

Лабораторная работа №5 “Метод опорных векторов”

Долатович Алина, 858461

Набор данных ex5data1.mat представляет собой файл формата *.mat (т.е. сохраненного из Matlab). Набор содержит три переменные X1 и X2 (независимые переменные) и y (метка класса). Данные являются линейно разделимыми.

Набор данных ex5data2.mat представляет собой файл формата *.mat (т.е. сохраненного из Matlab). Набор содержит три переменные X1 и X2 (независимые переменные) и y (метка класса). Данные являются нелинейно разделимыми.

Набор данных ex5data3.mat представляет собой файл формата *.mat (т.е. сохраненного из Matlab). Набор содержит три переменные X1 и X2 (независимые переменные) и y (метка класса). Данные разделены на две выборки: обучающая выборка (X, y), по которой определяются параметры модели; валидационная выборка (Xval, yval), на которой настраивается коэффициент регуляризации и параметры Гауссового ядра.

Набор данных spamTrain.mat представляет собой файл формата *.mat (т.е. сохраненного из Matlab). Набор содержит две переменные X - вектор, кодирующий отсутствие (0) или присутствие (1) слова из словаря vocab.txt в письме, и y - метка класса: 0 - не спам, 1 - спам. Набор используется для обучения классификатора.

Набор данных spamTest.mat представляет собой файл формата *.mat (т.е. сохраненного из Matlab). Набор содержит две переменные Xtest - вектор, кодирующий отсутствие (0) или присутствие (1) слова из словаря vocab.txt в письме, и ytest - метка класса: 0 - не спам, 1 - спам. Набор используется для проверки качества классификатора.

```
In [1]: 1 from scipy.io import loadmat
        2 import matplotlib.pyplot as pyplot
        3 from sklearn import svm
        4 import numpy as np
        5 import re
        6 from nltk.stem import PorterStemmer
        7 import pandas
        8 import os
```

Загрузите данные ex5data1.mat из файла.

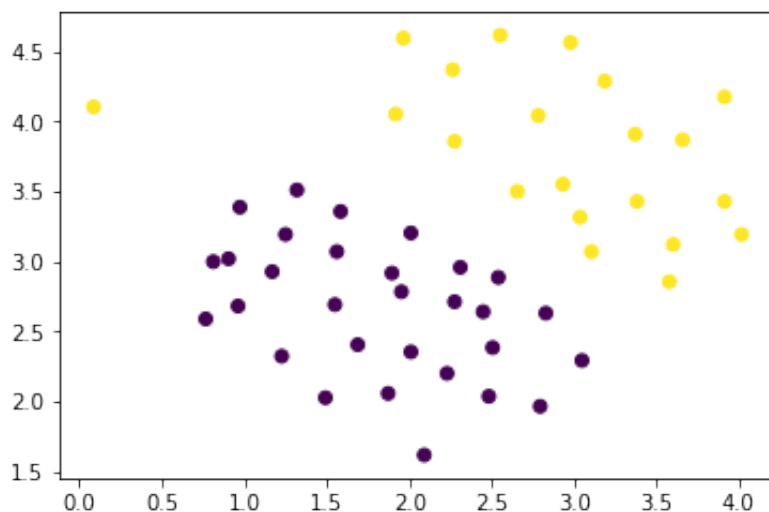
```
In [2]: 1 data = loadmat('ex5data1.mat')
        2
        3 x = data['X']
        4 y = data['y']
        5
        6 print(x.shape)
        7 print(y.shape)
```

```
(51, 2)
```

```
(51, 1)
```

Постройте график для загруженного набора данных: по осям - переменные X1, X2, а точки, принадлежащие различным классам должны быть обозначены различными маркерами.

```
In [3]: 1 def scatterPlot(x, y):
        2     pyplot.scatter(x[:, 0], x[:, 1], c=y.squeeze())
        3
        4 scatterPlot(x, y)
        5 pyplot.show()
```



Обучите классификатор с помощью библиотечной реализации SVM с линейным ядром на данном наборе.

```
In [4]: 1 def getSVMModel(x, y, c=1, gamma=10, kernel='linear'):
        2     model = svm.SVC(kernel=kernel, C=c, gamma=gamma)
        3     model.fit(x, y.ravel())
        4     model.score(x, y.ravel())
        5     return model
        6
        7 svmModel = getSVMModel(x, y)
```

