

sheet09

July 5, 2018

1 Part A: Kernels for Genes Sequences

In this first exercise, various *degree kernels* such as the weighted degree kernel (WDK) will be implemented. We will use Scikit-Learn (<http://scikit-learn.org/>) for training SVMs. The focus of this exercise is therefore on the computation of the kernels.

We consider a problem of binary classification of genes sequences. The training and test data is available in the folder `splices-data`. The following code reads the gene sequences data and stores it in numpy arrays.

```
In [82]: import numpy
Xtrain = numpy.array([numpy.array(list(l)) for l in open('splice-data/splice-train-data.txt')])
Xtest  = numpy.array([numpy.array(list(l)) for l in open('splice-data/splice-test-data.txt')])
Ttrain = numpy.array([int(l) for l in open('splice-data/splice-train-label.txt','r')])
Ttest  = numpy.array([int(l) for l in open('splice-data/splice-test-label.txt','r')])
```

1.1 Degree Kernels (20 P)

We consider the degree kernel of degree d applying to two genes sequences x and x' and defined as:

$$k_d(x, x') = \sum_{l=1}^{L-d+1} \mathbf{1}_{u_{l,d}(x)=u_{l,d}(x')}$$

where l iterates over the whole genes sequence, $u_{l,d}(x)$ is a subsequence of x starting at position l and of length d , and $\mathbf{1}_{\{ \}}$ is an indicator variable for the equality test given as argument. Given a training set and test set of genes sequences, implement a function that *efficiently* computes the kernel matrices for a certain degree $d \in \{1, 2, 3, 4\}$.

```
In [83]: def getdegreekernels(Xtrain, Xtest, degree):
```

```
    """ Replace by our own code

    # l_train: length of training dataset
    # l_test: length of test dataset
    # L: length of DNA
    # d: degree
    l_train = len(Xtrain)
    l_test = len(Xtest)
```

```

L = len(Xtrain[0]) - 1
d = degree

# Ktrain: kernel matrice of training data
# Ktest: kernel matrice of test data
Ktrain = numpy.zeros((l_train, l_train))
Ktest = numpy.zeros((l_test, l_train))

# calculating kernel matrices of training data
for i in range (l_train):

    # is_equal_1 is a l_train * l_train matrix, tracks where train data equals train data
    is_equal = Xtrain[i,:] == Xtrain

    sum_ = is_equal[:,L - d - 1]

    # finds where 'd'(degree) data in succession are identical
    for k in range(d):
        sum_ = sum_ & is_equal[:,k:L - d - 1 + k]

    # sum up
    Ktrain[i,:] = sum_.sum(axis = 1)

# calculating kernel matrices of test data
for i in range (l_test):

    # is_equal_1 is a l_test * l_train matrix, tracks where train data equals train data
    is_equal = Xtest[i,:] == Xtrain

    sum_ = is_equal[:,L - d - 1]

    # finds where 'd'(degree) data in succession are identical
    for k in range(d):
        sum_ = sum_ & is_equal[:,k:L - d - 1 + k]

    # sum up
    Ktest[i,:] = sum_.sum(axis = 1)

###

assert(Ktrain.shape==(len(Xtrain),len(Xtrain))) and Ktest.shape==(len(Xtest),len(Xtrain))
return Ktrain,Ktest

```

The code below calls the function you implemented for various degrees d , trains SVMs based on these kernels, and measures the prediction accuracy. It can be expected to run in less than 1 minute.

```
In [84]: from sklearn import svm
```

```

Ktrains,Ktests = [None]*4,[None]*4
import time

start_time = time.time()

for i in range(4):
    Ktrains[i],Ktests[i] = getdegreekernels(Xtrain,Xtest,i+1)
    mysvm = svm.SVC(kernel='precomputed').fit(Ktrains[i],Ttrain)
    Ytrain = mysvm.predict(Ktrains[i])
    Ytest = mysvm.predict(Ktests[i])
    print('degree: %d    training accuracy: %.3f    test accuracy: %.3f'% \
          (i+1,(Ytrain==Ttrain).mean(),(Ytest==Ttest).mean()))

running_time = round(time.time() - start_time,2)

print('running time is {} s'.format(running_time))

degree: 1    training accuracy: 0.995    test accuracy: 0.918
degree: 2    training accuracy: 1.000    test accuracy: 0.940
degree: 3    training accuracy: 1.000    test accuracy: 0.963
degree: 4    training accuracy: 1.000    test accuracy: 0.959
running time is 11.92 s

```

1.2 Weighted Degree Kernel (10 P)

We now consider a weighted degree kernel with uniform weights:

$$k(x, x') = \sum_{d=1}^4 k_d(x, x')$$

where $k_d(x, x')$ is the kernel with degree d that was implemented in the previous section. *Construct* the kernel matrices for the weighted degree kernel and *compute* the training and test accuracy of an SVM trained with this new kernel.

