

Assignment 5: Ensemble Methods

Copyright and Fair Use

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Automatic Testing Guidelines

Automatic unittesting requires you, as a student, to submit a notebook which contains strictly defined objects. Strictness of definition consists of unified shapes, dtypes, variable names and more.

Within the notebook, we provide detailed instruction which you should follow in order to maximise your final grade.

Name your notebook properly, follow the pattern in template name:

Assignment_N_NameSurname_matrnumber

1. N - number of assignment
2. NameSurname - your full name where every part of the name starts with a capital letter, no spaces
3. matrnumber - your 8-digit student number on ID card (without k)

Example:

- ✓ Assignment_0_RenéDescartes_12345678
- ✓ Assignment_0_SørenAabyeKierkegaard_12345678
- ✗ Assignment0_Peter_Pan_k12345678

Don't add any cells but use the ones provided by us. You may notice that most cells are tagged such that the unittest routine can recognise them.

We highly recommend you to develop your code within the provided cells. You can implement helper functions where needed unless you put them in the same cell they are actually called. Always make sure that implemented functions have the correct output and given variables contain the correct data type. Don't import any other packages than listed in the cell with the "imports" tag.

Note: Never use variables you defined in another cell in your functions directly; always pass them to the function as a parameter. In the unittest they won't be available either.

Good luck!

Task 1: AdaBoostM1 is an instance of forward stagewise modelling

In the lecture it was mentioned that one of the first boosting algorithms, i.e.

AdaBoostM1, is equivalent to forward stagewise modelling using the exponential loss $L(y, g(\mathbf{x})) = \exp(-yg(\mathbf{x}))$ for a binary classification problem ($y \in \{-1, 1\}$). In this task we intend to provide proof of this fact. We will guide you through the most important steps and you will have to add some details.

For AdaBoostM1, the basis functions at timestep n are the individual classifiers $b_n(\mathbf{x}) \in \{-1, 1\}$. We assume that all of them are slightly better than random guessing. Note that we use b_n here for the resulting classifier at timestep n , which differs slightly from the notation in the slides, mainly to not confuse it with the corresponding approximation from forward stagewise modelling, which is also called g_n there.

Using the exponential loss in each timestep n we have to solve

$$(\beta_n, b_n) = \arg \min_{\beta, b} \sum_{i=1}^l \exp(-y_i(g_{n-1}(\mathbf{x}_i) + \beta b(\mathbf{x}_i))),$$

for the classifier b_n and the coefficient β_n which are added at each step. This can be rewritten as

$$(\beta_n, b_n) = \arg \min_{\beta, b} \sum_{i=1}^l w_i^{(n)} \exp(-y_i \beta b(\mathbf{x}_i)), \quad (1)$$

with $w_i^{(n)} = \exp(-y_i g_{n-1}(\mathbf{x}_i))$. Since each $w_i^{(n)}$ depends neither on β nor b , it can be regarded as a weight that is applied to each observation. This weight depends on $g_{n-1}(\mathbf{x}_i)$, and so the individual weight values change with each iteration n . The solution of (1) can be found in two steps:

Calculation 1 (10 points):

In the **first step** fix $\beta \geq 0$ and show that in this case the solution to (1) is

$$b_n = \arg \min_b \sum_{i=1}^l w_i^{(n)} I(b(\mathbf{x}_i) \neq y_i). \text{ Hints:}$$

1. Try to write the expressions in (1) after $\arg \min$ in the form

$$\sum_{b(\mathbf{x}_i)=y_i} (\dots) + \sum_{b(\mathbf{x}_i) \neq y_i} (\dots).$$

Find the right expressions for (...), such that the product $b(\mathbf{x}_i)y_i$ doesn't appear there anymore.

2. Now show that this can be written as

$$(\exp(\beta) - \exp(-\beta)) \sum_{b(\mathbf{x}_i) \neq y_i} w_i^{(n)} + \exp(-\beta) \sum_{i=1}^l w_i^{(n)}. \quad (2)$$

The equation $\sum_i w_i^{(n)} = \sum_{b(\mathbf{x}_i) \neq y_i} w_i^{(n)} + \sum_{b(\mathbf{x}_i)=y_i} w_i^{(n)}$ might be helpful.

