

assignment3

October 21, 2019

1 Support Vector Machine Classifiers

In this assignment, we use the scikit-learn package to train an SVM classifier. To do so, we need to tune 2 hyperparameters: the cost C and precision γ (gamma). We are going to use K -fold cross-validation to determine the best combination of values for this pair.

- Question 0 (0%) Have a look at the 3 first cells. In the third one, take note of how the SVC object is instantiated and trained, how labels are predicted, and finally how the fitting error is computed. In this assignment, the prediction error after a given training is simply defined as the *number of misclassified labels*.
- Question 1 (60%) Using an SVM classifier with an RBF kernel, use 10-fold cross-validation to find the best cost and precision parameters. The range of test values for each parameter is provided.
 - a. First compute the cross-validation error matrix: for each parameter combination, instantiate an SVM classifier; for each split provided by the KFold object, re-train this classifier and compute the prediction error; the cross-validation error is the average of these errors over all splits.
 - b. Use the error matrix to select the best parameter combination.
 - c. Visualize the error matrix using `imshow` and the 'hot' colormap.
- Question 2 (30%) Plot the decision boundaries of this classifier, by appropriately modifying the code from the previous assignments. Display the support vectors on the same figure.
- Question 3 (10%) Evaluate and print the generalization error of this classifier, computed on the test set.

2 Code

2.1 Imports

```
[1]: import matplotlib.pyplot as plt
import numpy as np
from sklearn.model_selection import KFold
from sklearn.svm import SVC
%matplotlib inline
```

2.2 Load and display the training data

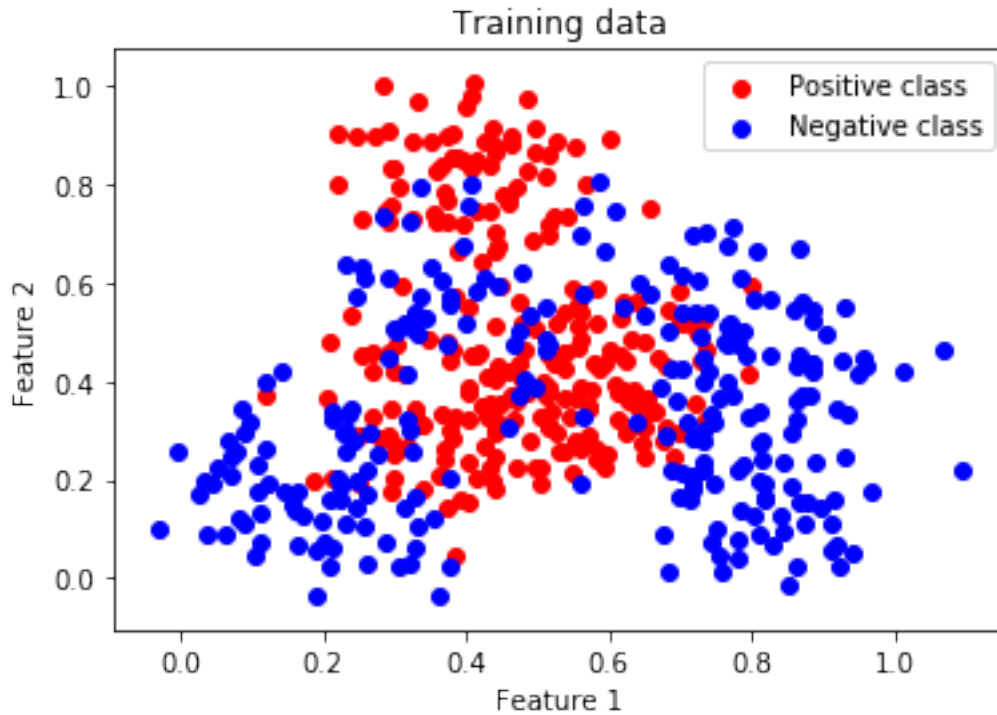
```
[2]: features = np.load("features.npy")
labels = np.load("labels.npy")
print("features size:", features.shape)
print("labels size:", labels.shape)

# Extract features for both classes
pos = labels == 1 # 1D array of booleans, with pos[i] = True if labels[i] == 1
features_pos = features[pos] # filter the array with the boolean array
neg = labels != 1
features_neg = features[neg]

# Display data
fig, ax = plt.subplots()
ax.scatter(features_pos[:, 0], features_pos[:, 1], c="red", label="Positive_
→class")
ax.scatter(features_neg[:, 0], features_neg[:, 1], c="blue", label="Negative_
→class")
ax.set_title("Training data")
ax.set_xlabel("Feature 1")
ax.set_ylabel("Feature 2")
ax.legend()

plt.show()
```

```
features size: (500, 2)
labels size: (500,)
```