In [85]: *### Replace by our own code*

```

weighted_Ktrain = numpy.zeros([len(Xtrain),len(Xtrain)])
weighted_Ktest = numpy.zeros([len(Xtest),len(Xtrain)])

for i in range(4):
    Ktrains,Ktests = getdegreekernels(Xtrain,Xtest,i+1)
    weighted_Ktrain += Ktrains
    weighted_Ktest += Ktests

mysvm = svm.SVC(kernel='precomputed').fit(weighted_Ktrain,Ttrain)
Yweighttrain = mysvm.predict(weighted_Ktrain)
Yweighttest = mysvm.predict(weighted_Ktest)
print('training accuracy: %.3f    test accuracy: %.3f'% \

```

```

        ((Yweighttrain==Ttrain).mean(),(Yweighttest==Ttest).mean()))
    ###

training accuracy: 1.000    test accuracy: 0.967

```

2 Part B: Kernels for Text

Structured kernels can also be used for classifying text data. In this exercise, we consider the classification of a subset of the 20-newsgroups data (available at <http://qwone.com/~jason/20Newsgroups/>). A subset of this data composed only of texts of class comp.graphics and sci.med is given in the folder newsgroup-data. The first class is assigned label -1 and the second class is assigned label +1. Furthermore, the beginning and the end of the newsgroup messages are removed as they typically contain information that makes the classification problem trivial. Like for the genes sequences dataset, data files are composed of multiple rows, where each row corresponds to one example. The code below extracts the fifth message of the training set and displays its 500 first characters.

```

In [86]: import textwrap
        text = list(open('newsgroup-data/newsgroup-train-data.txt','r'))[4]
        print(textwrap.fill(text[:500]+' [...]'))

```

```

hat is, >>center and radius, exactly fitting those points? I know how
to do it >>for a circle (from 3 points), but do not immediately see a
>>straightforward way to do it in 3-D. I have checked some >>geometry
books, Graphics Gems, and Farin, but am still at a loss? >>Please have
mercy on me and provide the solution? > >Wouldn't this require a
hyper-sphere. In 3-space, 4 points over specifies >a sphere as far as
I can see. Unless that is you can prove that a point >exists in
3-space that [...]

```

2.1 Creating Bag-Of-Words (10 P)

A convenient way of representing text data is as bag-of-words: a set composed of all the words occurring in the document. For the purpose of this exercise, we formally define a word as an isolated sequence of at least three consecutive alphabetical characters. Furthermore, a set of stopwords containing mostly uninformative words such as prepositions or conjunctions that should be excluded from the bag-of-word representation is provided in the file stopwords.txt. Create a function text2bow(text) that converts a text into a bag of words following the just described specifications.

```

In [87]: def text2bow(text):

        ### Replace by your own code
        import re
        # preprocessing of text
        text = re.sub('[^a-zA-Z]+', ' ', text) # replace non-alphabetical characters with

```

```

text = text.lower() #lower case
text = text.split(' ') #split
# initialize bow as a set
bow = set()

# read stopwords.txt
stopwords = open('stopwords.txt').read().splitlines()

# filter out words that are stopwords and words < 3 characters
for word in text:
    if (word not in stopwords) & (len(word)>=3):
        bow.add(word)
return bow

