
    testing = testing.reset_index(drop=True)
    validation = validation.reset_index(drop=True)
```

---

# Importing data

---

- Import using Pandas dataframe
- Split continuous data into 4 bins after checking if numerical attributes have more than 4 unique values
- Randomly divide the data into training, testing, validation samples

```

def decision_tree(df, tree=None, used=None):
    # Recursively builds a decision tree from a given dataframe until entropy(df)
    # is 0 or all attributes used

    truth_col = df.keys()[-1]

    if used is None:
        used = []
    # get best attribute based on information gain
    attr = find_best(df[df.columns.difference(used, sort=False)])
    used.append(attr)
    attributes = list(df)

    values = df[attr].unique()

    if tree is None:
        tree = {}
        tree[attr] = {}

    # loop through unique values of attr and create a branch for each value.
    # If entropy of sorting on value is 0, the branch is a lead
    for value in values:
        sub = df[df[attr] == value].reset_index(drop=True)

        if entropy(sub) == 0 or (len(attributes) - 1) == len(np.unique(used)):
            unique, counts = np.unique(sub[truth_col], return_counts=True)
            tree[attr][value] = unique[np.argmax(counts)]
        else:
            tree[attr][value] = decision_tree(sub, used=used)
    return tree

```

---

# Decision Tree algorithm

---

- Use attribute with best information gain as the root of the tree
- For each possible value of the attribute, create a new branch by recursively calling the decision\_tree function on the subset of the data with that value
- If the entropy of the subset is 0 or all attributes have been used, then the branch is a leaf and the tree stops growing

```

def pruning(df, tree):
    # prunes a decision tree by finding all nodes whose children are all leaves,
    # then checking to see if simplifying the node to a leaf gives a better accuracy,
    # and removes the node that gives the greatest gain in accuracy
    # essentially tries to find the shortest tree without loss of accuracy

    nodes = []
    # test on current tree for benchmark
    default_accuracy = test_all(df, tree)
    max_accuracy = 0
    # save current tree as best possible tree (that we know of)
    best_tree = tree
    # get all nodes with all leaf children
    get_nodes(tree, nodes)

    # test if simplifying each node produces a better or equal accuracy than the benchmark
    for i in range(len(nodes)):
        tcopy = copy.deepcopy(tree)
        tcopy = simplify_node(tcopy, nodes[i])
        accuracy = test_all(df, tcopy)
        if accuracy >= max_accuracy:
            best_tree = copy.deepcopy(tcopy)
            max_accuracy = accuracy
    # recursively prune the tree if the accuracy of the pruned tree is greater or equal
    # to the benchmark, otherwise return the current tree
    if (max_accuracy >= default_accuracy):
        return pruning(df, best_tree)
    else:
        return tree

```

---

# Post-pruning algorithm

- Get all nodes that have all-leaf-children
- For each node, test if removing it yields a greater or equal accuracy on the validation data
- Remove the node with the greatest increase in accuracy
- Recursively prune the tree until accuracy drops

```

def test(hypothesis, tree):
    # test a hypothesis by traversing the tree and comparing the leaf
    # of the tree to the label of the hypothesis

    if type(tree) is dict:
        try:
            return test(hypothesis, tree[list(tree)[0]][hypothesis[list(tree)[0]]])
        except:
            print("Label found in test set that does not exist in the decision tree.")
            return 0
    else:
        if(hypothesis[-1] == tree):
            return 1
        else:
            return 0

