

Hyperparameter Optimization

Introduction:

For this exercise I experimented with both the GridSearchCV and BayesSearchCV model selectors on a couple of ML algorithms. I decided on k-nearest neighbors and random forest and compared the accuracy of the models from the default results, GridSearchCV optimized and BayesSearchCV optimized.

Dataset Description:

I used the wine quality dataset^[1] for this and picked the white wine dataset because it has more entries than the red dataset. This data includes 11 features and a target value called “quality” that is between 0 and 10.

[1]: <https://archive-beta.ics.uci.edu/dataset/186/wine+quality>

Experimental Setup:

For the data setup, I went with the white wine as mentioned above, then split the test data to be 20% of the overall data and for the models I went with k nearest neighbors and random forest classifiers, knn was picked because I wanted to play around with the hyperparameters on this model more and random forest was picked because it was the most accurate in the previous exercise.

Accuracy this time is recorded using scikit-learn’s accuracy_score function that exists for classification models.

For this experiment I used scikit-learn’s GridSearchCV and BayesSearchCV model selectors which I gave a model, a list of hyperparameters, a scoring system and cross validation fold count to determine out of the list of parameters I give it, which ones are the best. BayesSearchCV I experimented with a little bit and reduced it’s number of iterations because it was testing with hyperparameters it had already tested on.

I later noticed that the Bayesian optimization on random forest was running faster so I decided to time all of the models and output that too.

Also both algorithms are set to doing 5 cross validation folds.

For the hyperparameters I optimized:

Random Forest:

Max_depth:

The default value for this is None, and I tested this with values 25, 50 and None, this parameter sets the max depth of the tree.

N_estimators:

The default value for this is 100, and I tested this with values 100, 125 and 150, this parameter sets how many random trees the algorithm will generate.

Min_samples_split:

The default value for this is 2, and I tested this with values 2, 4, 6, this parameter determines the minimum number of samples before a node can split.

Min_samples_leaf:

The default value for this is 1, and I tested this with 1, 3, 5, this parameter determines the minimum number of samples in a leaf.

K-Nearest-Neighbors:

n_neighbors:

The default value for this is 5, and I tested this with values 3, 4, 5, 6, 10, this parameter sets the number of neighbors.

leaf_size:

The default value for this is 30, and I tested this with values 10, 30, 50, 100, 200, this parameter sets the leaf size passed into the balltree and kdtree algorithms.

algorithm:

The default value for this is 'auto' and I tested this with 'ball_tree', 'kd_tree', 'brute', 'auto', this parameter changes which algorithm is being used.

p:

The default value for this is 2, and I tested this with both 1 and 2, it changes which distance formula is used.

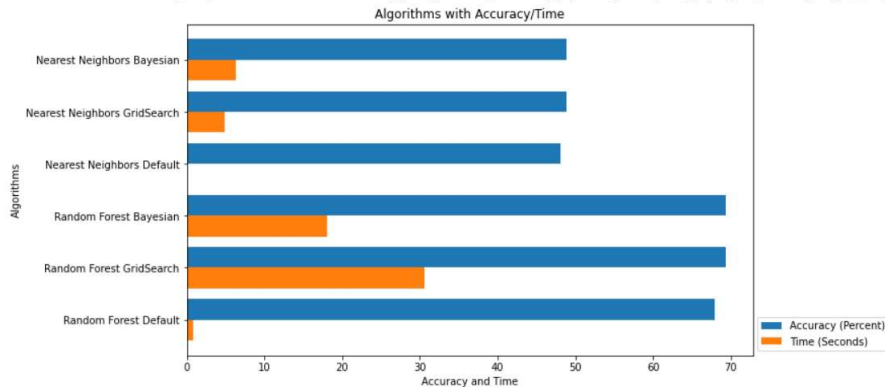
I then would graph the results and make a list of the algorithms sorted by accuracy^[2] then by time.

We also output the best hyperparameters found by each model selector.

[2]: code for that inspired by <https://stackoverflow.com/questions/48053979/print-2-lists-side-by-side> user SCB

Results:

```
Best Parameters found by gridsearch random forest: {'max_depth': 25, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 125}
Best Parameters found by bayesian random forest: OrderedDict([('max_depth', 25), ('min_samples_leaf', 1), ('min_samples_split', 2), ('n_estimators', 125)])
Best Parameters found by gridsearch knn: {'algorithm': 'ball_tree', 'leaf_size': 10, 'n_neighbors': 5, 'p': 1}
Best Parameters found by bayesian knn: OrderedDict([('algorithm', 'auto'), ('leaf_size', 10), ('n_neighbors', 5), ('p', 1)])
```



List of algorithms sorted best to worst:

```
Random Forest Bayesian: 69.38776% accuracy, 18.09990 seconds
Random Forest GridSearch: 69.38776% accuracy, 30.56478 seconds
Random Forest Default: 67.95918% accuracy, 0.89277 seconds
Nearest Neighbors GridSearch: 48.87755% accuracy, 4.90672 seconds
Nearest Neighbors Bayesian: 48.87755% accuracy, 6.33202 seconds
Nearest Neighbors Default: 48.06122% accuracy, 0.03653 seconds
Best Algorithm by accuracy and time is: Random Forest Bayesian
```

This image contains the output of my program, as we can see the gridsearch and bayesearch performed the same with what it was given and both outperformed default hyperparameters. Though one thing to note is that the bayesearch ran much faster than the gridsearch for random forest and about as fast for the nearest neighbors. From this, I can see with a more complex problem or larger dataset the Bayesian optimization may be better to use given how much faster it runs for a similar result.

Also the gridsearch and Bayesian optimization picked different algorithms, Bayesian picked auto and gridsearch picked 'ball_tree' but I am assuming that the 'auto' selected picked the same algorithm.

They only ended up with small improvements, likely due to how simple the problem was and how good the default values are for the problem.

Code:

In the repository under main.py.