

hw4 > src_and_data > hw4.ipynb

hw4.ipynb – hw4.ipynb

Project hw4.ipynb

Managed: http://localhost:8889 Python 3 (ipykernel) Trusted 4 18

(a) Compute the counting-based heuristic, and order the features by it.

```
In 3 1 from collections import OrderedDict
2
3 # TODO counting_heuristic()
4 my_dict = OrderedDict()
5
6 for feature_num in range(feature_len):
7
8     total_correct = counting_heuristic(x_inputs=x_set, y_outputs=y_set, feature_index=feature_num, classes=classes)
9     my_dict[feature_names[feature_num]] = total_correct
10
11 # TODO sort the feature names by their correct counts
12 # e.g., sort_correct = [best, second_best, ..., worst]
13 # e.g., sort_names = ["A", "B"] or ["B", "A"]
14
15 high_to_low = OrderedDict(sorted(my_dict.items(), key=lambda x: x[1]))
16 sort_correct = list(high_to_low.values()) #TODO: FIX ME
17 sort_names = list(high_to_low.keys()) #TODO: FIX ME
18
19 # print(f'{sort_correct = }')
20 # print(f'{sort_names = }')
21 # Print the sorted features along with their correct predictions count in the smaller dataset
22 longest = max(len(name) for name in sort_names)
23 for name, correct in zip(sort_names, sort_correct):
24     print("%s: %d/%d" % (longest, name, correct, len(x_set)))
25
26 A: 6/8
27 B: 6/8
```

(b) Compute the information-theoretic heuristic, and order the features by it.

```
In 4 1 # TODO information_remainder()
2 gains_dict = OrderedDict()
3
4 for feature_num in range(feature_len):
5     gain = information_remainder(x_inputs=x_set, y_outputs=y_set, feature_index=feature_num, classes=classes)
6     gains_dict[feature_names[feature_num]] = gain
7
8 # TODO sort the feature names by their gains
9 gains_sorted = OrderedDict(sorted(gains_dict.items(), key=lambda x: x[1], reverse=True))
10
11 sort_gains = list(gains_sorted.values()) #TODO: FIX ME
12 sort_names_by_gains = list(gains_sorted.keys()) #TODO: FIX ME
13
14 # print(f'{sort_gains = }')
```

File Project Bookmarks Structure Notebooks

File TODO Problems Terminal Python Packages Python Console Services Jupyter

Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)

2:1 LF UTF-8 4 spaces mL0135 886 of 2048M

hw4 > src_and_data > hw4.ipynb

hw4.ipynb – hw4.ipynb

Project hw4.ipynb

Managed: http://localhost:8889 Python 3 (ipykernel) Trusted 4 18

```

13 # print(f'{sort_gains = }')
14 # print(f'{sort_names_by_gains = }')
15
16 longest = max(len(name) for name in sort_names_by_gains)
17 for name, gain in zip(sort_names_by_gains, sort_gains):
18     print(f'{name: {longest}f}: {gain:.3f}')

```

A: 0.311
B: 0.189

(c) Discussion of results.

Using each of these heuristics, which feature do you choose to start the tree? Explain each choice of feature. How does heuristics affect the tree?

Since A and B are equivalent in num correct over the total in the counting-based heuristic and A has a better score (0.311) vs B (0.189) in the information-theoretic heuristic, the tie is broken and A is a better feature to start the tree. Starting the tree is done best by most important features at the top and less prominent features below. If you were to split on less relevant, niche, features at the top, then the decision tree would be more convoluted to view and the tree structure would not be used optimally. Much of its use is for visualization to even non computer science people.

If the counting heuristic yielded slightly better results for A and the information gain heuristic gave slightly better results for B this would be a tough decision to make on which feature to start the tree with. Information gain is a stronger way to classify features since the uncertainty level is measured and visible. Choosing one of the two features in each class iteration isn't as good at model building. Starting with feature B would be narrowing the tree too soon in optimal cases.

2 Compute both heuristics for simplified abalone data.

In 5

```

1 # ignore but helpful for inspecting data in pycharm
2 """
3 # https://bids.github.io/2015-06-04-berkeley/intermediate-python/03-sklearn-abalone.html
4 # TODO:fix the empty lists below
5 # full-feature abalone data
6 x_train_p = pd.read_csv('data_abalone/x_train.csv')
7 x_test_p = pd.read_csv('data_abalone/x_test.csv')
8 y_train_p = pd.read_csv('data_abalone/y_train.csv')
9 y_test_p = pd.read_csv('data_abalone/y_test.csv')
10 # np.loadtxt('data_abalone/x_train.csv')
11
12 # TODO:fix the empty lists below
13 # simplified version of the data (Restricted-feature)
14 simple_x_train_p = pd.read_csv('data_abalone/small_binary_x_train.csv')
15 simple_x_test_p = pd.read_csv('data_abalone/small_binary_x_test.csv')
16 simple_y_train_p = pd.read_csv('data_abalone/3class_y_train.csv')
17 simple_y_test_p = pd.read_csv('data_abalone/3class_y_test.csv')
18 """

```

Out 5

```

"\n# https://bids.github.io/2015-06-04-berkeley/intermediate-python/03-sklearn-abalone.html\n# TODO:fix the empty lists below\n# full-feature abalone data\nx_train_p = pd.read_csv('data_abalone/x_train.csv')\nx_test_p = pd.read_csv('data_abalone/x_test.csv')\ny_train_p = pd.read_csv('data_abalone/y_train.csv')\ny_test_p = pd.read_csv('data_abalone/y_test.csv')\n# np.loadtxt('data_abalone/x_train.csv')\n#\n# TODO:fix the empty lists below\n# simplified version of the data (Restricted-feature)\nsimple_x_train_p = pd.read_csv('data_abalone/small_binary_x_train.csv')\nsimple_x_test_p = pd.read_csv('data_abalone/small_binary_x_test.csv')\nsimple_y_train_p = pd.read_csv('data_abalone/3class_y_train.csv')\nsimple_y_test_p = pd.read_csv('data_abalone/3class_y_test.csv')\n"

```

File Structure Bookmarks Jupyter Server

File Edit TODO Problems Terminal Python Packages Python Console Services Jupyter

Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)

2:1 LF UTF-8 4 spaces mL0135 998 of 2048M

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** hw04.ipynb – hw04.ipynb
- Toolbar:** Includes icons for back, forward, search, and file operations.
- Project Bar:** Shows 'hw04' and 'src_and_data'.
- Code Editor:** The main area contains Python code for reading and processing the Abalone dataset. The code includes imports from pandas and numpy, and various steps for loading different versions of the dataset (full-feature, restricted-feature, simplified). It also prints feature names and class labels.
- Output:** Shows the output of the code execution, including the printed feature names and class labels.
- Bottom Navigation:** Includes links for Git, TODO, Problems, Terminal, Python Packages, Python Console, Services, and Jupyter.
- Status Bar:** Shows the file path as 'Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)', the file size as '376 of 2048M', and other system information.

hw4 > src_and_data > hw4.ipynb

(a) Compute the counting-based heuristic, and order the features by it.

```
In 7 1 # TODO counting_heuristic()
2
3 aba_dict = OrderedDict()
4 # for aba_feature in range(len(full_feature_names)):
5 ✓ for aba_feature in range(len(simple_feature_names)):
6
7     # total_correct_aba = counting_heuristic(x_inputs=x_train, y_outputs=y_train, feature_index=aba_feature, classes=classes_abalone)
8     total_correct_aba = counting_heuristic(x_inputs=simple_x_train, y_outputs=simple_y_train, feature_index=aba_feature, classes=classes_abalone)
9
10    aba_dict[full_feature_names[aba_feature]] = total_correct_aba
11
12    # TODO sort the feature names by their correct counts
13    high_to_low_aba = OrderedDict(sorted(aba_dict.items(), key=lambda x: x[1], reverse=True))
14    sort_correct_abalone = list(high_to_low_aba.values()) #TODO: FIX ME
15    sort_names_abalone = list(high_to_low_aba.keys()) #TODO: FIX ME
16
17    # Print the sorted features along with their correct predictions count in the smaller dataset
18    longest = max(len(name) for name in sort_names_abalone)
19    for name, correct in zip(sort_names_abalone, sort_correct_abalone):
20        print("%*s: %d/%d" % (longest, name, correct, len(simple_x_train)))
21
22    height_mm: 2316/3176
23    diam_mm: 2266/3176
24    length_mm: 2230/3176
25    is_male: 1864/3176
```

