# Report I - SF2957 Statistical Machine Learning

Anna Barysheva 960603-7423          Karin Larsson 980522-4384
Ruxue Zeng 981015-T549

2 December, 2021

## 1    Classification of MNIST

### 1.1    Objectives

In this project, we will use three different methods to design and implement a classifier for handwritten digit recognition : Stochastic subgradients, Gaussian process and Convolutional neural network. The dataset is made up of 1797 images of size $8 \times 8$. Each image is of a low-resolution hand-written digit.

### 1.2    Mathematical Background

In this section, we describe the general set up of our multi-class support vector machine.

#### 1.2.1    $R^{emp}$ convex

We consider a classification problem with $K = 10$ classes $\{0, 1, ..., 9\}$ and $d = 8 \times 8 = 64$. The classifier is represented using matrix

$$\Theta = [\theta_1 \theta_2 \ldots \theta_K] \in \mathbb{R}^{d \times K}$$

Given $\Theta$ the predicted class for a data vector image $x \in \mathbb{R}^d$ is the label

$$\hat{y} = \operatorname*{argmax}_{j=1,...,K}[\Theta^T x]_j$$

To train the classifier, we minimize the empirical risk using a multi-class hinge-loss

$$L(\Theta, (x, y)) = \max_{j \neq y} \left[1 + \sum_{i=1}^{d} x_i(\Theta_{ij} - \Theta_{iy})\right]_+,$$

where $t_+ = max(t, 0)$ is the positive part. More precisely, we aim to minimize

$$R^{emp}(\Theta) = \frac{1}{n} \sum_{k=1}^{n} L(\Theta, (x_k, y_k))$$

We want to give a subgradient G of $R^{emp}(\Theta)$. To do so, firstly, we must to show that $R^{emp}(\Theta)$ is convex.

Let $\lambda \in [0, 1]$ and $(\Theta, \Psi) \in \mathbb{R}^{d \times K} \times \mathbb{R}^{d \times K}$, we have:

$$R^{emp}(\lambda\Theta + (1-\lambda)\Psi) = \frac{1}{n} \sum_{k=1}^{n} \max_{j \neq y_k} \left(1 + \sum_{i=1}^{d} x_{k,i}(\lambda\Theta_{ij} + (1-\lambda)\Psi_{ij} - \lambda\Theta_{iy_k} - (1-\lambda)\Psi_{iy_k})\right)$$

We rewrite $1 = \lambda + (1 - \lambda)$ and separate the sum in two parts, we obtain:

$$R^{emp}(\lambda\Theta + (1-\lambda)\Psi) = \frac{1}{n} \sum_{k=1}^{n} \max_{j \neq y_k} \left(\lambda + \sum_{i=1}^{d} x_{k,i}(\lambda\Theta_{ij} - \lambda\Theta_{iy_k}) + (1-\lambda) + \sum_{i=1}^{d} x_{k,i}((1-\lambda)\Psi_{ij} - (1-\lambda)\Psi_{iy_k})\right)$$

Since the maximum of a sum is less or equal to the sum of the maximum, we obtain the inequality :

$$R^{emp}(\lambda\Theta + (1-\lambda)\Psi) \leq \frac{1}{n}\left[\sum_{k=1}^{n} \max_{j \neq y_k} \left(\lambda + \sum_{i=1}^{d} x_{k,i}(\lambda\Theta_{ij} - \lambda\Theta_{iy_k})\right)\right.$$
$$\left. + \sum_{k=1}^{n} \max_{j \neq y_k} \left((1-\lambda) + \sum_{i=1}^{d} x_{k,i}((1-\lambda)\Psi_{ij} - (1-\lambda)\Psi_{iy_k})\right)\right]$$

Thus, we have :

$$\boxed{R^{emp}(\lambda\Theta + (1-\lambda)\Psi) \leq \lambda R^{emp}(\Theta) + (1-\lambda)R^{emp}(\Psi)}$$

So $R^{emp}$ is a convex function, now we can calculate its subgradient.

### 1.2.2  Subgradient $G$ of $R^{emp}$

We assume we have only one data set $(x, y)$, thus $L(\Theta, (x, y)) = R^{emp}(\Theta)$. We consider

$$j^* = \max_{j \neq y} \left[1 + \sum_{i=1}^{d} x_i(\Theta_{ij} - \Theta_{iy})\right]_+ , \quad L(\Theta, (x, y)) = \left[1 + \sum_{i=1}^{d} x_i(\Theta_{ij^*} - \Theta_{iy})\right]_+$$

We have two cases:

- $L(\Theta, (x, y)) = 0$ : which means $1 + \sum_{i=1}^{d} x_i(\Theta_{ij^*} - \Theta_{iy}) \leq 0 \ \forall j$. If the inequality is strict, the expression is differentiable, we have gradient equals to zero matrix. If we have the equality, we can take the simplest subgradient zero matrix.Thus, we have

$$G = 0 \in \mathbb{R}^{d \times k}$$

- $L(\Theta, (x, y)) > 0$ : which means $L(\Theta, (x, y)) = 1 + \sum_{i=1}^{d} x_i(\Theta_{ij^*} - \Theta_{iy})$. We take the derivative of this expression with regard to $j^{th}$ column, we obtain :

$$G = \begin{cases} x & \text{if } j = j^* \\ -x & \text{if } j = y \\ 0 & \text{otherwise} \end{cases}$$

## 1.3 Results

### 1.3.1 Stochastic Subgradient

Using the subgradient calculated of $R^{emp}$, a projected stochastic gradient descent classifier was implemented to classify the MNIST dataset.

Initially, we took the stepsize $\epsilon_k$ proportional to $k^{-\alpha}$ for $\alpha \in (0, 1]$, and project $\Theta$ to the ball

$$B_r = \{\Theta \in \mathbb{R}^{dk} : \sum_{i,j} \Theta_{ij}^2 \leq r^2\}$$

We took $\alpha = 0.5$ and evaluated the performance of the classifier using training sets of size 20, 50, 100, 500, 1000 and 1500, and a test set of size 100. We plot the error and the standard deviation over 30 trials in figure 1a. We observed the classifier gets a better performance when increasing the number of training sets.

Then, with 1697 training sets and 100 test sets, we investigated the performance of the classifier with different choices of $\alpha$. We observe in figure 1b that the performance of the classifier gets worse when increasing the $\alpha$. We remark also that the standard derivation become extremely large for $\alpha \in (0.6, 1]$.