def test_all(df, tree):
    # tests all hypotheses in the set against the tree

    count = 0
    for i in range(len(df[df.keys()[-1]])):
        count += test(df.loc[i, :], tree)
    return count / len(df[df.keys()[-1]])

```

---

# Testing the trees

---

- For each instance in the given dataframe
  - Traverse the tree using the values in the instance
  - Return 1 if the classification given by the tree matches that of the instance
  - Return 0 if the classification given by the tree doesn't match, or if a certain attribute value is missing from the tree

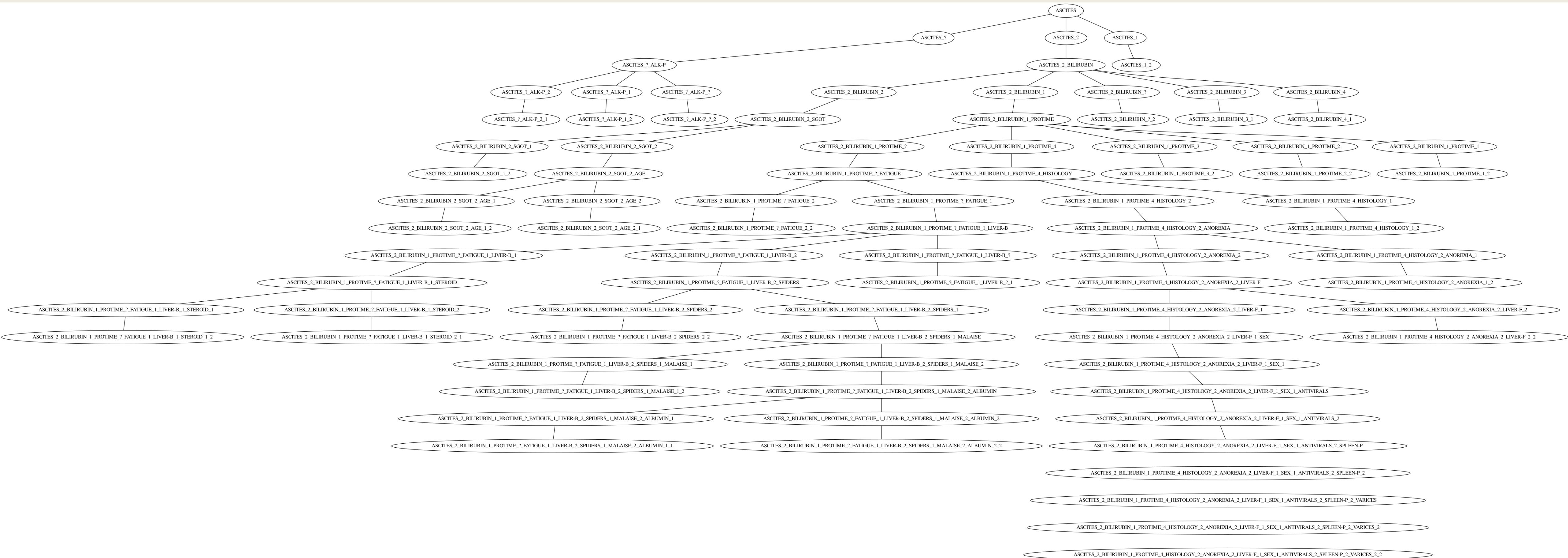
---

# Results

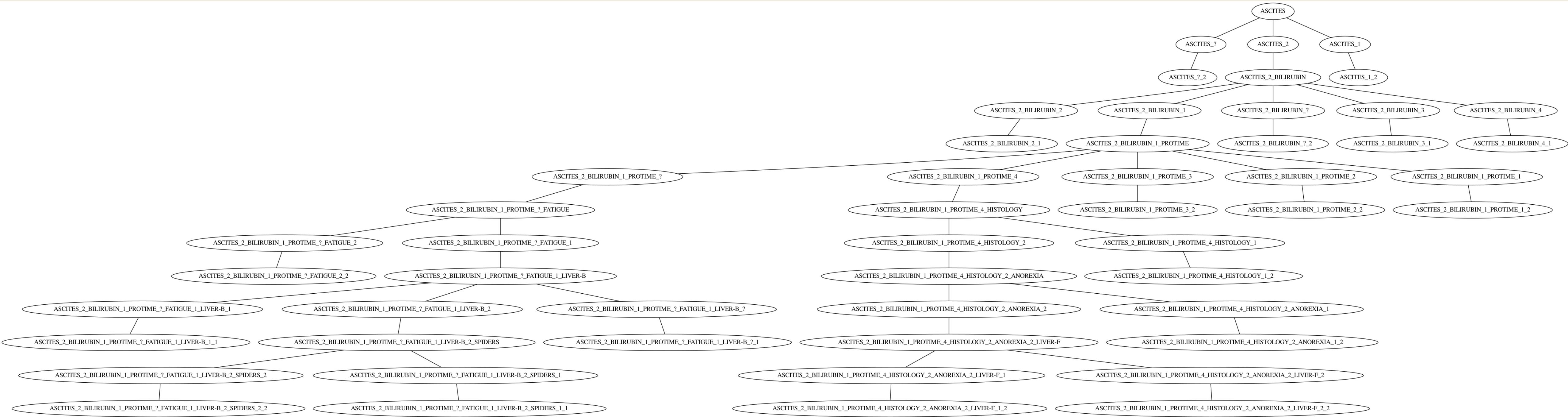
# Hepatitis Dataset

- Testing on training data:
  - Original Accuracy: 0.8467741935483871
  - Pruned Accuracy: 0.75
- Testing on test data:
  - Original Accuracy: 0.8125
  - Pruned Accuracy: 0.8125
- Testing on validation data:
