

# Práctica 7. Detección de Spam

## Parte A. Support Vector Machines

### 1. Kernel lineal

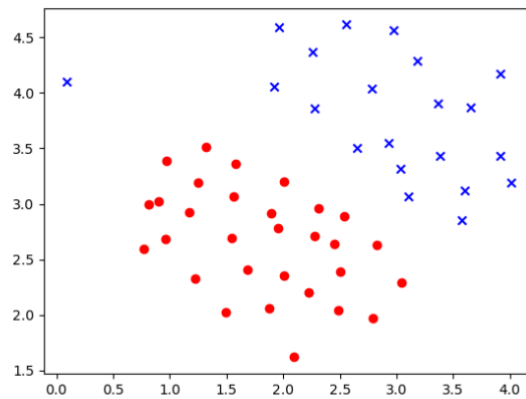
Para la carga del conjunto de datos se ha utilizado el método *load\_data()*:

```
1 def load_data(file):
2     m = loadmat(file)
3     X = m['X']
4     y = m['y'].ravel()
5
6     return X, y
```

Para empezar, se muestra la gráfica del primer conjunto de datos del fichero *ex6data1.mat*. En esta gráfica se aprecia una separación considerable en el conjunto de datos. Para mostrar la gráfica se ha utilizado el método *print\_data()*:

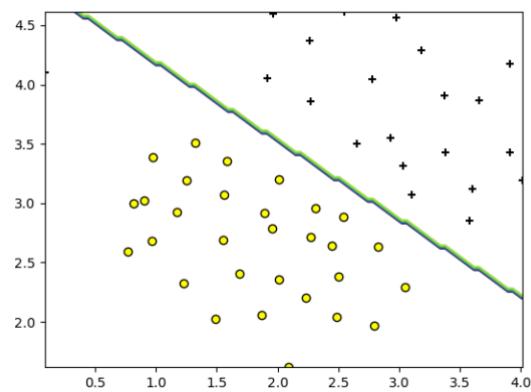
```
1 def print_data(X, y):
2     x0 = (X[:, 0].min(), X[:, 0].max())
3     x1 = (X[:, 1].min(), X[:, 1].max())
4     p = np.linspace(x0, x1, 100)
5
6     x1, x2 = np.meshgrid(p, p)
7     zeros = (y == 0).ravel()
8     ones = (y == 1).ravel()
9
10    plt.figure()
11    plt.scatter(X[ones, 0], X[ones, 1], color='b', marker='x')
12    plt.scatter(X[zeros, 0], X[zeros, 1], color='r', marker='o')
13
14    plt.show()
```

La gráfica es la que se muestra a continuación:

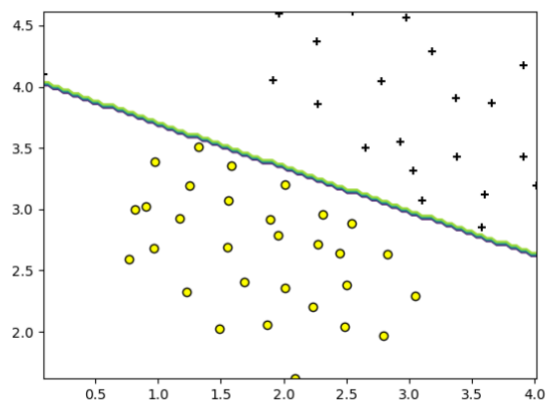


Utilizando la clase *SVC* proporcionada por la librería *sklearn* y especificando el kernel como lineal se obtienen, dependiendo del parámetro *C* las siguientes gráficas:

- Con  $C = 1$ :



- Con  $C = 100$ :



Como vemos en las gráficas, aunque ambas son capaces de diferenciar entre ambos tipos de manera perfecta, para  $C=1$  podemos ver una separación con una mayor distancia entre la línea y cada uno de los puntos. Esto permite una mayor distribución a la hora de introducir nuevos datos.

El código para obtener este resultado es el siguiente:

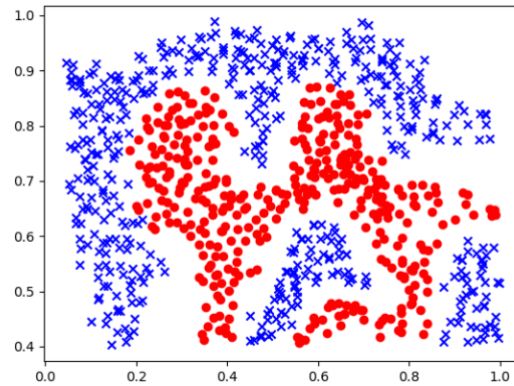
```
1 X, y = load_data(file='data/ex6data1.mat')
2 svm = skl.SVC(kernel='linear', C=1) # C = 100
3 svm.fit(X, y)
4 border_data(X, y, svm)
```

Y la definición del método `border_data()` es el siguiente:

```
1 def border_data(X, y, svm):
2     x1 = np.linspace(X[:, 0].min(), X[:, 0].max(), 100)
3     x2 = np.linspace(X[:, 1].min(), X[:, 1].max(), 100)
4
5     x1, x2 = np.meshgrid(x1, x2)
6     yp = svm.predict(np.array([x1.ravel(), x2.ravel()]).T).reshape(x1.shape)
7     pos = (y == 1).ravel()
8     neg = (y == 0).ravel()
9
10    plt.figure()
11    plt.scatter(X[pos, 0], X[pos, 1], color='black', marker='+')
12    plt.scatter(X[neg, 0], X[neg, 1], color='yellow', edgecolors='black', marker='o')
13    plt.contour(x1, x2, yp)
14    plt.show()
15    plt.close()
```

## 2. Kernel Gaussiano

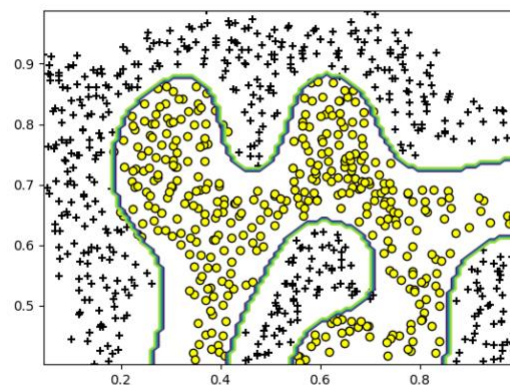
A continuación, se muestra la gráfica obtenida a partir del conjunto de datos proporcionado por el fichero *ex6data2.mat*.



El código para mostrar esta gráfica es el mismo utilizado en el apartado anterior. Para mostrar la gráfica donde se aprecia el uso del kernel gaussiano a partir de la clase SVC proporcionada por la librería *sklearn* es el siguiente:

```
1  c = 1
2  sigma = 0.1
3  svm = skl.SVC(kernel='linear', C = c, gamma = (1/(2 * sigma**2)))
4  svm.fit(X, y)
5  border_data(X, y, svm)
```

Y la gráfica resultante del método *border\_data()* es:



Si se aplicase el kernel lineal el resultado sería una línea recta que no separaría de manera correcta los diferentes conjuntos de datos. Con el kernel gaussiano se puede apreciar que el modelo lo separa casi a la perfección.

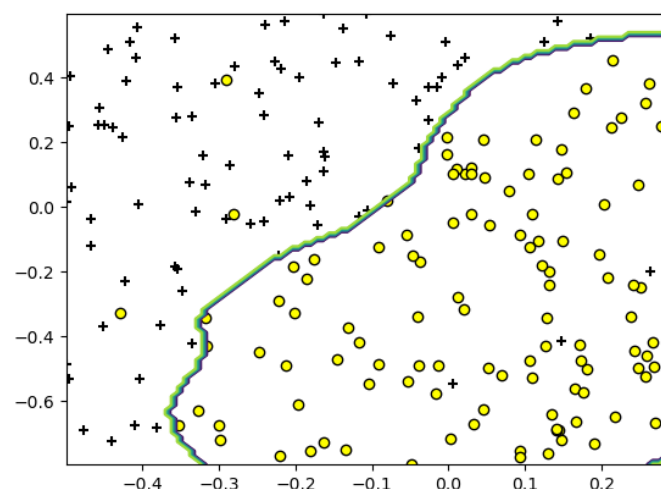
### 3. Elección de los parámetros C y $\sigma$

A continuación, se seleccionarán los valores de C y  $\sigma$  para un modelo de SVM con kernel gaussiano que clasifique el tercer conjunto de datos *ex6data3.mat*. En este conjunto de datos, además de los datos de entrenamiento X e y, se incluyen datos de validación Xval e yval que servirán para evaluar el modelo aprendido.

Se ha definido el método *main\_chooseCAndSigma()*:

```
1  def main_chooseCAndSigma():
2      data = loadmat('data/ex6data3.mat')
3      X = data['X']
4      y = data['y'].ravel()
5      x_val = data['Xval']
6      y_val = data['yval'].ravel()
7      scores = chooseCAndSigma(X, y, x_val, y_val)
8      print(scores)
9      ind = np.where(scores == scores.max())
10     params = [0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30]
11     indp1 = ind[0][0]
12     indp2 = ind[1][0]
13
14     svm = skl.SVC(kernel='rbf', C = params[indp1], gamma = 1 / (2*params[indp2] ** 2))
15     svm.fit(X, y)
16     border_data(x_val, y_val, svm)
17     print("C = {}, sigma = {}".format(params[indp1], params[indp2]))
18
19     print(ind)
```

Como resultado se obtiene que el C óptimo es 1 y sigma = 0.1. La gráfica resultante es la siguiente:



El método para hallar estos parámetros es el siguiente:

```
1 def chooseCAndSigma(X, y, x_val, y_val):
2     params = [0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30]
3     scores = np.zeros((len(params), len(params)))
4     res = []
5
6     for v in params:
7
8         for sigma in params:
9             svm = skl.SVC(kernel='rbf', C = v, gamma = 1 / (2*sigma ** 2))
10            svm.fit(X, y)
11            x_pred = svm.predict(x_val)
12            correct = sum(x_pred == y_val) / x_pred.shape[0] * 100
13            res.append(correct)
14            scores[params.index(v), params.index(sigma)] = correct
15
16     return scores
```

## Detección de Spam

### 1. Lectura y procesamiento de los datos

Se ha procesado los conjuntos de datos conjuntos de datos de correo spam (spam.zip) y no spam (easy\_ham.zip, más fáciles de identificar como correo no spam, y hard\_ham.zip más fáciles de confundir con spam).

