

## 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 [1]: import numpy
Xtrain = numpy.array([numpy.array(list(l)) for l in open('splice-data/splice-train-data.txt', 'r')])
Xtest  = numpy.array([numpy.array(list(l)) 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')])
```

### 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 [2]: def getdegreekernels(Xtrain,Xtest,degree):

    ### Replace by our own code
    import solutions
    Ktrain,Ktest = solutions.getdegreekernels(Xtrain,Xtest,degree)
    ###

    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 [3]: from sklearn import svm
Ktrains,Ktests = [None]*4,[None]*4

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()))

degree: 1    training accuracy: 0.994    test accuracy: 0.916
degree: 2    training accuracy: 1.000    test accuracy: 0.934
degree: 3    training accuracy: 1.000    test accuracy: 0.964
degree: 4    training accuracy: 1.000    test accuracy: 0.956
```

## 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 [4]: ### Replace by our own code
import solutions
solutions.wdk(Ktrains, Ktests, Ttrain, Ttest)
###
```

training accuracy: 1.000    test accuracy: 0.967

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

```
count, I think. Most of them have >said either "You have mono" or
"You're full of it; there's nothing wrong >with you." One has
admitted to having no idea what was wrong with her, >and one has
claimed that it is Epstein-Barr syndrome. > >Now, what she told me
about EBS is that very few doctors even believe that >it exists.
(Obviously, this has been her experience.) So, what's the >story?
Is it real? Does the medical profession believe it to be real? > >Has
anyone had success is [...]
```

## Creating Bag-Of-Words (15 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 [6]: def text2bow(text):

    ### Replace by your own code
    import solutions
```

```

bow = solutions.text2bow(text)
###

return bow

```

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

```

In [7]: print(textwrap.fill(str(text2bow(text))))

set(['saying', 'all', 'help', 'just', 'entity', 'discovered',
'thanks', 'yes', 'still', 'probs', 'whose', 'immune', 'had', 'either',
'jan', 'annals', 'combo', 'has', 'worth', 'renaming', 'real', 'info',
'them', 'very', 'docs', 'disagreement', 'not', 'candida', 'now',
'massively', 'treating', 'astonishly', 'outright', 'like', 'success',
'disregulation', 'did', 'retrovirus', 'admitted', 'try', 'these',
'she', 'told', 'barr', 'prompted', 'exists', 'people', 'some', 'idea',
'ebv', 'are', 'ebs', 'bacteria', 'even', 'what', 'said', 'for',
'fevers', 'assistance', 'since', 'centers', 'does', 'medicine',
'newly', 'cause', 'epstein', 'full', 'genes', 'cfids', 'stress',
'same', 'herpes', 'about', 'mono', 'wrong', 'experience', 'etc',
'place', 'massive', 'think', 'among', 'seems', 'profession',
'thereof', 'one', 'psychological', 'chronic', 'because', 'viruses',
'least', 'your', 'dysfunction', 'considering', 'story', 'her',
'outbreaks', 'nightsweats', 'there', 'been', 'anyone', 'few', 'live',
'was', 'tell', 'sort', 'partly', 'knows', 'that', 'took', 'but',
'moment', 'environmental', 'doctors', 'believe', 'with', 'count',
'include', 'toxins', 'this', 'originally', 'official', 'theories',
'were', 'called', 'and', 'say', 'cure', 'something', 'have',
'claimed', 'seem', 'different', 'syndrome', 'thing', 'amounts',
'things', 'recurrent', 'elevated', 'also', 'causing', 'internal',
'which', 'you', 'clear', 'may', 'glands', 'who', 'variant', 'most',
'virus', 'levels', 'fatigue', 'nothing', 'thoughts', 'antibodies',
'maybe', 'medical', 'disease', 'obviously', 'swollen', 'the',
'having'])

```

## 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 [8]: def kernel_naive(bow1,bow2):

        ### Replace by your own code
        import solutions
        return solutions.kernel_naive(bow1,bow2)
        ###
```

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 [9]: import utils
        utils.analyze_worstcase_performance(text2bow, kernel_naive)

kernel score: 827.000 , computation time: 0.815
```

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 [10]: def kernel_sorted(sbow1,sbow2):

        ### Replace by your own code
        import solutions
        return solutions.kernel_sorted(sbow1,sbow2)
        ###
```

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 [11]: def text2sbow(text): return sorted(list(text2bow(text)))

        import utils
        utils.analyze_worstcase_performance(text2sbow, kernel_sorted)

kernel score: 827.000 , computation time: 0.002
```

## Classifying Documents with a Kernel SVM (20 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 [12]: 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(map(int, open('newsgroup-data/newsgroup-train-label.txt', 'r')))
        Ttest  = numpy.array(map(int, open('newsgroup-data/newsgroup-test-label.txt', 'r')))
        print(len(Xtrain), len(Xtest))
```

(134, 106)

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 [13]: def build_kernels(Xtrain,Xtest):

    ### Replace by your own code
    import solutions
    Ktrain,Ktest = solutions.build_kernels(Xtrain,Xtest)
    ###

    assert(Ktrain.shape==(len(Xtrain),len(Xtrain)) and Ktest.shape==(len(Xtest),len(Xtrain)))
    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 [14]: def get_svm_prediction(Ktrain,Ttrain,Ktest):

    ### Replace by your own code
    import solutions
    Ytrain,Ytest = solutions.get_svm_prediction(Ktrain,Ttrain,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 [15]: 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