(b) Compute the information-theoretic heuristic, and order the features by it.

```
In 8 1 # TODO information_remainder()
2 info_gains_dict = OrderedDict()
3 ✓ for aba_feature in range(len(simple_feature_names)):
4     gain = information_remainder(x_inputs=simple_x_train, y_outputs=simple_y_train, feature_index=aba_feature, classes=classes_abalone)
5     info_gains_dict[simple_feature_names[aba_feature]] = gain
6
7    # TODO sort the feature names by their gains
8    info_gains_sorted = OrderedDict(sorted(info_gains_dict.items(), key=lambda x: x[1], reverse=True))
9
10   sort_gains_abalone = list(info_gains_sorted.values()) #TODO: FIX ME
11   sort_names_by_gains_abalone = list(info_gains_sorted.keys()) #TODO: FIX ME
12
13   # longest = max(len(name) for name in sort_names_by_gains_abalone)
14   for name, gain in zip(sort_names_by_gains_abalone, sort_gains_abalone):
15       print("%*s: %.3f" % (longest, name, gain))
```

File Requests Project hw4.ipynb Current File Git Help

Managed: http://localhost:8889 Python 3 (ipykernel) Trusted 4 18 18 Scala

Structure Bookmarks Notebooks

File Edit View Insert Cell Kernel Help

Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)

2:1 LF UTF-8 4 spaces mLO135 main 422 of 2048M

hw4 > src_and_data > hw4.ipynb

hw4.ipynb – hw4.ipynb

Project hw4.ipynb ×

File Requests

Markdown

Managed: http://localhost:8889

Python 3 (ipykernel)

Trusted

4 18

(b) Compute the information-theoretic heuristic, and order the features by it.

```
In 8 1 # TODO: information_remainder()
2 info_gains_dict = OrderedDict()
3 for aba_feature in range(len(simple_feature_names)):
4     gain = information_remainder(x_inputs=simple_x_train, y_outputs=simple_y_train, feature_index=aba_feature, classes=classes_abalone)
5     info_gains_dict[simple_feature_names[aba_feature]] = gain
6
7 # TODO: sort the feature names by their gains
8 info_gains_sorted = OrderedDict(sorted(info_gains_dict.items(), key=lambda x: x[1], reverse=True))
9
10 sort_gains_abalone = list(info_gains_sorted.values()) #TODO: FIX ME
11 sort_names_by_gains_abalone = list(info_gains_sorted.keys()) #TODO: FIX ME
12
13 # longest = max(len(name) for name in sort_names_by_gains_abalone)
14 for name, gain in zip(sort_names_by_gains_abalone, sort_gains_abalone):
15     print("%*s: %.3f" % (longest, name, gain))

    height_mm: 0.173
    diam_mm: 0.150
    length_mm: 0.135
    is_male: 0.025
```

3) Generate decision trees (criterion='entropy', random_state=42) for full- and simple-feature data

(a) Train and eval on entire train and test sets. Print accuracy values and generate tree images.

Render the tree diagram, naming it "full." A text file and PDF should be created and saved (i.e., `full` and `full.pdf`) – include both in submission.

```
In 9 1 from sklearn.tree import export_graphviz
2 from sklearn.tree import DecisionTreeClassifier
3 import graphviz
4
5 full_model = DecisionTreeClassifier(criterion='entropy', random_state=42)
6 full_model.fit(x=x_train, y=y_train)
7
8 train_accuracy = full_model.score(X=x_train, y=y_train)
9 test_accuracy = full_model.score(X=x_test, y=y_test)
10
11 # TODO calculate accuracies
12 print(f"Accuracy (train): {train_accuracy:.3f}")
13 print(f"Accuracy (test): {test_accuracy:.3f}")
```

File Structure Bookmarks Notebooks

File TODO Problems Terminal Python Packages Python Console Services Jupyter

Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)

2:1 LF UTF-8 4 spaces mL0135 main 477 of 2048M

hw4 > src_and_data > hw4.ipynb

hw4.ipynb – hw4.ipynb

Project hw4.ipynb

Managed: http://localhost:8889 Python 3 (ipykernel) Trusted 4 18

3) Generate decision trees (criterion='entropy', random_state=42) for full- and simple-feature data

(a) Train and eval on entire train and test sets. Print accuracy values and generate tree images.

Render the tree diagram, naming it "full." A text file and PDF should be created and saved (i.e., `full` and `full.pdf`) – include both in submission.

```
In 9 1 from sklearn.tree import export_graphviz
2 from sklearn.tree import DecisionTreeClassifier
3 import graphviz
4
5 full_model = DecisionTreeClassifier(criterion='entropy', random_state=42)
6 full_model.fit(X=x_train, y=y_train)
7
8 train_accuracy = full_model.score(X=x_train, y=y_train)
9 test_accuracy = full_model.score(X=x_test, y=y_test)
10
11 # TODO calculate accuracies
12 print(f"Accuracy (train): {train_accuracy:.3f}")
13 print(f"Accuracy (test): {test_accuracy:.3f}")
14
15 # TODO generate tree image
16 full_tree = export_graphviz(decision_tree=full_model, feature_names=full_feature_names)
17 generated_graph_full = graphviz.Source(full_tree)
18 generated_graph_full.render('full')

Accuracy (train): 1.000
Accuracy (test): 0.202

Out 9 'full.pdf'
```

(b) Restricted-feature (aka simple) data.

