

1 Part 1

Mini-exam question:

Question 1: Write a proper introductory chapter for your report:

- Explain the task and the objectives **0.1p**
- Describe, what kind of data were used? **0.1p** Where did it originate? **0.1p**
Give the correct reference. **0.1p** Explain shortly, how the data has been prepared in the original article. **0.2p** Explain, how you deployed the data **0.2p**
- Report shortly the process used in this task (preprocessing **0.2p**, feature derivation **0.2p**, methods **0.2p**, training **0.2p**, evaluation **0.2p**)
- Describe the results **0.2p**

Introduction

This is an exercise for Machine Learning and Pattern Recognition course. The purpose of this task was to learn how to train and evaluate machine learning (ML) models, which classify images into classes. The exercise had different tasks: to import the data, to derive essential features, to analyse the gathered data, to build and select best ML models using hyperparameter selection, and finally to evaluate the performance of the selected models.

The exercise was based on the work of Çınar and Koklu [1]. The task was to classify rice species. The data consisted of images (size: 250 x 250 pixels) including one rice grain. The rice images were taken with a camera placed in a closed box to ensure the same lighting environment for all samples. The rices were captured on a black background. The original images were color images, which then were converted to grayscale images. The data is openly available at <https://www.muratkoklu.com/datasets>. The original data includes images for five rice species, 75000 images in total. In this exercise, we used only a small random sample of the original data (100 images for each rice species).

Different features were derived from the rice images. First, a contour around the rice grain was detected for each image. Then, the detected contour was used to calculate color features using the pixels within the contour, such as mean, variance, skewness, and kurtosis of the YCbCr-color channel values. Morphological features such as area, perimeter, major and minor axis length, equivalent diameter, compactness and shape factors were determined using an ellipse fitted to the detected contour. The gathered data was explored using histograms, pair plots and Principal Component Analysis (PCA).

The feature data was standardized. Three classifiers were used for classification: Random Forest (RF), Support Vector Machine (SVM) and Multi-Layer Perceptron (MLP). The best models were selected using hyperparameter selection and 5-fold 3 times repeated Kfold cross validation (cv). The model performance was estimated using nested cv (stratified 10-fold cv in the outer loop and 3 times repeated 5-fold in the inner loop). The performance of the models was evaluated using accuracy.

The estimated mean accuracies were 0.994/0.998/0.968 for RF/SVM/MLP, respectively. Morphological and shape features seemed to be the most important features according to RF feature importance.

Reference

- [1] İ. Çınar and M. Koklu. Identification of rice varieties using machine learning algorithms. Journal of Agricultural Sciences, 28(2):307–325, 2022. doi: 10.15832/ankutbd.862482.

2 Part 2

3 Part 3

In [1]: ►

```
1 import numpy as np
2 import pandas as pd
3 from matplotlib import pyplot as plt
4 # import cv2 as cv
5 # import glob, os
6 import itertools
7 # from random import sample
8 # from random import seed
9 # from scipy.stats import skew
10 # from scipy.stats import kurtosis
11 # from scipy.stats import entropy
12 import scipy.stats
13 from sklearn.decomposition import PCA
14 from sklearn.model_selection import RepeatedKFold
15 from sklearn.model_selection import StratifiedKFold
16 from sklearn.model_selection import GridSearchCV
17 from sklearn.ensemble import RandomForestClassifier
18 from sklearn.svm import SVC
19 from sklearn.neural_network import MLPClassifier
20 from sklearn.metrics import confusion_matrix
21 from sklearn.metrics import accuracy_score
22 from sklearn.model_selection import cross_val_predict
23
24 from IPython.display import display, HTML
25 display(HTML("<style>.container { width:90% !important; }</style>"))
26
27 import time
```

3.1 Performance estimation

3.1.1 Random forest

Let's estimate the performance of each model using nested cross-validation. Again, we'll use 5-fold repeated cross validation with 3 repetitions and the same parameter ranges as in Part 2 for the inner loop. For the outer loop we'll use 10-fold Stratified Kfold cross-validation.

Study the following implementation of performance estimation for Random Forest model.

