

Classification of Human Brain Image Data using Support Vector Machines

Internship Report

Philipp Glock

February 2, 2015

Acknowledgments

I would like to thank the JSC for its support during my internship. A special thanks goes to Markus Axer and Stefan Köhnen for providing the brain data used in this internship.

1 Introduction

The Research Centre Jülich investigates key technologies for the 21st century. Different tools are needed for all kinds of experiments. Because these tools are often very specialized, the Research Centre often develops them itself. Apart from the experiment related software, there are many computational simulation and data problems that have to be solved.

The Jülich Super Computing Centre (JSC) provides the resources for these computations such as the supercomputing capacity and capability needed for a fast calculation. The JSC has several high performance clusters that can be used for compute- and/or data-intensive problems [1]. Besides hardware the institute investigates, explores, and develops different software solutions for parallel computing.

The department “Federated Systems and Data” (FSD) of the Jülich Super Computing Centre provides access to distributed systems that consist on a wide variety of different resources (e.g. HPC systems, parallel filesystem environments, high throughput computing resources, or clouds). This is done by implementing open standards and simplifying usage and administration of these services. Furthermore the division provides in particular, middleware services, simple upload / download services, replication services, or data management know-how in general.

The research group on High Productivity Data Processing works on solutions to overcome problems and challenges of applications that specifically require big data analytics solutions. This can be splitted into three categories.

1. **Investigate Generic Data Methods:** How to overcome limitations of processing and analyzing large amounts of data (e.g. in-memory databases, data privacy methods and query processing).
2. **Machine Learning:** Explore, develop and tune serial or parallel machine learning techniques, like classification and clustering, in order to enable solutions that work with large quantities of scientific data or high dimensional datasets.
3. **Smart Data Analytics Applications:** Find and develop solutions that are specifically applicable in real-world applications for general data analysis and statistical data analysis.

During the internship in this group the intention is to analyse, evaluate and build classification models for data sets that show problems when working with ‘big data’ . As one of the best out-of-the-box solutions preferred by application domain scientists, Support Vector Machines (SVMs) are the focused method, but other methods are considered

as well. This report highlights selected different problems and their solutions based on several approaches and implementations by taking advantage of the scikit-learn module.

After the Introduction, this report will provide the necessary background in Section 2, while Section 3 highlights one application domain that represents the origin of the data-sets analysed. In Section 4 the process of developing a classification model is described and in Section 5 the results are evaluated. Section 6 finishes the report with a summary and Conclusion.

2 Background

Given the tremendous uptake of python in scientific working environments, the primary programming language used during this internship is python. It is a dynamic scripting language and provides many nice modules for scientific use cases like numpy [19] and scipy[10]. There are also some machine learning modules and the scikit-learn module [11] is one of the most used and stable ones. While python alone may not have the best performance, it is easy to improve it by using c or fortran code. This is done by most libraries. The serial SVM algorithms of scikit-learn are based on the well known libsvm library [5]. The adoption of this stable library makes python and the scikit-learn module a good choice for classification tasks that play an instrumental part in the application domain problems in the internship.

Some tasks are executed on a super computer at the Juelich Supercomputing Centre. It uses a Portable Batch System (PBS) [9] based batch system to handle the different jobs of users and distribute the available nodes. A popular method for distributed memory multiprocessing is the Message Passing Interface (MPI) [7]. The processes do not share any memory and therefore data has to be communicated over messages. This ranges from simple send and receive functions to more complex ones like broadcast or scatter and gather functions. The left side of figure 2.1 shows an illustration of MPI multiprocessing. The red arrows are the messages that are passed between the processes, which are represented by the blue arrows.

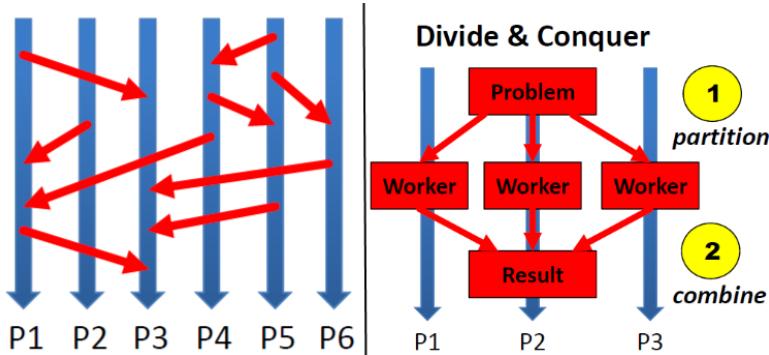


Figure 2.1: Multiprocessing on P1-P6 processors using MPI (left) and on P1-P3 processors using map/reduce (right).

An alternative to MPI based multiprocessing is the more recent map/reduce approach [6] used by Apache Hadoop. That also takes advantage of a distributed filesystem like the Hadoop Distributed File System (HDFS) [16]. The data is split into subsets and distributed on the single machines that often also include a replication strategy (e.g. 3 times the same data in HDFS at different sites). On each machine the same map job is executed on the local data set. The fundamental idea of Hadoop and HDFS is therefore to 'transfer the job to the data' in contrast to transfer the data to the job resources as typically done by traditionally MPI approaches. The results of the map jobs are saved as key value pairs, so that the reducer jobs can further evaluate them after

Listing 1: pbs.controller.template

```
#PBS -l nodes={n/4}:ppn=4
#PBS -l walltime=24:00:00
#PBS -l pvmem=1GB
#PBS -N ipython-engine-glock
#PBS -m abe
#PBS -j oe
#PBS -o ipc.engine-$PBS_JOBID.out

module load python
module load intel
module load mvapich2

cd $PBS_O_WORKDIR
which ipengine
mpiexec -np {n} ipengine --profile=pbs
```

a sort/shuffle/group process of keys performed by the map-reduce run-time itself. The number of algorithms that can profit from this method are limited because they need to be seperatable into a map and a reduce job. There are some frameworks that can handle more general algorithms, including data mining [14, p. 1-18].

