

Amin Fadaee

Machine Learning HW4

Python 3.5, Libraries: LIBSVM, Numpy, Matplotlib, Pandas, Math, random, sklearn, Documentation: HTML, Bootstrap 4, jQuery, Graphic: Xkcd

Problem 1

a. Here is a plot depicting the data:



b. To transform the data into a 1-D linearly separable space, we can use the following transformation function:

$$f(x_1, x_2) = x_1^2 + x_2^2$$

The reason for choosing this transformation is the circular shape of the data in each class and the fact that the data in different classes have a distinguishable distance from the origin. Although for simplicity we do not perform the square operation but this will not cause any harm.

Here is the transformed data:

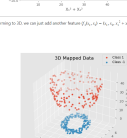


As can be seen, the data is now perfectly separable.

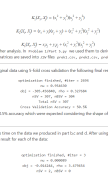
c. While preserving the features we created in the previous part, we can add another feature that simply sums the two features:

$$f(x_1, x_2) = (x_1 + x_2, x_1^2 + x_2^2)$$

Doing so would result in the following transformed data:



d. As for the choice of transforming to 3D, we can just add another feature $f(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2)$ to the original data and deriving the following new data:



e. We can derive the kernels of the mentioned transformation functions in the following way $(x = x_1, y = x_2, X = [x_1, x_2]^T)$ and $K(X_i, X_j) = X_i^T X_j$:

$$1. f(x_1, x_2) = x_1^2 + x_2^2$$

$$K_1(X_i, X_j) = (x_1^2 + x_2^2)(x_1^2 + x_2^2)$$

$$2. f(x_1, x_2) = (x_1 + x_2, x_1^2 + x_2^2)$$

$$K_2(X_i, X_j) = (x_1 + x_2)(x_1 + x_2) + (x_1^2 + x_2^2)(x_1^2 + x_2^2)$$

$$3. f(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2)$$

$$K_3(X_i, X_j) = (x_1 + x_2)(x_1 + x_2) + (x_1^2 + x_2^2)(x_1^2 + x_2^2)$$

Now we can use these kernels for further analysis. In `model.LinearSVC` we used them to derive the Kernel Matrix of data for each of the transformations (these kernel matrices are saved into `svm_fit_1.py`, `svm_fit_2.py`, `svm_fit_3.py` for later use).

f. After training a linear SVM on the original data using 5-fold cross validation the following final result were obtained:

```
optimization finished, #iter = 1000
nu = 0.000000
obj = -99.450000, smo = 0.507500
nSV = 101, nSV2 = 100
Total nSV = 101
Cross Validation Accuracy = 50.0%
```

As can be seen the model achieved 50.0% accuracy which were expected considering the shape of the data and trying to separate it using a simple line.

g. Now let's use the same model but this time on the data we produced in part bc and d. After using the same options for LIBSVM (5-fold CV, linear kernel) here is the final result for each of the data:

1-D Data (b):

```
optimization finished, #iter = 5
nu = 0.000000
obj = -0.004444, smo = 1.175000
nSV = 1, nSV2 = 0
Total nSV = 1
Cross Validation Accuracy = 100%
```

2-D Data (c):

```
optimization finished, #iter = 50
nu = 0.000000
obj = -0.004444, smo = 1.175100
nSV = 1, nSV2 = 0
Total nSV = 1
Cross Validation Accuracy = 100%
```

3-D Data (d):

```
optimization finished, #iter = 75
nu = 0.000000
obj = -0.011750, smo = 1.450000
nSV = 4, nSV2 = 0
Total nSV = 4
Cross Validation Accuracy = 100%
```

As can be seen the result for all 3 data is perfect (100% accuracy). This is due to the fact that using the provided mapping function we were able to make the data linearly separable and therefore a linear separator can perform well with them.

h. Using the kernel matrix we saved in e, we can use LIBSVM feature namely precomputed kernels and use these matrices as the training data. Doing so with (again) 5-fold CV will result in:

Kernel Matrix 1:

```
optimization finished, #iter = 5
nu = 0.000000
obj = -0.004444, smo = 1.175000
nSV = 1, nSV2 = 0
Total nSV = 1
Cross Validation Accuracy = 100%
```

Kernel Matrix 2:

```
optimization finished, #iter = 50
nu = 0.000000
obj = -0.004444, smo = 1.175100
nSV = 1, nSV2 = 0
Total nSV = 1
Cross Validation Accuracy = 100%
```

Kernel Matrix 3:

```
optimization finished, #iter = 75
nu = 0.000000
obj = -0.011750, smo = 1.450000
nSV = 4, nSV2 = 0
Total nSV = 4
Cross Validation Accuracy = 100%
```

As expected, the result are completely correspondent with that of the previous section.