3. Argue why this already implies the claim, i.e., the solution to (1) is

$$b_n = \arg \min_b \sum_{i=1}^l w_i^{(n)} I(b(\mathbf{x}_i) \neq y_i).$$

Please provide reasoning and explanations in full sentences. Grading of the task will heavily depend on it.

Write your first step* calculation here in **this cell***

1. Hint 1 ...
2. Hint 2 ...
3. Hint 3 ...

Calculation 2 (20 points):

In the **second step** you need to optimize the following expression with respect to β :

$$(\exp(\beta) - \exp(-\beta)) \sum_{b_n(\mathbf{x}_i) \neq y_i} w_i^{(n)} + \exp(-\beta) \sum_{i=1}^l w_i^{(n)}$$

Hint: do it in the usual way (i.e. by differentiating the expression and setting it to 0).

Using the abbreviation $\text{err}_n = \frac{\sum_{i=1}^l w_i^{(n)} I(y_i \neq b_n(\mathbf{x}_i))}{\sum_{i=1}^l w_i^{(n)}}$ show that the obtained expression is

$$\beta_n = \frac{1}{2} \ln \frac{1 - \text{err}_n}{\text{err}_n}.$$

Note that $\beta_n \geq 0$, by our assumption that all classifiers are better than random guessing, i.e. $\text{err}_n \leq \frac{1}{2}$, so the result is in accordance with the previous subtask.

Please provide reasoning and explanations in full sentences. Grading of the task will heavily depend on it.

Write your second step* calculation here in **this cell***

1. Derivate ...
2. Using given abbreviation show ...

Calculation 3 (10 points):

In the **final step** we can update the approximation as follows:

$g_n(\mathbf{x}) = g_{n-1}(\mathbf{x}) + \beta_n b_n(\mathbf{x})$. To finish the proof proceed by deriving the following relations:

1. The weights for the next generation can be computed as follows:

$$w_i^{(n+1)} = w_i^{(n)} \exp(-y_i \beta_n b_n(\mathbf{x})).$$
2. $-y_i b_n(\mathbf{x}) = 2I(y_i \neq b_n(\mathbf{x}_i)) - 1$
3. Use these two relations to show that:

$$w_i^{(n+1)} = w_i^{(n)} \exp(-\beta_n) \exp(\alpha_n I(y_i \neq b_n(\mathbf{x}_i))) \quad (2)$$

where $\alpha_n = 2\beta_n$ is the α_n from the AdaBoostM1 algorithm from the lecture.

Please provide reasoning and explanations in full sentences. Grading of the task will heavily depend on it.

Write your *final step** calculation here in **this cell***

1. Calculate weights ...
2. ...
3. ...

Task 2: Random Forests and Feature importance

In this task you will train a Random Forest (RF) Classifier on a subset of fashionMNIST. You should observe how these models can immediately give you useful information about feature importance, which is a very convenient property of RFs.

First you should re-use the given code from the previous assignment to load the whole data set (the procedure is completely analogous).

- Next implement code that filters the data set for the classes with labels 3 (dresses) and 8 (bags). You should create the filtered train data set from `x_train` and `y_train` and the test data set from `x_test` and `y_test`.
Hint: Masks provide a convenient solution to this task.
- After the filtering procedure the data samples corresponding to dresses should be labelled as 1 and the bags as 0. Perform this step on the test and train data set.
- To accomplish this task, implement a function `_filter_()`.

Code 2.1 (15 points):

```
In [1]: #NOTE#####  
#Please add all your imports in this cell only  
#####
```

```
#Nothing to do here  
import numpy as np  
import pandas as pd  
import sys  
import time  
import numpy as np  
from sklearn.ensemble import RandomForestClassifier  
from mnist_loader import MNIST  
import matplotlib.pyplot as plt  
from matplotlib import style  
import matplotlib as mpl  
import seaborn as sns  
from sklearn.metrics import accuracy_score, confusion_matrix  
# Set random seed to ensure reproducible runs  
RSEED = 10
```

```
In [2]: #Load training and test data (routine from previous assignment)  
data = MNIST('./dataset/')  
img_train, labels_train = data.load_training()
```

```

x_train = np.array(img_train)
y_train = np.array(labels_train)
x_test,y_test = data.load_testing()
x_test = np.array(x_test)
y_test = np.array(y_test)
print(y_train)
print(y_test)

