

Classification Of Visualization In Scientific Literature

Masterarbeit zur Erlangung des akademischen Grades
Master of Science (M.Sc.)

Lehrstuhl für Intelligent Systems und Lehrstuhl für Data Science
der Fakultät für Informatik und Mathematik
der Universität Passau

| | |
|-----------------|-----------------------------|
| Name: | Arnold Azeem |
| Matrikelnummer: | 79176 |
| Fachbereich: | Informatik |
| Studiengang: | Master Informatik |
| Erstprüfer: | Prof. Dr. Christin Siefert |
| Zweitprüfer: | Prof. Dr. Michael Granitzer |
| Date: | June 7, 2018 |

0.1 Abstract

Distinct visualization techniques are used in scientific research publications to summarize large amount of data and also represent a variety of data. These visualizations help to communicate complex information and support the arguments being presented in the publication in a way that is easy to understand and follow. These figures tend to reveal trends, patterns or relations that might have otherwise been difficult to grasp using only text. It is therefore relevant that we extract the data from these visualizations since the extracted data can be used for validating the publication or presenting the data in another form for a different audience. In this context, classifying these visualizations is the initial step since, there is a variety of visualizations and each one is processed in a specific way. It is only after classification that extracting of raw data from these visualizations can be acquired for other tasks. This thesis presents an approach whereby real world data is used to create four types of plots (scatter plots, bar charts, line charts, and box-plots) and random plots also of the same kind from the Internet are added together and used to train and evaluate a CovNet model to be able to classify these plots.

Contents

| | | |
|----------|--|-----------|
| 0.1 | Abstract | 2 |
| | List of Figures | 3 |
| | List of Tables | 4 |
| 1 | Introduction | 5 |
| 1.1 | MOTIVATION | 6 |
| 1.2 | OBJECTIVE | 6 |
| 2 | Background and Related Work | 8 |
| 2.1 | Model Based Approaches | 8 |
| 2.2 | Machine Learning Based Approaches | 9 |
| 2.2.1 | Convolutional Neural Network | 10 |
| | Convolutional Layer | 10 |
| | ReLU Layer | 10 |
| | Pooling Layer | 11 |
| | Fully Connected Layer | 11 |
| | Machine Learning Classification Algorithms to Recognize Chart Types in Portable Document Format(PDF) Files | 11 |
| | Architecture proposal for data extraction of chart images using Convolutional Neural Network | 12 |
| | Chart classification by combining deep convolutional networks and deep belief networks | 13 |
| 3 | Approach | 15 |
| 3.1 | Dataset | 15 |
| 3.1.1 | Dataset for Matlab | 16 |
| 3.1.2 | Dataset for R | 16 |
| 3.1.3 | Dataset for Python | 16 |
| 3.1.4 | Dataset for Java | 17 |
| 3.2 | Creating Plots | 18 |
| 3.3 | Training, Validation and Test set | 21 |

| | | |
|----------|-------------------------------|-----------|
| 3.4 | The Architecture | 22 |
| 3.5 | Image Preprocessing | 22 |
| 4 | Evaluation | 23 |
| 5 | Summary | 24 |
| 6 | Bibliography | 25 |

List of Figures

| | | |
|-----|---|----|
| 1.1 | What we hope to achieve in the thesis | 7 |
| 3.1 | Scatter plots | 18 |
| 3.2 | Example Bar Charts | 20 |
| 3.3 | Example Line Charts | 20 |
| 3.4 | Example Box Plots | 20 |
| 3.5 | Steps involved in building the Classification model | 22 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Details on Dataset from [1] | 12 |
| 2.2 | Number of Train and Test Dataset collected | 13 |
| 2.3 | Comparing Results of Proposed Framework from (Liu et al. [2]) | 14 |
| 3.1 | Names of datasets used in each plotting program | 17 |
| 3.2 | Overview of the varied parameters for creating plots in the different plotting programs | 19 |

