

Part A: Kernels for DNA Sequences (20 P) ¶

In this first part the weighted degree kernel (WDK) will be implemented for the purpose of classifying DNA sequences. We will use Scikit-Learn (<http://scikit-learn.org/>) for training SVM classifiers. The focus of this exercise is therefore on the computation of the kernels. The training and test data is available in the folder `splices-data`. The following code reads the DNA sequence data and stores it in numpy arrays of characters.

In [1]:

```
import numpy
Xtrain = numpy.array([numpy.array(list(l.rstrip('\r\n'))) for l in open('splice-data/splice-train-data.txt', 'r')])
Xtest  = numpy.array([numpy.array(list(l.rstrip('\r\n'))) for l in open('splice-data/splice-test-data.txt', 'r')])
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')])
```

We consider the weighted degree kernel described in the lecture. It applies to two genes sequences $x, z \in \{A, T, G, C\}^L$ and is defined as:

$$k(x, z) = \sum_{m=1}^M \beta_m \sum_{l=1}^{L-m+1} I(u_{l,m}(x) = u_{l,m}(z))$$

where l iterates over the whole genes sequence, $u_{l,m}(x)$ is a subsequence of x starting at position l and of length m , and $I(\cdot)$ is an indicator function that returns 1 if the argument is true and 0 otherwise. We would like to implement a function that is capable of *efficiently* computing this weighted degree kernel for any degree M . For this, we will make use of the block method presented in the lecture.

As a first step, we would like to implement a function `size2contrib`, which builds a mapping from a given block size to the kernel contribution, i.e. the sum of beta values associated to all substrings contained in this block. The relation between block size and contribution to the kernel score is as follows:

- Block size 1: contribution = β_1
- Block size 2: contribution = $2\beta_1 + \beta_2$
- Block size 3: contribution = $3\beta_1 + 2\beta_2 + \beta_3$
- etc.

The function should return an integer array of size 101 containing the contribution of blocks of size zero, one, two, up to 100.

In [2]:

```
def size2contrib(beta):
    s2ctable_length = 101
    beta_length = beta.shape[0]
    beta = numpy.flip(beta)
    s2ctable = numpy.zeros(s2ctable_length)
    contribs = numpy.concatenate( (numpy.zeros((beta_length,)), numpy.arange(1,101)) )
    for i in range(s2ctable_length):
        contribs = numpy.array([])
        s2ctable[i] = beta.dot(contribs[i:i+beta_length])
    return s2ctable
```

The function can be tested on a simple weighted degree kernel of degree 3 where beta coefficients are given by $\beta_1 = 1, \beta_2 = 3, \beta_3 = 9$.

In [3]:

```
size2contrib(numpy.array([1.0,3.0,9.0]))[:20]
```

Out[3]:

```
array([ 0.,  1.,  5., 18., 31., 44., 57., 70., 83., 96., 109.,
       122., 135., 148., 161., 174., 187., 200., 213., 226.])
```

Having implemented this index, we now focus on implementing the weighted degree kernel. Here, the function `wdk` we would like to implement receives two data arrays `X` and `Z`, and some parameter vector β . The function should return the kernel Gram matrix associated to `X` and `Z`, and run *efficiently*, i.e. as much as possible performing operation over several data points simultaneously by means of array computations. (Hint: An array of block sizes can be transformed to an array of kernel contributions by using the indexing operation `blockcontribs = s2ctable[blocksizes]`.)

In [4]:

```
def wdk(X,Z,beta):

    K = numpy.zeros((X.shape[0],Z.shape[0]))
    weighted_contribs = size2contrib(beta)

    for i in range(X.shape[0]):
        total_score_row = numpy.zeros((Z.shape[0],))
        block_size = numpy.zeros((Z.shape[0],))
        x = X[i]
        for j,c in enumerate(x):
            matches = ( Z[:,j] == c)
            mismatches_index = numpy.where(matches == 0)
            total_score_row[mismatches_index] += weighted_contribs[block_size[mismatches_index].astype(int)]
            block_size[mismatches_index] = 0
            block_size += matches
        total_score_row += weighted_contribs[block_size.astype(int)]
        K[i] = total_score_row
    return K
```

Once the weighted degree kernel has been implemented, the code below trains SVMs on the classification task of interest for different choices of parameters β .