```

```

[2 9 6 ... 8 8 7]
[0 1 2 ... 8 8 1]

```

```

In [3]: def _filter_(
        x_train: np.ndarray,
        y_train: np.ndarray,
        x_test: np.ndarray,
        y_test: np.ndarray,
        labels_list: list[int, int]
    ) -> tuple[np.ndarray, np.ndarray, np.ndarray, np.ndarray]:
        """This function filters the datasets w.r.t to given labels_list.
        So, in the end x_train and x_test only contain the samples with labels that

        Parameters
        -----
        x_train : np.ndarray
            Training data matrix.
        y_train : np.ndarray
            Training labels vector.
        x_test : np.ndarray
            Test data matrix.
        y_test : np.ndarray
            Test labels vector.
        labels_list : list[int, int]
            list of length 2 which consists of integer labels.

        Returns
        -----
        tuple(np.ndarray, np.ndarray, np.ndarray, np.ndarray)
            Returns the filtered train and test data matrix and labels vector (which
            Check the return statement to see the actual order of returned arrays.
        """
        #Your code goes here ↓↓↓
        mask = [True if y in labels_list else False for y in y_train]
        x_train_filtered = x_train[mask]
        y_train_filtered = y_train[mask]
        y_train_filtered = np.array([1 if y == 3 else 0 for y in y_train_filtered])

        mask = [True if y in labels_list else False for y in y_test]
        x_test_filtered = x_test[mask]
        y_test_filtered = y_test[mask]
        y_test_filtered = np.array([1 if y == 3 else 0 for y in y_test_filtered])
        #Your code ends here _____

        return x_train_filtered, y_train_filtered, x_test_filtered, y_test_filtered

```

```

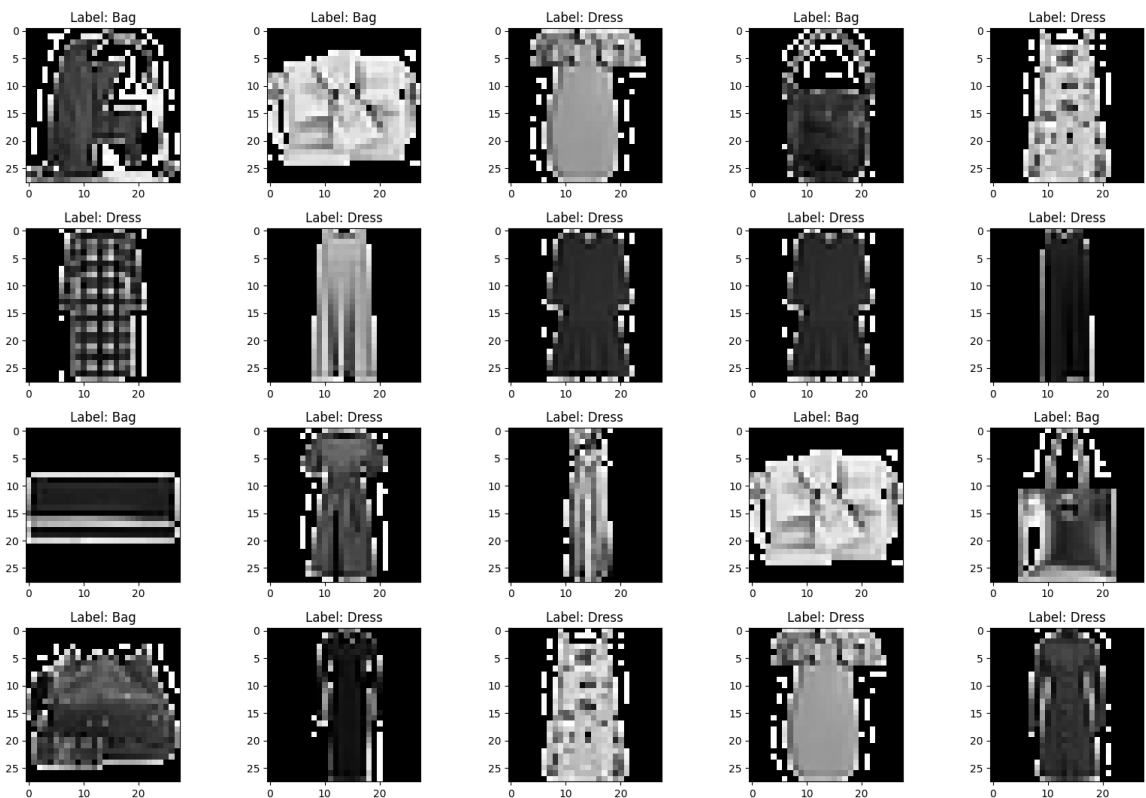
In [4]: #NOTHING TO CHANGE HERE
        x_train, y_train, x_test, y_test = _filter_(x_train, y_train, x_test, y_test, [3
        print(y_train)
        print(y_test)