1 Introduction

"A picture is worth a thousand words" even though a widely used phrase stands to be very true especially when complex data is visualized and presented in scientific research publications. Data is ever growing and sometimes complex, using figures and diagrams to interpret and represent this data cannot be undermined since, they provide a way to easily give insight into the research findings, which would have otherwise been more complex relying on only textual data. For this reason, there has been a growing interest in the chart analysis area and quite a number of techniques have been developed. In spite of this growing interest, there has been little groundbreaking results achieved due to different variations in appearance of plots [2]. For example, Manollis Savva, et al [3] proposed a model to classify charts using extracted low-level features and textual features. After extracting the features, a Support Vector Machines (SVMs) classifier is used for the classification step. This method was limited since most charts contain the same type of features like axes, grid lines, and legends. In V Shiv Naga Prasad's work [4] classification was based on using features based on the shape and spatial relationships of their primitives. This work was limited due to the inconstancy in which data in most charts can be depicted. The process of extracting data already visualized as figures can be done relatively easier manually but becomes more complicated if done automatically. This process can be divided into two main steps [3]. The first step which our work focuses on, classifying the chart and the second step which involves extracting the data from the classified chart. To achieve the classification step, this paper presents an approach where charts are created with real-world datasets, different plotting programs (Python, Matlab, R, and Java) and different libraries supported by these plotting programs were used together with downloaded chart images from the Internet. We then use these images as input to Convolutional Neural Network (CNN). CNN was used instead of primitive approaches because it has achieved ground breaking results in the area of image classification [5]. Our model can identify four classes of plots namely Box-plots, Line Charts, Scatter-plots and, Bar-charts. The other parts of this paper are organized in the follows: In the next Sections, we present the motivation behind this work. Other works related to this, our pro-

posed method is described, Experimental evaluation and results are reported, and finally, the conclusion and the way to approach this work in the future.

1.1 MOTIVATION

Complex data is better explained in scientific papers with the aid of visualizations. These plots present complex data in an easy to understand way compared to textual representation. The data which these visualizations contain when extracted play an important role in events where another researcher wants to verify the work of the publisher, this data can also be used to develop other visualizations in situations where the paper needs to be presented to a different audience with a different background as opposed to the audience which the visualizations were created for, Also when comparing two plots the raw data helps make a better decision than just the figures. Since each plot will be processed differently to extract the raw data, it very relevant that we can distinguish one plot from another and this is the main aim of this thesis.

1.2 OBJECTIVE

The purpose of this thesis is to answer the question:

How Well Can We Classify the Four Different Types of Plots (Line Chart, Bar Chart, Scatter Plot and Boxplot) in Scientific Literature ?

In this work we focus on only four plots. These plots are scatter plots, bar charts, line charts and Box plots. Figure 1.1 shows the vision of this work, The first part of the diagram involves extracting or obtaining the four different types of plots mentioned earlier, after which we then label our plots and train a neural network model to be able to classify with high accuracy any of the four plots if shown to our model, then finally the raw data can be extracted from the detected plot. But this work mainly focuses on the red dotted lines shown below in the diagram which which is getting the plots, labeling them, training the model and classifying the plots.