(a) With different size of training set          (b) With different stepsize parameter $\alpha$
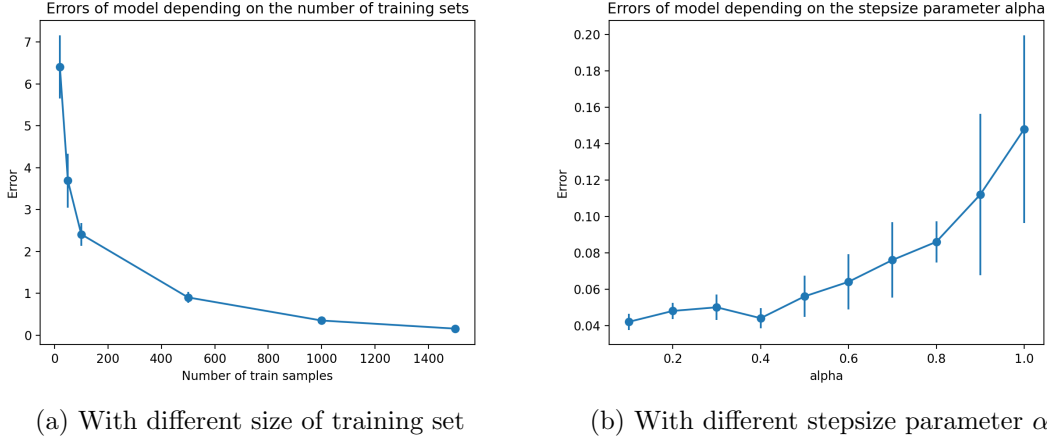
Figure 1: Error rate and SGD on the test set of size 100

### 1.3.2 Gaussian process

A Gaussian process classifier was implemented to classify the MNIST dataset using two different methods. Firstly the scikit learn method GaussianProcessClassifier (GPC) was used. Secondly the scikit learn GaussianProcessRegressor (GPR) was used with $y$ as a one hot vector. The performance was evaluated by analysing the mean of the relative error of the test set. The mean was computed over 3 tries. This was done for multiple kernels. In order to be able to compare the same kernel on GPC and GPR, GPC (optimize = None) was used.

When using a radial basis function (RBF) as the kernel, the length scale $L$ was of high importance to the error rate. This can be seen for some parameters for the GPC in figure 2a and for the GPR in figure 2b. GPR seems to be less affected by the choice of $L$, and gave lower error rates for most points, see figure 3
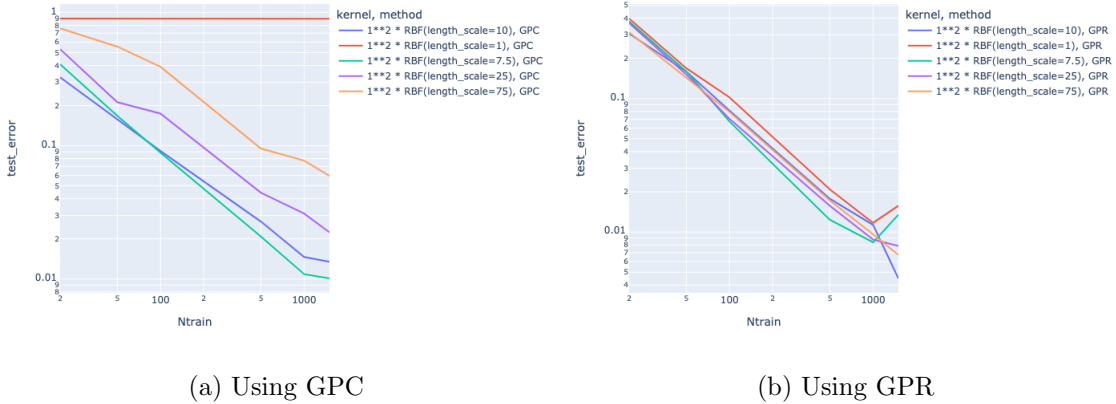


(a) Using GPC  (b) Using GPR

Figure 2: Error rate on the test set using RBF with some different parameters

Since RBF gave a reasonably good error rate, the more general Matern kernel was also tested with different values of $\nu$. $\nu$ can be seen as a smoothness parameter where the larger $\nu$, the smoother the resulting function has to be. As $\nu \to \infty$ Matern()$\to$ RBF(), therefore it is reasonable that the error values of Matern() approaches the values of RBF() as $\nu$ increases. The Matern kernel was applied with different $\nu$ and compared to the corresponding RBF function, as can be seen in figure 4.

When the dot product was used as the kernel, there seemed to be very little difference depending on the sigma parameter. Some results can be observed in figure 5. There is however a big difference between the different methods. GPC (optimize = None) gives a lower error rate than GPR for most of the kernels and number of training elements. For GPR the DotProduct() behaves unexpectedly an increases with some sizes of the training
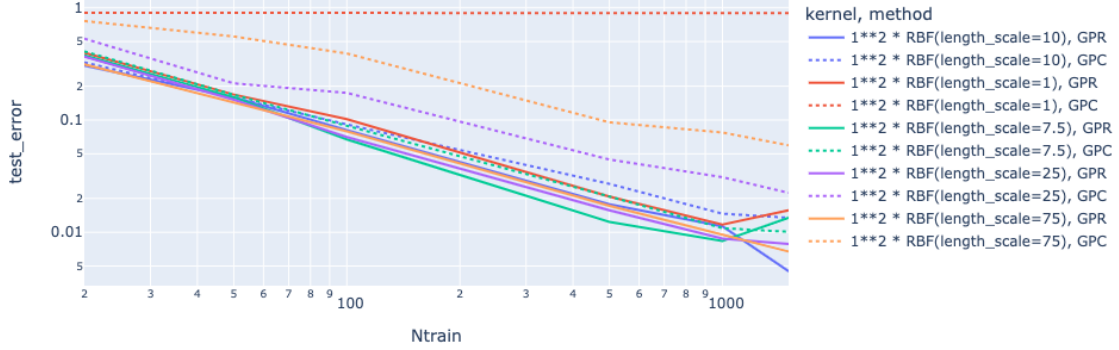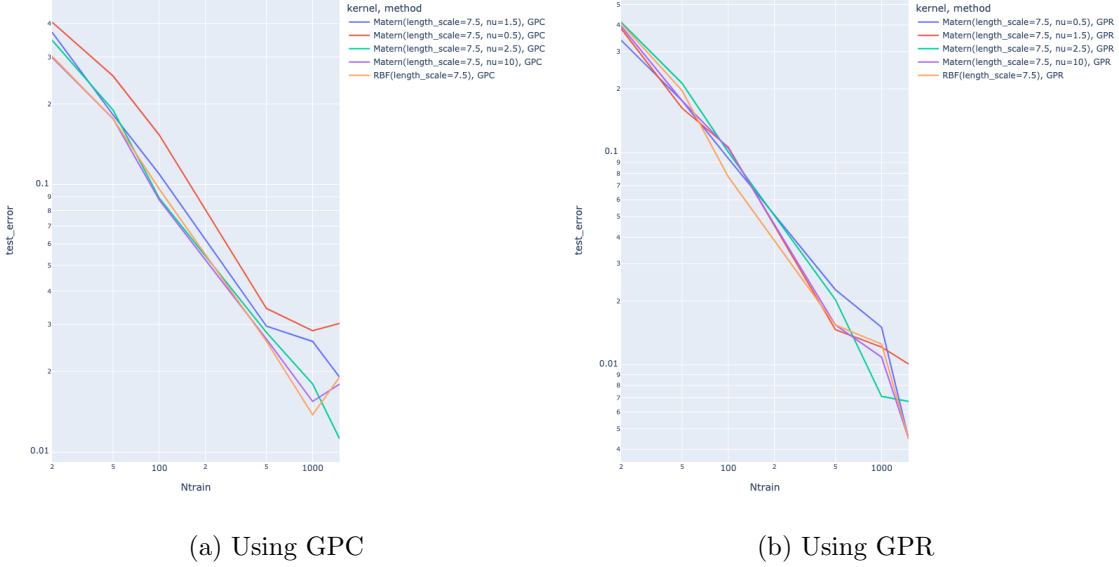
Figure 3: Error rate on RBF