Problem 2

a. 1. We split the Professor's data set into training (60%), validation (20%) and test set (20%). Using RBF kernel and the validation set, different values for C and Gamma (10^{-3} to 10^3) has been tested. Here is the final result:

```
Best Validation Result: C=100.0, gamma=0.05, Accuracy=91.5%
```

Using Gamma = $1e-3$ and C = 100 as the final hyperparameters, here is the accuracy obtained on classifying the test set:

```
Test Accuracy: 84.3%
```

2. Also 20 random generated values for the parameters (from the range 10^{-3} to 10^3) have been used in the same way which resulted in:

```
Best Validation Result: C=0.000000, gamma=0.0001, Accuracy=91.1%
```

And plugging C = 100000000 and Gamma = 0.0001, the test accuracy is obtained:

```
Test Accuracy: 84.3%
```

3. Using grid search to find the parameters is exhaustive and time consuming, on the other hand random search might help a good result and map as well fast. Using an optimization procedure like Bayesian optimization or gradient descent might help finding good parameters faster.

b. For the degree kernel, after scaling the data we used the following parameters $d \in \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ and Gamma in range 10^{-3} to 10^3 in a grid search to find the best ones via validation analysis. Here is the result:

```
Best Validation Result: C=10.0, gamma=0.01, coeff=0.1, Accuracy=96.8%
```

As can be seen the accuracy is 100% for $D = 4$ and Gamma = 1, using these parameters we can obtain the test error:

```
Test Accuracy: 97.9%
```

Now using 20 random generated set of parameters for C, gamma, coeff in range 10^{-3} to 10^3 and $d \in \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ the validation accuracy:

```
Best Validation Result: C=0.000000, gamma=0.0, degree=1, coeff=1.0, Accuracy=98.8%
```

And using the obtained parameters we get the test accuracy:

```
Test Accuracy: 97.1%
```

c. For the degree kernel, after again scaling the data we used the following parameters C, coeff and Gamma in range 10^{-3} to 10^3 in a grid search to find the best ones via validation analysis. Here is the result:

```
Best Validation Result: C=10.0, gamma=0.01, coeff=0.1, Accuracy=96.8%
```

Using these best parameters we can obtain the test error:

```
Test Accuracy: 97.1%
```

Now using 20 random generated set of parameters for C, gamma, coeff in range 10^{-3} to 10^3 the validation accuracy:

```
Best Validation Result: C=0.01, gamma=0.0000000000, coeff=1.000000, gamma=0.0, Accuracy=97.4%
```

And using the obtained parameters we get the test accuracy:

```
Test Accuracy: 96.9%
```

d. The C parameter tells the SVM optimization how much we want to avoid misclassifying each training example. For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly. Conversely, a very small value of C will cause the optimizer to look for a large-margin separating hyperplane, even if that hyperplane misclassifies more points. For very tiny values of C, we get misclassified examples, often even if the training data is linearly separable.

We just try to the test by training a polynomial data with gamma=1, C=1 and C=10e-30 and 10e+20.

C=10e-30

```
Test Accuracy: 96.9%
```

C=10e+20

```
Test Accuracy: 96.9%
```

e. And finally here is the number of support vectors for each of the procedures we did in section a through c:

1. Number of Support Vectors: 40
2. Number of Support Vectors: 40
3. Number of Support Vectors: 60
4. Number of Support Vectors: 50
5. Number of Support Vectors: 50
6. Number of Support Vectors: 40

Problem 3

a. Adaboost is sensitive to noise. This is due to the fact that if something is deeply in the positive class region but actually labeled as being in the negative class it would suffer a large loss/penalty since the penalty is exponential, so this one outlier/misclassified point could end up having a very strong influence on the final model learned.

b. Given the data, in the first step we choose the best base classifier which is:



Using this classifier would result in 2 misclassification which is highlighted in the table below:

ID	A	B	Class	w	w*
x1	0	0	+	0.6	0.00000000
x2	0.5	0.5	+	0.6	0.00000000
x3	0	1	-	0.6	0.00000000
x4	0.5	1	-	0.6	0.00000000
x5	1	0	+	0.6	0.00000000
x6	1	1	+	0.6	0.00000000
x7	0	1	-	0.6	0.00000000
x8	0	0	-	0.6	0.00000000
C = 0.000000000000					

In the table above, w is the weights associated with each datum at the beginning of the iteration and w* is the unnormalized weight after we apply the classifier. w* is computed in the following way:

$$w^* = \begin{cases} w e^{-\gamma \text{ misclassification}} & \text{if } \text{misclassification} < 0.5 \\ w e^{\gamma \text{ misclassification}} & \text{if } \text{misclassification} > 0.5 \end{cases}$$

Where $\gamma = \frac{1}{2} \ln \left(\frac{1 - \epsilon}{\epsilon} \right)$ which in this iteration would be $\gamma_1 = 0.6931471805599453$ because $\epsilon_1 = 0.25$. Using $\gamma_1 = 0.6931471805599453$ as a normalizing factor for the weights and using the next base classifier:



We would obtain the following:

ID	A	B	Class	w	w*
x1	0	0	+	0.25	0.00000000
x2	0.5	0.5	+	0.25	0.00000000
x3	0	1	-	0.25	0.00000000
x4	0.5	1	-	0.25	0.00000000
x5	1	0	+	0.25	0.00000000
x6	1	1	+	0.25	0.00000000
x7	0	1	-	0.25	0.00000000
x8	0	0	-	0.25	0.00000000
C = 0.0					

With $\epsilon_1 = 0.0041189502195004$, $\epsilon_2 = 0.000000000000$ and $\epsilon_3 = 0.000000000000$

In the final iteration, using this classifier:



The following table is obtained:

ID	A	B	Class	w	w*
x1	0	0	+	0.25	0.00000000
x2	0.5	0.5	+	0.25	0.00000000
x3	0	1	-	0.25	0.00000000
x4	0.5	1	-	0.25	0.00000000
x5	1	0	+	0.25	0.00000000
x6	1	1	+	0.25	0.00000000
x7	0	1	-	0.25	0.00000000
x8	0	0	-	0.25	0.00000000
C = 0.0					

With $\epsilon_1 = 1.000000000000$, $\epsilon_2 = 0.0$ and $\epsilon_3 = 0.0$