Se ha creado los siguientes métodos para conseguir esto:

```
1 def transformAVector(vector, dictionary):
2     email = np.zeros(1899)
3     for word in vector:
4         if word in dictionary.keys():
5             email[dictionary.get(word) - 1] += 1
6     return email
7
8 def load_examples(file, num_examples):
9     response = np.zeros((num_examples, 1899))
10    dictionary = getVocabDict()
11    for i in range(1, num_examples + 1):
12        content = codecs.open('{}/{}.txt'.format(file, str(i).zfill(4)), encoding = 'utf', errors = 'ignore').read()
13        vector = email2TokenList(content)
14        transform = transformAVector(vector, dictionary)
15        response[i - 1] = transform
16    return response
```

Este código es parte de un proceso para cargar y transformar datos de correo electrónico en vectores numéricos, con el objetivo de analizar o entrenar modelos de aprendizaje automático. El proceso convierte el contenido de archivos de texto en una representación vectorial, usando un diccionario para mapear palabras a índices específicos y generar vectores de conteo. Estos vectores luego se pueden usar para análisis posteriores, como el entrenamiento de modelos de clasificación para detección de spam.

## 2. Entrenamiento y validación

El proceso para encontrar los valores óptimos de C y sigma se divide en dos funciones principales: *findOptimalCAndSigma* y *train*.

- *findOptimalCAndSigma*

Esta función busca los valores óptimos de C y sigma evaluando un conjunto de parámetros y seleccionando los que producen la mayor precisión en un conjunto de validación para luego utilizarlos en el modelo SVM.

```
1 def findOptimalCAndSigma(x_train, y_train, x_val, y_val):
2     params = [ 0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30 ]
3     scores = np.zeros((len(params), len(params)))
4
5     for C in params:
6         print('Comprobación de C en {}'.format(C))
7         for sigma in params:
8             print('Comprobación de Sigma en {}'.format(sigma))
9             svm = skl.SVC(kernel='rbf', C = C, gamma = 1/(2*sigma**2))
10            svm.fit(x_train, y_train)
11            accuracy = accuracy_score(y_val, svm.predict(x_val))
12            print('Precisión con estos parámetros: {}'.format(accuracy))
13            scores[params.index(C), params.index(sigma)] = accuracy
14
15     optimalC = params[np.where(scores == scores.max())[0][0]]
16     optimalSigma = params[np.where(scores == scores.max())[1][0]]
17     print('Se ha determinado que el C óptimo es {}'.format(optimalC))
18     print('Se ha determinado que el sigma óptimo es {}'.format(optimalSigma))
19     return optimalC, optimalSigma
```

- *train*

Esta función entrena un modelo SVM con los datos de correo electrónico proporcionados y utiliza *findOptimalCAndSigma* para determinar los mejores valores de C y sigma.

```
1 def train(eham, hham, spam):
2     yeham = np.zeros(2551)
3     yham = np.zeros(250)
4     yspam = np.ones(500)
5     random_state = 231052021
6
7     eham_train, eham_test, yeham_train, yeham_test = train_test_split(eham, yeham, test_size = 0.25, random_state = random_state)
8     hham_train, hham_test, yham_train, yham_test = train_test_split(hham, yham, test_size = 0.25, random_state = random_state)
9     spam_train, spam_test, yspam_train, yspam_test = train_test_split(spam, yspam, test_size = 0.25, random_state = random_state)
10
11     x_train1 = np.concatenate((eham_train, hham_train, spam_train))
12     y_train1 = np.concatenate((yeham_train, yham_train, yspam_train))
13     x_test = np.concatenate((eham_test, hham_test, spam_test))
14     y_test = np.concatenate((yeham_test, yham_test, yspam_test))
15
16     xtrain, xval, ytrain, yval = train_test_split(x_train1, y_train1, test_size = 0.30, random_state = random_state)
17     optimalC, optimalSigma = findOptimalCAndSigma(x_train = xtrain, y_train = ytrain, x_val = xval, y_val = yval)
18
19     svmop = skl.SVC(kernel = 'rbf', C = optimalC, gamma = 1/(2*optimalSigma**2))
20     svmop.fit(x_train1, y_train1)
21     print(accuracy_score(y_test, svmop.predict(x_test)))
```

### 3. Resultados

Ejecutando estos métodos en la función *main* se ha llegado a la conclusión que los parámetros óptimos son  $C = 30$  y  $\sigma = 30$ , proporcionando una precisión del 0.9576271186440678