(a) Using GPC



(b) Using GPR

Figure 4: Error rate on the test set using Matern($l = 7.5$) with some different values of $\nu$

set. For large sets of training data and some values $\sigma_0$, the model can not be fitted to the data since the dot product is not returning a positive definite matrix.

Different kernels (with parameters that gave a decent error rate on the test set) are compared in figure 6. Form this it can be observed that RBF ($L = 10$) gives a lower error rate than the dot product and the dot product with an additional noise term. By adding
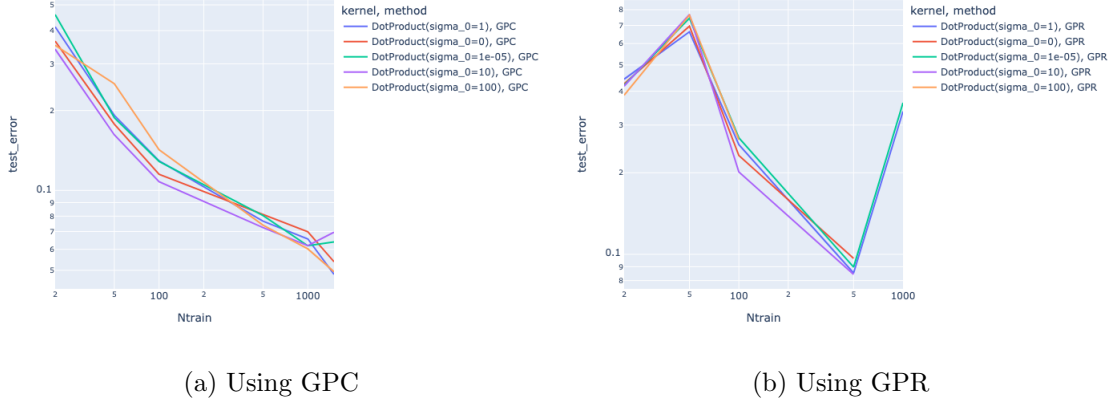
(a) Using GPC  (b) Using GPR

Figure 5: Error rate on the test set using DotProduct with some different parameters

the noise term the error is slightly lowered, which is to be expected.

When comparing GPC and GPR for these methods (figure 6), it can be observed that the GPC gives a lower error rate on the test set when using the dot product kernel, but GPR gives a lower error rate for the RBF and Matern kernels.
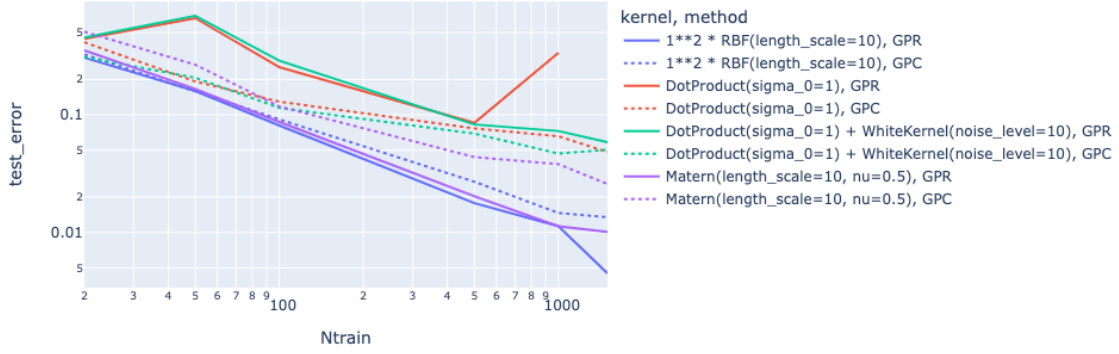


Figure 6: Error rate on test set

As described above the best test error rate was obtained by using the RBF ($L = 10$) kernel closely followed by the Matern kernels ($L = 10$) with the GaussianProcessRegressor being slightly more accurate than the GaussianProcessClassifier. Both methods are pretty computational heavy, however a high precision of classification was achieved. The computation time increased significantly with the size of the training data.

### 1.3.3 Convolutional Neural Network

In the third part of the project, a simple convolutional neural network (CNN) is used to perform the multi-classification of the digits data set using `Tensorflow` and `Keras`. The task is to classify a given image of a handwritten digit into one of 10 classes $(0, 1, ...9)$.

**Digits classification data set**

The `digits` dataset from `sklearn` consists of 1797 small square $8 \times 8$ pixel grayscale images of handwritten single digits. A quarter of data (450 examples) is held as a test data set, the other 75% of data is used to train models.
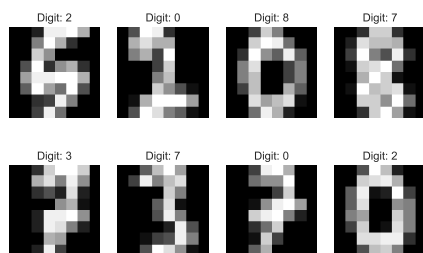


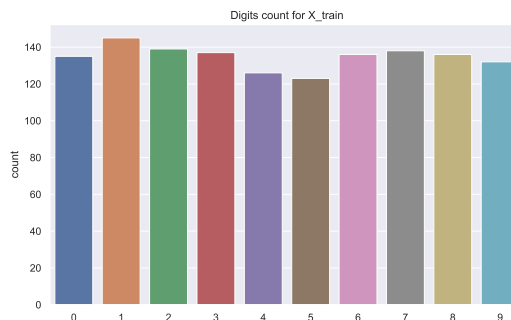Figure 7: Subset of images from `digits` data set



Figure 8: Occurrence of each label in a train data set

Figure 7 shows the first eight images in the training data set. It is important to check if the train data set is balanced, namely the frequency of digits occurrence. From Figure 8, the training set is balanced and thus it can be assumed there is an equal distribution of classes in the training data set.

