

---

## Assignment 4

---

### 1. **\*\*Support Vector Machines with Synthetic Data\*\***, 50 points. ¶

For this problem, we will generate synthetic data for a nonlinear binary classification problem and partition it into training, validation and test sets. Our goal is to understand the behavior of SVMs with Radial-Basis Function (RBF) kernels with different values of  $C$  and  $\gamma$ .

```

In [1]: # DO NOT EDIT THIS FUNCTION; IF YOU WANT TO PLAY AROUND WITH DATA GENERATION,
# MAKE A COPY OF THIS FUNCTION AND THEN EDIT
#
import numpy as np
from sklearn.datasets import make_moons
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap

def generate_data(n_samples, tst_frac=0.2, val_frac=0.2):
    # Generate a non-linear data set
    X, y = make_moons(n_samples=n_samples, noise=0.25, random_state=42)

    # Take a small subset of the data and make it VERY noisy; that is, generate outliers
    m = 30
    np.random.seed(30) # Deliberately use a different seed
    ind = np.random.permutation(n_samples)[:m]
    X[ind, :] += np.random.multivariate_normal([0, 0], np.eye(2), (m, ))
    y[ind] = 1 - y[ind]

    # Plot this data
    cmap = ListedColormap(['#b30065', '#178000'])
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap, edgecolors='k')

    # First, we use train_test_split to partition (X, y) into training and test sets
    X_trn, X_tst, y_trn, y_tst = train_test_split(X, y, test_size=tst_frac,
                                                    random_state=42)

    # Next, we use train_test_split to further partition (X_trn, y_trn) into training and validation sets
    X_trn, X_val, y_trn, y_val = train_test_split(X_trn, y_trn, test_size=val_frac,
                                                    random_state=42)

    return (X_trn, y_trn), (X_val, y_val), (X_tst, y_tst)

```

```

In [2]: #
# DO NOT EDIT THIS FUNCTION; IF YOU WANT TO PLAY AROUND WITH VISUALIZATION,
# MAKE A COPY OF THIS FUNCTION AND THEN EDIT
#

def visualize(models, param, X, y):
    # Initialize plotting
    if len(models) % 3 == 0:
        nrows = len(models) // 3
    else:
        nrows = len(models) // 3 + 1

    fig, axes = plt.subplots(nrows=nrows, ncols=3, figsize=(15, 5.0 * nrows))
    cmap = ListedColormap(['#b30065', '#178000'])

    # Create a mesh
    xMin, xMax = X[:, 0].min() - 1, X[:, 0].max() + 1
    yMin, yMax = X[:, 1].min() - 1, X[:, 1].max() + 1
    xMesh, yMesh = np.meshgrid(np.arange(xMin, xMax, 0.01),
                                np.arange(yMin, yMax, 0.01))

    for i, (p, clf) in enumerate(models.items()):
        # if i > 0:
        #     break
        r, c = np.divmod(i, 3)
        ax = axes[r, c]

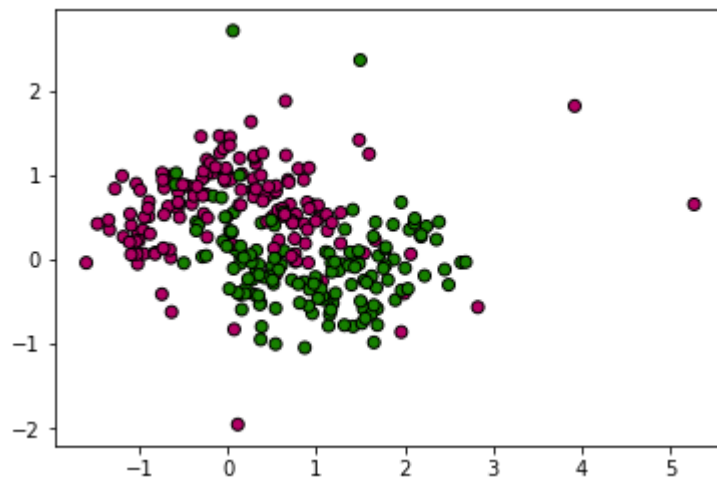
        # Plot contours
        zMesh = clf.decision_function(np.c_[xMesh.ravel(), yMesh.ravel()])
        zMesh = zMesh.reshape(xMesh.shape)
        ax.contourf(xMesh, yMesh, zMesh, cmap=plt.cm.PiYG, alpha=0.6)

        if (param == 'C' and p > 0.0) or (param == 'gamma'):
            ax.contour(xMesh, yMesh, zMesh, colors='k', levels=[-1, 0, 1],
                       alpha=0.5, linestyles=['--', '-', '--'])

        # Plot data
        ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap, edgecolors='k')
        ax.set_title('{0} = {1}'.format(param, p))

```

```
In [3]: # Generate the data
n_samples = 300 # Total size of data set
(X_trn, y_trn), (X_val, y_val), (X_tst, y_tst) = generate_data(n_samples)
```



### a. (25 points) The effect of the regularization parameter, $C$

Complete the Python code snippet below that takes the generated synthetic 2-d data as input and learns non-linear SVMs. Use scikit-learn's [SVC](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>) function to learn SVM models with **radial-basis kernels** for fixed  $\gamma$  and various choices of  $C \in \{10^{-3}, 10^{-2}, \dots, 1, \dots, 10^5\}$ . The value of  $\gamma$  is fixed to  $\gamma = \frac{1}{d \cdot \sigma_X}$ , where  $d$  is the data dimension and  $\sigma_X$  is the standard deviation of the data set  $X$ . SVC can automatically use these setting for  $\gamma$  if you pass the argument `gamma = 'scale'` (see documentation for more details).