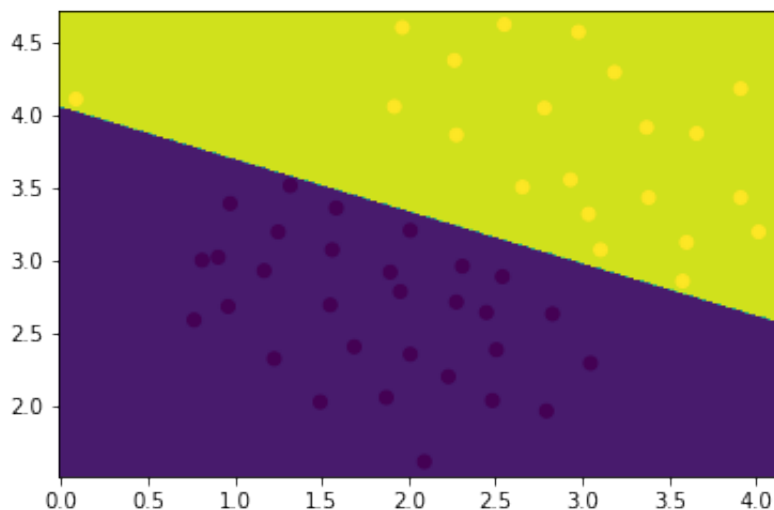
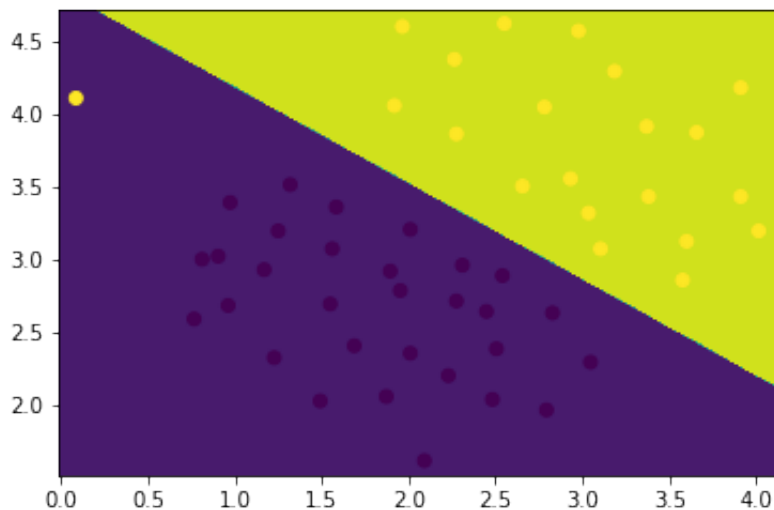
Постройте разделяющую прямую для классификаторов с различными параметрами $C = 1$, $C = 100$ (совместно с графиком из пункта 2). Объясните различия в полученных прямых?

In [5]:

```

1  def contourfPlot(x, y, svmModel, step=0.01):
2      xMin, xMax = min(x[:, 0]) - 0.1, max(x[:, 0]) + 0.1
3      yMin, yMax = min(x[:, 1]) - 0.1, max(x[:, 1]) + 0.1
4      xx, yy = np.meshgrid(np.arange(xMin, xMax, step), np.arange(yMin, yMax, step))
5
6      Z = svmModel.predict(np.c_[xx.ravel(), yy.ravel()])
7      Z = Z.reshape(xx.shape)
8      pyplot.contourf(xx, yy, Z)
9
10 def plotData(x, y, cs, kernel, gamma=10):
11     for c in cs:
12         svmModel = getSVMModel(x, y, c, gamma=gamma, kernel=kernel)
13         contourfPlot(x, y, svmModel)
14         scatterPlot(x, y)
15         pyplot.show()
16
17 cs = [1, 100]
18 plotData(x, y, cs, 'linear')

```



Реализуйте функцию вычисления Гауссового ядра для алгоритма SVM.

```
In [6]: 1 #k(x1,x2) = exp(-q * ||x1-x2||^2), q>0 \\ q=1/(2*sigm^2)
2
3 def gausKernel(x, gamma=0.5): #q>0
4     x0, x1 = x[:, 0], x[:, 1]
5     return np.exp(-gamma * max(x0 - x1) ** 2)
6
7 gausKernel(x)
```