```

```
[1 0 0 ... 0 0 0]
[1 0 1 ... 0 0 0]
```

In the following we visualize a few randomly selected samples from our training data:

```
In [5]: #A routine that you can use for plotting some of the data.
arr = ['Bag', 'Dress'] # 0: Bag, 1: Dress
a = np.random.randint(1,40,20)
plt.figure(figsize=(20, 13))
for n,i in enumerate(a):
    plt.subplot(4, 5, n+1)
    two_d = (np.reshape(x_train[i], (28, 28)) * 255).astype(np.uint8)
    plt.title('Label: {0}'.format(arr[y_train[i]]))
    plt.imshow(two_d, interpolation='nearest', cmap='gray')
plt.subplots_adjust(hspace = 0.3)
```



Code 2.2 (5 points):

Your task now is to train a sklearn [RandomForestClassifier](#) with the default parameters on the training data set.

Then get the model's predictions for the test data set. Use **RSEED** as random_seed for the RandomForestClassifier.

For this, we ask you to implement a function fit_predict.

```
In [6]: def fit_predict(
    x_train: np.ndarray,
    y_train: np.ndarray,
    x_test: np.ndarray,
    y_test: np.ndarray,
    rseed: int
) -> tuple[RandomForestClassifier, np.ndarray]:
```

```

"""Function fits a RandomForestClassifier on the training data and returns p

Parameters
-----
x_train : np.ndarray
    Training data matrix.
y_train : np.ndarray
    Training labels vector.
x_test : np.ndarray
    Test data matrix.
y_test : np.ndarray
    Test labels vector.
rseed : int
    Random Seed used for initializing the Classifier.

Returns
-----
tuple(RandomForestClassifier, np.ndarray)
    Where the classifier is the already trained classifier and the array is
"""
#Your code goes here ↓↓↓
clf = RandomForestClassifier(random_state=rseed)
model = clf.fit(x_train, y_train)

prediction = model.predict(x_test)
#Your code ends here _____

return model, prediction

```

```
In [7]: model, prediction = fit_predict(x_train,y_train,x_test,y_test, RSEED)
```

Code 2.3 (15 points):

Now, within the function `get_n_items_wrong` compute and return the size of the test set and number of misclassified test samples. Furthermore, retrieve the misclassified samples, their actual labels and the predicted labels from the classifier.

Use your predictions to plot up to 20 test data sample(s) that were **misclassified** (there can be fewer than 20 misclassified samples as well):

```
In [8]: def get_n_items_wrong(
        x_test:
            np.ndarray, y_test:
            np.ndarray,
            prediction: np.ndarray
    ) -> tuple[int, int, np.ndarray, np.ndarray, np.ndarray]:
    """Calculates the size of the test set, number of misclassified samples, the
    the actual labels that are misclassified in the prediction, and the labels t

    Parameters
    -----
    x_test : np.ndarray
        Test data matrix.
    y_test : np.ndarray
        Test labels vector.
    prediction : np.ndarray
        The predicted labels vector from the model.