You do not need to run the cell below, you can import the results and start from there.

```

In [ ]: █ 1 df = pd.read_parquet('../training_data/rice_feature_data.parquet')
2
3 feats = ['mean_y', 'var_y', 'skew_y', 'kurt_y',
4           'mean_cr', 'var_cr', 'skew_cr', 'kurt_cr',
5           'mean_cb', 'var_cb', 'skew_cb', 'kurt_cb',
6           'major_axis_length', 'minor_axis_length',
7           'area', 'perimeter',
8           'equivalent_diameter', 'compactness',
9           'shape_factor1', 'shape_factor2']
10
11 # Z-Score using pandas
12 for feat in feats:
13     df['{}_Z'.format(feat)] = (
14         df[feat] - df[feat].mean()) / df[feat].std()
15
16 # Z-scored feature names
17 feats_Z = [feat + '_Z' for feat in feats]
18
19 y = df['class']
20 X = df[feats_Z]
21
22 kf_outer = StratifiedKFold(
23     n_splits=10,
24     random_state=10,
25     shuffle=True)
26 kf_inner = RepeatedKFold(
27     n_splits=5,
28     n_repeats=3,
29     random_state=5)
30
31 n_estimators = range(100, 350, 50)
32 max_features = ['sqrt', 'log2', None]
33 bootstrap = [True, False]
34
35 parameters = {
36     'n_estimators': n_estimators,
37     'max_features': max_features,
38     'bootstrap': bootstrap}
39
40 rf = RandomForestClassifier(random_state=20)
41
42 results = []
43
44 start_time = time.time()
45
46 for fold, (train_index, test_index) in enumerate(
47     kf_outer.split(X, y)):
48     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
49     y_train, y_test = y.iloc[train_index], y.iloc[test_index]
50
51     # temporary dataframe for saving the results
52     temp_res = df.iloc[test_index].copy()
53     temp_res['fold'] = [fold] * temp_res.shape[0]
54
55     # select the best model
56     gscv = GridSearchCV(
57         rf,
58         parameters,
59         cv=kf_inner,
60         return_train_score=False,
61         n_jobs=-1)

```

```

62     gscv.fit(X_train, y_train)
63
64     # use the selected model to predict the classes for the
65     temp_res['pred_class'] = gscv.predict(X_test)
66
67     # save the best hyperparameter combination
68     temp_res['best_params'] = [
69         gscv.best_params_] * temp_res.shape[0]
70
71     results.append(temp_res)
72
73 results = pd.concat(results, axis=0)
74 results['accuracy'] = results.apply(
75     lambda i: 1 if i['class'] == i['pred_class'] else 0, axis=1)
76
77 results.index.name = 'orig_index'
78 # sort the results in the same order as in df
79 results = results.sort_values(
80     'orig_index').reset_index(drop=True)
81 results.to_parquet('../results/rf_results.parquet')
82
83 print("--- %s seconds ---" % (time.time() - start_time))

```

In [3]: ►

```

1 # results.index.name = 'orig_index'
2 # results = results.sort_values('orig_index').reset_index()
3 # results.to_parquet('../training_data/rf_results.parquet')

```

Mini-exam question:

Question 2. Nested cross validation

- is used to estimate the performance of a model selected in a particular way

Save the supplied (can be found from Moodle) rf_results.parquet-file in a folder named "results".

```
In [2]: ► 1 # start from here!
      2 results = pd.read_parquet('../results/rf_results.parquet')
```

Answer the questions/tasks:

- Which parameter combination is selected the best most often?
- Calculate the mean accuracy of each fold and the mean accuracy over all folds
- Calculate the confusion matrix over all samples

```
3 1 # just check that there are 50 test samples in each fold  
2 results.groupby('fold').size()
```

```
Out[3]: fold  
0    50  
1    50  
2    50  
3    50  
4    50  
5    50  
6    50  
7    50  
8    50  
9    50  
dtype: int64
```

```
In [4]: ► 1 results['param_comb'] = results.apply(  
2     lambda i: str(i.best_params['bootstrap']) +  
3         '_' + str(i.best_params['max_features']) +  
4         '_' + str(i.best_params['n_estimators']), axis=1)  
5 results.groupby(['param_comb']).size() / int(results.shape[0])
```

```
Out[4]: param_comb  
False_sqrt_100      4.0  
False_sqrt_150      1.0  
True_sqrt_100       1.0  
True_sqrt_150       2.0  
True_sqrt_250       1.0  
True_sqrt_300       1.0  
dtype: float64
```