Out[6]: 0.70966828666490156

Загрузите данные ex5data2.mat из файла.

```
In [7]: 1 data2 = loadmat('ex5data2.mat')
2
3 y = data2['y']
4 x = data2['X']
5
6 print(x.shape, y.shape)
```

((863, 2), (863, 1))

Обработайте данные с помощью функции Гауссового ядра.

```
In [8]: 1 width = gausKernel(x)
2 print(width)
```

0.844384016754

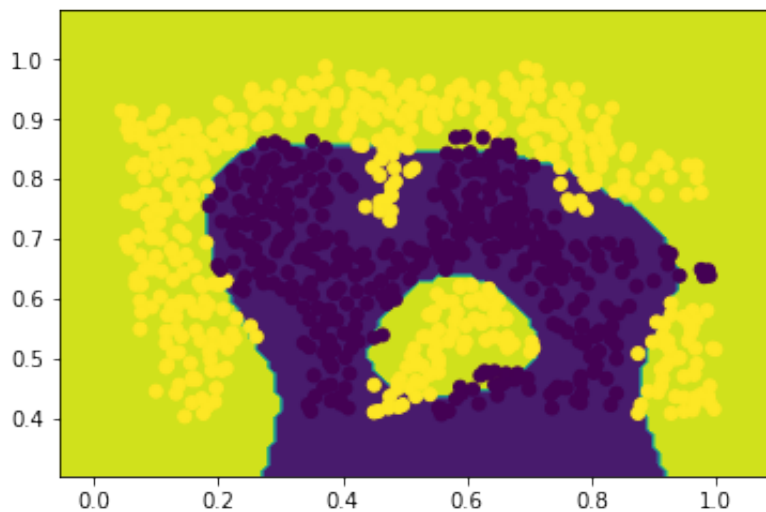
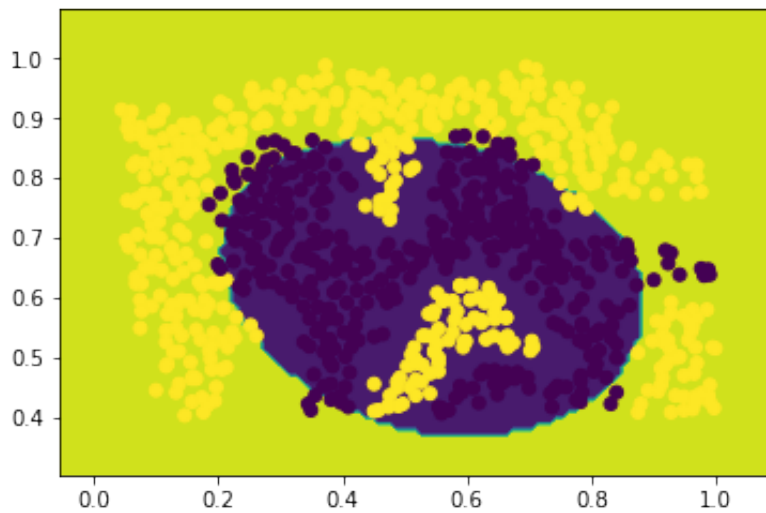
Обучите классификатор SVM.

```
In [9]: 1 svmModel = getSVModel(x, y, gamma=width, kernel='rbf')
2 print(svmModel)

SVC(C=1, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma=0.84438401675358987,
    kernel='rbf', max_iter=-1, probability=False, random_state=None,
    shrinking=True, tol=0.001, verbose=False)
```

Визуализируйте данные вместе с разделяющей кривой (аналогично пункту 4).

```
In [10]: 1 cs = [10, 100000]
2 plotData(x, y, cs, 'rbf', gamma=width)
```



Загрузите данные ex5data3.mat из файла.

```
In [11]: 1 data3 = loadmat('ex5data3.mat')
2
3 yVal = data3['yval']
4 xVal = data3['Xval']
5
6 y = data3['y']
7 x = data3['X']
8
9 print(xVal.shape, yVal.shape, x.shape, y.shape)
```

```
((200, 2), (200, 1), (211, 2), (211, 1))
```