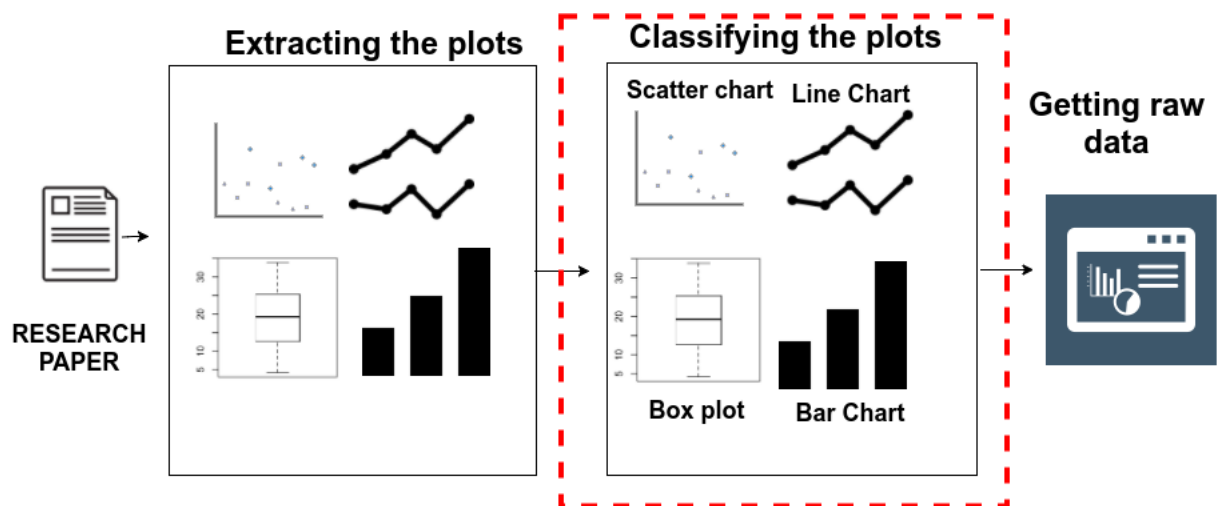


Figure 1.1: What we hope to achieve in the thesis

2 Background and Related Work

Classification involving charts has been a topic of huge interest in recent times. This is due to the fact that most publication are embedded with these charts as a way to convey complex research findings in more visualizable and easy to understand form. As a result of its popularity different techniques have over time been used for classification of charts and these techniques have evolved over the years. The techniques can be put into two categories; Model Based Approaches and Machine Learning Based Approaches [5]. Next the approaches will be explained into details.

2.1 Model Based Approaches

The following is written with inspiration from Boyle et al. work [6]. Many chart are very different in sense of their context structure. This is as a result of the variability of positions of the context structure of a chart. Content like; legends, axes, grids etc have no fixed positions and sometimes are not even present in some charts, as a result classifying becomes a difficult task. The model based approach is divided into two steps. Firstly, a number of predefined object classes are created with abstract models and then an image is matched with the created model. For example, a model for a 2D/3D pie chart consists of line segments (radii and circular/elliptical arcs), and these line segments are used to create the model, however, just using these line segments wont be enough, so some constraints are introduced into the model, example of constraints of a pie chart include; the center of the pie chart is where all radii meet; all 2D pie charts have radii of equal length; arcs mainly form parts of the same 2D pie chart circle, or if its a 3D then the arcs form part of the ellipse. After this a goodness-of-fit is introduced to help measure the discrepancy between the image and the conditions for the model. For example, a goodness-of-fit of a pie chart is as follows: difference in length of radii; difference in how the arcs are curved; distance between a center of a circle and that of a radii. When a chart is presented for classification, the edges are detected, thinned, linked, vectorized, extracted and compared to the various edge models created and the best match is selected. After this, the

good-to-fit criteria is used to measure the difference between the edges of the model and the image. Finally, a voting is performed base on the value of the good-to-fit results. The model based approach has some draw backs though. The main drawback is that human intervention is needed to develop the various models. This could be a cumbersome and time consuming procedure. Due to the limitation of this approach a machine learning approach was introduced to handle some of the model based approach drawbacks.

2.2 Machine Learning Based Approaches

This approach also has different method. Some methods involve using hand-crafted features from the charts for classifying the chart. In Zhou and Tan’s [7] work, features like legend, x-y-axis title, chart title and value of the bar were extracted for purposes of classifying a bar-chart. In Shao and Futrelle’s [8] work, graphical elements of the charts like colors,tick marks on an axis and data point markers are used. Another instance where handcrafted features were used was in Inokuchi et al. work [9]. In this work regularly appearing substructures of chart are extracted and used as features for classification. The next step is the classification step, after the features are obtained, they are then represented in vector form for the classification stage. The features are stored in a matrix vector like structure. Algorithm 1 shows an example structure. This structure has columns containing the features and one column indicating the target or label of the chart. The structure is then fed as input into a machine learning technique for the model training step. In Karthikeyani and Nagarajan’s paper [1] after the features were obtained SVM Classification¹ MLP Classification² and KNN Classification³ were used to create a function that maps the set of features to a predefined label. This function is known as the classification model. This technique also achieves very good results but, the drawback is the reliance on handcrafted features such as edges and corners and therefore, does not handle a large amount of data with variable context properly. Recently however, CovNets (this abbreviation stands for Convolutional Neural Networks and will be used a lot through this work)⁴, a machine learning technique which does not need to be supplied with handcrafted features, since they learn the features and classify them on their own. Hence, CovNets are used in our work to automatically classify charts and avoid the cumbersome task of having to extract handcrafted features.