The following parameter combination is chosen most often as the best (4 times out of 10):

- bootstrap: False
- max_features: sqrt
- n_estimators: 100

As we can see, different hyperparameter combinations can be selected as the best in different cv rounds.

```

5   1 display(results.groupby('fold').accuracy.mean())
  2 print('Mean accuracy over all folds:', 
  3       results.groupby(
  4           'fold').accuracy.mean().mean().round(3))

fold
0    0.98
1    1.00
2    1.00
3    1.00
4    1.00
5    1.00
6    0.96
7    1.00
8    1.00
9    1.00
Name: accuracy, dtype: float64

Mean accuracy over all folds: 0.994

```

Mini-exam question:

Question 3. The smallest fold-wise accuracy for Random Forest is

Answer: 0.96

Note: The following answers were accepted: 0.96, 0.96, 96, 96%

Here's another way to implement nested cross-validation which uses *cross_val_predict* function from scikit-learn.

Modify it to make the performance estimation for MLP and SVM (again, use the parameter ranges given in the Part 2).

For both models:

- Calculate the accuracy over all samples
- Show the confusion matrix (you can combine all samples in one confusion matrix as all folds are stratified and include the same number of samples)

Compare the three models. Which one performs the best?

```

In [6]: ► 1 feats_Z = ['mean_y_Z', 'var_y_Z',
  2 'skew_y_Z', 'kurt_y_Z',
  3 'mean_cr_Z', 'var_cr_Z',
  4 'skew_cr_Z', 'kurt_cr_Z',
  5 'mean_cb_Z', 'var_cb_Z',
  6 'skew_cb_Z', 'kurt_cb_Z',
  7 'major_axis_length_Z',
  8 'minor_axis_length_Z',
  9 'area_Z', 'perimeter_Z',
 10 'equivalent_diameter_Z',
 11 'compactness_Z',
 12 'shape_factor1_Z',
 13 'shape_factor2_Z']

```

```

7 # another way
8 start_time = time.time()
9 y = results['class']
10 X = results[feats_Z]
11
12 rf = RandomForestClassifier(random_state=20)
13
14 kf_outer = StratifiedKFold(
15     n_splits=10,
16     random_state=10,
17     shuffle=True)
18 kf_inner = RepeatedKFold(
19     n_splits=5,
20     n_repeats=3,
21     random_state=5)
22
23 n_estimators = range(100, 350, 50)
24 max_features = ['sqrt', 'log2', None]
25 bootstrap = [True, False]
26
27 parameters = {
28     'n_estimators': n_estimators,
29     'max_features': max_features,
30     'bootstrap': bootstrap}
31
32 gscv = GridSearchCV(
33     rf,
34     parameters,
35     cv=kf_inner,
36     return_train_score=False,
37     n_jobs=-1)
38 nested_pred = cross_val_predict(
39     gscv, X, y, cv=kf_outer)
40 print("--- %s seconds ---" % (time.time() - start_time))

```

--- 169.75516724586487 seconds ---

In [8]: ►

```

1 # Predicted values on columns, true values on rows
2 display(pd.DataFrame(
3     confusion_matrix(y, nested_pred)).rename(
4         mapper=dict(
5             zip(range(5),
6                 ['Arb', 'Bas', 'Ips', 'Jas', 'Kar'])),
7         axis=0).rename(
8             mapper=dict(
9                 zip(range(5),
10                    ['Arb', 'Bas', 'Ips', 'Jas', 'Kar'])),
11             axis=1))

```

	Arb	Bas	Ips	Jas	Kar
Arb	99	0	0	0	1
Bas	0	100	0	0	0
Ips	0	0	99	1	0
Jas	0	0	0	100	0
Kar	1	0	0	0	99

```
9 1 accuracy_score(y, nested_pred)
```

```
Out[9]: 0.994
```

```
In [10]: ► 1 # number of hyperparameter combinations
          2 a = [item[1] for item in parameters.items()]
          3 len(list(itertools.product(*a)))
```

```
Out[10]: 30
```