```

```

Returns
-----
tuple(int, int, np.ndarray, np.ndarray, np.ndarray)
    Where the first two ints are the size_test and num_wrong and the arrays
    """
# Your code goes here ↓↓↓
size_test = len(x_test)
# create mask via np.not_equal
mask = np.not_equal(y_test, prediction)
num_wrong = len(prediction[mask])
samples_wrong = x_test[mask]
labels_wrong = np.array(y_test)[mask]
predictions_wrong = prediction[mask]
# Your code ends here _____

return size_test, num_wrong, samples_wrong, labels_wrong, predictions_wrong

```

```

In [9]: #SOLUTION TO NUMBER OF WRONG PREDICTIONS
#Your code goes here ↓↓↓
size_test, num_wrong, samples_wrong, labels_wrong, predictions_wrong = get_n_ite
#Your code ends here _____

#Following print statement might be evaluated
print("Number of test samples: {0}\nNumber of misclassified samples: {1}".format

Number of test samples: 2000
Number of misclassified samples: 11

```

```

In [10]: def plot_wrong_predictions(samples_wrong: np.ndarray, labels_wrong: np.ndarray,
    """Function creates a figure that shows misclassified samples.

Parameters
-----
samples_wrong : np.ndarray
    The samples that got misclassified by the classifier.
labels_wrong : np.ndarray
    The actual true labels of the misclassified samples.
predictions_wrong : np.ndarray
    The predicted labels of the misclassified samples.

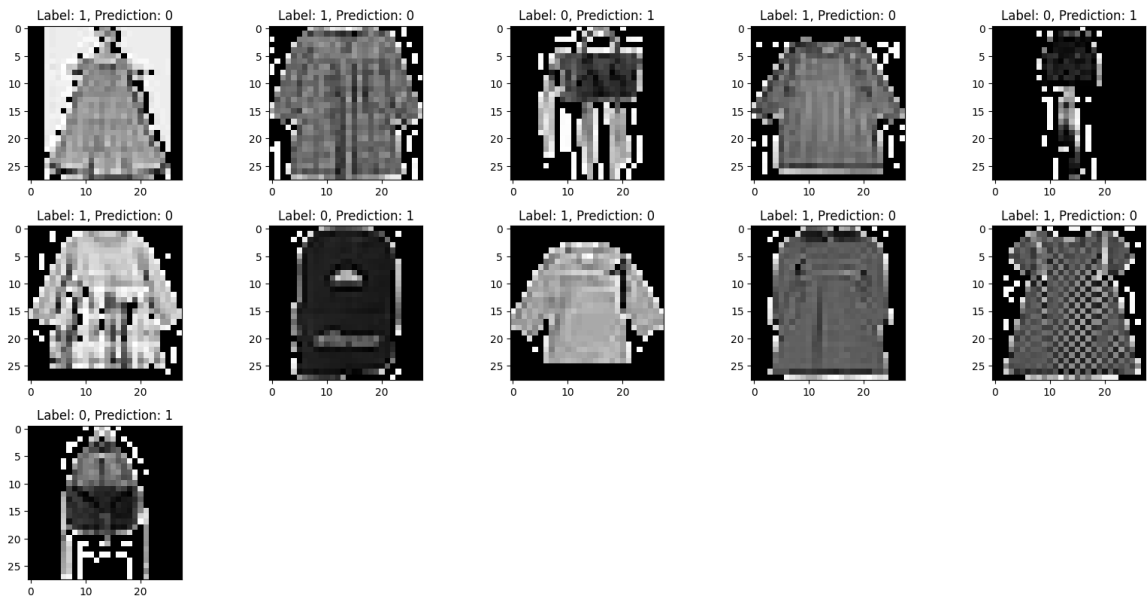
Returns
-----
plt.Figure
    A matplotlib.pyplot figure object (i.e. the misclassified samples).
    """
# Your code goes here ↓↓↓
fig = plt.figure(figsize=(20, 13))
for i in range(len(samples_wrong)):
    plt.subplot(4, 5, i + 1)
    two_d = (np.reshape(samples_wrong[i], (28, 28)) * 255).astype(np.uint8)
    plt.title('Label: {0}, Prediction: {1}'.format(labels_wrong[i], predicti
    plt.imshow(two_d, interpolation='nearest', cmap='gray')
plt.subplots_adjust(hspace=0.3)
plt.show()
# Your code ends here _____

return fig

```



```
In [11]: fig = plot_wrong_predictions(samples_wrong, labels_wrong, predictions_wrong)
```



Code 2.4 (20 points):

Within this part we will try to see the decision-making incentives of Random Forest.