IPython [12] is an interactive python shell with many functionalities. One of its most known features is the ipython notebook. It is browser based and supports code, text, mathematical expressions, plots and other media. Because of this it is often used in teaching or science to save code and calculations for a better reproducibility . Another feature is its capabilty of parallel computing. It supports different styles of parallelism like task or data parallelismwell suited for High Throughput Computing (HTC) resources and MPI for HPC resources. This allows parallelizing python code in an interactive python session and has a robust error handling. Furthermore the ipython clusters can be started with a PBS scheduler.

After setting up a profile and creating PBS templates (e.g. listing 1), you can start a cluster easily using:

```
ipcluster start --profile = pbs - n 8
```

More information on ipython in general and the setup for parallel usage can be found in [2].

3 Human Brain Project - Brain Analytics Case Study

Simulating the human brain in order to get a better understanding of it, is the intention of the human brain project [3]. Scientists from many countries create an infrastructure to link brain research and information technology. The project unites neuro and computer scientists as well as mathematicians, physicists and medical experts as well as scientists all over Europe.

The Research Centre Jülich researches basics about structure and functioning of the brain. Apart from that the JSC develops new data processors and software that can handle the huge amount of data about the brain.

One sub task of this project is to build a 3D model of the human brain. The brain analytics data set created by the institute of neuroscience and medicine (INM) provides images of different cross sections of a brain. The cross sections of one brain were manually labelled in a time-consuming process, so that each pixel either belongs to the cross section or not. Labeling a whole brain takes a huge amount of time and resources and as such an automation process that supports this process is one of the research goals of this case study. So instead of labeling new images by hand, this data is used to build a model that automatically classifies a pixel. If this classification is good enough, the predicted 2D brain parts should be used to rebuild a 3D model of a human brain. Figure 3.1 shows such a cross section. The trained classifier then in turn can be applied to unseen new future brain data images when they become available in the project.

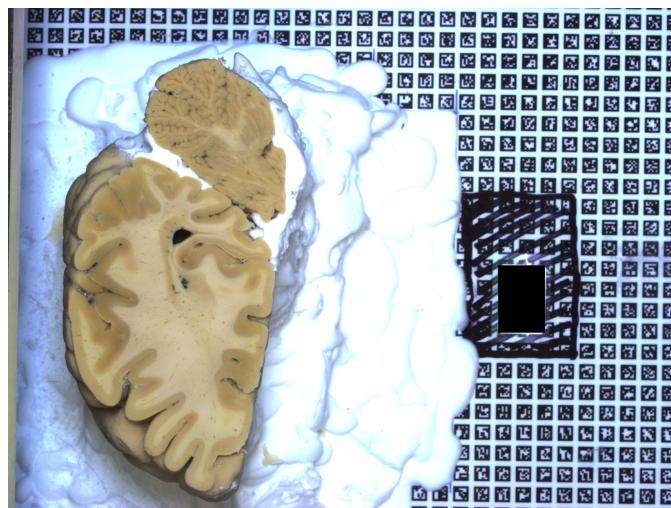


Figure 3.1: An image of a cross section of a human brain

The reconstructed 3D brain and its related research contributes to the understanding of creating computational simulations of the human brain, it can be used by medical scientists to better understand healthy and diseased brains and to create new drugs. But it can also support the development of high performance computers which are also energy efficient, whereby the modelling is inspired from the human brain.

4 Scientific Data Analysis

This section starts by giving an overview over the used scientific dataset in Section 4.1 and the data preparation described in Section 4.2. This includes sampling and feature extraction. This is followed by the Data Modelling in Section 4.3 and the Deployment with an IPython notebook in Section 4.4.

4.1 Scientific Dataset

The data is not given as a table with numeric values and is an image-based dataset and therefore not directly useable for classification. It consists of two sets of images. The first one contains the original images of each cross section as RGB images (Figure 3.1). The second type is a mask for each of the different RGB images thus providing labelled data (Figure 4.1), a gray scale image with only black or white pixels. A pixel is white if it belongs to the cross section and black if it does not. There are 762 images (and additional 762 mask images) with a resolution of 3272×2469 .



Figure 4.1: Mask of a cross section

As stated above there are two classes in total, so that the task is a binary classification problem. Before a classifier can be trained, the images have to be converted into a attribute value format and saved in a file. This could be a simple csv file or any other format that can handle numeric like data. In this case the libsvm format was chosen, because it is supported by many svm libraries, e.g. scikit-learn. Every row represents one instance, which is in this case one pixel and its features. Without applying any feature extraction method, the dataset (called A) has only the RGB values of a pixel from the original image and the class label from the mask. Listing 2 shows the libsvm format schematically.