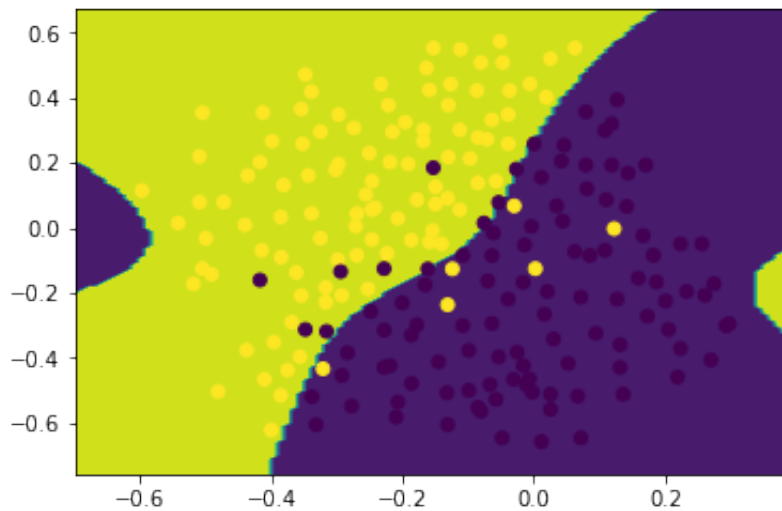
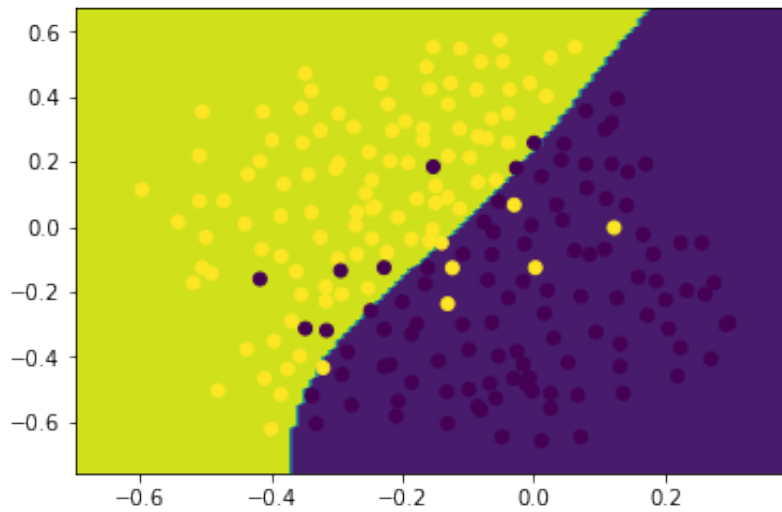
Вычислите параметры классификатора SVM на обучающей выборке, а также подберите параметры C и σ^2 на валидационной выборке.

```
In [12]: 1 width = gausKernel(x)
          2
          3 model = getSVMModel(x, y, c=10000, gamma=width, kernel='rbf')
          4 model.score(xVal, yVal)
          5 params = model.get_params()
          6 print(params)
          7
          8 predicted = model.predict(xVal)
          9
         10 errorCount = 0
         11 for p, val in zip(predicted, yVal.squeeze()):
         12     if p != val:
         13         errorCount += 1
         14
         15 print(float(errorCount)/len(predicted))
         16
```

```
{'kernel': 'rbf', 'C': 10000, 'verbose': False, 'probability': False,
 'degree': 3, 'shrinking': True, 'max_iter': -1, 'decision_function_shape': 'ovr', 'random_state': None, 'tol': 0.001, 'cache_size': 200, 'coef0': 0.0, 'gamma': 0.77298573383986124, 'class_weight': None}
0.065
```

Визуализируйте данные вместе с разделяющей кривой (аналогично пункту 4).

```
In [13]: 1 cs = [100, 10000]
          2 plotData(x, y, cs, 'rbf', gamma=width)
          3
```