3.1.2 SVM

```
In [11]: ► 1 y = results['class']
          2 X = results[feats_Z]
```

```
In [12]: ► 1 start_time = time.time()
          2 kf_outer = StratifiedKFold(
          3     n_splits=10,
          4     random_state=10,
          5     shuffle=True)
          6 kf_inner = RepeatedKFold(
          7     n_splits=5,
          8     n_repeats=3,
          9     random_state=5)
         10
         11 gamma = ['scale', 'auto']
         12 C = [0.1, 1, 10, 100]
         13 kernel = ['linear', 'rbf', 'poly']
         14
         15 parameters = {'C': C,
         16                 'gamma': gamma,
         17                 'kernel': kernel}
         18
         19 svm = SVC()
         20
         21 gscv = GridSearchCV(
         22     svm,
         23     parameters,
         24     cv=kf_inner,
         25     return_train_score=False,
         26     n_jobs=-1)
         27 nested_pred = cross_val_predict(
         28     gscv, X, y, cv=kf_outer)
         29 print("--- %s seconds ---" % (time.time() - start_time))
```

--- 2.7521047592163086 seconds ---

In [13]

```
1 # Predicted values on columns, trues on rows
2 display(pd.DataFrame(
3     confusion_matrix(y, nested_pred)).rename(
4         mapper=dict(
5             zip(range(5),
6                 ['Arb', 'Bas', 'Ips', 'Jas', 'Kar'])),
7         axis=0).rename(
8             mapper=dict(
9                 zip(range(5),
10                    ['Arb', 'Bas', 'Ips', 'Jas', 'Kar']))),
11         axis=1))
```

	Arb	Bas	Ips	Jas	Kar
Arb	99	0	0	1	0
Bas	0	100	0	0	0
Ips	0	0	100	0	0
Jas	0	0	0	100	0
Kar	0	0	0	0	100

In [14]: ► 1 accuracy_score(y, nested_pred)

Out[14]: 0.998

In [15]: ► 1 # number of hyperparameter combinations
2 a = [item[1] for item in parameters.items()]
3 len(list(itertools.product(*a)))

Out[15]: 24

3.1.3 MLP

```
In [16]: ► 1 start_time = time.time()
  2 kf_outer = StratifiedKFold(
  3     n_splits=10,
  4     random_state=10,
  5     shuffle=True)
  6 kf_inner = RepeatedKFold(
  7     n_splits=5,
  8     n_repeats=3,
  9     random_state=5)
10
11 mlp = MLPClassifier(
12     max_iter=500,
13     early_stopping=True,
14     random_state=20)
15
16 hidden_layer_sizes = [(i,) for i in range(15, 45, 5)]
17 activation = ['tanh', 'relu']
18 solver = ['sgd', 'adam']
19 val_fr = [0.1, 0.3]
20 alpha = [0.01, 0.1, 1]
21
22 parameters = {
23     'hidden_layer_sizes': hidden_layer_sizes,
24     'activation': activation,
25     'solver': solver,
26     'validation_fraction': val_fr
27 }
28
29 gscv = GridSearchCV(
30     mlp,
31     parameters,
32     cv=kf_inner,
33     return_train_score=False,
34     n_jobs=-1)
35 nested_pred = cross_val_predict(
36     gscv, X, y, cv=kf_outer)
37 print("--- %s seconds ---" % (time.time() - start_time))
```

--- 43.2096381187439 seconds ---

```
In [17]: # Predicted values on columns, trues on rows
          1 display(pd.DataFrame(
          2     confusion_matrix(y, nested_pred)).rename(
          3         mapper=dict(
          4             zip(range(5),
          5                 ['Arb', 'Bas', 'Ips', 'Jas', 'Kar'])),
          6         axis=0).rename(
          7             mapper=dict(
          8                 zip(range(5),
          9                     ['Arb', 'Bas', 'Ips', 'Jas', 'Kar'])),
          10            axis=1))
```

	Arb	Bas	Ips	Jas	Kar
Arb	93	0	0	1	6
Bas	0	100	0	0	0
Ips	0	0	100	0	0
Jas	4	2	0	93	1
Kar	1	0	0	1	98

```
In [18]: accuracy_score(y, nested_pred)
```

Out[18]: 0.968

```
In [19]: # number of hyperparameter combinations
          1 a = [item[1] for item in parameters.items()]
          2
          3 len(list(itertools.product(*a)))
```

Out[19]: 48

Mini-exam question:

Question 4. Analyze the confusion matrices. Were the acquired results as expected? Why?