Listing 2: libsvm format

```
<class_label> <feature_id>:value ... <feature_id>:value  
<class_label> <feature_id>:value ... <feature_id>:value ...
```

4.2 Data Preparation

In the first preprocessing step the image-based data is transformed into the libsvm format. The easiest method saves for each pixel the RGB values and the class label. This is done for every image resulting in 762 files. Every file has 8.078.568 instances, so there are 6.155.868.816 pixels in total. The data preparation process in itself is already a challenging task for serial approaches and a time-consuming process. A first trivial approach took about 3 days to generate all attribute value files. The RGB values can be read directly from the original images. The class label is set to 1 of the corresponding pixel value in the mask is greater than 0. If this is not the case it is set to 0. This was done with a python script, which was executed on Judge, a supercomputer of the JSC with 2472 nodes and 412 graphic processors. The python script is started with a msub job submission script reserving one core with one node. To speed up the calculation time, the python script was changed so that it expects a range as parameters that says which image files it should convert. With this several independent jobs can be started with a different range argument resulting in a faster computation due to parallelization.

Building a model on the whole set would take a huge amount of time and is therefore unreasonable for creating a classifier. Instead of using the complete data set samples have to be used. As figure 4.1 shows there are more instances of class 0 (black pixels) than of class 1 (white pixels). To weight both classes equally, the sampled set is balanced.

The first approach is a balanced random sampling method. For a given percentage p and a data set with N instances the method draws $p \cdot N$ samples from a file. If the sample is balanced, $p \cdot N \cdot 0.5$ instances of each class are drawn randomly. This can easily be parallelized since the files can be sampled independently.

An alternative method is to resize the images to a much smaller resolution (e.g. 256×256). This is done for a number of images that are uniformly distributed among the different layers of the brain. An image has several different sections like the cross section, the brain in the background or the ice. If random sampling is used, it may happen that some of these sections are not represented by the sample. The lower the percentage p is, the higher the chance gets that some sections are missed.

Because of this later samples were made by resizing the images.

This yields a first sample of dataset A. To improve the classifier the number of features has to be increased.

4.2.1 Hue-Saturation-Value color space

A feature extraction method is used in order to provide more meaningful data to be used in learning the classifier. To increase the number of features another color space

in addition to the given RGB values is used. In this case the HSV color space was chosen, because it is similar to the human color vision. Furthermore it is used for object detection [18] because methods used on grayscale images can easily be used on images in the HSV space by using them on each color component.

The HSV color space has three components hue, saturation and value. The hue (H) is the angle on the chromatic circle and is therefore between 0 and 360. The second and third components saturation (S) and lightness (V) are given in percentage, so their value is between 0 and 1 (or 0 – 100%). If the saturation is 1, the color is clear. For a low saturation the color becomes gray. If the lightness is 0, the color is black.

The values for a pixel can easily be transformed from RGB to HSV with the following formula:

$$\begin{aligned}
 Max &= \max(R, G, B) \\
 Min &= \min(R, G, B) \\
 H &= \begin{cases} 0, & \text{if } Max = Min \Leftrightarrow R = G = B \\ 60^\circ \cdot \left(0 + \frac{G-B}{Max-Min}\right) & \text{if } Max = R \\ 60^\circ \cdot \left(2 + \frac{B-R}{Max-Min}\right) & \text{if } Max = G \\ 60^\circ \cdot \left(4 + \frac{R-G}{Max-Min}\right) & \text{if } Max = B \end{cases} \\
 S &= \begin{cases} 0 & \text{if } Max = 0 \Leftrightarrow R=G=B=0 \\ \frac{Max-Min}{Max} & \text{else} \end{cases} \\
 V &= Max
 \end{aligned}$$

Since this computation can be done independently for every pixel, the features can easily be computed in parallel and be added to existing data. While computing the HSV values is simple, they also add only few additional information. Dataset B has HSV features additional to the RGB features.

4.2.2 Local Features

Despite adding the HSV color space, there is still no information on the neighbourhood of a pixel that in many cases brings enormous results in classification accuracy of images (e.g. land-cover classification of remote sensing images [4]). A first level approach to get information on the neighbourhood of a pixel are statistical values like mean or standard deviation. The statistical values are calculated on an $N \times N$ array around a pixel and assigned to the centered pixel. This is done for every pixel of the image. Listing 3 shows how the standard deviation is calculated for an image. This is done with the numpy and scipy modules, which are optimized.

The *generic_filter* function generates a window with the given size for every pixel and calls the given function (e.g. `numpy.std`). However, the parameter of the called function is an one dimensional array. If the function needs a two dimensional one, it has to know the original shape of the window and reshape it. For most statistical measures this is

not needed. The standard deviation is an indication for the homogeneity of the pixel's neighbourhood.

Listing 3: Calculate standard deviation

```

1 def calc_std(img, size, use_mean=False):
2     """Calculate the local standard deviation for each pixel of an image.
3     Arguments:
4         img : rgb image
5         size : size of the window of neighbouring pixels
6         use_mean : boolean, if False the standard deviation of R, G and B
7             value are calculated and return as a 3 dimensional array.
8             if True the standard deviation of the mean is returned.
9     Return:
10        array_like, 3d or 2d
11    """
12    if not use_mean:
13        std = np.zeros(img.shape)
14        for i in range(img.shape[2]):
15            std[:, :, i] = filters.generic_filter(img[:, :, i], np.std, size)
16    else:
17        mean_img = np.mean(img, axis=2)
18        std = filters.generic_filter(mean_img, np.std, size)
19    return std

```

By using the local standard deviation as a feature some first information of the neighbourhood is available. This can be further improved by other features.

4.2.3 Image Segmentation

Watershed [15] is an image segmentation algorithm that is used on grayscale images. The gray level of a pixel is interpreted as its altitude and the image is seen as a ground. Water, that is filled into this ground from different markers, flows to a local minimum.

The algorithm floods basins from markers, which were set by the user. When two basins from different markers meet a watershed is drawn. The markers can be set in an easy way by using thresholds. While there are more complex ways using image processing methods, the threshold approach is used to keep the feature generation as simple as possible.

