

This is the template for the image recognition exercise.

Some general instructions, read these carefully:

- The final assignment is returned as a clear and understandable report
 - define shortly the concepts and explain the phases you use
 - use the Markdown features of Jupyter notebooks for large explanations
- return your output as a working Jupyter notebook
- name your file as `Exercise_MLPR2023_PartX_uuid.jsonnb`
 - use the uuid code determined below
 - use the standard naming convention for assignments
- write your code with comments
- if you copy/paste code from web, provide a reference
- it is ok to discuss with a friend about the assignment, but it is not ok to copy someone's work. Everyone should submit their own implementation
- in case of identical submissions, both submissions are failed

Deadlines:

- Part 1: Mon 6.2 at 23:59
- Part 2: Mon 20.2 at 23:59
- Part 3: Mon 6.3 at 23:59

No extensions for the deadlines

- after each deadline, example results are given, and it is not possible to submit anymore

If you encounter problems, Google first and if you can't find an answer, ask for help

- Moodle area for questions
 - pekav@utu.fi
- teacher available for questions
 - Monday 30.1 at 14:00-15:00 room 407B Honka (Agora 4th floor)
 - Monday 13.2 at 14:00-15:00 room 407B Honka (Agora 4th floor)
 - Thursday 2.3 at lecture 10:15-12:00

Grading

The exercise covers a part of the grading in this course. The course exam has 5 questions, 6 points of each. Exercise gives 6 points, i.e. the total score is 36 points.

From the template below, you can see how many exercise points can be acquired from each task. Exam points are given according to the table below:

Exercise Points	1 point	2 points	3 points	4 points	5 points	6 points
7 exercise points	1	0	0	0	0	0
9 exercise points	0	1	0	0	0	0
10 exercise points	0	0	1	0	0	0
11 exercise points	0	0	0	1	0	0
12 exercise points	0	0	0	0	1	0

To pass the exercise, you need at least 7 exercise points, and at least 1 exercise point from each Part.

Each student will grade one submission from a peer and their own submission. After each Part deadline, example results are given. Study them carefully and perform the grading according to the given instructions. Mean value from the peer grading and self-grading is used for the final points.

```
In [1]: # Import used
# Run this cell only once and save the code. Use the same id code for each Part.
# Printing random seed using uuid()
# print(uuid.uuid1())
# print(uuid.uuid4())

The id code is: 1Nw3Elgo7HhYDZPZh7HsLH7QHlpVW
```

Part 1

Read the original research article:

I. Çınar and M. Koklu. Identification of rice varieties using machine learning algorithms. Journal of Agricultural Sciences, 28(2):307–325, 2022. doi: 10.1583/ankutbd.862482.
<https://engpark.org.tr/en/download/article-file/1519632>

Introduction (1 p)

Will be written in Part 3

Preparations of the data (1 p)

Make three folders in your working folder: "notebooks", "data" and "training_data". Save this notebook in "notebooks" folder.

Perform preparations for the data

- import all the packages needed for this notebook in one cell
- import the images. Data can be found (downloading starts as you press the link) <https://www.muratkoklu.com/datasets/vtdhd9.php>
- save the data folders "Arborio", "Basmati" and "Jasmine" in "data" folder
- take a random sample of 100 images from Arborio, Basmati and Jasmine rice species (i.e. 300 images in total)
- determine the contour of each rice (you can use e.g. `findContours` from OpenCV)
- plot one example image of each rice species, including the contour

Feature extraction (2 p)

Gather the feature data

Color features (15)

- Calculate the following color features for each image, including only the pixels within the contour (you can use e.g. `pointPolygonTest` from OpenCV)
- Mean for each RGB color channel
- Variance for each RGB color channel
- Skewness for each RGB color channel
- Kurtosis for each RGB color channel
- Entropy for each RGB color channel

Dimension features (6)

- Fit an ellipse to the contour points (you can use e.g. `fitEllipse` from OpenCV)
- Plot one example image of each rice species including the fitted ellipse
- Calculate the following features for each image (for details, see the original article)
 - the major axis length of the ellipse
 - the minor axis length of the ellipse
 - area inside the contour (you can use e.g. `contourArea` from OpenCV)
 - perimeter of the contour (you can use e.g. `arcLength` from OpenCV)
 - roundness
 - aspect ratio

Gather all the features in one array or datframe: one data point in one row, including all feature values in columns.

For each data point, include also information of the original image and the label (rice species). Save the data in "training_data" folder.

Part 2

Data exploration (2 p)

Standardize the data,

- Plot a boxplot of each feature.
- Plot histogram for each feature using a different color for each class
- Plot each feature against each feature and the label against each feature

Discuss your findings from the above figures, e.g. can you spot features which might be very useful in predicting the correct class?

- Fit PCA using two components.
- Plot the PCA figure with two components, color the data points according to their species
- Can you see any clusters in PCA? Does this figure give you any clues, how well you will be able to classify the image types? Explain.
- How many PCA components are needed to cover 99% of the variance?
- Make clear figures, use titles and legends for clarification.

```
In [7]: # Reading data from the provided solution of part 1 in moodle.
data_path = os.path.join('..', 'training_data', 'rice_feature_data.parquet')
df = pd.read_parquet(data_path)
df.head(10)
```