Since the pixel values for each grayscale image in the data set are unsigned integers in the range from 0 to 255, normalisation of the pixel values is needed. This involves dividing the pixel values by 255 which gives values from 0 to 1.

**Model validation**

In order to estimate the performance of developed models, both 5-fold cross-validation and evaluation on a test data set generated at random are used. k-Fold cross-validation is under consideration since it results in a less unbiased estimate of the models' skill in comparison with a simple train/test split.
A sparse categorical cross-entropy loss and the accuracy as a metric are shown for implemented models.

7

**Simple NN**

To perform mini-batch gradient descent, batch_size is fixed to be 32 for all models. The number of simulation times (epoch) is set to be 100. Moreover, given that the problem is a multi-class classification and for each of 10 output nodes, it is required to assign a probability, a softmax activation function is always used in the output layer.

At first, a one-layer NN is designed as shown below. A hidden layer with 1000 nodes uses ReLU as an activation function. This model has 75,010 trainable parameters and achieves already a good fit with the accuracy of 0.913 on the test set and 0.829 when averaging over 5 folds for cross-validation. Here SGD with the learning rate of 0.01 is used as an optimizer.

```
1 inputs=keras.Input(shape=(pic_size,1))
2 fl=Flatten()(inputs)
3 hidden1=Dense(1000, activation='relu')(fl)
4 outputs=Dense(10, activation='softmax')(hidden1)
5 opt=SGD(learning_rate=0.01)
```

The loss and accuracy functions over episodes for a simple one-layer NN are shown on Figure 9. It can be noted that for such a simple neural network, the performance is already outstanding. The loss function rapidly decreases for both train and test data sets and there is no overfitting.
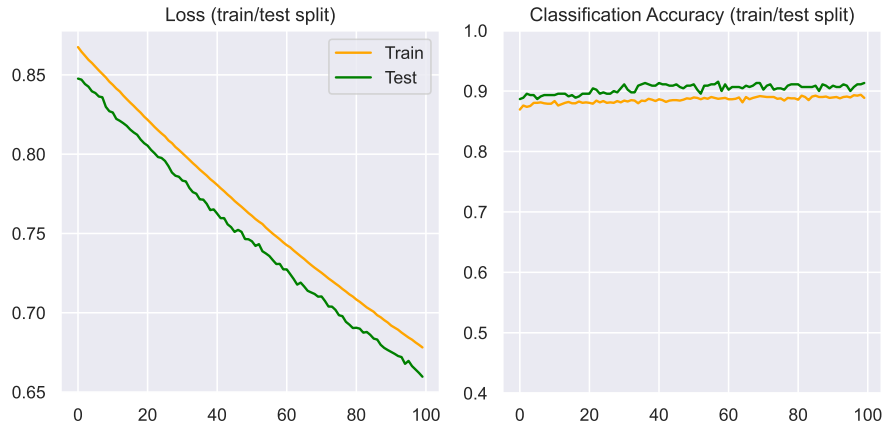


Figure 9: Loss and accuracy for a one-layer CNN, train/test split

To improve it further, a kernel for the activation function is set by the He weight scheme and the momentum method is added to the optimizer.

```
1  inputs=keras.Input(shape=(pic_size,1))
2  fl= latten()(inputs)
3  hidden1=Dense(1000, activation='relu', kernel_initializer='he_uniform')(fl)
4  outputs=Dense(10, activation='softmax')(hidden1)
5  opt=SGD(learning_rate=0.01, momentum=0.9)
```

In this case, the average accuracy score on 5-fold cross-validation reaches 0.964 and 0.982 on a test data set which is much better compared to the fisrt model having the same number of parameters. However, the loss function on the test set tends to increase over epoch and diverge from the loss function on a train test. While the training error goes down, the test loss goes up (Figure 10). Thus other models should be studied.
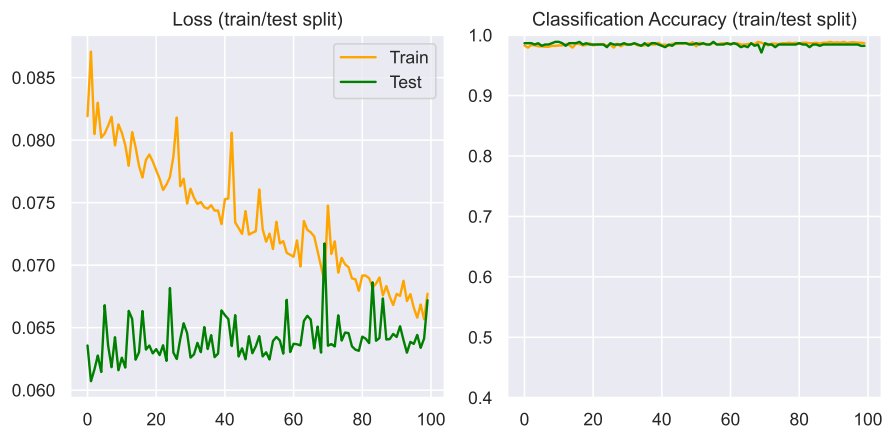


Figure 10: Loss and accuracy for an ungraded one-layer CNN, train/test split

**Multi-layer neural network**

Consider a deep neural network with four hidden layers where ReLU is used as an activation function. This NN has 621,160 training parameters and its architecture looks as follows:

```
1  inputs_deep=keras.Input(shape=(pic_size,1))
2  fl=Flatten()(inputs_deep)
3  hidden1=Dense(1000,activation='relu')(fl)
4  hidden2=Dense(500, activation='relu')(hidden1)
5  hidden3=Dense(100, activation='relu')(hidden2)
6  hidden4=Dense(50, activation='relu')(hidden3)
7  outputs_deep = Dense(10, activation = 'softmax')(hidden4)
8  opt=SGD(learning_rate=0.01)
```