  - Original Accuracy: 0.666666666666
  - Pruned Accuracy: 0.8

# Original Tree



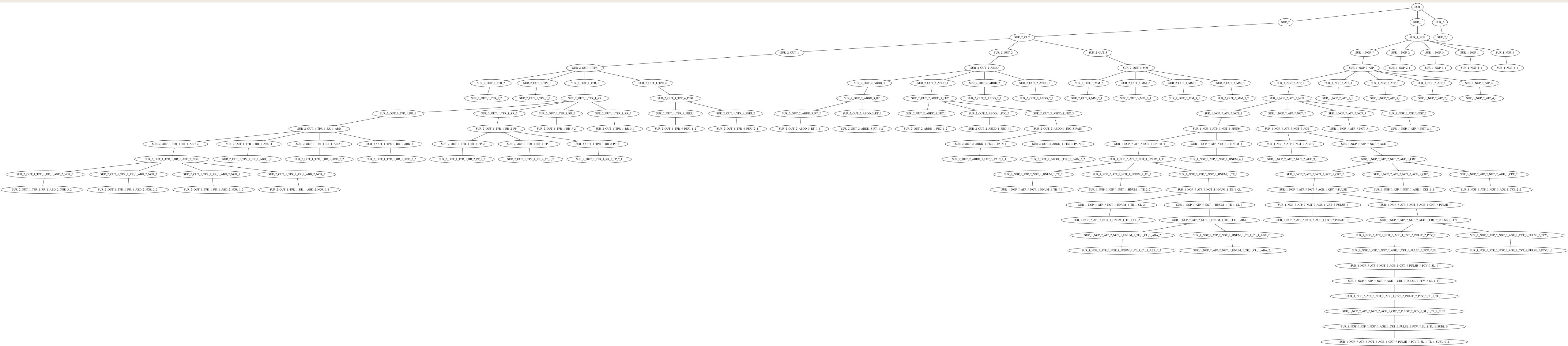
# Pruned Tree



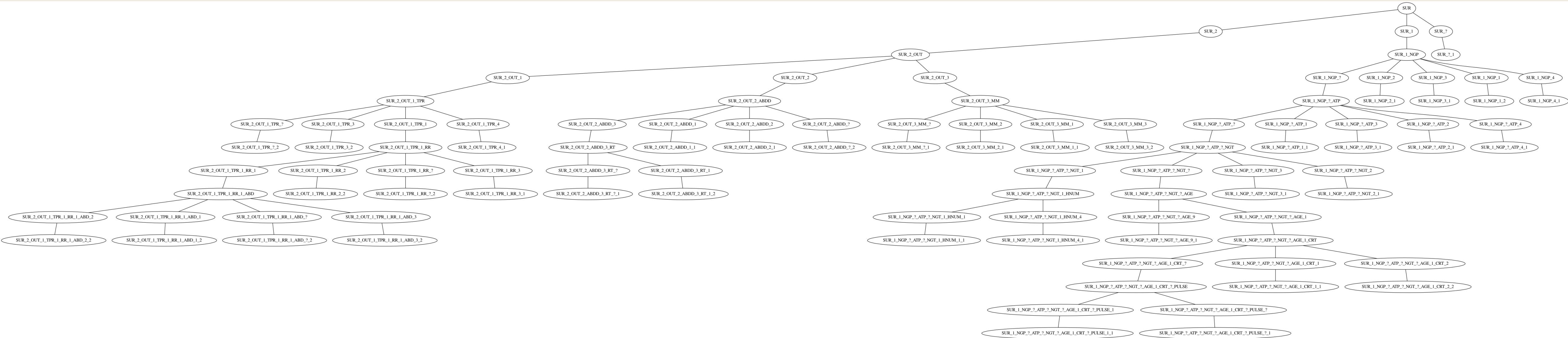
# Horse Colic Dataset

- Testing on training data:
  - Original Accuracy:  
0.8916666666666667
  - Pruned Accuracy:  
0.858333333333333
- Testing on test data:
  - Original Accuracy:  
0.8
  - Pruned Accuracy:  
0.833333333333334
- Testing on validation data:
  - Original Accuracy:  
0.7666666666666667
  - Pruned Accuracy:  
0.833333333333334

# Original Tree



# Pruned Tree

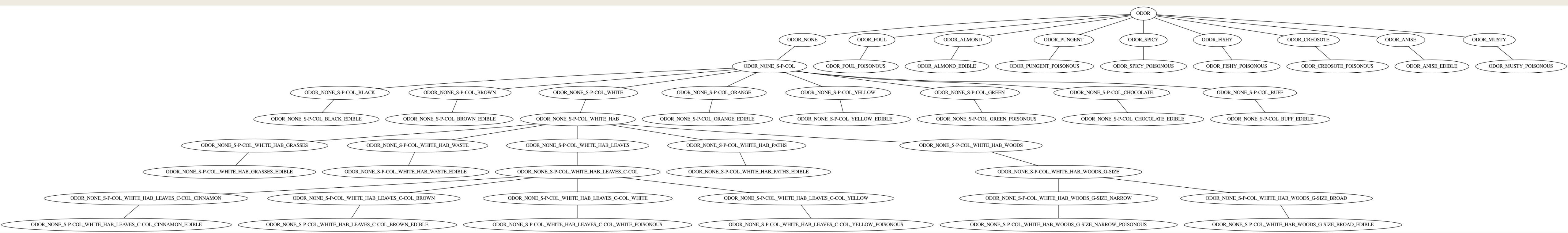


---

# Mushroom Dataset

- Testing on training data:
  - Original Accuracy: 1.0
  - Pruned Accuracy: 1.0
- Testing on test data:
  - Original Accuracy: 1.0
  - Pruned Accuracy: 1.0
- Testing on validation data:
  - Original Accuracy: 1.0
  - Pruned Accuracy: 1.0

# Original Tree



# Pruned Tree