Train and eval on simple train and test sets. Same as above, accept this time use the `simple` set. Render the tree diagram, naming it "simple." A text file and PDF should be created and saved (i.e., `simple` and `simple.pdf`) – include both in submission.

```
In 10 1 from sklearn.tree import export_graphviz
2 from sklearn.tree import DecisionTreeClassifier
3 import graphviz
4
5 simple_model = DecisionTreeClassifier(criterion='entropy', random_state=42)
6 simple_model.fit(X=simple_x_train, y=simple_y_train)
7
8 # TODO calculate out accuracies
9 simple_train_accuracy = simple_model.score(X=simple_x_train, y=simple_y_train) # Fix me
```

File Structure Bookmarks Help Notebooks

File Edit View Insert Cell Kernel Help

File Git Current File Terminal Services Jupyter

Jupyter Server started at http://localhost:8889 // Open in Browser (2 minutes ago)

2:1 LF UTF-8 4 spaces mL0135 566 of 2048M

hw4 > src_and_data > hw4.ipynb

hw4.ipynb – hw4.ipynb

Project Pull Requests

hw4.ipynb Current File Git

(b) Restricted-feature (aka simple) data.

Train and eval on simple train and test sets. Same as above, accept this time use the `simple` set. Render the tree diagram, naming it "simple." A text file and PDF should be created and saved (i.e., `simple` and `simple.pdf`) – include both in submission.

```
In 10 from sklearn.tree import export_graphviz
from sklearn.tree import DecisionTreeClassifier
import graphviz

simple_model = DecisionTreeClassifier(criterion='entropy', random_state=42)
simple_model.fit(X=simple_x_train, y=simple_y_train)

# TODO calculate out accuracies
simple_train_accuracy = simple_model.score(X=simple_x_train, y=simple_y_train) # Fix me
simple_test_accuracy = simple_model.score(X=simple_x_test, y=simple_y_test) # Fix me

print(f"Accuracy (train): {simple_train_accuracy:.3f}")
print(f"Accuracy (test): {simple_test_accuracy:.3f}")

# TODO generate tree image
simple_tree = export_graphviz(decision_tree=simple_model, feature_names=simple_feature_names, class_names=class_names)
generated_graph_simple = graphviz.Source(simple_tree)
generated_graph_simple.render('simple')

Accuracy (train): 0.733
Accuracy (test): 0.722

Out 10 'simple.pdf'
```

(c) Discuss the results seen for the two trees

Discuss the results you have just seen. What do the various accuracy-score values tell you? How do the two trees that are produced differ? Looking at the outputs (leaves) of the simplified-data tree, what errors does that tree make?

The entire sets gave an excellent 1.000 accuracy rating on training data, but on test it was 0.202. This suggests overfitting when the train accuracy is very good but test is poor. The model gets so good at fitting the training data that adapting to unseen data is not successful. The model could be overly complex and it would be better potentially to have a lower-degree polynomial equation for lesser complexity and avoiding overfitting. There are 8 input features for the full feature data but only 4 for the simple data.

The simple sets had very comparable accuracy for training and testing at 0.733 and 0.722, respectively. The almost identical accuracies show consistency. The model is predictable. Although the train accuracy is not as high as the full feature data, it is better than a nearly 80% accuracy gap. Low 70s percent accuracy leaves something to be desired, as 90+ percent is ideal; however. The complexity of this model was appropriate given the similar training and testing fits. A reason why the accuracy wasn't as good as it perhaps could have been is the unused "Large" class. If we just had "Small" and "Medium" classes like the simple.pdf shows, then I expect accuracy would be better. Complexity is only good if it improves model performance.

Poor accuracy on simple and overfitting on full suggests that we might want to use a different classifier and training method instead of Decision Trees. Alternatively we could try using different criterion types to see if that would boost accuracy.

Seeing entropy values greater than 1.0 suggest high levels of disorder and uncertainty. Those leaves on the simple and full graphs alike mean low confidence in the classified results. On the other hand, the many entropy values close to 0.0 or 0.0 in the full graph show high accuracy, which is reflected from the training results. A strategy to build a better full test accuracy would be to reduce the number of full feature train columns to target high 80's accuracy or low 90's accuracy on train and test.