In [6]:

```

from sklearn import svm

for beta in [
    numpy.array([1]),
    numpy.array([1,3,9]),
    numpy.array([0,0,0,1,3,9]),
]:
    Ktrain = wdk(Xtrain,Xtrain,beta)
    Ktest = wdk(Xtest,Xtrain,beta)
    mysvm = svm.SVC(kernel='precomputed').fit(Ktrain,Ttrain)
    print('beta = %-20s training: %.3f test: %.3f'%(beta,mysvm.score(Ktrain,Ttrain),
    mysvm.score(Ktest,Ttest)))

```

```

beta = [1]                training: 0.994 test: 0.916
beta = [1 3 9]            training: 1.000 test: 0.963
beta = [0 0 0 1 3 9]      training: 1.000 test: 0.933

```

We observe that it is necessary to include non-unigram terms in the kernel computation to achieve good prediction performance. If however, we rely mainly on long substrings, e.g. 4-, 5-, and 6-grams, there are not sufficiently many matchings in the data to obtain reliable similarity scores, and the prediction performance decreases as a result.

Part B: Kernels for Text (20 P)

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 composed only of texts of classes `comp.graphics` and `sci.med`. 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 [7]:

```

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 [...]

```

Converting Texts to a Set-Of-Words

A convenient way of representing text data is as "set-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 set-of-words representation is provided in the file `stopwords.txt`. Create a function `text2sow(text)` that converts a text into a set of words following the just described specifications.

In [10]:

```
stopwords_newline = list(open('stopwords.txt','r'))
stopwords = []
for word in stopwords_newline:
    stopwords.append(word.strip())

import re
def text2sow(text):
    words= text.split()
    clean_words = [clean_word.lower() for word in words for clean_word in re.findall(
'[a-zA-Z]{3,}', word) ]
    textSet=set(clean_words)
    stopwordSet=set(stopwords)
    return textSet.difference(stopwordSet)
```

The set-of-words implementation is then tested for the same text shown above:

In [11]:

```
print(textwrap.fill(str(text2sow(text))))
```

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

Implementing the Set-Of-Words Kernel

The set-of-words kernels between two documents x and z is defined as

$$k(x, z) = \sum_{w \in \mathcal{L}} I(w \in x \wedge w \in z)$$

where $I(w \in x \wedge w \in z)$ is an indicator function testing membership of a word to both sets of words. As for the DNA classification exercise, it is important to implement the kernel in an efficient manner.

The function `benchmark(text2sow, 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 implementation. Here, the function is tested on some naive implementation of the set-of-words kernel available in `utils.py`.

In [12]:

```
import utils
utils.benchmark(text2sow, utils.naivekernel)
```

kernel score: 761.000 , computation time: 0.605

The goal of this exercise is to accelerate the procedure by sorting the words in the set-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 `sow1` is called `ssow1`. *Implement* a function `sortedkernel(ssow1, ssow2)` that takes as input two sets of words (sorted in alphabetic order) and that computes the kernel score in an efficient manner, by taking advantage of the sorting structure.

In [13]:

```
def sortedkernel(ssow1, ssow2):

    total_score = 0
    sow1_size = len(ssow1)
    i = 0
    words_dict = {}
    while i in range(sow1_size):
        match_word = ssow1[i]
        matches1 = 1
        i+=1
        while i < sow1_size and ssow1[i] == match_word :
            matches1+= 1
            i+=1
        words_dict[match_word] = matches1

    for w2 in ssow2:
        if w2 in words_dict.keys():
            total_score += words_dict[w2]

    ## -----

    return total_score
```

This efficient implementation of the set-of-words kernel can be tested for worst case performance by running the code below. Here, we define an additional method `text2ssow(text)` for computing the sorted set-of-words.

In [16]:

```
def text2ssow(text): return sorted(list(text2sow(text)))

import utils
utils.benchmark(text2ssow,sortedkernel)
```

kernel score: 761.000 , computation time: 0.001

The kernel score remains the same, showing that our new sorted implementation still produces the same function, however, the computation time has dropped drastically.

Classifying Documents with a Kernel SVM

The set-of-words kernel implemented above can be used to build a SVM-based text classifier. Here, we would like to separate our two classes `comp.graphics` and `sci.med`. The code below reads the whole dataset and stores it in a sorted set-of-words format.

In [17]:

```
import numpy
Xtrain = list(map(text2ssow,open('newsgroup-data/newsgroup-train-data.txt','r')))
Xtest  = list(map(text2ssow,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')
))))
```

Kernel matrices are then produced using our efficient kernel implementation with pre-sorting, and a SVM can be trained to predict the document class.

In [18]:

```
Ktrain = numpy.array([[sortedkernel(ssow1,ssow2) for ssow2 in Xtrain] for ssow1 in Xtrain])
Ktest  = numpy.array([[sortedkernel(ssow1,ssow2) for ssow2 in Xtrain] for ssow1 in Xtest])
mysvm = svm.SVC(kernel='precomputed').fit(Ktrain,Ttrain)
print('training: %.3f  test: %.3f'% (mysvm.score(Ktrain,Ttrain),mysvm.score(Ktest,Ttest)))
```

training: 1.000 test: 0.962

In []: