

ECE 532 Final Project

The MNIST is a “Hello World” project in machine learning (ML). At the same time, it is one of the most appropriate databases for machine learning beginners who want to practice their machine learning techniques. This MNIST dataset contains real-world hand-written images which makes it become popular for machine learning projects. During these years, more and more machine learning algorithms researchers improve this MNIST project in order to help ML beginners. In this project, we study several ML solvers using **sk.learn** and applied those on the MNIST dataset. These ML solvers include: **linear ridge regression**, **nearest neighbor classification** and **kernel based support vector machine(SVM)**.

1 Downloading the dataset

1.1 Downloading the MNIST

The MNIST is a large database of handwritten digits. There are 60,000 images of 28x28 pixel in the training set, and the testing set is composed of 10,000 patterns of 28x28 pixel. Each pixel forms a feature, so one image has 784 features, and the value range of each pixel is [0, 255]. The training label set has 60,000 numbers, and the testing label set contains 10,000 results. Each label is an actual number between 0-9. We can use `mnist = fetch_openml("mnist_784")` to download the MNIST.

```
1 #1 Downloading the Data (MNIST)
2 import numpy as np
3 from sklearn.datasets import fetch_openml
4 mnist = fetch_openml("mnist_784")
```

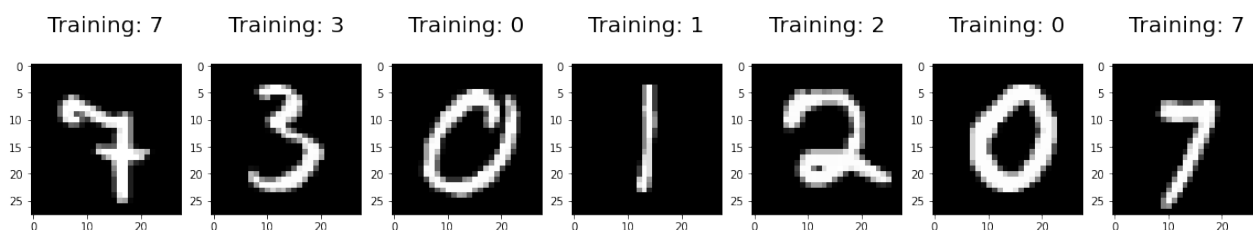
1.2 Splitting the data into Training and Test Sets

Using `train_test_split` separates the MNIST into the training subset and the test subset. The training set has 60,000 images and corresponding labels. The test set involves 10,000 images and corresponding labels.

```
1 #2 Splitting the MNIST
2 from sklearn.model_selection import train_test_split
3 train_img, test_img, train_lbl, test_lbl = train_test_split(mnist.data,
4 mnist.target, test_size=1/7.0, random_state=0)
5 train_lbl = np.asarray(train_lbl, 'float64')
6 test_lbl = np.asarray(test_lbl, 'float64')
7 train_img.shape
8 test_img.shape
```