**Plot:** For each classifier, compute **both** the **training error** and the **validation error**. Plot them together, making sure to label the axes and each curve clearly.

**Discussion:** How do the training error and the validation error change with  $C$ ? Based on the visualization of the models and their resulting classifiers, how does changing  $C$  change the models? Explain in terms of minimizing the SVM's objective function  $\frac{1}{2} \mathbf{w}' \mathbf{w} + C \sum_{i=1}^n \ell(\mathbf{w} \mid \mathbf{x}_i, y_i)$ , where  $\ell$  is the hinge loss for each training example  $(\mathbf{x}_i, y_i)$ .

**Final Model Selection:** Use the validation set to select the best the classifier corresponding to the best value,  $C_{best}$ . Report the accuracy on the **test set** for this selected best SVM model. *Note: You should report a single number, your final test set accuracy on the model corresponding to  $C_{best}$ .*

```
In [4]: # Learn support vector classifiers with a radial-basis function kernel with
# fixed gamma = 1 / (n_features * X.std()) and different values of C
C_range = np.arange(-3.0, 6.0, 1.0)
C_values = np.power(10.0, C_range)

models = dict()
trnErr = dict()
valErr = dict()

for C in C_values:
    #
    #
    # Insert your code here to learn SVM models
    #
    #

visualize(models, 'C', X_trn, y_trn)

#
#
# Insert your code here to perform model selection
#
#
```

```
File "<ipython-input-4-8875a1448a41>", line 17
    visualize(models, 'C', X_trn, y_trn)
    ^
```

**IndentationError:** expected an indented block

## b. (25 points) The effect of the RBF kernel parameter, $\gamma$

Complete the Python code snippet below that takes the generated synthetic 2-d data as input and learns various non-linear SVMs. Use scikit-learn's [SVC](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>) function to learn SVM models with **radial-basis kernels** for fixed  $C$  and various choices of  $\gamma \in \{10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3\}$ . The value of  $C$  is fixed to  $C = 10$ .

**Plot:** For each classifier, compute **both** the **training error** and the **validation error**. Plot them together, making sure to label the axes and each curve clearly.

**Discussion:** How do the training error and the validation error change with  $\gamma$ ? Based on the visualization of the models and their resulting classifiers, how does changing  $\gamma$  change the models? Explain in terms of the functional form of the RBF kernel,  $\kappa(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \cdot \|\mathbf{x} - \mathbf{z}\|^2)$

**Final Model Selection:** Use the validation set to select the best the classifier corresponding to the best value,  $\gamma_{best}$ . Report the accuracy on the **test set** for this selected best SVM model. *Note: You should report a single number, your final test set accuracy on the model corresponding to  $\gamma_{best}$ .*

```

In [ ]: # Learn support vector classifiers with a radial-basis function kernel with
# fixed C = 10.0 and different values of gamma
gamma_range = np.arange(-2.0, 4.0, 1.0)
gamma_values = np.power(10.0, gamma_range)

models = dict()
trnErr = dict()
valErr = dict()

for G in gamma_values:
    #
    #
    # Insert your code here to learn SVM models
    #
    #

visualize(models, 'gamma', X_trn, y_trn)

#
#
# Insert your code here to perform model selection
#
#

```

## 2. **\*\*Breast Cancer Diagnosis with Support Vector Machines\*\***, 25 points.

For this problem, we will use the [Wisconsin Breast Cancer](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)) ([https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))) data set, which has already been pre-processed and partitioned into training, validation and test sets. Numpy's [loadtxt](https://docs.scipy.org/doc/numpy-1.13.0/reference/generated/numpy.loadtxt.html) (<https://docs.scipy.org/doc/numpy-1.13.0/reference/generated/numpy.loadtxt.html>) command can be used to load CSV files.

```

In [ ]: # Load the Breast Cancer Diagnosis data set; download the files from eLearning
# CSV files can be read easily using np.loadtxt()
#
# Insert your code here.
#

```

Use scikit-learn's [SVC](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>) function to learn SVM models with **radial-basis kernels** for **each combination** of  $C \in \{10^{-2}, 10^{-1}, 1, 10^1, \dots 10^4\}$  and  $\gamma \in \{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$ . Print the tables corresponding to the training and validation errors.

**Final Model Selection:** Use the validation set to select the best classifier corresponding to the best parameter values,  $C_{best}$  and  $\gamma_{best}$ . Report the accuracy on the **test set** for this selected best SVM model. *Note: You should report a single number, your final test set accuracy on the model corresponding to  $C_{best}$  and  $\gamma_{best}$ .*

```
In [ ]: #  
#  
# Insert your code here to perform model selection  
#  
#
```

### 3. **\*\*Breast Cancer Diagnosis with $k$ -Nearest Neighbors\*\***, 25 points.

Use scikit-learn's [k-nearest neighbor](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>) classifier to learn models for Breast Cancer Diagnosis with  $k \in \{1, 5, 11, 15, 21\}$ , with the kd-tree algorithm.

**Plot:** For each classifier, compute **both** the **training error** and the **validation error**. Plot them together, making sure to label the axes and each curve clearly.

**Final Model Selection:** Use the validation set to select the best classifier corresponding to the best parameter value,  $k_{best}$ . Report the accuracy on the **test set** for this selected best kNN model. *Note: You should report a single number, your final test set accuracy on the model corresponding to  $k_{best}$ .*

```
In [ ]: #  
#  
# Insert your code here to perform model selection  
#  
#
```

**Discussion:** Which of these two approaches, SVMs or kNN, would you prefer for this classification task? Explain.