Загрузите данные spamTrain.mat из файла.


```
In [14]: 1 data4 = loadmat('spamTrain.mat')
          2
          3 trainX = data4['X']
          4 trainY = data4['y']
          5
          6 print(trainX, trainY.shape)

(array([[0, 0, 0, ..., 0, 0, 0],
        [0, 0, 0, ..., 0, 0, 0],
        [0, 0, 0, ..., 0, 0, 0],
        ...,
        [0, 0, 0, ..., 0, 0, 0],
        [0, 0, 1, ..., 0, 0, 0],
        [0, 0, 0, ..., 0, 0, 0]], dtype=uint8), (4000, 1))
```

Обучите классификатор SVM.

```
In [15]: 1 model = getSVMModel(trainX, trainY, c=100, gamma=1, kernel='rbf')
          2 print(model)

SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma=1, kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

Загрузите данные spamTest.mat из файла.

```
In [16]: 1 data5 = loadmat('spamTest.mat')
          2
          3 xTest = data5['Xtest']
          4 yTest = data5['ytest']
          5
          6 print(xTest.shape, yTest.shape)

((1000, 1899), (1000, 1))
```

Подберите параметры C и σ^2 .

```
In [17]: 1 model.score(xTest, yTest)
```

```
Out[17]: 0.80900000000000005
```

In [18]:

```
1 emailText = open("emailSample1.txt", "r").read()  
2 print(emailText)
```

> Anyone knows how much it costs to host a web portal ?

>

Well, it depends on how many visitors you're expecting.

This can be anywhere from less than 10 bucks a month to a couple of \$100.

You should checkout <http://www.rackspace.com/>

(<http://www.rackspace.com/>) or perhaps Amazon EC2

if youre running something big..

To unsubscribe yourself from this mailing list, send an email to:
groupname-unsubscribe@egroups.com

Реализуйте функцию предобработки текста письма, включающую в себя: -перевод в нижний регистр; -удаление HTML тэгов; -замена URL на одно слово (например, "httpaddr"); -замена email-адресов на одно слово (например, "emailaddr"); -замена чисел на одно слово (например, "number"); -замена знаков доллара (\$) на слово "dollar"; -замена форм слов на исходное слово (например, слова "discount", "discounts", "discounted", "discounting" должны быть заменены на слово "discount"). Такой подход называется stemming; -остальные символы должны быть удалены и заменены на пробелы, т.е. в результате получится текст, состоящий из слов, разделенных пробелами.

In [19]:

```

1  def processText(text):
2      text = text.lower()
3
4      htmlReg = r'<.*?>'
5      text = re.sub(htmlReg, '', text)
6
7      urlReg = r'[a-z]*[.:]+\S+'
8      text = re.sub(urlReg, ' httpaddr ', text)
9
10     emailReg = r'[\w\.-]+@[\w\.-]+'
11     text = re.sub(emailReg, ' emailaddr ', text)
12
13     numberReg = r'\d+'
14     text = re.sub(numberReg, ' number ', text)
15
16     text = text.replace('$', ' dollar ')
17
18     # ps = PorterStemmer()
19     # wordList = text.split(' ')
20     # wordList = [ps.stem(w) for w in wordList]
21     # text = ' '.join(wordList)
22
23     symbolsReg = r'[?_()^@%*&*+>\/=<\t!@,.\~\'|&\"';\\n#$-]'
24     text = re.sub(symbolsReg, '', text)
25     return text
26
27 text = processText(emailText)
28 print(text)

```

anyone knows how much it costs to host a web portal well it depends on how many visitors youre expectingthis can be anywhere from less than number bucks a month to a couple of dollar number you should checkout httpaddr or perhaps amazon ec number if youre running something httpaddr to unsubscribe yourself from this mailing list send an email togroupnameunsubscribe httpaddr

Загрузите коды слов из словаря vocab.txt.

In [20]:

```

1 vocab = pandas.read_csv('vocab.txt', sep = "\t", header=None)
2 print(vocab.shape)

```

```
(1899, 2)
```

Реализуйте функцию замены слов в тексте письма после предобработки на их соответствующие коды.

In [21]:

```
1 def replaceWordToCode(text, vocab):
2     wordsData = vocab[1]
3     codesData = vocab[0]
4
5     words = text.split(' ')
6     codes = []
7     for word in words:
8         index = wordsData[wordsData == word].index.tolist()
9         if index:
10             code = codesData[index].values[0]
11             codes.append(code)
12     return codes
13
14 codes = replaceWordToCode(text, vocab)
15 print(codes)
```