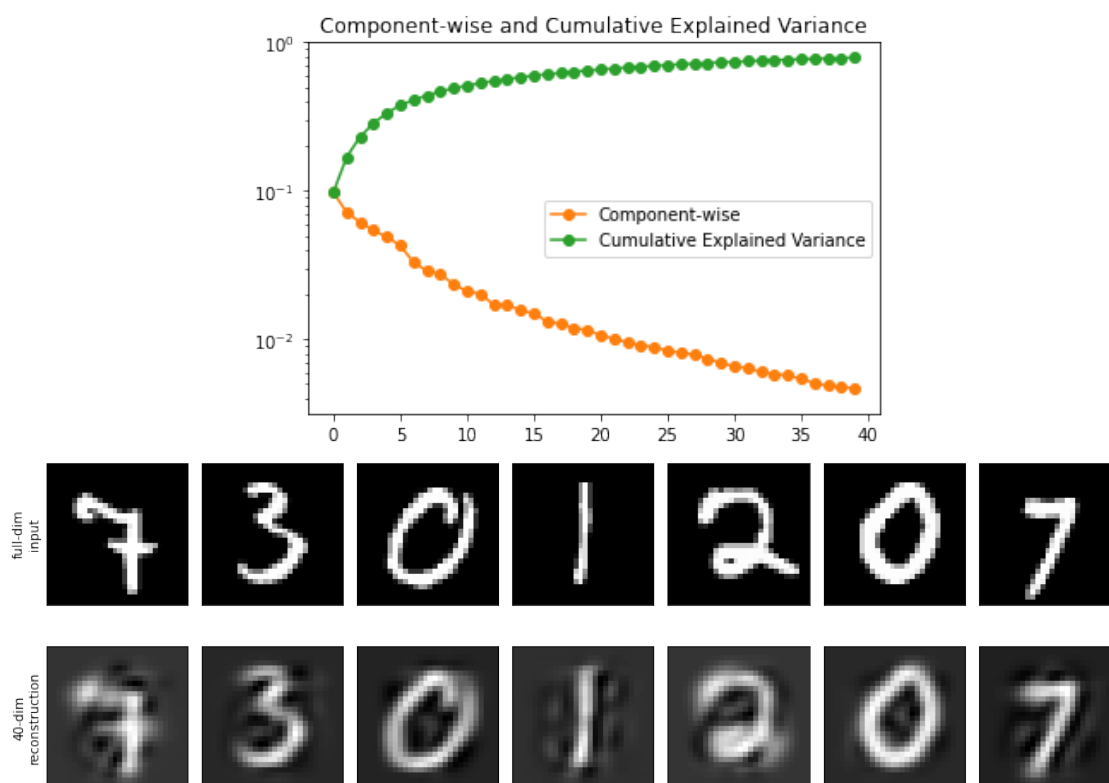
1.3 Training data display

```
1 # Data display
2 import matplotlib.pyplot as plt
3 plt.figure(figsize=(20,4))
4 for index, (image, label) in enumerate(zip(train_img[0:7], train_lbl[0:7])):
5     plt.subplot(1, 7, index + 1)
6     plt.imshow(np.reshape(image, (28,28)), cmap=plt.cm.gray)
7     plt.title('Training: %i\n' % label, fontsize = 20)
```



2 Dimensionality reduction

The MNIST has a huge computation in a learning process, so we need to use principal component analysis (PCA) to reduce its dimensionality. For PCA, the higher cumulative which explained variance of principal components is chosen, the more information in the original features can be retained. we choose the first 40 components as the new features that collects 78.7% of variance. Also, we can directly find the difference between the original image and the reconstruction image using the first 40 components.



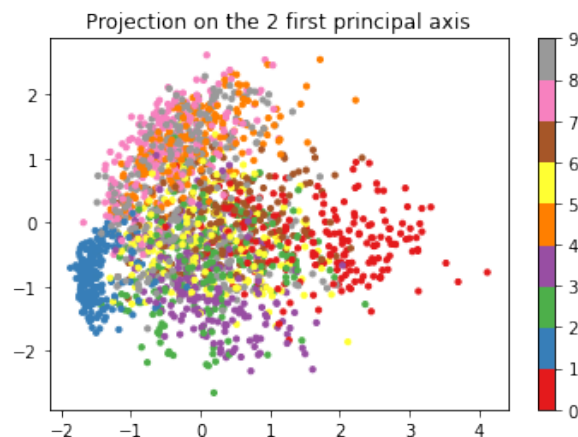
```
1 # PCA with the first 40 components
2 from sklearn.decomposition import PCA
```

```

3  from time import time
4  n_components = 40
5  t0 = time()
6  pca = PCA(n_components=n_components, svd_solver='randomized',
7           whiten=True).fit(train_img)
8  print("done in %0.3fs" % (time() - t0))
9  train_img_pca = pca.transform(train_img)
10
11 # Component-wise and Cumulative Explained Variance
12 plt.hist(pca.explained_variance_ratio_, bins=n_components, log=True)
13 sum_var = pca.explained_variance_ratio_.sum()
14 plt.plot(range(40), pca.explained_variance_ratio_, 'o-', label='Component-
15 wise')
16 plt.plot(range(40), np.cumsum(pca.explained_variance_ratio_), 'o-',
17 label='Cumulative Explained Variance')
18 plt.title("Component-wise and Cumulative Explained Variance")
19 plt.legend()
20 plt.ylim(0, 1)
21 print(sum_var)
22
23 # Reconstruct the original image using the new features.
24 projected_pca = pca.inverse_transform(train_img_pca)
25 fig, ax = plt.subplots(2, 7, figsize=(15, 5),
26 subplot_kw={'xticks':[], 'yticks':[]},
27 gridspec_kw=dict(hspace=0.1, wspace=0.1))
28
29 for i in range(7):
30     ax[0, i].imshow(train_img[i].reshape(28, 28), cmap=plt.cm.gray)
31     ax[1, i].imshow(projected_pca[i].reshape(28, 28), cmap=plt.cm.gray)
32
33 ax[0, 0].set_ylabel('full-dim\ninput')
34 ax[1, 0].set_ylabel('40-dim\nreconstruction');

```

The visualization of data using PCA that shows that the multi-class classification is in the 2-dimensions.

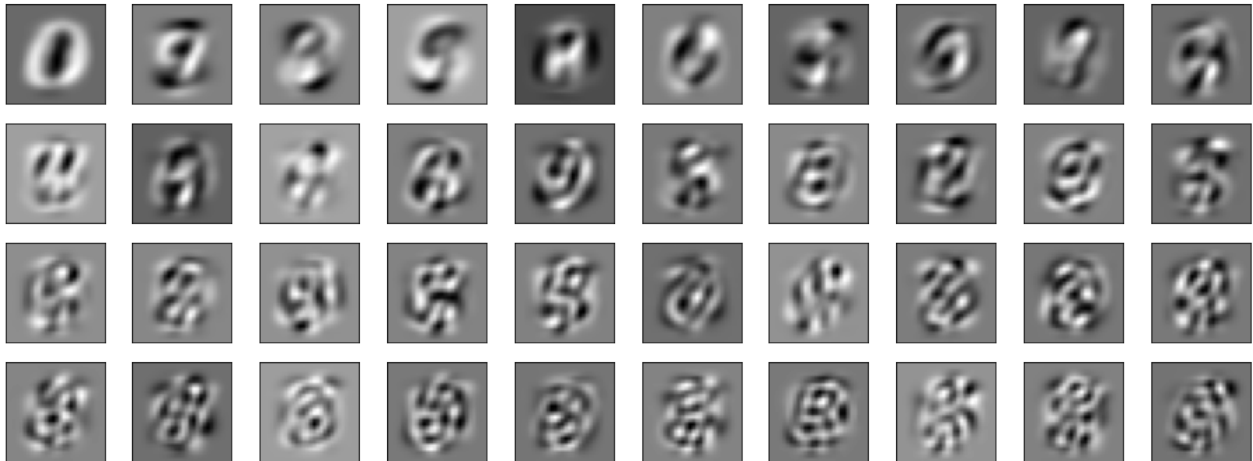


```

1 # Visualization of the MNIST dataset in 2D
2 plt.scatter(train_img_pca[0:2000, 0], train_img_pca[0:2000, 1],
3             c=train_lbl[0:2000], s=10, cmap='Set1')
4 plt.colorbar()
5 plt.title('Projection on the 2 first principal axis')

```

The components are ordered by their importance from top-left to bottom-right in the following figure. We see that the first few components seem to depict the outline, and the remaining components highlight more details.



```

1 # The ordered components display
2 fig = plt.figure(figsize=(16, 6))
3 for i in range(40):
4     ax = fig.add_subplot(4, 10, i + 1, xticks=[], yticks=[])
5     ax.imshow(pca.components_[i].reshape(np.reshape(train_img[0],
6             (28,28)).shape), cmap=plt.cm.gray)

```

3 Three algorithms and results

3.1 Linear ridge regression classifier

Ridge regression imposing a penalty on the size of the coefficients based on the ordinary Least Squares problem to control the amount of shrinkage. The ridge coefficients minimize a penalized residual sum of squares. Its cost function can be expressed as

$$\min_w ||Xw - y||_2^2 + \alpha ||w||_2^2, \quad \alpha \geq 0$$

The first term is the loss function of least squares, and the second term is a penalty term for the ridge regression. The key parameter α is a trade off between variance and bias. Thus, it is necessary to use cross-validation to choose an appropriate α . The Ridge regressor has a classifier variant: *RidgeClassifier*, which can achieve the multiclass classification. Besides, we need to use cross-validation to determine the best α . We set the range of α from 10^{-7} to 10^6 . Finally, we find that indicates Various α have the similar performance except for $\alpha=10^6$.