¹<http://www.statsoft.com/Textbook/Support-Vector-Machines>

²http://www.iro.umontreal.ca/~bengioy/ift6266/H12/html/mlp_en.html

³<https://kevinzakka.github.io/2016/07/13/k-nearest-neighbor/>

⁴<https://deeplearning4j.org/convolutionalnetwork>

2.2.1 Convolutional Neural Network

The following is inspired by the following blogs by [ujjwalkarn \[10\]](#), [Ifu Aniemeka\[11\]](#) and [Adit Deshpande \[12\]](#). CovNets are a variety of Artificial Neural Network which have given results near to human performance. They especially perform well in the areas of image recognition and classification. All CovNets are divided into steps called layers, the layers are Convolution, Pooling or Sub Sampling, Classification (Fully Connected Layer). When an image is presented to a computer, the image is seen as an array of pixel values. For instance if a color image of size 48x48 is presented the computer sees it as 48x48x3 where the 3 represents the RGB values.

Convolutional Layer

The convolutional layer will be explained using the following example, let consider the image as a magical cake of height and width 48x48, and there is a cake cutter of size say 5x5. The cutter is used to cut the cake from the top left. The cutter is referred to as a filter in machine learning. The cutter in this case is also an array of numbers called weights and for the maths to work the cutter must have the same depth as our cake (5x5x3). Lets first consider the cutter being in the first position ie. the top left of the cake. The values in the cutter are multiplied with the values of the cake (element wise multiplications). These multiplications are all summed up and result in a single number. This single number is only for the first part and this process is repeated throughout the whole cake (Next step would be moving the cutter to the right by 1 unit, then right again by 1, and so on) keep in mind its a magical cake so you can cut a part more than once. After using the cutter on the whole of the magical cake we result in a new cake of size 28x28x1, the results is called a feature map. The filters are low level feature identifiers (straight edges, simple colors, and curves).

ReLU Layer

The next layer after the convolution phase is the Rectified Linear Units (ReLU). This layers helps to handle situations where the relation between the input values and the CovNet output is non-linear. The ReLU has a function $f(x) = \max(0, x)$ which means if you give it a value x, it will return 0 if x is negative and will return the value itself if its positive. There are also other functions which can used n place of the ReLU; that is the tanh or the sigmoid function but the ReLU mostly works well.

Pooling Layer

The pooling layers basically have two main functions, first one is that it helps the CovNet to locate features regardless of which part of the image it is located. This results in the model being robust against small changes in position of the features of the images. The second function is that it also helps to reduce the size of the feature map. Therefore computations in the futures layers are relatively less complex. One way of performing pooling is using max pooling technique, thats is sliding a window through the feature map, the windows fills a number of arrays in the feature map therefore, you pick the largest among all the numbers and disregard the rest of the numbers.

Fully Connected Layer

This is the last layer and where the decision about the image is made. i.e a probability saying which class the image belongs to. This layer takes the output of the previous layers, which are the high level features and check which class these features strongly correlate to. All the above mentioned layers can be multiply stacked on each other. In the rest of this chapter, three related works that our work builds on are described into details.