This yields a binary feature that is either one or zero and is used for training the classifier. Dataset C has RGB,HSV, standard deviation of RGB and watershed as features.

4.3 Data Modeling

4.3.1 Support Vector Machines

Support vector machines (SVMs) are one of the preferred classification methods lately, because tools are stable and widely available and it is one of the best out-of-the-box methods scientists can use that are not particularly trained in machine learning algorithms.

They have a high accuracy, but their training time is quite long. A model can easily be described by the found support vectors. Figure 4.2 shows a basic visualization of an SVM. On the left side a line is drawn that separates both classes. This is the decision boundary. The data in the right side is not linear separable. A kernel is used to push the data into a higher dimensional space, so that the classes are linear separable in this space.

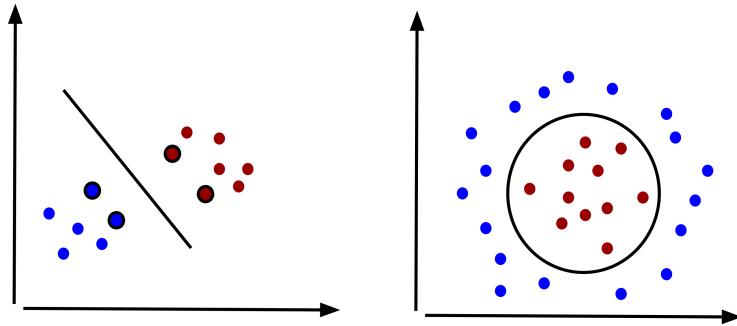


Figure 4.2: Classification using Support Vector Machines. Linear example on the left side (support vectors have a black border). Non linear example on the right side.

A kernel defines the used scalar product. The most popular kernels are the following:

- linear: $K(x, y) = \sum_{i=0}^n x_i \cdot y_i$
- polynomial: $K(x, y) = (c + \sum_{i=0}^n x_i \cdot y_i)^d$
- rbf: $\exp(-\frac{\|x-y\|^2}{2\sigma^2})$

Unlike other classification methods (e.g. neural networks, naive bayes) SVMs need only a small training set in theory. Additional instances only affect the classifier if they contain new support vectors. However depending on the chosen kernel cross-validation runs are necessary to determine the parameter for the kernels (e.g. rbf works with parameter gamma).

Figure 4.3 shows the predicted labels with a SVM on dataset A and with a RBF kernel. More information on SVMs in general can be found at [8].

Because support vector machines maximize the margin between the decision boundary and the support vectors, they do not overfit as easily as other classifiers like decision trees for example. However, solving the quadratic programming problem can have a complexity between $O(n^2)$ and $O(n^3)$ depending on the used algorithm and implementation. As figure 5.1 shows training an svm takes much more time to be trained than other classifiers. This is especially true if a non linear kernel is used, e.g. rbf kernel. A linear svm is much faster but this results in a drop in accuracy. To get both a feasible training time and a good accuracy, a kernel approximation is used together with a linear SVM. A kernel



Figure 4.3: predicted image using svm with rbf kernel

approximation is used before training a SVM (see listing 4). It adds additional random features to the data. For the rbf kernel this is done by a Monte Carlo approximation of its Fourier transform. More information on random fourier features and random binning features can be found at [13].

Listing 4: kernel approximation using scikti learn

```
from sklearn.kernel_approximation import RBFSampler
rbf = RBFSampler(gamma=2)
X_features = rbf.fit_transform(X_train)
X_test_features = rbf.transform(X_test)
```

A second approach is the Nyström method [20] which can approximate every kernel and not only the rbf one. It uses a subsample of the data set to approximate a kernel. As [21] shows the nyström method can achieve a better generalization in some cases.

4.3.2 Other Classifiers

There are some other classifiers besides support vector machines. While the focus is on SVMs and they were optimized the most, some models with other classifiers were tested too. Naive Bayes is one of them. Due to its simplicity it is quite fast. Unfortunately the accuracy score is not good since Naive Bayes assumes independency of features which is not the case for RGB values of an image.

While KNearrestNeighbour considers feature dependencies and therefore has a higher accuracy, it's training time is also higher than the one of naive bayes. Apart from that

table 4.1 shows that the prediction takes much longer since the distances to all instances have to be calculated.

Decision trees without any pruning or ensemble methods have a better accuracy than naive bayes as well but not as good as KNN. While their training time is larger than the one of naive bayes, it is way faster than KNN. For the random forest classifier, which uses ten randomized decision trees, the accuracy as well as the training and testing time increase.

Classifier	Accuracy	training time in s	test time in s
Decision Tree	0.9569552165	0.563767	0.057934
Random Forest	0.9637635858	1.698716	0.399092
KNN	0.9677196684	1.57	5.29
GaussianNB	0.9475908597	0.042224	0.075098
SVM	0.9579808675	45.242423	66.846724

Table 4.1: Different classifiers used on Dataset A, 0.1% sample.

4.3.3 Multiple Models

With only one model for all the different cross section the prediction is not very good for some cross sections. While it is not reasonable to use the x and y coordinates of a pixel as features because the brains position changes between the images, it is possible to use the z axis as a feature. The z position of an image is contained in its file name, since all the images are numbered. Naturally, for this to work future cross sections have to be numbered in the same way. The alternative to use the z position as a feature with one model is to divide the brain in several layers and create a model for every layer. This however could have impacts for the classification of future brains that might differ in volume, orientation, and location of brain elements.