```
[794, 1077, 883, 1699, 790, 1822, 1831, 883, 1171, 794, 238, 162, 68
8, 945, 1663, 1120, 1062, 1699, 1162, 477, 1120, 1893, 1510, 799, 11
82, 512, 1120, 810, 799, 1699, 1896, 688, 961, 1477, 71, 530, 799]
```

Реализуйте функцию преобразования текста письма в вектор признаков (в таком же формате как в файлах spamTrain.mat и spamTest.mat).

In [22]:

```
1 def conversion(codes, vocab):
2     signs = np.zeros((len(vocab), len(codes)))
3     for codeIndex in range(len(codes)):
4         index = codes[codeIndex]
5         signs[index-1, codeIndex] = 1
6     return signs.T
7
8 conversionX = conversion(codes, vocab)
9 print(conversionX.shape)
```

```
(37, 1899)
```

Проверьте работу классификатора на письмах из файлов emailSample1.txt, emailSample2.txt, spamSample1.txt и spamSample2.txt.

In [23]:

```
1 def checkPredictedValue(vocab):
2     mails = ['emailSample1.txt', 'emailSample2.txt', 'spamSample1.
3     for mail in mails:
4         emailText = open(mail, "r").read()
5         text = processText(emailText)
6         codes = replaceWordToCode(text, vocab)
7         conversionX = conversion(codes, vocab)
8         #         print(conversionX)
9         predicted = model.predict(conversionX)
10        print(predicted)
11
12 checkPredictedValue(vocab)
```

[illegible]

Создайте свой набор данных из оригинального корпуса текстов -

<http://spamassassin.apache.org/old/publiccorpus/>

(<http://spamassassin.apache.org/old/publiccorpus/>).

In [24]:

```
1 def generateDict():
2     paths = ['hard_ham/', 'spam/']
3     keys = [0, 1]
4     count = 0
5
6     returnedText = ''
7     dictValue = dict()
8     for key, path in zip(keys, paths):
9         fullText = ''
10        for filename in os.listdir(path):
11            fullText += open(path + filename, "r").read()
12
13        fullText = processText(fullText)
14        returnedText += fullText
15        setWord = set(fullText.split(' '))
16        for word in setWord:
17            dictValue[count] = [count+1, word, key]
18            count += 1
19
20    return dictValue, returnedText
21
22 values, text = generateDict()
23 print(values)
```

IOPub data rate exceeded.

The notebook server will temporarily stop sending output to the client in order to avoid crashing it.

To change this limit, set the config variable

`--NotebookApp.iopub_data_rate_limit`.

Постройте собственный словарь.

In [25]:

```
1 fullDict = pandas.DataFrame.from_dict(values, orient='index')
2 newVocab = fullDict.loc[:, fullDict.columns.isin([0, 1])]
3 newY = np.array(fullDict.loc[:, fullDict.columns.isin([2])])
4 xxx = list(fullDict.loc[:, fullDict.columns.isin([0])].values.squeeze())
5 newX = conversion(xxx, newVocab)
```

In []:

```
1 # codes = replaceWordToCode(text, newVocab)
```

In []:

```
1 model = getSVMModel(newX, newY, c=100, gamma=1, kernel='rbf')
2 checkPredictedValue(newVocab)
```

Вывод

Метод опорных векторов — набор схожих алгоритмов обучения с учителем, использующихся для задач классификации и регрессионного анализа. Принадлежит семейству линейных классификаторов. Особым свойством метода опорных векторов является непрерывное уменьшение эмпирической ошибки классификации и увеличение зазора, поэтому метод также известен как метод классификатора с максимальным зазором.

Основная идея метода — перевод исходных векторов в пространство более высокой размерности и поиск разделяющей гиперплоскости с максимальным зазором в этом пространстве. Две параллельных гиперплоскости строятся по обеим сторонам гиперплоскости, разделяющей классы. Разделяющей гиперплоскостью будет гиперплоскость, максимизирующая расстояние до двух параллельных гиперплоскостей. Алгоритм работает в предположении, что чем больше разница или расстояние между этими параллельными гиперплоскостями, тем меньше будет средняя ошибка классификатора.