Table 1: Performance summary of different α

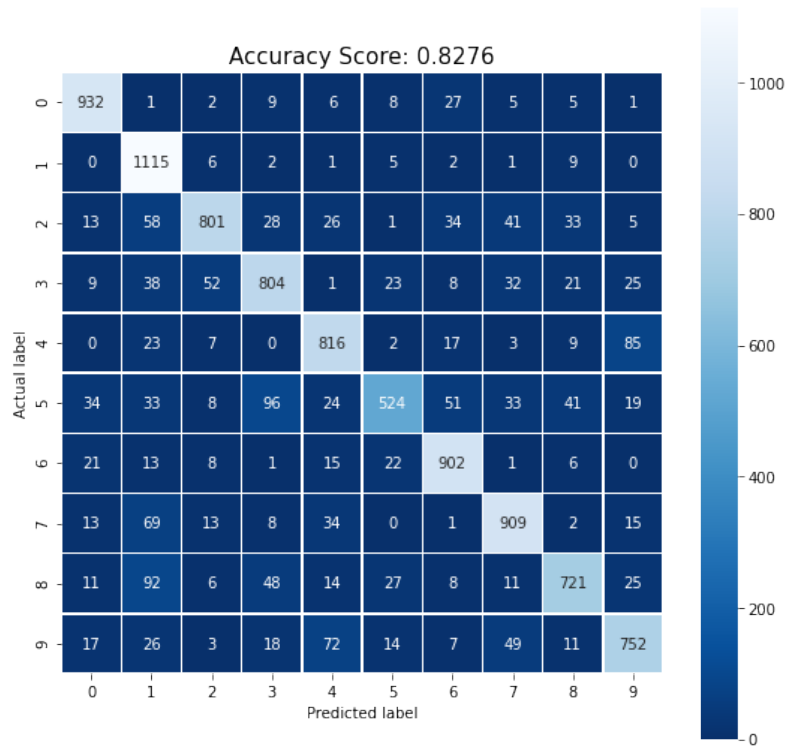
| α | 10^{-7} | 10^{-5} | 10^{-3} | 10^{-1} | 10^2 | 10^4 | 10^6 |
|----------|---------------|-----------|-----------|-----------|--------|--------|--------|
| Score | 0.8346 | 0.8346 | 0.8346 | 0.8346 | 0.8346 | 0.8341 | 0.7447 |

```

1  # Learning using the best parameters and prediction_ridge regression
2  from sklearn.linear_model import RidgeClassifier
3  from sklearn.model_selection import GridSearchCV
4  param_grid = [
5      {'alpha':[1e-7, 4e-7, 8e-7, 1e-6, 4e-6, 8e-6, 1e-5, 4e-5, 8e-5, 1e-4,
6          4e-4, 8e-4, 1e-4, 1e-2, 1e-1, 1e2, 1e4, 1e6]}
7  ]
8  t0 = time()
9  reg_clsf = RidgeClassifier()
10 reg_grid_clsf = GridSearchCV(reg_clsf,param_grid,n_jobs=1, verbose=2,
11     cv=10)
12 reg_grid_clsf.fit(train_img_pca, train_lbl)
13 print("done in %0.3fs" % (time() - t0))
14 reg_classifier = reg_grid_clsf.best_estimator_
15 reg_params = reg_grid_clsf.best_params_
16 # results of cross-validation
17 pd.pivot_table(pd.DataFrame(reg_grid_clsf.cv_results_),
18     values='mean_test_score', index='param_alpha')
19 # prediction
20 reg_predictions = reg_classifier.predict(pca.transform(test_img))
21 # evaluation
22 reg_score = reg_classifier.score(pca.transform(test_img), test_lbl)
23 print(reg_score)

```

Using the best α of 10^{-7} fit the model on the entire training data(60000), and predict new labels on the entire testing data, then score this method. The final score is **0.8276**.



```

1  # Confusion matrix
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  from sklearn import metrics
5  cm = metrics.confusion_matrix(test_lbl, reg_predictions)
6  print(cm)
7
8  plt.figure(figsize=(9,9))
9  sns.heatmap(cm, annot=True, fmt="1d", linewidths=.5, square = True, cmap =
    'Blues_r');
10 plt.ylabel('Actual label');
11 plt.xlabel('Predicted label');
12 all_sample_title = 'Accuracy Score: {0}'.format(reg_score)
13 plt.title(all_sample_title, size = 15);

```

3.2 KNN classifier

Nearest neighbor classification is a non-linear instance-based learning classifier. It simply stores instances from the training data and compare distances between unknown data with the training data. Then based on distances, it predicts labels for the unknown data.