2.3 Training the SVM classifier with arbitrary hyperparameters

```
[3]: cost = 1
gamma = 1

# Train the SVM classifier.
svm = SVC(C=cost, kernel='rbf', gamma=gamma)
svm.fit(features, labels)

# Predict labels.
# Note that here we use the same set for training and testing,
# which is not the case in the remainder of the assignment.
predicted_labels = svm.predict(features)

# Compute the error.
# Note: since in Python, True and False are equivalent to 1 and 0, we can
# directly sum over the boolean array returned by the comparison operator.
error = sum(labels != predicted_labels)
print("Prediction error:", error)
```

Prediction error: 98

2.4 Training with K-fold cross-validation

2.4.1 Define test values for the cost and precision parameters

```
[4]: def logsample(start, end, num):  
      return np.logspace(np.log10(start), np.log10(end), num, base=10.0)  
  
num_gammas = 20  
num_costs = 20  
gamma_range = logsample(1e-1, 1e3, num_gammas)  
cost_range = logsample(1e-1, 1e3, num_costs)
```

2.4.2 Compute the cross-validation error for each parameter combination

The KFold class from scikit-learn is a “cross-validation” object, initialized with a number of folds. For each fold, it randomly partitions the input data into a training set and a validation set. The documentation provides an example of use.

```
[5]: K = 10 # number of folds for cross validation  
kf = KFold(n_splits=K)  
cv_error = np.zeros((num_gammas, num_costs)) # error matrix  
  
# TODO (Question 1)  
  
# /TODO (Question 1)
```

2.4.3 Train the classifier with the best parameter combination

```
[6]: # Find gamma and cost giving the smallest error  
# TODO (Question 1)  
  
# /TODO (Question 1)  
  
# Train the SVM classifier using these parameters  
svm = SVC(C=cost, kernel='rbf', gamma=gamma)  
svm.fit(features, labels)  
support_vectors = svm.support_vectors_
```

2.4.4 Display cross-validation results and decision function

```
[7]: # Sample points on a grid
num_points = 100
x_rng = np.linspace(0, 1, num_points)
y_rng = np.linspace(0, 1, num_points)
grid_x, grid_y = np.meshgrid(x_rng, y_rng)

# Evaluate decision function for each point
xy_list = np.column_stack((grid_x.flat, grid_y.flat))
values = svm.decision_function(xy_list)
values = values.reshape((num_points, num_points))

# Display
fig = plt.figure(figsize=plt.figaspect(0.25))

ax = fig.add_subplot(1, 3, 1)
ax.set_title("Cross-validation error")
ax.set_xlabel("Log10 of the cost parameter")
ax.set_ylabel("Log10 of the precision parameter")
# TODO (Question 1)

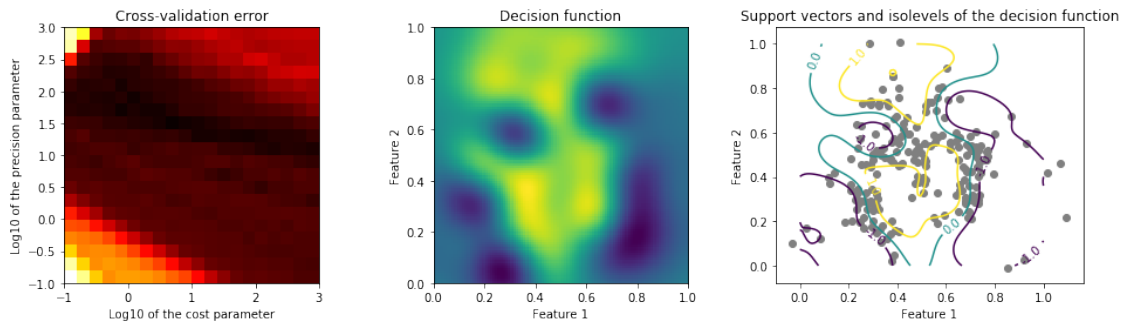
# /TODO (Question 1)

ax = fig.add_subplot(1, 3, 2)
ax.set_title("Decision function")
ax.set_xlabel("Feature 1")
ax.set_ylabel("Feature 2")
ax.imshow(values, extent=[0, 1, 0, 1], origin='lower')

ax = fig.add_subplot(1, 3, 3)
ax.set_title("Support vectors and isolevels of the decision function")
ax.set_xlabel("Feature 1")
ax.set_ylabel("Feature 2")
# TODO (Question 2)

# /TODO (Question 2)

plt.show()
```



2.5 Generalization error

2.5.1 Load the test data

```
[8]: # Load the training data
test_features = np.load("test_features.npy")
test_labels = np.load("test_labels.npy")
print(test_features.shape)
print(test_labels.shape)
```

(500, 2)

(500,)

2.5.2 Print the number of misclassified points in the test set

```
[9]: # TODO (Question 3)

# /TODO (Question 3)
```

90