```

Your bag-of-words implementation can be tested for the same text shown above by running the code below.

```

In [88]: print(textwrap.fill(str(text2bow(text))))

{'immediately', 'center', 'meet', 'know', 'defined', 'coplaner',
'circle', 'surface', 'wrong', 'cannot', 'infinity', 'circumference',
'angles', 'non', 'choose', 'bisector', 'graphics', 'yes',
'equidistant', 'distant', 'define', 'abc', 'collinear', 'one', 'well',
'lie', 'coincident', 'possibly', 'fitting', 'correct', 'solution',
'geometry', 'say', 'either', 'otherwise', 'exactly', 'unless',
'points', 'point', 'best', 'centre', 'algorithm', 'loss', 'radius',
'provide', 'books', 'close', 'pictures', 'specifies',
'straightforward', 'containing', 'least', 'must', 'far', 'quite',
'let', 'normally', 'plane', 'diameter', 'numerically', 'consider',
'exists', 'space', 'passing', 'four', 'two', 'see', 'hat', 'happen',
'farin', 'gems', 'line', 'way', 'relative', 'right', 'checked',
'still', 'please', 'equi', 'desired', 'check', 'sorry', 'could',
'error', 'subject', 'take', 'perpendicular', 'find', 'call',
'require', 'lies', 'fact', 'sphere', 'hyper', 'failure',
'intersection', 'mercy', 'three', 'may', 'necessarily', 'normal',
'need', 'since', 'steve', 'prove', 'bisectors'}

```

2.2 Implementing Bag-Of-Words Kernels (15 P)

In the following, your task is to implement a simple kernel over bag-of-words. The kernel between two bag-of-words \mathcal{X} and \mathcal{Y} is defined as

$$k(\mathcal{X}, \mathcal{Y}) = \sum_{w \in \mathcal{L}} 1_{w \in \mathcal{X} \wedge w \in \mathcal{Y}}$$

where $1_{w \in \mathcal{X} \wedge w \in \mathcal{Y}}$ is an indicator function testing membership to both bags of words. The language \mathcal{L} (set of all existing words) is typically unknown and very large. However, it is computationally equivalent to reduce the language \mathcal{L} to the union $\mathcal{X} \cup \mathcal{Y}$ of the two considered bag-of-words. Thus, we can rewrite the kernel as:

$$k(\mathcal{X}, \mathcal{Y}) = \sum_{w \in (\mathcal{X} \cup \mathcal{Y})} 1_{w \in \mathcal{X} \wedge w \in \mathcal{Y}}$$

Create a kernel method that implements this kernel function in a *naive* way. Your naive implementation will then be compared to an optimized one. The naive implementation can be summarized as follows:

- Iterate over all possible words w in $\mathcal{X} \cup \mathcal{Y}$.
- At each iteration test membership of w to \mathcal{X} and \mathcal{Y} .
- If both memberships are satisfied, increment the kernel score by 1. If not, leave it to its current value.

Remark: To test the membership of w to \mathcal{X} and \mathcal{Y} , do *not* use the operator "in" in Python, as it makes use of special data structures behind the scenes. Instead, iterate over all elements of \mathcal{X} and \mathcal{Y} using a for loop, and test membership using "==".

In [89]: `def kernel_naive(bow1,bow2):`

```

    ### Replace by your own code

    # to make this function in a naive way, convert the bows from set to str
    bow1 = list(bow1)
    bow2 = list(bow2)
    w = []

    # calculate the union set w of X and Y
    for i in range(len(bow1)):
        w.append(bow1[i])
    for i in range(len(bow2)):
        flag = True
        for j in range(len(w)): #without using operator "in"
            if bow2[i] == w[j]:
                flag = False
                break #break when found to save computation
        if flag:
            w.append(bow2[i]) #if bow2[i] not found in w, append w with it

    #Iterate over all possible words in w
    count = 0
    for i in range(len(w)):
        for j in range(len(bow1)):
            if w[i] == bow1[j]:
                for k in range(len(bow2)):
                    if w[i] == bow2[k]:
                        count += 1
                        break #break when found to save computation

    return count
    ###

```

The method `analyze_worstcase_performance(text2bow, kernel)` in `utils.py` computes the worst-case performance (i.e. when applied to the two longest texts in the dataset) of a specific kernel. Run the code below to test the performance of your implementation of the naive kernel.

```
In [90]: import utils
         utils.analyze_worstcase_performance(text2bow, kernel_naive)
```

kernel score: 761.000 , computation time: 1.510

This baseline implementation can be greatly accelerated (by a factor more than 100) by sorting the words in the bag-of-words in alphabetic order, and making use of the new sorted structure in the kernel implementation. In the code below, the sorted list associated to `bow1` is called `sbow1`. *Implement* a function `kernel_sorted(sbow1, sbow2)` that takes as input two lists of words (sorted in alphabetic order) and computes the kernel value in a more efficient manner. Like for the naive implementation, do *not* use the Python operator "in".

```
In [92]: def kernel_sorted(sbow1, sbow2):

         ### Replace by your own code
         bow1 = set(sbow1)
         bow2 = set(sbow2)
         w = bow1.intersection(bow2)
         return len(w)
         ###
```