To do this we ask you to implement the following tasks:

Hint: In the following you will have to plot heatmaps, to do this you might want to take a look at seaborn (sns): [Seaborn heatmap](#)

1. Take your training dataset and split it into 2: dresses and bags.
Calculate the average of the features. Reshape the averages to 2D arrays of shape 28*28 and plot them as heatmaps. Save the results under variables **dress_av** and **bag_av**
2. From the average of dresses subtract the average of bags. Save the result as in the variable **diff**. Plot it as a heatmap.
3. Define the feature importance of the previously trained RF classifier within the variable **importances**. Visualize it as a heatmap and don't forget to reshape it to 28*28.

Hint: Check scikit-learn documentation to access feature importance.

The evaluation of the following code will be done by viewing your plots.

For those who are curious: run RF under different seeds, and look how plots are changing.

BEFORE SUBMISSION RETURN TO ORIGINAL SEED = 10

```
In [12]: #PLOTING HEATMAPS
def plot_heatmaps(x_train: np.ndarray, y_train: np.ndarray, model: RandomForestC
    """Create plots as described in the task description above.

    Parameters
    -----
    x_train : np.ndarray
```

```

    Training data matrix.
y_train : np.ndarray
    Training labels vector.
model : RandomForestClassifier
    Your already trained classifier.

Returns
-----
plt.Figure
    A matplotlib.pyplot figure object (i.e. the heatmaps).
"""
#Your code goes here ↓↓↓

#Use your freestlye plotting

#Step 1 Split and calculate averages
dress_av = np.reshape(np.mean(x_train[y_train == 1], axis=0), (28, 28))
bag_av = np.reshape(np.mean(x_train[y_train == 0], axis=0), (28, 28))

# Step 2 Subtract
diff = dress_av - bag_av

# Step 3 Extract feature importance
importances = np.reshape(model.feature_importances_, (28, 28))

# Step 4 plot everything together
fig, axs = plt.subplots(2,2,figsize=(14, 14), gridspec_kw = {'wspace':0.15, '
# bag
sns.heatmap(bag_av, ax=axs[0, 0])
axs[0, 0].set_title('Average bags')

# dress
sns.heatmap(dress_av, ax=axs[0, 1])
axs[0, 1].set_title('Average dresses')

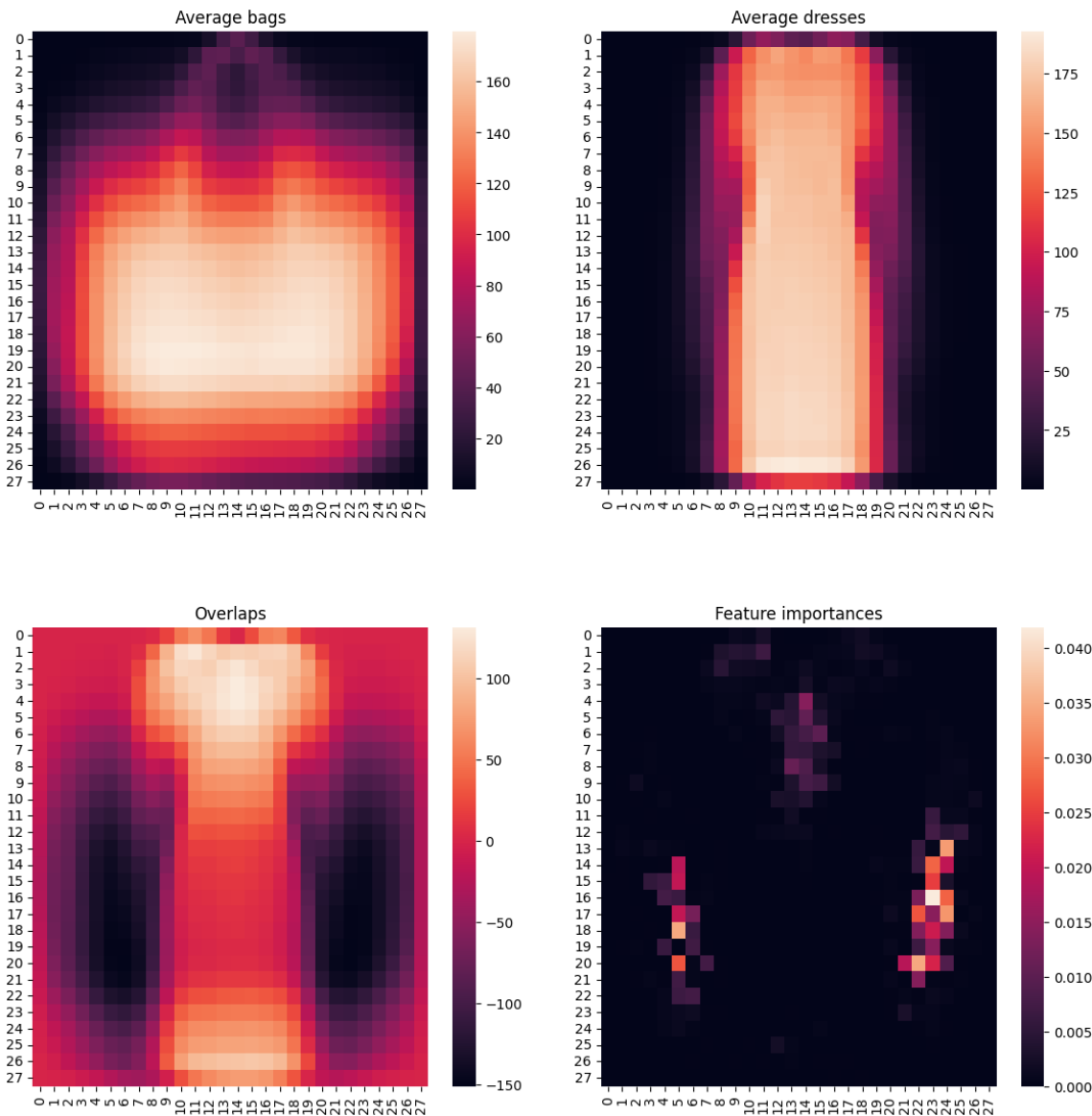
# diff
sns.heatmap(diff, ax=axs[1, 0])
axs[1, 0].set_title('Overlaps')

# importance
sns.heatmap(importances, ax=axs[1, 1])
axs[1, 1].set_title('Feature importances')
plt.show()
#Your code ends here _____

return fig

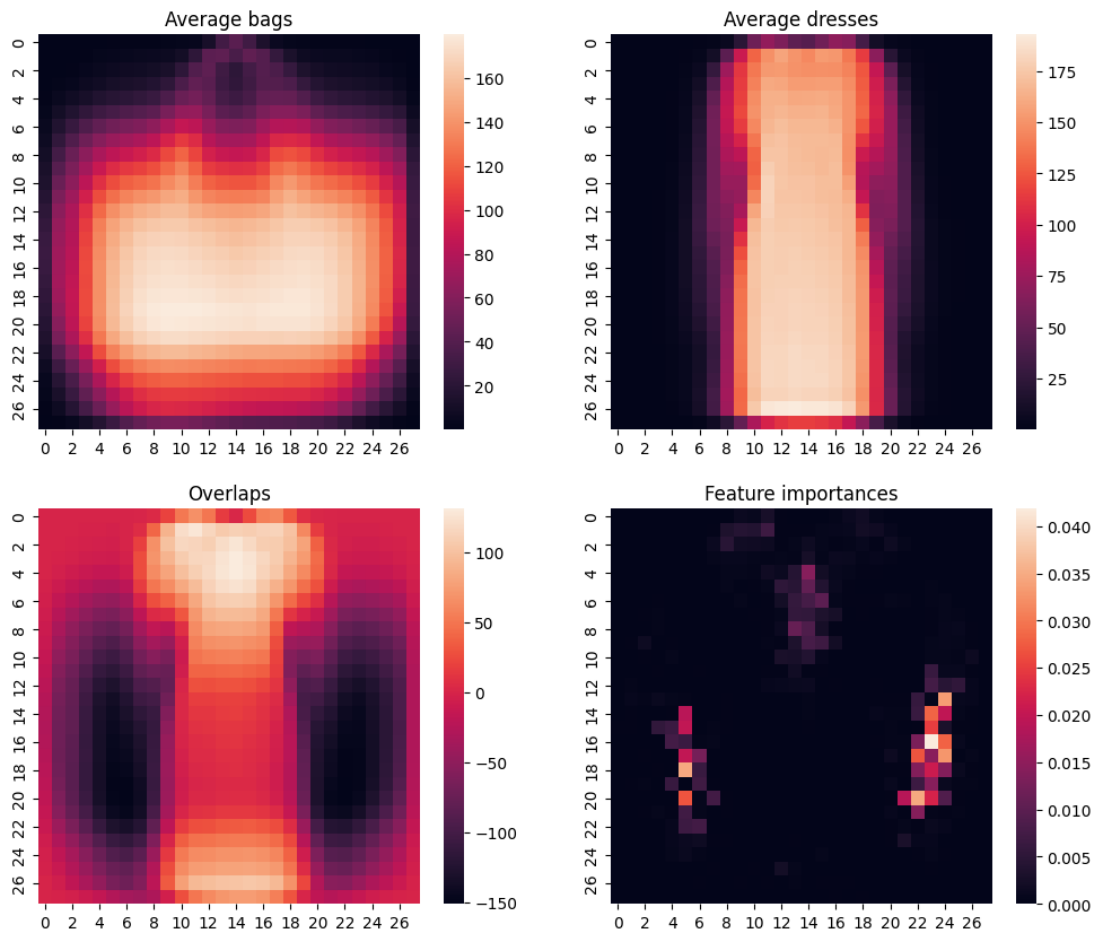
```

```
In [13]: fig = plot_heatmaps(x_train, y_train, model)
```