9

From Figure 11, it can be concluded this model does not perform correctly and heavily overfits the given data. Since we work with image recognition, it might be a good idea to add convolutional layers and reduce the number of parameters in a model to prevent overfitting.
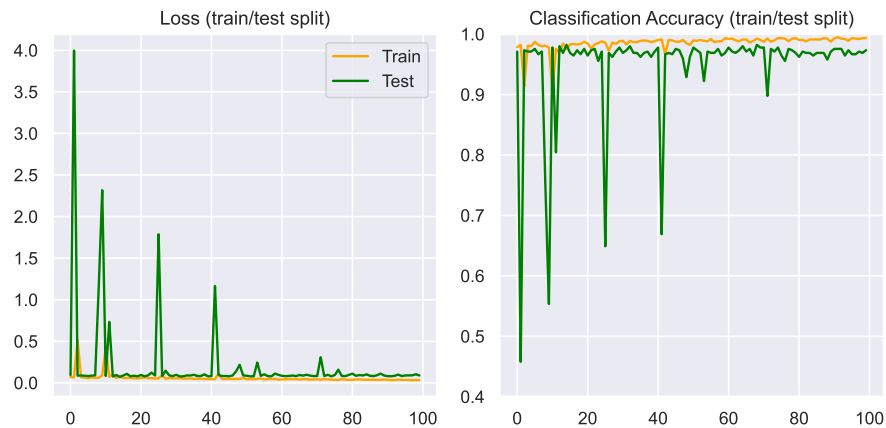


Figure 11: Loss and accuracy for a multi-layer NN, train/test split

**Convolutional neuron network**

In this section, we illustrate a CNN with a smaller number of training parameters that is enough to show better performance then before.

To do so, we add a conv2D layer where convolutional layers learn filters that are set up as parameters. To reduce the number of parameters in a model, max pooling operation for 2D spatial data is used. To prevent overfitting and regularizing deep neural networks, dropout is implemented.

Let us look at the final model we chose, its structure is shown below. The average accuracy score on 5-fold cross validation achieves 0.992 and around 1.000 on the test set. This model has 32,764 training parameters, half as many as we used in a simple NN.

```
1  inputs_cnn=Input(shape=(8,8,1))
2  cnn_1=Conv2D(filters=4, padding='same',kernel_size=3, strides=(1,1),
       activation='relu')(inputs_cnn)
3  dr_1=Dropout(0.1)(cnn_1)
4  cnn_2=Conv2D(filters=8, padding='same',kernel_size=3, strides=(1,1),
       activation='relu')(dr_1)
5  mp_1=MaxPool2D(pool_size=2)(cnn_2)
6  dr_2=Dropout(0.1)(mp_1)
7
```

```
8  cnn_3=Conv2D(filters=16, padding='same',kernel_size=3, strides=(1,1),
       activation='relu')(dr_2)
9  dr_3=Dropout(0.1)(cnn_3)
10 fl=Flatten()(dr_3)
11 hidden_1=Dense(100, activation='relu', kernel_initializer='he_uniform')(fl)
12 hidden_2=Dense(50, activation='relu', kernel_initializer='he_uniform')(
       hidden_1)
13 output_cnn=Dense(10, activation='softmax')(hidden_2)
14 opt = Adam()
```

As it can be seen from Figure 13, the test loss converges to the train loss for all 5 folders as well as the accuracy score, thus it makes the model reliable enough to use.
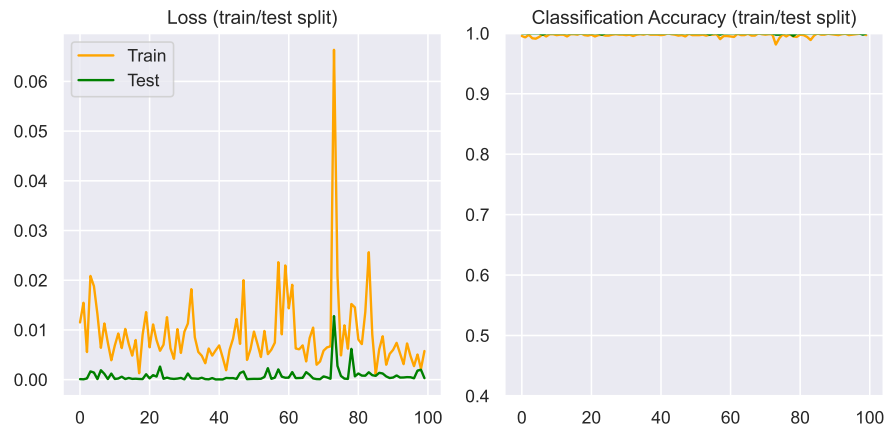


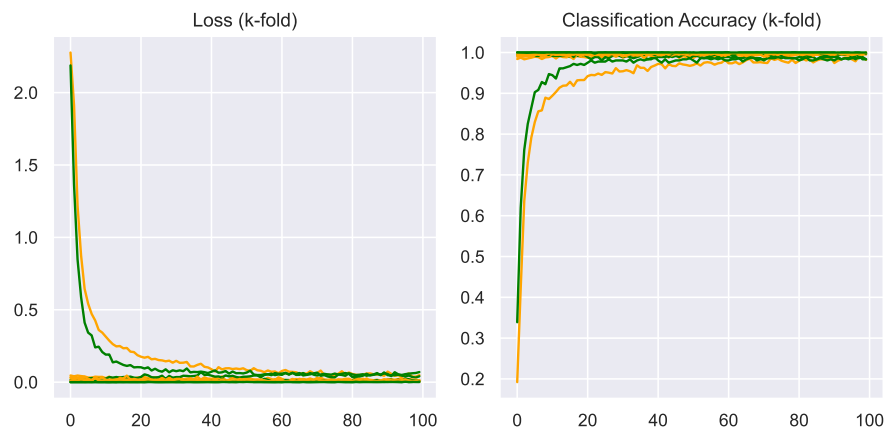Figure 12: Loss and accuracy for a chosen CNN, train/test split



Figure 13: Loss and accuracy for a chosen CNN, 5-fold CV

11

## Splitting proportion

So far we worked with the data splitted into train and test sets with proportion 3:1. Here we study how this proportion results on the accuracy and loss functions. For six different split proportions taking values from 0.1 to 0.6, the accuracy for the chosen model is plotted on Figure 14.



Figure 14: Accuracy for a chosen CNN for different proportion splits value

The minimum value of the accuracy of the test set is shown when a training set consists of 90% of the data, for other values of proportion the accuracy is relatively large. The first thought here is that training a model on a large data set while having a small one to validate the results can lead to overfitting. It should be noted that the results might be different for a data set with a higher resolution.

In the matter of the loss functions, as it was expected, having a small test set leads to overfit (large overfit gap on Figure 15 for 9:1 proportion).



Figure 15: Loss for a chosen CNN for different proportion splits value

12

## 1.4   Summary

To summarise, three different methods have been applied to classify hand written digits form the MNIST dataset.

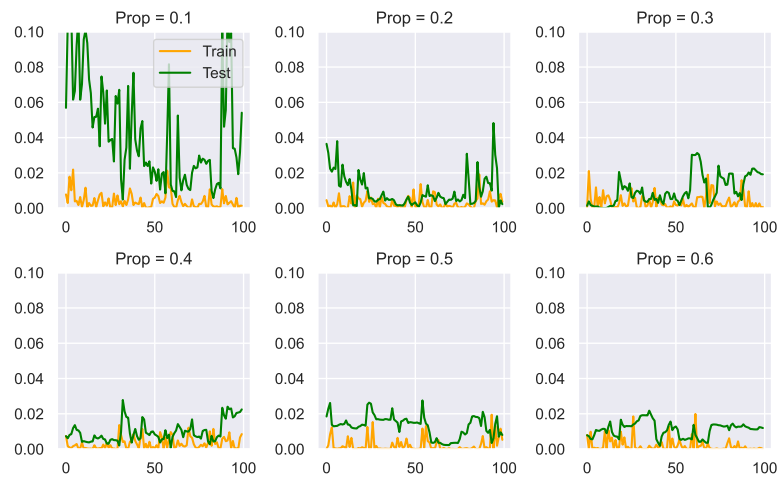The stochastic subgradient projection method is suitable for large-scale problems with decomposition techniques. It is efficient in terms of processing time and precision of classification.

The Gaussian Process Classifier achieved a high precision in the classification, however the computational time increased significantly when the size of the training data increased.

In usual practice, convolutional neural networks achieve high performance as it was shown in this example. However, there are some cons to take into account. First, determining a proper network architecture is a challenging and time consuming process. Secondly, small amount of training data can lead to poor performance and thus some neural networks are pre-trained on different data sets. Also, there is a urgent problem with unpredictability and explainability of neural networks.

## 2 Appendix

### 2.1 Stochastic subgradient classifier

#### 2.1.1 Implementation

```
1  def loss_function(Theta, x, y):
2  # Returns the loss_value vector and the index j_star where j_star is the
       argmax of loss vector
3      d, K = Theta.shape
4      loss_vector = np.maximum(np.zeros([K, 1]), (np.ones([K, 1]) + x.
       transpose().dot(Theta - np.reshape(Theta[:,y], [-1, 1])).transpose()))
5      loss_vector[y] = -1
6      j_star = np.argmax(loss_vector)
7      loss_value = loss_vector[j_star][0]
8      return loss_value, j_star
9
10 def subgradient(Theta, x, y):
11 # Computes a subgrafient of the objective empirical hinge loss
12 # one pair of (x, y) given, x of size 1, and y an integer in {0, 1, ...,
       9}.
13     g = np.zeros(Theta.shape)
14     x = np.reshape(x, [-1, 1])
15     y = int(y)
16     d, K = Theta.shape
17     loss, j_star = loss_function(Theta, x, y)
18     if loss > 0:
19         g[:, j_star] = x.transpose()
20         g[:, y] = -x.transpose()
21     return(g)
22
23 def svmsubgradient(Theta, x, y):
24 #  Returns a subgradient of the objective empirical hinge loss
25 # The inputs are Theta, of size n-by-K, where K is the number of classes,
26 # x of size n, and y an integer in {0, 1, ..., 9}.
27     G = np.zeros(Theta.shape)
28     for x_k, y_k in zip(x, y):
29         G += subgradient(Theta, x_k, y_k)
30     G /= x.shape[0]
31     return(G)
32
33 def sgd(Xtrain, ytrain, maxiter = 10, init_stepsize = 1.0, l2_radius =
       10000, alpha = 0.5):
34 # Performs maxiter iterations of projected stochastic gradient descent
35 # on the data contained in the matrix Xtrain, of size n-by-d, where n
36 # is the sample size and d is the dimension, and the label vector
37 # ytrain of integers in {0, 1, ..., 9}. Returns two d-by-10
38 # classification matrices Theta and mean_Theta, where the first is the
       final
39 # point of SGD and the second is the mean of all the iterates of SGD.
```

```
40  # Each iteration consists of choosing a random index from n and the
41  # associated data point in X, taking a subgradient step for the
42  # multiclass SVM objective, and projecting onto the Euclidean ball
43  # The stepsize is init_stepsize / sqrt(iteration).
44      K = 10
45      NN, dd = Xtrain.shape
46      print(NN)
47      Theta = np.zeros(dd*K)
48      Theta.shape = dd,K
49      mean_Theta = np.zeros(dd*K)
50      mean_Theta.shape = dd,K
51      ## YOUR CODE HERE -- IMPLEMENT PROJECTED STOCHASTIC GRADIENT DESCENT
52      for i in range(maxiter):
53          index = np.random.randint(0, NN)
54          stepsize = init_stepsize / ((i+1)**alpha)
55          Theta -= stepsize * subgradient(Theta, Xtrain[index,:], ytrain[
    index])
56          radius = np.linalg.norm(Theta.flatten())
57          if radius > l2_radius:
58              Theta = Theta / radius * l2_radius
59          mean_Theta += Theta
60      mean_Theta /= maxiter
61      return Theta, mean_Theta
```

## 2.1.2 Evaluation

```
1  #choose a seed
2  seed = 1
3  np.random.seed(seed)
4  # Load data into train set and test set
5  digits = datasets.load_digits()
6  X = digits.data
7  y = np.array(digits.target, dtype = int)
8  X,y = shuffle(X,y)
9  N,d = X.shape
10 Ntest = np.int(100)
11
12 ###### Question 1d : performance vs size of training sets
13 error_vector = np.zeros([6, 5])
14 Ntrains = [20, 50, 100, 500, 1000, 1500]
15 for i in range(len(Ntrains)):
16     Xtrain = X[0:Ntrains[i],:]
17     ytrain = y[0:Ntrains[i]]
18     Xtest = X[Ntrains[i]:N,:]
19     ytest = y[Ntrains[i]:N]
20
21     alpha = 0.5
22     l2_radius = 40.0
23     M_raw = np.sqrt(np.mean(np.sum(np.square(Xtrain))))
24     init_stepsize = l2_radius/M_raw
```

```
25      maxiter = 40000
26      for j in range(5):
27          Theta, mean_Theta = sgd(Xtrain, ytrain, maxiter, init_stepsize,
        l2_radius, alpha)
28          error_vector[i, j] = np.sum(np.not_equal(Classify(Xtest, mean_Theta
        ),ytest)/Ntest)
29
30  std = np.zeros([6, 1])
31  avg = np.zeros([6, 1])
32  for i in range(len(Ntrains)):
33      std[i] = stdev(error_vector[i])
34      avg[i] = mean(error_vector[i])
35
36  fig, ax = plt.subplots()
37  ax.errorbar(Ntrains, avg,
38              yerr=std,
39              fmt='-o')
40  ax.set_xlabel('Number of train samples')
41  ax.set_ylabel('Error')
42  ax.set_title('Errors of model depending on the number of training sets')
43  plt.show()
44
45  ###### Question 1e : Performance vs learning rate
46  Ntrain = np.int(1697)
47  Xtrain = X[0:Ntrain,:]
48  ytrain = y[0:Ntrain]
49  Xtest = X[Ntrain:N,:]
50  ytest = y[Ntrain:N]
51
52  error_alpha_vector = np.zeros([10, 5])
53
54  alphas = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1]
55  for i in range(len(alphas)):
56      alpha = alphas[i]
57      l2_radius = 40.0
58      M_raw = np.sqrt(np.mean(np.sum(np.square(Xtrain))))
59      init_stepsize = l2_radius/M_raw
60      maxiter = 40000
61      for j in range(5):
62          Theta, mean_Theta = sgd(Xtrain, ytrain, maxiter, init_stepsize,
        l2_radius, alpha)
63          error_alpha_vector[i, j] = np.sum(np.not_equal(Classify(Xtest,
        mean_Theta),ytest)/Ntest)
64
65  std = np.zeros([10, 1])
66  avg = np.zeros([10, 1])
67  for i in range(len(alphas)):
68      std[i] = stdev(error_alpha_vector[i])
69      avg[i] = mean(error_alpha_vector[i])
70
```

```python
71 fig, ax = plt.subplots()
72 ax.errorbar(alphas, avg,
73               yerr=std,
74               fmt='-o')
75 ax.set_xlabel('alpha')
76 ax.set_ylabel('Error')
77 ax.set_title('Errors of model depending on the stepsize parameter alpha')
78 plt.show()
```

## 2.2 Gaussian Process Classifier

```python
 1 """
 2 =====================================================
 3 Gaussian process classification (GPC)
 4 =====================================================
 5
 6 """
 7
 8 import numpy as np
 9 import matplotlib.pyplot as plt
10 import plotly.express as px
11 import pandas as pd
12 from sklearn import datasets
13 from sklearn.utils import shuffle
14 from sklearn.gaussian_process import GaussianProcessClassifier,
      GaussianProcessRegressor
15 from sklearn.gaussian_process.kernels import RBF, DotProduct,
      ConstantKernel, Matern, WhiteKernel, RationalQuadratic, ExpSineSquared
16 from sklearn.preprocessing import LabelBinarizer
17
18 def return_onehot(y):
19     #Note, use LabelBinarizer instead of OneHotEncoder
20     lb = LabelBinarizer()
21     lb.fit(range(0,10,1))
22     yonehot=lb.transform(y)
23     #print(np.size(yonehot, 0),np.size(yonehot, 1) )
24     return yonehot
25
26 def train_model_GPClass(Xtrain, ytrain, kernel, Xtest, ytest):
27     #gpc_rbf = GaussianProcessClassifier(kernel=kernel).fit(Xtrain, ytrain)
28     gpc_rbf = GaussianProcessClassifier(kernel=kernel, optimizer = None).
      fit(Xtrain, ytrain)
29     #GaussianProcessClassifier performs hyperparameter estimation, this
      means that the the value specified above may not be the final
      hyperparameters
30     #If you dont want it to do hyperparameter optimization set optimizer =
      None like this GaussianProcessClassifier(kernel=kernel,optimizer = None)
      .fit(Xtrain, ytrain)
31
32
```

```python
33     yp_train = gpc_rbf.predict(Xtrain)
34     train_error_rate = np.mean(np.not_equal(yp_train,ytrain))
35
36     yp_test = gpc_rbf.predict(Xtest)
37     test_error_rate = np.mean(np.not_equal(yp_test,ytest))
38
39     kernel_params = gpc_rbf.kernel_.get_params()['kernels']
40
41     return train_error_rate, test_error_rate, str(kernel_params[0])
42
43 def train_model_GPRegr(Xtrain, ytrain, kernel, Xtest, ytest):
44     #transform to onehot
45     yonehot = return_onehot(ytrain)
46
47     gpr = GaussianProcessRegressor(kernel=kernel).fit(Xtrain, yonehot)
48
49     yp_train = gpr.predict(Xtrain)
50     yp_train_indx = np.argmax(yp_train, axis = 1) #Find max of each row, i.
    e find class
51
52     train_error_rate = np.mean(np.not_equal(yp_train_indx,ytrain))
53
54     yp_test = gpr.predict(Xtest)
55     yp_test_indx = np.argmax(yp_test, axis = 1) #Find class
56     test_error_rate = np.mean(np.not_equal(yp_test_indx,ytest))
57
58     #print(gpr.get_params())
59     kernel_param = gpr.get_params()['kernel']
60
61     return train_error_rate, test_error_rate, str(kernel_param)
62
63 def calculate_new_values(error_df, methods, training_sizes, kernels, iters
    = 3):
64
65     # import data
66     digits = datasets.load_digits()
67
68     X = digits.data
69     y = np.array(digits.target, dtype = int)
70     N,d = X.shape
71
72     for method in methods:
73         for kernel_ind in range(len(kernels)):
74             kernel = kernels[kernel_ind]
75
76             for size_ind in range(len(training_sizes)):
77                 Ntrain = training_sizes[size_ind]
78                 temporary_df = error_df[(error_df['method'] == method) & (
    error_df['Ntrain'] == Ntrain) &(error_df['kernel'] == str(kernel) )]
79                 #Check that the value is not already calculated
```

```python
                    if temporary_df.empty:

                        print('method = %s, Ntrain = %d, kernel = %s has not
    been prev. calc.' %(method, Ntrain,str(kernel)) )
                        try:
                            mean_test_error = 0
                            mean_train_error = 0
                            for i in range(iters): #do the errorcalculation
    multiple times, take average
                                X,y = shuffle(X,y)
                                Xtrain = X[0:Ntrain-1,:]
                                ytrain = y[0:Ntrain-1]
                                Xtest = X[Ntrain:N,:]
                                ytest = y[Ntrain:N]

                                if method == 'GPC':
                                    train_error_rate, test_error_rate,
    kernel_params = train_model_GPClass(Xtrain, ytrain, kernel, Xtest, ytest
    ) #a)
                                elif method == 'GPR':
                                    train_error_rate, test_error_rate,
    kernel_params = train_model_GPRegr(Xtrain, ytrain, kernel, Xtest, ytest)
     #b)
                                else:
                                    raise ValueError("%s is not an acceptable
    method" %method)

                                mean_train_error += train_error_rate / iters
                                mean_test_error += test_error_rate/iters
                        #      print(test_error_rate)
                        # print(mean_test_error)
                            print("Success when computing model")
                            row_df = {'method': method, 'Ntrain': Ntrain, '
    kernel':kernel_params, 'train_error': mean_train_error , 'test_error':
    mean_test_error}
                            error_df = error_df.append(row_df, ignore_index=
    True)

                        except:
                            print("Error when computing model")
                    else:
                        pass

        return error_df

def plot(error_df, methods, training_sizes, kernels):
    #Get the subdataframe to be plotted
    kernel_names = [str(kernels[i]) for i in range(len(kernels))]
    sub_df = error_df[(error_df['method'].isin(methods)) & (error_df['
    Ntrain'].isin(training_sizes)) &(error_df['kernel'].isin (kernel_names))
```

```python
      ]
119
120     #plot it
121     fig = px.line(sub_df, x = 'Ntrain', y = 'test_error', color = 'kernel',
         line_dash = 'method', log_x = True, log_y = True)
122
123     img_name = "images/" + "".join(methods) + str(kernels[0]) + ".png"
124     fig.write_image(img_name)
125     fig.show()
126     return
127
128 def main():
129     #choose a seed
130     seed = 1
131     np.random.seed(seed)
132
133     training_sizes = [20, 50, 100, 500, 1000, 1500]
134     #training_sizes = [20, 50, 100, 500]
135
136     ##SPECIFY WHICH KERNELS TO USE
137     kernels = []
138     #kernels = [DotProduct(1.0), 1.0 * RBF([10]),  DotProduct() +
         WhiteKernel(10), Matern(7.5, nu = 0.5)];
139     #kernels = [DotProduct() + ConstantKernel() +WhiteKernel(10), Matern(),
         RationalQuadratic(), ExpSineSquared()]
140
141     params = [1,7.5,10,25,75]
142     nu_vals = [ 0.5, 1.5, 2.5, 10]
143     dot_params = [0,1e-5,1,10,100]
144     for param in params:
145         kernels.append(1.0*RBF(param))
146
147     # for nu_val in nu_vals:
148     #     kernels.append(Matern(7.5, nu = nu_val))
149     # kernels.append(RBF(7.5))
150     #
151     # for param in dot_params:
152     # #    kernels.append(DotProduct(param))
153     #     kernels.append(DotProduct(10) + WhiteKernel(param))
154
155     ##SPECIFY WHICH METHODS TO USE
156     #methods = ['GPR']
157     methods = ['GPR', 'GPC']
158
159     ##USE DATAFRAME TO SAVE CALCULATIONS
160     #error_df = pd.DataFrame(columns=['method', 'Ntrain', 'kernel', '
         train_error', 'test_error']) #Use if wanted to start from beginning
161     error_df = pd.read_csv('error.csv', index_col = 0)
162
163     error_df = calculate_new_values(error_df, methods, training_sizes,
```

```
        kernels)
164     error_df.to_csv("error.csv")
165     # print(error_df)
166
167     plot(error_df, methods, training_sizes, kernels)
168
169     return
170
171 main()
```