Answer: SVM had the highest accuracy, only once Arborio was predicted as Jasmine. RF gave the second best result: Arborio and Karacadag were once mislabeled as each other and Ipsala once as Jasmine. MLP made more errors than the other two classifiers: it struggled with Arborio and Jasmine the most. According to the PCA figure from part 2, it was expected that it is possible to train well-behaving classifiers, however more mislabeling might have been expected between Karacadag and Arborio. We used only two PCA components, probably we could have seen better discrimination between Karacadag and Arborio already with 3-dimensional PCA.

3.2 Discussion

Mini-exam question:

Question 5. Discuss your results:

- Ponder the limitations and generalization of the models. How well will they perform for data outside this data set? Explain. **1p**
- Compare your results with the original article. Are they comparable? **1p**

- Ponder applications for these kind of models, who could benefit from them?
Ponder also what would be interesting to study more on this area. **1b**

SVM gave the best results. With this data set size, it was very fast to train and evaluate. RF gave almost as good results as SVM, and it also offers an interpretation, how important each feature is. However, it took rather long time to evaluate its performance. MLP reached a bit lower accuracy than the other classifiers, but it would probably work better with a larger data set and when using a more extensive hyperparameter selection.

The sample images were prepared with a special preprocessing procedure. We can expect that the model produces similar results for rices from the same distribution (similar growing environment, same image preprocessing pipeline). The generalization outside this dataset, i.e. for the same rice species but from a different growing environment and conditions and with another preprocessing pipeline (differing lighting, resolution, cropping) is something which we cannot determine directly from the acquired results but should be tested separately. It would be interesting to test the models just by taking random rice samples available e.g. from our own kitchen.

Our results are rather comparable with the original article, however some differences exist: The original article compared six classifiers, including the three studied classifiers here. They used more features than we did. Performance of the classifiers were estimated using 10-fold cross validation. It seems they did not perform model selection, but the model parameters were fixed, e.g. they used four hidden layers and sigmoid activation for MLP. In this exercise we used only a small sample, 500 images in total. We made model selection by selecting the best model from several hyperparameter combinations for each studied classifier and estimated their performance using nested cross validation. The accuracies in the original article were 99.91/99.88/99.85% for MLP/RF/SVM, respectively. The results for RF and SVM classifiers were very close in our subset, whereas the accuracy for MLP was somewhat lower. This was probably due to the fact that we had only one hidden layer. In small data sets outliers may have a large impact. Generally, it is important to use cross validation or repeated splits, single splits to training and test set would not give a reliable estimation of the performance. In this subset, the performance evaluation results were very close to the accuracies obtained in the model selection phase, which shows that the selected models are performing consistently across different partitions of the dataset, so overfitting was not a problem here.

It is interesting, how these models could be deployed. The article states that human decisions are inconsistent, subjective and slow, and machine vision systems could be cost-effective, fast and accurate technique. The machine learning algorithm requires carefully prepared samples with only one rice grain present, which must limit the otherwise fast process of making the classification. The same type of machine learning classification scheme can be (and has been) used for other plant classification as well. It could also be used for detecting rotten or broken individuals from good ones. The results of these kind of models could benefit biodiversity researchers, and ultimately farmers and sellers. Personally, I'd like to see a research which studies more in detail the variation (due to different growing conditions and production stages) within one species, and how to tackle that (e.g. the bimodal feature values for Jasmine).

3.3 Reflection

Mini-exam question:

Question 6. Reflect your work: (0.5 points per answered question)

- Did you have any (mis)conceptions which changed while doing this work?
- What did you learn?
- What was difficult?
- In which aspects of work did you succeed well?
- How could you improve your working process?
- How can you deploy the learnings from this exercise in future?