If you have solved the previous task correctly, the resulting plot should look close to this:

Main Task



In [14]: *# You can use this cell to run RF under different seeds, if you are curious (wil*

Question (5 points):

What observations can you make?

(Multiple answers might be correct)

To answer the question assign to variables in the nex cell **True** or **False** boolean values.

To earn points **assign values to all variables**.

NOTE Do not reuse these variable names. They are used for testing.

- a_) RF achieves an accuracy higher than 90% on the test set.
- b_) The most important features are located in image regions that primarily show either only bags or dresses as depicted by the plots of average bags and dresses.
- c_) If one pixel would always be bright for bags but never for dresses, the RF classifier would certainly learn to use only this pixel as a feature.
- d_) Comparing the misclassified bags to other bags suggests that they might be outliers as they are rather tall, but not as long as most other bags.

In [15]: *#examples for you*
 example_of_true_variable = **True**
 example_of_false_variable = **False**

```
#your answers go here ↓↓↓
a_=True
b_=True
c_=False
d_=True
```

```
In [16]: # this cell and the next one can help you as a sanity check of your implementati
def filter_tl(xt,yt,xte,yte,labels_list):
    try:
        xtf, ytf, xtef, ytef = _filter_(xt,yt,xte,yte,labels_list)
        res = (xtf.tolist(), ytf.tolist(), xtef.tolist(), ytef.tolist())
        print('Filtering works')
        return res
    except Exception as e:
        print("Execution of _filter_, or conversion from np.array to list failed")
        print(e)
        raise ValueError()

def model_param(xt,yt,xte,yte):
    try:
        model, pred = fit_predict(xt,yt,xte,yte,RSEED)
        nfeat = model.n_features_in_
        nclass = model.n_classes_
        print("Model is valid object")
        return nfeat,nclass
    except Exception as e:
        print("Execution of fit_predict, or model parameters extraction failed")
        print(e)
        raise ValueError()

def testoptions(options):
    for elem in options:
        if elem!=True and elem!=False and elem!=None:
            raise ValueError(f"Check answers for questions again")
    print("Test questions answers are ok")
```

```
In [17]: #these are dummy inputs, your code should be runnable on such
xt,yt,xte,yte = np.array([[1,1,1],[9,9,9],[5,5,5]]),np.array([1,9,5]),np.array([

xt,yt,xte,yte = filter_tl(xt,yt,xte,yte,[1,9])
model_param(xt,yt,xte,yte)
testoptions(np.array([a_,b_,c_,d_]))
print('Executable')
```

```
Filtering works
Model is valid object
Test questions answers are ok
Executable
```