mean_b	var_b	skew_b	kurt_b	entr_b	mean_g	var_g	skew_g	kurt_g	entr_g	mean_r	var_r	skew_r	kurt_r	entr_r	major_axis_length	minor_axis_length	area	perimeter	roundness	aspect_ratio	class
0.33487800	0.354786	0.009677	0.716162	0.322982	1.612640	0.421272	0.367326	-0.806359	0.203586	-0.925652	0.307681	0.839521	0.941032	0.36762	0.533615	0.994183	-1.046037	Arb			
0.971039	0.562107	-0.17077	0.403674	0.566321	5.930377	1.205033	2.041398	0.042097	0.598818	5.506168	0.407264	0.882907	-1.113487	0.622454	0.589942	1.186738	1.49388	-1.009049	Arb		
0.228070	-0.187133	0.403340	0.536525	0.538448	-0.186878	-0.270346	-0.702446	-0.585377	-0.201177	-0.107171	-0.207167	-0.110771	-0.107695	0.452525	-0.621208	1.025026	1.27009	1.077009	Arb		
0.318534	0.136491	-0.896994	0.164531	0.274261	1.277638	1.58599	1.581537	0.025030	0.329043	0.153871	0.533894	0.982067	1.164369	0.12230	0.628372	1.167169	1.187910	Arb			
0.1704427	0.136049	-0.310771	0.545430	0.027434	-0.919257	0.916752	-1.061714	0.145295	0.372460	-0.416213	0.375261	-0.872028	1.238729	0.36698	0.520795	1.277529	1.156844	Arb			
1.207410	0.299632	0.201206	0.051133	1.223305	1.193072	0.294386	0.509605	0.457637	0.562891	-0.443309	0.509779	0.707508	1.297360	0.499928	0.392800	1.125222	1.125222	Arb			
0.014069	0.569314	0.093162	-0.366289	0.140851	-0.170144	0.042043	-0.493130	-0.193444	0.709855	-0.094949	0.729935	-0.425598	2.074187	0.287426	1.198620	1.189615	1.189615	Arb			
0.1448078	0.857563	0.437674	0.751124	0.476165	5.701034	2.087674	2.104520	0.049483	0.110714	0.070424	0.677668	-0.118269	0.556969	145.034149	0.884285	7.7405	374.007140	0.695375	1.205475	Arb	
0.22763776	0.20527217	-0.518637	0.134448	0.572952	2.1669795	2.1512185	0.856129	0.945222	0.571909	1.023406	0.518659	0.118209	0.556969	155.566954	70.409998	76.225	369.800129	0.700421	1.974105	Arb	
0.21148584	0.20036110	-0.015902	0.519230	0.310937	5.786502	2.21587139	2.04195499	-0.405279	0.129115	0.578623	0.140829	0.547217	145.594498	79.186180	8657.0	309.64749	0.732085	1.836053	Arb		
0.0140772	0.385757	0.092628	0.519230	0.310937	5.786502	2.21587139	2.04195499	-0.405279	0.129115	0.578623	0.140829	0.547217	145.594498	79.186180	8657.0	309.64749	0.732085	1.836053	Arb		

10 rows × 23 columns

```
In [10]: # Standardization the data
labels = data.loc[:, data.columns[-2:]]
numerical_features = data.loc[:, data.columns[:-2]]
```

```
scaler = StandardScaler()
standardized_numerical_features = scaler.fit_transform(numerical_features)

# Standardized features and labels into a new DataFrame
standardized_df = pd.concat([standardized_numerical_features, labels], axis=1)

standardized_df.head(10)
```

mean_b	var_b	skew_b	kurt_b	entr_b	mean_g	var_g	skew_g	kurt_g	entr_g	mean_r	var_r	skew_r	kurt_r	entr_r	major_axis_length	minor_axis_length	area	perimeter	roundness	aspect_ratio	class
0.036251	0.281051	0.009577	0.716162	0.322982	1.612640	0.421272	0.367326	-0.806359	0.203586	-0.925652	0.307681	0.839521	0.941032	0.36762	0.533615	0.994183	-1.046037	Arb			
0.971039	0.562107	-0.17077	0.403674	0.566321	5.930377	1.205033	2.041398	-0.406297	0.598818	5.506168	0.407264	0.882907	-1.113487	0.622454	0.589942	1.186738	1.49388	-1.009049	Arb		
0.228070	-0.187133	0.403340	0.536525	0.538448	-0.186878	-0.270346	-0.702446	-0.585377	-0.201177	-0.107171	-0.207167	-0.110771	-0.107695	0.452525	-0.621208	1.025026	1.27009	1.077009	Arb		
0.318534	0.136491	-0.896994	0.164531	0.274261	1.277638	1.58599	1.581537	0.025030	0.329043	0.153871	0.533894	0.982067	1.164369	0.12230	0.628372	1.167169	1.187910	Arb			
0.1704427	0.136049	-0.310771	0.545430	0.027434	-0.919257	0.916752	-1.061714	0.145295	0.372460	-0.416213	0.375261	-0.872028	1.238729	0.36698	0.520795	1.277529	1.156844	Arb			
1.207410	0.299632	0.201206	0.051133	1.223305	1.193072	0.294386	0.509605	0.457637	0.562891	-0.443309	0.509779	0.707508	1.297360	0.499928	0.392800	1.125222	1.125222	Arb			
0.014069	0.569314	0.093162	-0.366289	0.140851	-0.170144	0.042043	-0.493130	-0.193444	0.709855	-0.094949	0.729935	-0.425598	2.074187	0.287426	1.198620	1.189615	1.189615	Arb			
0.1448078	0.857563	0.437674	0.751124	0.476165	5.701034	2.087674	2.104520	0.049483	0.110714	0.070424	0.677668	-0.118269	0.556969	145.034149	0.884285	7.7405	374.007140	0.695375	1.205475	Arb	
0.22763776	0.20527217	-0.518637	0.134448	0.572952	2.1669795	2.1512185	0.856129	0.945222	0.571909	1.023406	0.518659	0.118209	0.556969	155.566954	70.409998	76.225	369.80012				