The optimized kernel can be tested for worst case performance by running the code below. Here, we define an additional method `text2sbow(text)` for computing the sorted bag-of-words. Verify that the kernel score remains the same as with the naive implementation. The computation time is expected to drop drastically.

```
In [93]: def text2sbow(text): return sorted(list(text2bow(text)))

         import utils
         utils.analyze_worstcase_performance(text2sbow, kernel_sorted)
```

kernel score: 761.000 , computation time: 0.000

2.3 Classifying Documents with a Kernel SVM (15 P)

The kernel function between two text documents can be used to build a SVM-based text classifier. Here, we would like to discriminate between the two classes `comp.graphics` and `sci.med` present in the dataset. The code below reads the whole dataset and stores input (mapped to sorted bag-of-words) and labels in the appropriate data structures.

```
In [94]: import numpy
         Xtrain = map(text2sbow, open('newsgroup-data/newsgroup-train-data.txt', 'r'))
         Xtest  = map(text2sbow, open('newsgroup-data/newsgroup-test-data.txt', 'r'))
```

```

Ttrain = numpy.array(list(map(int,open('newsgroup-data/newsgroup-train-label.txt','r')
Ttest  = numpy.array(list(map(int,open('newsgroup-data/newsgroup-test-label.txt','r')

#converting map to list (because I am using python3)
Xtrain = list(Xtrain)
Xtest  = list(Xtest)

```

As a first step, one needs to build the kernel matrices between pairs of training examples and between training and test examples. After evaluating whether building such matrices is computationally feasible given the performance of your optimized bag-of-words kernel implementation, write the function `build_kernels(Xtrain,Xtest)` for constructing these matrices.

```

In [95]: def build_kernels(Xtrain,Xtest):

    ### Replace by your own code
    Ktrain = numpy.zeros([len(Xtrain),len(Xtrain)])
    Ktest  = numpy.zeros([len(Xtest),len(Xtrain)])
    for i in range(len(Xtrain)):
        for j in range(len(Xtrain)):
            Ktrain[i,j] = kernel_sorted(Xtrain[i],Xtrain[j])

    Ktest = numpy.zeros([len(Xtest),len(Xtrain)])
    for i in range(len(Xtest)):
        for j in range(len(Xtrain)):
            Ktest[i,j] = kernel_sorted(Xtest[i],Xtrain[j])
    ###

    assert(Ktrain.shape==(len(Xtrain),len(Xtrain)) and Ktest.shape==(len(Xtest),len(X
    return Ktrain,Ktest

```

These kernel matrices along with the vector of training labels `Ttrain` can be used to train an SVM in the same way as in the previous exercise on genes sequences classification. Write a function that trains an SVM (using scikit-learn with default parameters) and computes the predictions on the training and test data.

```

In [96]: def get_svm_prediction(Ktrain,Ttrain,Ktest):

    ### Replace by your own code
    mysvm = svm.SVC(kernel='precomputed').fit(Ktrain,Ttrain)
    Ytrain = mysvm.predict(Ktrain)
    Ytest  = mysvm.predict(Ktest)
    ###

    assert(Ytrain.shape==Ttrain.shape and Ytest.shape==Ttest.shape)
    return Ytrain,Ytest

```

Finally, the functions that you have implemented for classifying the texts can be tested by measuring the training and test accuracy.


```
In [97]: Ktrain,Ktest = build_kernels(Xtrain,Xtest)
        Ytrain,Ytest = get_svm_prediction(Ktrain,Ttrain,Ktest)

        print('training accuracy: %.3f    test accuracy: %.3f'% \
              ((Ytrain==Ttrain).mean(),(Ytest==Ttest).mean()))

training accuracy: 1.000    test accuracy: 0.962
```