There are also advantages in comparison to using one model. By increasing the number of models each single model is simpler. It also decreases the training time. This has two reasons. The first is that the amount of training data for each model is less because the data is sliced into different layers. In addition the training of the models can easily be parallelized because they are completely independent from each other.

4.3.4 Grid Search

While scikit-learn already implements a function that performs grid search on a given set of parameters and a classifier, this function can not be used if kernel approximation is used as well. The parameters that have to be optimized are the penalty parameter C and RBF Kernel coefficient γ . If kernel approximation is used, the linear SVM does not expect an argument γ . Instead γ is set in the RBFSampler. Therefore the scikit-learn function can not be used and a modified version was implemented. The new version generates a new data set C^* with the additional features and the current γ . Thereafter a linear SVM is trained with the given C and cross validated. All results are saved in a

list and returned. To speed up the grid search, it is parallelized with IPython using the LoadBalancedView. Unlike the DirectView interface this one does not allow direct access to the individual engines. Instead of that the IPython scheduler assigns the tasks to the engines minimizing the idle time of the engines. Because of this the grid search can not only be used on a single machine but on clusters, if the IPython cluster is setup.

$g = 0.1, C = 1$	$g = 0.1, C = 10$	$g = 0.1, C = 100$	$g = 0.1, C = 1000$
$g = 0.25, C = 1$	$g = 0.25, C = 10$	$g = 0.25, C = 100$	$g = 0.25, C = 1000$
$g = 0.5, C = 1$	$g = 0.5, C = 10$	$g = 0.5, C = 100$	$g = 0.5, C = 1000$
$g = 1, C = 1$	$g = 1, C = 10$	$g = 1, C = 100$	$g = 1, C = 1000$

Table 4.2: Parameter space for a grid search on set C

4.4 Deployment

In this section the deployment is described as a so called IPython notebook as introduced in Section 2. It contains the complete workflow from data preprocessing to model evaluation as well as some explanations and comments. This allows a scientist to easily change and use the functions he needs.

The following part can also be watched as a html file or opened using the *ipython notebook FullRun.ipynb* command.

4.4.1 Setup

This notebook shows the different steps that are needed to classify the brain cross sections. It starts with the data preprocessing and ends with a final model. Several different functions are used by the imported custom modules. This is done, so that the notebook stays clear. If you want to take a detailed look at the used functions, feel free to use the different source files.

At the beginning the client is set up, so that the engines (processes) can be accessed and used.

Listing 5: Client Setup

```

1 # Import the client and the imp module to load source files on all engines
2 from IPython.parallel import Client
3 c = Client() dview = c[:]
4 dview.block = True
5 with dview.sync_imports():
6     import imp

```

The training data can be sampled. In this example this is done by rescaling some images to a size of 256x256.

Listing 6: Generating samples

```
1 # resizing the images
2 import glob
masks = glob.glob("../classification/data/rescaled256x256/masks/MSA*.tif")
3 images = glob.glob("../classification/data/rescaled256x256/images/MSA*.tif")
4 from PIL import Image
5 size=(256,256)
6 resized_images = [Image.open(i).resize(size) for i in images]
7 resized_masks = [Image.open(m).resize(size) for m in images]
```

4.4.2 Feature Extraction and Sampling

The following code block shows how to generate libsvm formated data from an original image and the hand labeled mask. This is done on multiple engines started by ipython. In this example the used features are:

$$R, G, B, \text{std}(R), \text{std}(G), \text{std}(B), \text{segmentationbit}, H, S, V$$

The window size for the image is 16x16 pixel. The thresholds for watershed segmentation are set as $p=0.2$ and $q=0.8$. Custom features can be added easily. A feature that needs the original image needs to be set in the new_features list. If only the pixel is needed, the online_feature list is sufficient.

Because callables are passed, the user can add any new function. The only condition is that the first argument is the image/pixel. If the function has more arguments, the *partial* function can be used to create a new one with only one argument.

Listing 7: Converting and Adding Features

```
1 # load data conversion and feature extraction modules locally and on the engines
2 dc = imp.load_source("data_conversion","../parallel/data_generation/data_conversion.py")
3 fe = imp.load_source("feature_extraction","../classification/feature_extraction.py")
4 # use the direct view to load the modules
5 dview.execute('dc = imp.load_source("data_conversion","../parallel/data_generation/data_conversion.py")')
6 dview.execute('fe = imp.load_source("feature_extraction","../classification/feature_extraction.py")')
7 # set masks and images, that are to be converted, as lists
8 # be aware that masks[0] must belong to images[0]. This is done by sorting the lists
9 import glob
10 masks = sorted(glob.glob("../classification/data/rescaled256x256/masks/MSA*.tif"))
11 images = sorted(glob.glob("../classification/data/rescaled256x256/images/MSA*.tif"))
12 # scatter the images and mask to the engines
13 dview.scatter("images",images)
14 dview.scatter("masks",masks)
15 with dview.sync_imports():
16     from functools import partial
17     # add features and convert the images
18     out_dir = "../classification/data/rescaled256x256/test"
19     dview[out_dir] = out_dir
20     cmd = 'dc.convert_and_save(images,masks,out_dir,new_features=[partial(fe.calc_std,'
21     cmd += 'size=16),partial(fe.findSegmentation,p=0.2)],online_features=[fe.add_hsv])'
21 dview.execute(cmd)
```

The data from set C can be combined into training and testing sets. The data is shuffled and the sets are class balanced by default. The *combine_files* method can be used to generate training sets for the different layers, if multiple models are used.

Listing 8: Create training and testing sets

```

1 sampling = imp.load_source("sample_files","../data_generation/random_sampling/
    sample_files.py")
2 sampling.OUT = "../classification/data/rescaled256x256/test"
3 sampling.combine_files("../classification/data/rescaled256x256/test/","ALL.svm")
4 sampling.split_train_test("../classification/data/rescaled256x256/test/","ALL.svm",
    test_size=0.5,shuffle=True)

```

4.4.3 Creating a Model

To build a classifier, first the training and test data have to be loaded. In this case it is a single cross section, but any combination of the different images can be used. This can be done with *sample.combine_files*.

The data is normalized and additional features are added. Instead of the RBFSampler the Nystroem method can be used to generate different features or approximate another kernel.

At the end the classifier is trained with parameters optimized with a grid search. Instead of building one classifier, grid search may be used to test different parameters on the current data set. After this the model can easily be saved with the pickle module which is part of the python standard library. It can also be used to save the RBFSampler which is needed if new data should be predicted.

Listing 9: Build a SVM classifier

```

1 from skimage.io import imshow,imread
2 import numpy as np
3 from sklearn.kernel_approximation import RBFSampler
4 from sklearn.svm import LinearSVC
5 from sklearn.datasets import load_svmlight_file
6 # load training set
7 X_train,y_train = load_svmlight_file("../classification/data/rescaled256x256/MSA_03-
    -2009_dXXXX-XX-XX_s0110.svm")
8 X_train = X_train.toarray()
9 # normalize data
10 max_val = X_train.max(axis=0)
11 min_val = X_train.min(axis=0)
12 X_train = (X_train-min_val)/(max_val-min_val)
13 # load test image
14 X_test,y_test = load_svmlight_file("../classification/data/segmentation/MSA_03-2009-
    _dXXXX-XX-XX_s0200.svm")
15 X_test = X_test.toarray()
16 X_test = (X_test-min_val)/(max_val-min_val)
17 # generate additional features
18 rbf = RBFSampler(gamma=2)
19 X_features = rbf.fit_transform(X_train)
20 X_test_features = rbf.transform(X_test)
21 # train linear SVM on highdimensional data
22 linearSVM = LinearSVC(dual=False,C=1000000)
23 linearSVM.fit(X_features,y_train)

```

```
24 # save model
25 import pickle with open("SVMClassifier.pkl", "w") as f:
26     pickle.dump(linearSVM,f)
```

4.4.4 Predicting new Data

If the classifier is not trained in the current session, it can be loaded with the pickle module. After adding the random features the class can be predicted using the *predict* method of the classifier. The predicted labels can be reshaped, so that the new mask can be printed as seen in figure 4.4.

Listing 10: Classify test image

```
1 # load model
2 import pickle
3 with open("SVMClassifier.pkl", "r") as f:
4     linearSVM = pickle.load(f)
5 # load hand labeled mask
6 mask = imread("../classification/data/segmentation/masks/MSA_03-2009_dXXXX-XX-
    XX_s0200.tif")
7 labels = linearSVM.predict(X_test_features)
8 labels = labels.reshape(mask.shape)
9 %matplotlib inline
10 imshow(labels)
```

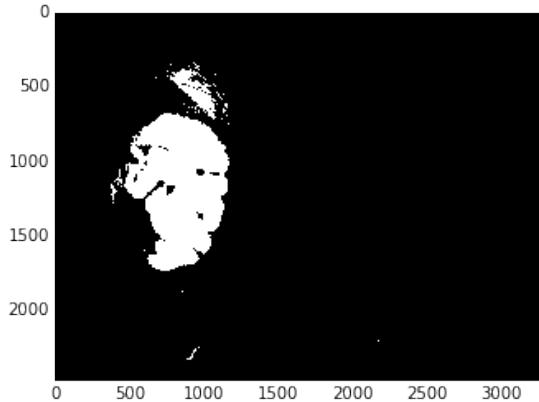


Figure 4.4: predicted mask

4.4.5 Postprocessing the results

If the data is also labeled by hand, the quality of the model can be measured. The following function calculates the accuracy as well as the f-score. It also returns the confusion matrix.

Because the classified data is an image, the confusion matrix is visualized as an image. All true positives are colored white and all true negatives black. With 100% accuracy

this would result in the hand labeled mask. Since this is most often not the case, all false positives are marked red and all false negatives blue.

This gives a good first impression on the quality of the classifier. If the predicted labels are post processed by hand, it furthermore gives an indication which regions cause problems.

Listing 11: Postprocessing

```
1 from sklearn.metrics import confusion_matrix,f1_score,accuracy_score
2 from collections import namedtuple
3 AccuracyMetrics = namedtuple('AccuracyMetrics',['accuracy','fscore',
4     'confusion_matrix','image'])
5 def compare_to_mask(pred,mask):
6     size = pred.shape
7     img = np.zeros((mask.shape[0],mask.shape[1],3),dtype='uint8')
8     tp = np.dstack(((pred>0).reshape(size),(mask>0).reshape(size)))
9     fp = np.dstack(((pred>0).reshape(size),(mask==0).reshape(size)))
10    fn = np.dstack(((pred==0).reshape(size),(mask>0).reshape(size)))
11    img[np.all(tp,axis=2),:] = 255
12    img[np.all(fp,axis=2),0] = 255
13    img[np.all(fn,axis=2),2] = 255
14    acc = accuracy_score(mask.flatten(),pred.flatten(),normalize=True)
15    fscore = f1_score(mask.flatten(),pred.flatten())
16    cm = confusion_matrix(mask.flatten(),pred.flatten())
17    return AccuracyMetrics(acc,fscore,cm,img)
18
19 metrics = compare_to_mask(labels,y_test.reshape(labels.shape))
20 imshow(metrics[3])
```

5 Evaluation On Dataset A And C

5.1 Comparison between Scikit-Learn and Twister SVM

Both scikit-learn and the Twister SVM [17] are based on the libsvm library. Twister is a map/reduce runtime that enables iterative map/reduce models. The SVM version for twister is based on the cascade SVM approach. The training data is split into several sets and a SVM is built for every subset. The support vectors of the models are used to create a new SVM model. If there is only one model left, the first iteration is finished. The support vectors of the last model are added to the subsets and the iteration step is repeated. This is done until the change of the final support vectors is less than a given threshold. In other words via iterations the not support vectors are filtered out.

To test the runtime of the two implementations, classifiers were built for training data of different size. Table 5.1 shows the results for 20.000, 40.000 and 100.000 instances. It can be seen that scikit-learn is faster. This is due to the overhead caused by the message passing and setting up the different jobs. With an increasing training set the twister SVM implementation promises to become faster than the serial one. Unfortunately, it was not possible to run more tests with larger training sets, because the twister implementations is still quite unstable due to its many dependencies on a map/reduce runtime, a messaging system for the iterative map/reduce technique and an on-top implementation of SVMs that was never released as an official software package. Hence, it is also rather a matter of an ongoing research project.

Classifier	Accuracy	training time in s	test time in s
scikit-learn, 20k instances	0.957980867538	45.242423	66.846724
scikit-learn, 40k instances	0.962536060588	254.560878	112.238295
scikit-learn, 100k instances	0.965938617227	1145.255996	209.801505
libSVM twister, 20k instances	0.959293049667559	445.867	-
libSVM twister, 40k instances	0.963317509002937	934.912	-
libSVM twister, 100k instances	0.966521447502947	3077.239	-

Table 5.1: Training time of scikit-learn and twister libsvm on dataset A

Figure 5.1 shows the accuracy and training time with an increasing number of samples of dataset A. The improvement of accuracy is low compared to the additional runtime that is needed. Additional samples do not improve the accuracy much. Therefore, it is more reasonable to add additional features to the data set, that improve the accuracy instead of increasing the number of instances. The training time of a SVM classifier is much more dependent on the number of instances than it is on the number of features, so that the training time can be reduced by using a small sample of the data with added features.

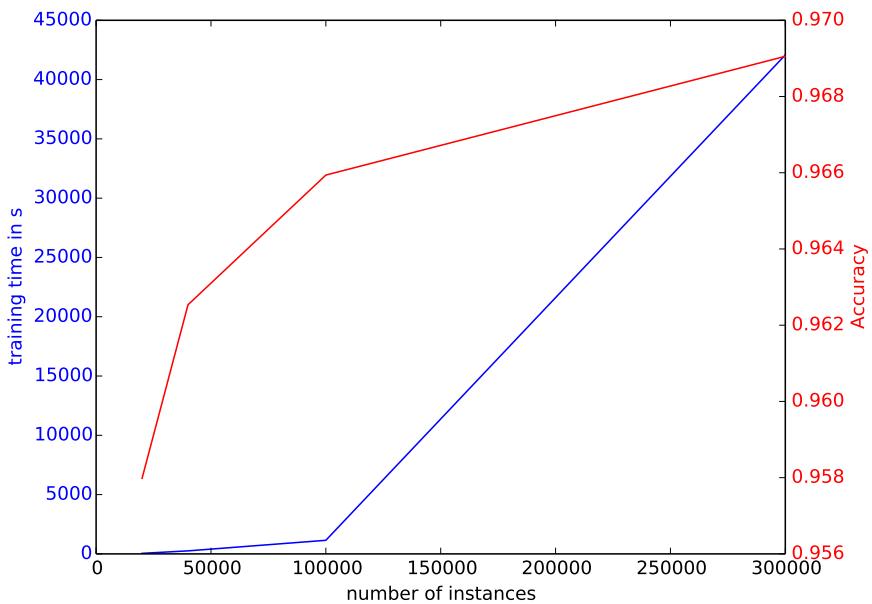


Figure 5.1: SVM training time with increasing sample

5.2 Improvements with Feature Extraction

Table 5.2 shows the accuracy on Dataset C and a 1% sample. Because of the additional features compared to dataset A and a larger sample the accuracy has increased. The performance dropped because of the larger training set. The effect of the new features can be seen best on the SVM, because the model built on dataset C as well as the one built on set A use 20000 instances.

Classifier	Accuracy	training time in s	test time in s
Decision Tree	0.9803088445	4.650993	0.048027
Random Forest	0.9867766293	9.456317	0.309398
KNN	0.9858822801	68.91499	31.868704
GaussianNB	0.9594231602	0.093499	0.106665
SVM 20k instances, no grid search	0.9732314167	2.342334	19.371704

Table 5.2: Performance and Accuracy of several Classifiers on Dataset C, sample 1%

5.3 Data Post Processing

After the model is created, an image can be predicted. During the preprocessing the image is flattened and several features are added, so that every row contains one pixel. Since the order of the pixel is contained, the predicted labels can be reshaped to form the outcome of the prediction as a mask for the given image. This gives a much better

impression of the model's performance than the accuracy score because the data is not balanced. As the mask of the image is given, the predicted mask and the original mask can easily be compared by opening both masks or calling the `compare_to_mask` function. The latter calculates the accuracy, f-score and the confusion matrix. Furthermore it also visualizes the confusion matrix as an image. An example of this image is shown in figure 5.2.

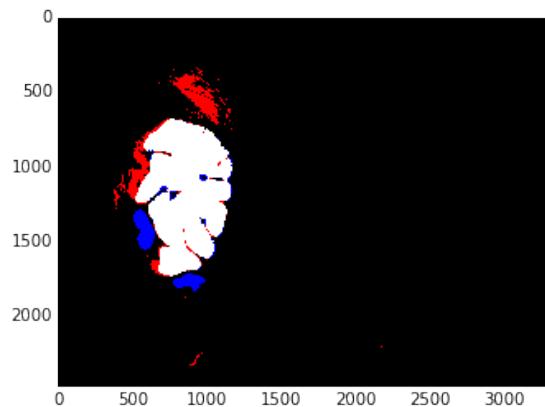


Figure 5.2: Confusion matrix visualized

6 Conclusion

The intention of this internship was to create a classifier for a given real-world application domain problem. It has to predict whether a pixel belongs to a cross section of a human brain or not. For this several problems had to be solved.

Because of the huge amount of data the preprocessing had to be done in parallel on a super computer. For this IPython was configured so that it is able to be started with the scheduler and create several engines. While IPython supports MPI, which is the most used framework on the super computers in Jülich, it is also possible to access the engines using the DirectView or LoadBalancedView classes of IPython.

In order to create a better classifier, feature extraction methods have been applied taking related work from image segmentation methods into account. This contains pixel based features and local neighbours.

While Support Vector Machines were the focused method for classification during this internship some other popular approaches were used as well for comparisons with the SVM approach. The one with the best result is the Random Forest method, which is an ensemble method based on several randomized Decision Trees that are used for majority voting. Both Support Vector Machines and Random Forests can handle feature dependency and do not overfit as easily as other classifiers. A negative example would be Naive Bayes which despite its good training speed has a bad accuracy because it assumes feature independency.

After optimizing the model with grid search, it is able to help the scientist by automatically predicting an approximated cross section which can be manually corrected. It can be easily used with the given IPython notebook.

While the notebook is based on scikit-learn and therefore a serial SVM implementation, a parallel approach was covered by testing the twister svm implementation. Although the cascade svm promises to increase the training time, this was not the case for the particular dataset used in the internship. Emerging Brain data is expected with new scanners with order of magnitude of better resolutions (i.e. bigger data) and as such this internship was a good proof-of-concept of evaluating different options. The used data sets are not large enough to show the advantages of a cascade svm and as such parallelization needs, except the preprocessing stage. The second reason is that the twister implementation as well as the SVM implementation for twister are still a matter of ongoing research and not particularly stable and optimized versions need to be developed. Because of this it was also not possible to test it on a larger training set.

References

- [1] http://www.fz-juelich.de/ias/jsc/EN/Expertise/Supercomputers/supercomputers_node.html.
- [2] <http://ipython.org/documentation.html>. [Online; accessed 20-January-2015].
- [3] <https://www.humanbrainproject.eu/de/home>.
- [4] G. Cavallaro, M. Riedel, J.A. Benediktsson, M. Goetz, T. Runarsson, K. Jonasson, and T. Lippert. Smart data analytics methods for remote sensing applications. In *Geoscience and Remote Sensing Symposium (IGARSS), 2014 IEEE International*, pages 1405–1408, July 2014.
- [5] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [6] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: Simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, January 2008.
- [7] Message P Forum. Mpi: A message-passing interface standard. Technical report, Knoxville, TN, USA, 1994.
- [8] J. Han, M. Kamber, and J. Pei. *Data Mining Concepts and Techniques*. Morgan Kaufmann, 2011.
- [9] Robert Henderson. Job scheduling under the Portable Batch System. In Dror Feitelson and Larry Rudolph, editors, *Job Scheduling Strategies for Parallel Processing*, volume 949 of *Lecture Notes in Computer Science*, pages 279–294. Springer Berlin / Heidelberg, 1995.
- [10] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–. [Online; accessed 2015-01-13].
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [12] Fernando Pérez and Brian E. Granger. IPython: a system for interactive scientific computing. *Computing in Science and Engineering*, 9(3):21–29, May 2007.
- [13] A. Rahimi and B. Recht. Random features for large-scale kernel machines. <http://www.eecs.berkeley.edu/~brecht/papers/07.rah.rec.nips.pdf>.
- [14] Anand Rajaraman and Jeffrey David Ullman. *Mining of Massive Datasets*. Cambridge University Press, New York, NY, USA, 2011.

- [15] Jos B.T.M. Roerdink and Arnold Meijster. The watershed transform: Definitions, algorithms and parallelization strategies. *Fundam. Inf.*, 41(1,2):187–228, April 2000.
- [16] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The hadoop distributed file system. In *Proceedings of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST)*, MSST ’10, pages 1–10, Washington, DC, USA, 2010. IEEE Computer Society.
- [17] Zhanquan Sun and Geoffrey Fox. Study on parallel svm based on mapreduce. In *International Conference on Parallel and Distributed Processing Techniques and Applications*, pages 16–19. Citeseer, 2012.
- [18] Shailaja Sunkutlawar, Prof Ramesh, and K Kulkarni. Shadow suppression using rgb and hsv color space in moving object detection.
- [19] Stefan van der Walt, S. Chris Colbert, and Gael Varoquaux. The numpy array: A structure for efficient numerical computation. *Computing in Science & Engineering* 13, 22-30, 2011.
- [20] Christopher Williams and Matthias Seeger. Using the nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 13*, pages 682–688. MIT Press, 2001.
- [21] Tianbao Yang, Yu-feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method vs random fourier features: A theoretical and empirical comparison. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 476–484. Curran Associates, Inc., 2012.