KNN requires a total of X training data points and M classes, it predicts an unobserved training point X_{new} as the mean of the closes k neighbours to X_{new} . Usually, the standard Euclidean metric $d(x, x_i) = \sqrt{\sum_{i=1}^n (x - x_i)^2}$ is used as the the distance matrix.

$$\hat{y} = \frac{1}{k} \sum_{x_i \in X} d(x_i, x)$$

There is a nearest neighbor classifier named *KNeighborsClassifier* that implements learning based on a k nearest neighbors of the new data. When k is too small, the model becomes susceptible to noise and outlier data points; however, if k is too large, it poses a risk over-smoothing the classification results and increasing bias. It need to learn the training data using different values k, then chooses the best according to their performance of cross-validation.

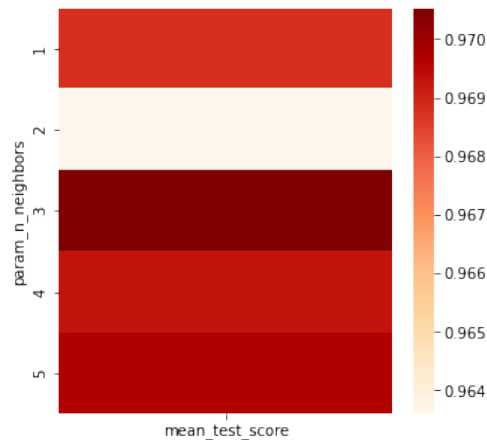


Table 2: Performance summary of different k

| k | 1 | 2 | 3 | 4 | 5 |
|-------|--------|--------|---------------|--------|--------|
| Score | 0.9688 | 0.9636 | 0.9725 | 0.9693 | 0.9697 |

```

1  # Learning using the best parameters and prediction_KNN
2  from sklearn import neighbors
3  param_grid = [
4      {'n_neighbors':[1, 2, 3, 4, 5]}
5  ]
6  t0 = time()
7  knn_clsif = neighbors.KNeighborsClassifier()
8  knn_grid_clsif = GridSearchCV(knn_clsif,param_grid,n_jobs=1, verbose=2,
9                                cv=10)
10 knn_grid_clsif.fit(train_img_pca, train_lbl)
11 print("done in %0.3fs" % (time() - t0))
12 knn_classifier = knn_grid_clsif.best_estimator_
13 knn_params = knn_grid_clsif.best_params_
14 # results of cross-validation
15 pd.pivot_table(pd.DataFrame(knn_grid_clsif.cv_results_),
16                 values='mean_test_score', index='param_n_neighbors')
17
18 # heatmap to display the results of cross-validation.
19 import seaborn as sns
20 import pandas as pd
21 pvt = pd.pivot_table(pd.DataFrame(knn_grid_clsif.cv_results_),
22                       values='mean_test_score', index='param_n_neighbors')
23 plt.figure(figsize=(5, 5))
24 sns.heatmap(pvt,cmap = 'OrRd')

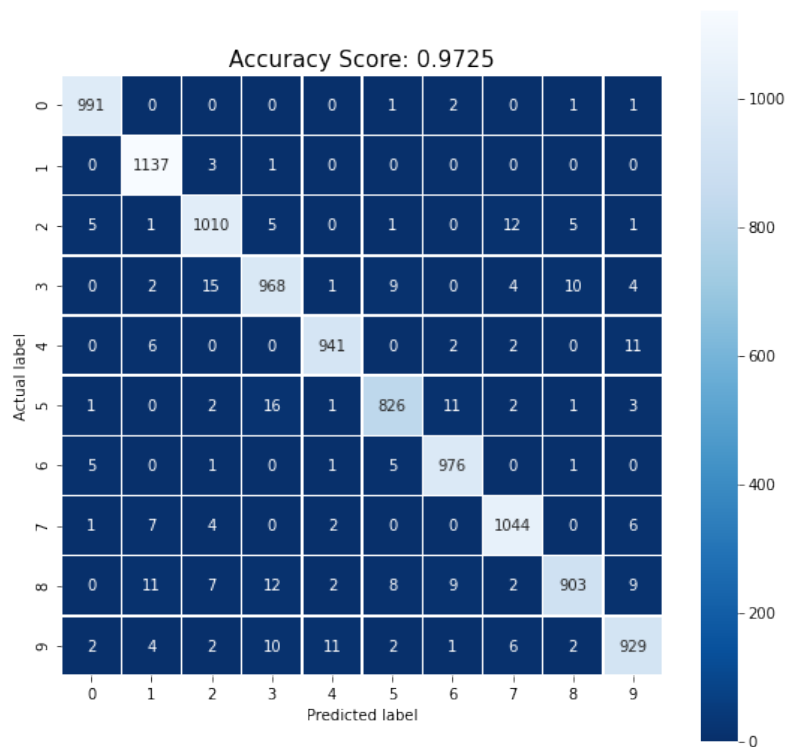
```

```

22
23 # prediction
24 knn_predictions = knn_classifier.predict(pca.transform(test_img))
25 # evaluation
26 knn_score = knn_classifier.score(pca.transform(test_img), test_lbl)
27 print(knn_score)