## 2.3 Evaluation functions for CNN

### 2.3.1 Implementation

```
 1 def evaluate(model, epochs, batch_size, X, y, n_folds=5):
 2   ''' Fuction for a model evaluation using 5-fold cross-validation '''
 3   scores, histories = [], []
 4   # Prepare cross-validation
 5   kfold = KFold(n_folds, shuffle=True, random_state=1)
 6   # Enumerate splits
 7   for train_ix, test_ix in kfold.split(X):
 8     # Select rows for train and test
 9     X_train, y_train, X_test, y_test = X[train_ix], y[train_ix], X[test_ix
      ], y[test_ix]
10     # Fit model
11     history = model.fit(X_train, y_train, epochs=epochs, batch_size=
      batch_size, validation_data=(X_test, y_test), verbose = 0)
12     # Evaluate model
13     _, acc = model.evaluate(X_test, y_test, verbose=0)
14     # Store scores
15     scores.append(acc)
16     histories.append(history)
17   return scores, histories
18
19
20 def diagnostics(fold_histories, histories):
21   fig, (ax1, ax2) = plt.subplots(1, 2, figsize = (8,4))
22   ''' Diagnostics of the model learning behavior during training and the
      estimation of the model performance '''
23   # Plotting a loss-finction
24   ax1.set_title('Loss (train/test split)')
25   ax1.plot(histories.history['loss'], color='orange', label='Train')
26   ax1.plot(histories.history['val_loss'], color='green', label='Test')
27   ax1.legend()
28   # Plotting accuracy
29   ax2.set_title('Classification Accuracy (train/test split)')
30   ax2.plot(histories.history['accuracy'], color='orange')
31   ax2.plot(histories.history['val_accuracy'], color='green')
32   ax2.set_ylim(0.4, 1.)
33   fig.tight_layout()
34   # fig.savefig('finaltest.pdf')
```

```
35    fig.show()
36
37    plt.figure(figsize=(8,8))
38    for i in range(len(fold_histories)):
39      # Plotting a loss-finction
40      plt.subplot(2, 2, 1)
41      plt.title('Loss (k-fold)')
42      plt.plot(fold_histories[i].history['loss'], color='orange', label='
        Train')
43      plt.plot(fold_histories[i].history['val_loss'], color='green', label='
        Test')
44      # Plotting accuracy
45      plt.subplot(2, 2, 2)
46      plt.title('Classification Accuracy (k-fold)')
47      plt.plot(fold_histories[i].history['accuracy'], color='orange')
48      plt.plot(fold_histories[i].history['val_accuracy'], color='green')
49    plt.tight_layout()
50    # fig.savefig('finaltestfold.pdf')
51    plt.show()
52
53    return
54
55  def scoreTest(kFold_accScores,acc_score):
56    print('Average Accuracy Score on k-Fold Cross Validation:  %.3f' % (np.
        mean(kFold_accScores)))
57    print('Accuracy Score on k-Fold Cross Validation:')
58    for i in range(len(kFold_accScores)):
59      print(f'{i+1}-Fold:  %.3f' % (kFold_accScores[i]))
60    print('Accuracy Score on a Test Dataset:  %.3f' % (acc_score))
61    return
```