```

Using the best k of **3** fit the model on the entire training data(60000), and predict new labels on the entire testing data. The final score is **0.9725**.



```

1 # confusion matrix
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn import metrics
5 cm = metrics.confusion_matrix(test_lbl, knn_predictions)
6 print(cm)
7
8 plt.figure(figsize=(9,9))
9 sns.heatmap(cm, annot=True, fmt="1d", linewidths=.5, square = True, cmap =
  'Blues_r');
10 plt.ylabel('Actual label');
11 plt.xlabel('Predicted label');
12 all_sample_title = 'Accuracy Score: {0}'.format(knn_score)
13 plt.title(all_sample_title, size = 15);

```

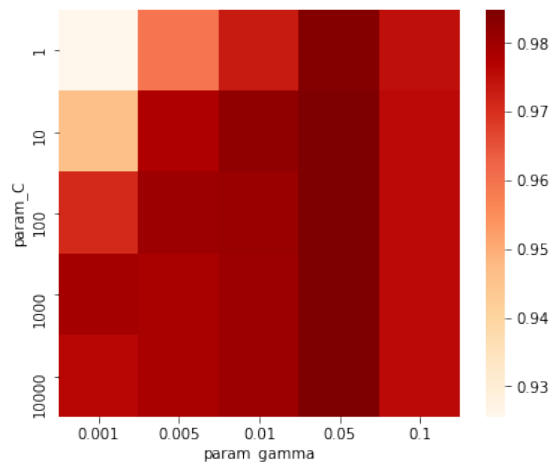

3.3 Kernel based SVM

Support Vector Machine is capable of performing multi-class classification on a large dataset, which maximizes the margin around the separating hyperplane and only depends on the support vectors. Kernel based support vector machines can be written as following. These kernel functions

has various forms, such as polynomial $K(u, v) = (u^T v + 1)^q$ or Gaussian $K(u, v) = e^{-\frac{\|u-v\|_2^2}{2\sigma^2}}$.

$$\min_{\alpha} (1 - y^i \sum_{j=1}^N \alpha_j K(x^i x^j))_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(x^i x^j)$$

The first term is the hinge loss function and the second term is the misclassification penalty. The penalty receives increasing emphasis as λ increases. Gamma is the Kernel coefficient. We used cross-validation to find the best parameters. There is a classifier named *sklearn.svm.SVC* that implements learning based on support vectors.



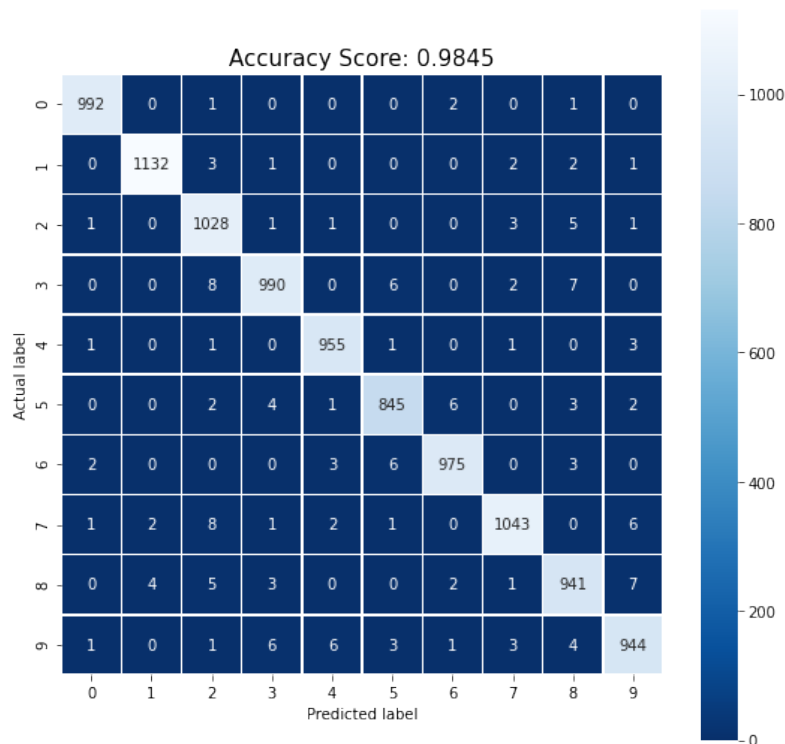
```
1 # Learning using the best parameters and prediction_SVM
2 from sklearn import svm
3 from sklearn.model_selection import GridSearchCV
4 param_grid = [
5     # {'C': [10, 100], 'kernel': ['linear']},
6     {'C': [1, 10, 100, 1000, 10000], 'gamma': [0.1, 0.05, 0.01, 0.005,
7         0.001], 'kernel': ['rbf']}
8 ]
9 t0 = time()
10 svm_cls = svm.SVC()
11 grid_cls = GridSearchCV(svm_cls, param_grid, n_jobs=1, verbose=2, cv=10)
12 grid_cls.fit(train_img_pca, train_lbl)
13 print("done in %0.3fs" % (time() - t0))
14 classifier = grid_cls.best_estimator_
15 params = grid_cls.best_params_
16 # results of cross-validation
17 pd.pivot_table(pd.DataFrame(grid_cls.cv_results_),
18     values='mean_test_score', index='param_C', columns='param_gamma')
19
20 # heatmap to display the results of cross-validation.
21 import seaborn as sns
22 import pandas as pd
```

```

21 pvt = pd.pivot_table(pd.DataFrame(grid_clsfc.cv_results_),
22 values='mean_test_score', index='param_C', columns='param_gamma')
23 plt.figure(figsize=(6,5))
24 sns.heatmap(pvt, cmap = 'OrRd')
25
26 # prediction
27 predictions = classifier.predict(pca.transform(test_img))
28 # evaluation
29 score = classifier.score(pca.transform(test_img), test_lbl)
30 print(score)

```

Using the best C of **100** and gamma of **0.05** fit the model, and get the final score is **0.9845**.



```

1 # confusion matrix
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn import metrics
5 cm = metrics.confusion_matrix(test_lbl, predictions)
6 print(cm)
7
8 plt.figure(figsize=(9,9))
9 sns.heatmap(cm, annot=True, fmt="1d", linewidths=.5, square = True, cmap =
10 'Blues_r');
11 plt.ylabel('Actual label');
12 plt.xlabel('Predicted label');
13 all_sample_title = 'Accuracy Score: {0}'.format(score)
14 plt.title(all_sample_title, size = 15);

```

4 Discussion

The linear ridge regression classifier is time-consuming, but when the dataset is non-linear, the limitation of the linear classifier results in poor performance because it is unable to capture more complex patterns. For KNN, it has no training period which means adding new data will not impact the accuracy of the algorithm; however, KNN doesn't work well with high dimensional data because high dimensionality increases the difficulty in calculating the distance in each dimension. In this dataset, the score of performance is 0.97 when I used the first 40 components, but this score is 0.95 when I used the first 100 components. Besides, KNN is sensitive to noise in the dataset because it has no penalty. Last, SVM is more effective in high dimensional spaces due to the tricky kernel; however, it can be computationally expensive.