Machine Learning Classification Algorithms to Recognize Chart Types in Portable Document Format(PDF) Files

Karthikeyani and Nagarajan's work [1] focuses on classifying charts found in pdf files. Th steps involve extracting texture features from the charts and then, using a machine learning technique for the classifying phase. The dataset consists of images extracted from various pdf's. The dataset consists of 155 images in total and images are 256*256 RGB. Table 2.1 shows the details of the dataset used. In this work, handcrafted features are extracted from the image and then used in the classification phase to train and test a model. For the feature extraction phase, Gray Level Co-Occurrence Matrix (GLCM) is employed. GLCM is a technique which uses co-occurrence matrix to extract texture features of an image with the use of statistical equations. GLCM was used to extract eleven features, these included: area, median, minimum and maximum intensity, contrast, homogeneity, energy, entropy, mean, variance, standard deviation and correlation. These extracted features are correlated to the pixels of the image. These extracted features are then stored using a 2-dimensional matrix vector data structure. This structure has thirteen (13) columns and 'x' rows, where x is the size of the dataset. The features are stored in twelve (12) columns and the label is stored in the thirteenth column. The features and labels were stored in the below structure 1.

| Chart Type | No of Charts | Chart Type | No of Charts |
|--------------|--------------|-------------|--------------|
| 2D Bar chart | 40 | Doughnut 2D | 7 |
| 3D Bar chart | 16 | Doughnut 3D | 11 |
| 2D Pie chart | 13 | Line | 35 |
| 3D Pie Chart | 20 | Mixed Chart | 13 |

Table 2.1: Details on Dataset from [1]

Algorithm 1 This is the structure used to store the features

```

Struct FeatureVector {
float feature1; float feature2;
float feature3; float feature4;
float feature5; float feature6;
float feature7; float feature8;
float feature9; float feature10;
float feature11; float feature12;
int target;
}

```

This vector like structure is then fed as input into three classifiers SVM, MLP neural network and K-Nearest Neighbour. These classifiers are then trained and a classification model is formed for the recognition step. After this, to see if the model works well a test set consisting of new records is fed into the model for the model to predict their labels. Three metrics were used to check the performance of the model; error rate, classification accuracy and speed of classification. The error rates are MLP 0.30, K-NN 0.22 and SVM 0.23. For the accuracy metric; KNN (78.06%), MLP (69.68%) and SVM (76.77%). Finally the speed metric, which sum of training and test time results in MPL being the slowest with 8.38, followed by SVM with 0.31secs and KNN with 0.26secs. Even though these results were very good the paper proposed, extracting features which are relate to shape and curves.

Architecture proposal for data extraction of chart images using Convolutional Neural Network

Our work is inspired by De Freitas et al. [13], proposed a way to extract the wealth of information contained in different visualization techniques. The paper talks about two main stages of accomplishing this task. Firstly, classification of the charts is done since it allows a different variety of chart to be detected automatically allowing the next step, which is the extraction of data from the classified plots. The paper, however, focuses on the first step, classification of charts. In this paper, a Convolutional Neural Network is used for the classification task. The Convolutional neural network encapsulates the

characterization and classification processes during its learning process, unlike other techniques. The dataset used for this task were searched for and downloaded from Google image search. Table 2.2 shows the chart types which were collected and the number of train and test sets which the respective charts were divided into.

| Chart Type | Test | Train |
|---------------|------|-------|
| Area Chart | 50 | 555 |
| Bar Chart | 50 | 657 |
| Line Chart | 50 | 489 |
| Map | 50 | 476 |
| Pareto Chart | 50 | 261 |
| Pie Chart | 50 | 361 |
| Radar Chart | 50 | 454 |
| Scatter Chart | 50 | 552 |
| Table | 44 | 236 |
| Venn Diagram | 48 | 304 |
| Total | 498 | 4345 |

Table 2.2: Number of Train and Test Dataset collected

For the classification, a variant of CovNet called LeNet-based CNN model is used. The model was implemented using Tensorflow⁵, LeNet-based CovNet has an architecture which is comprised of 3 convolutional layers, followed by a fully connected layer. The model is trained in a way that the dataset is divided into mini-batches, samples of fixed sizes(100) are selected and fed into the CNN, as a result of this process the model becomes robust since it learns to generalize from the different min-batches which are fed into the model. Also, all the images are converted to JPG and resized to 224x224x3, that is, 224 pixels of height, 224 of width and 3 layers of output. The other parameters used were 1000 epochs and a learning rate of 0.003. The accuracy at the end of the training process was 70%.

Chart classification by combining deep convolutional networks and deep belief networks

In another paper by Liu et al. [2], a new approach was proposed for the process of chart classification. The process involves using CovNet to extract deep hidden features of charts and then deep belief networks then use the extracted

⁵<https://www.tensorflow.org/>

features to predict the labels of the charts. Due to a difficulty in acquiring a large number of charts as training data, natural images were first used to train the model and later the model was fine-tuned with just over 5,000 collected charts. The types of charts collected were pie charts, scatter charts, line charts, bar charts, and flowcharts. The architecture of the CovNet is made up of five Convolutional layers and two fully-connected layers and then an output layer. The preprocessing steps for the images involve down-sampling them to $256 \times 256 \times 3$, after which each is cropped to a size 227×227 from the center and its horizontal flip are extracted as the input of the CovNet, other parameters used for the CovNet include a learning rate that starts with 0.01 initially and is then decreased by a factor of 0.1 after every 100k iterations, the weight decay parameter was set at 0.0005 and a dropout rate of 0.5. This results in an output of a 5-way softmax which produces the distribution over the 5 class labels and this is used as input for the deep belief network. The deep belief network architecture has three hidden layers, whose dimensions are 5000, 500 and 2000. This results in a softmax predicting the probability distribution over the 5 categories of charts as output. The training process was done with 4000 randomly selected images and the rest were used as test set. The accuracy of the model after the evaluation was 75.4%. Table 2.3 shows the results after the training was done without deep belief networks but pre-trained with the natural images and finally the training done with only the chart dataset but with deep belief networks.

| Chart | CovNets | CovNets+DBN without pre- training | CovNets+DBN |
|---------------|---------|---|-------------|
| Bar Chart | 75.6% | 45.6% | 74.2% |
| Flow Chart | 88.3% | 56.8% | 91.3% |
| Line Chart | 71.2% | 22.3% | 67.9% |
| Scatter Chart | 69.8% | 44.5% | 84.2% |
| Pie Chart | 58.1% | 50.1% | 59.4% |
| Ave. Accuracy | 72.6% | 43.9% | 75.4% |

Table 2.3: Comparing Results of Proposed Framework from (Liu et al. [2])

3 Approach

As seen in the previous chapter, different approaches and methods have been described for chart image recognition. Due to the complexity of this task, a CovNet approach is employed for our work. Reasons for selecting this approach include; automatic capturing, learning and extraction of features; parameter sharing which allows the model to learn a single set of weights once, rather than a different set of weights every time [14] and more importantly a performance accuracy nearing human capability. In spite of all the advantages of CovNets, they are very computationally expensive to apply on high resolution images. Fortunately, current GPU capabilities with the help of good optimization techniques can handle this issue [15].

The specific contributions of this work is as follows: to create a dataset of chart images (scatter-plot, bar chart, line-chart and box-plot) using real world data. This chart image dataset tries to capture all variety of structure in the various chart images to be handled. The scripts, libraries and parameters used in creating the chart images will readily available for recreation of the dataset. The second contribution, is training a CovNet model for recognizing these types of chart images when presented to it.

The rest of this chapter describes into details how the dataset was created, it also describes the CovNet architecture that was used for the classification task.

3.1 Dataset

The following paragraph is inspired by [16]. The saying 'Garbage in Garbage out' is a valid statement when it comes to creating a dataset for machine learning. The machine learning technique will learn from whatever data fed to it. So if a dataset of good quality is fed into the algorithm, then the model created will also be of good quality. The dataset creation stage is therefore a very important stage. In most approaches that worked on chart images classification, the dataset used consisted of images downloaded from Google and a few others were obtained from extracting them from pdf's. This approach we believe

is bit limited since, we don't know the programming languages, the libraries used and parameters used in creating these charts. This information is very relevant since it tells us how diverse our resultant dataset is. For example, how are we sure that all the charts that were downloaded from Google, were not only created with python or java?, and in such case how well will a new chart created with matlab or R be classified. For this reason the dataset comprised of image charts created with different libraries and in different programming languages. In the next sections, the various datasets and the languages used in plotting are described. To make the dataset as diverse as possible, charts created in each programing language used a different set of CSV files.

3.1.1 Dataset for Matlab

The Data used for creating the plots in Matlab were randomly chosen from Project Dataset [17], a free CSV data repository, DatPlot [18] and Plotly CSV repository in github [19]. The datasets are multidimensional and compiled from normal day to day activities like dating, what makes people happy etc, and objects like cameras and cars. On the average the datasets used contain about 500 instance and 5 different columns. The biggest dataset is called Speed dating data. It is made up of over 8,000 observations of answers to survey questions about how people rate themselves and how they rate others on several dimensions. The smallest dataset used has 33 instances and 12 columns. It contains information about cars. The number of gears and speed, just to name a few attributes.

3.1.2 Dataset for R

For the plots in R, 13 random CSV files where downloaded from an archive of datasets distributed with R called Rdatasets [20]. Rdatasets is a collection of dataset distributed with R. On the average there are 80 instances and 5 columns in each dataset. The biggest CSV file is the Australian athletes dataset. It's made of 203 instances and 14 columns and contains attributes like sex,height,weight and sports. The smallest dataset is the Canadian Women's Labour-Force Participation. This dataset has 30 rows and 7 columns. It contains information like average wages of women, percent of adult women in the workforce etc.

3.1.3 Dataset for Python

The data used for creating the plots in Python were 15 randomly seleted csv files also from Rdatasets [20]. The biggest dataset among the 15 is the Monoclonal gammopathy data, it contains natural history patients with monoclonal

gammopathy of undetermined significance. The dataset is made up of 1384 observations with 10 columns, it has attributes like age, sex, time of death and last contact in months. On the average each dataset contains about 200 instances and 7 columns of multi-dimensional data. The smallest dataset however contains only 33 instances with 11 columns and is called the Nuclear Power Station Construction Data. The data relate to the construction of 32 light water reactor (LWR) plants constructed in the U.S.A in the late 1960's and early 1970's.

3.1.4 Dataset for Java

For the plots created in java, I used the dataset made available by Plotly [19], a github repository of CSV datasets used in the Plotly API examples. 14 random CSV files were downloaded, the biggest file has 1002 instances and 9 columns, and on the average each file contains about 100 instances and 9 columns. The smallest file however is made of 33 instances and 12 columns called the mtcars file. It contains information about a variety of different car models like the number of gears, speed etc. The table 3.1 contains the names of all CSV files that were used in the different languages with the different plotting programs.

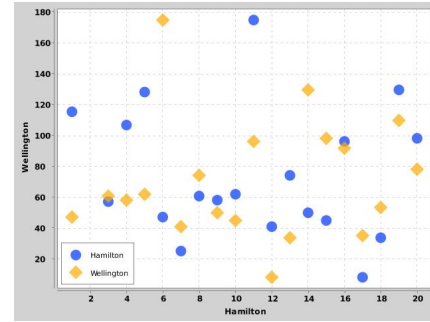
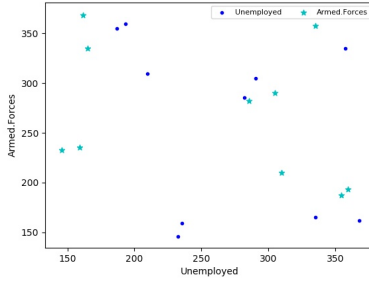
| Datasets | | | |
|--|--|--|---|
| Python | Matlab | R | Java |
| 3d_line_sample_data.csv LightForwardFlapStall.csv line_3d_dataset.csv longley.csv loti.csv lung.csv nuclear.csv timeseries.csv USJudgeRatings WVSCulturalMap.csv wind_rose.csv volcano.csv uspop2.csvm tips | Camera.csv Cars.csv speedDating.csv Cereal.csv happiness.csv TestData1.csv TestData2.csv mpg.csv okcupid- religion.csv spectral.csv stockdata.csv subplots.csv | ais.csv Angell.csv Baumann.csv Bfox.csv cane.csv carprice.csv Chirot.csv Davis.csv Ericksen.csv Florida.csv Highway1.csv Pottery.csv Prestige.csv salinity.csv urine.csv | 3d-line-plot.csv 3d-scatter.csv 2011_flight_paths.csv 2011_us_exports.csv auto-mpg.csv candlestick_dataset.csv finance-charts-apple.csv globe_contours.csv hobbs-pearson- trials.csv motor_trend_tests.csv nz_weather.csv volcano.csv iris.csv mtcars.csv |

Table 3.1: Names of datasets used in each plotting program

3.2 Creating Plots

The inspiration for creating a variety of plots to capture all type of plots used in scientific papers was gotten by inspecting the dataset of Architecture proposal for data extraction of chart images using CovNet paper [13] and Viziometrics: Analyzing visual information in the scientific literature [21] dataset. Scripts in various languages were written to handle the plotting and labeling process automatically. All datasets for a particular plot (example scatter plot for python) are put into one folder. The scripts reads each CSV file column by column while creating the plots.

Table 3.2 describe how the plots where created in each language. The type column describes the different variety of a particular plot, for example bar charts can be of type stacked, grouped, vertical and horizontal bar charts, also scatter plots types can be a scatter plot consisting of one type of marker, one scatter plot with multiple markers and finally a scatter plot with a line showing the correlation between the plots. Figure 3.1 shows two different types of bar charts. Figure 3.1a is a stacked bar chart and Figure 3.1b a normal vertical bar chart. The Library column shows the different plotting libraries used, the parameter column describes parameters that were changed and finally the number of plots created were also added. The images below the tables are sample images that exist in our dataset of created plots for each language.

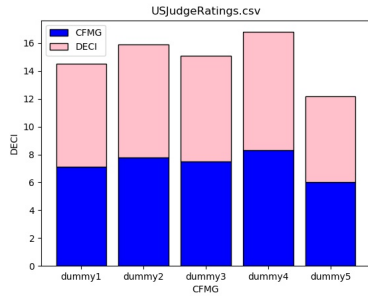


(a) Matplotlib scatter plot with star and circular markers (b) Java scatter plot with circular and diamond markers

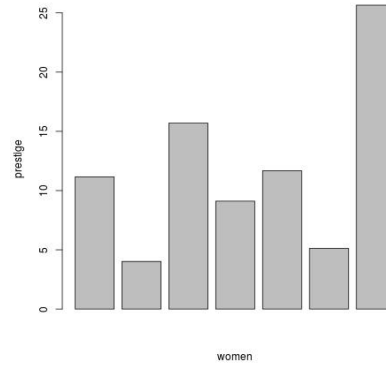
Figure 3.1: Scatter plots

Table 3.2: Overview of the varied parameters for creating plots in the different plotting programs

| SCATTER PLOTS | | | | |
|---------------|--|--|-----------------|---|
| Language | Library | Parameters | Number of plots | Type |
| Python | Matplotlib v2.1.2 Plotly v2.5.1 Seaborn v0.8.1 | MarkerStyle ['o', '*', '.', '+', 'x'] | 1020 | Unique markers, With legends, multiple markers regplot |
| MATLAB | Default Plotly | MarkerStyle ['o', '*', '+', 'x', 's'] | 1044 | |
| R | Plotly Lattice Ggplot2 | MarkerStyle ['o', '*', '+', 'x', 's'] | 1644 | |
| JAVA | XChart 3.5.1 jfreechart 1.0.1 | MarkerSize (15 - 18) | 1644 | |
| BAR CHARTS | | | | |
| Language | Library | Parameters | Number of plots | Types(bar) |
| Python | Matplotlib v2.1.2 Plotly v2.5.1 Seaborn v0.8.1 | | 1000 | Horizontal and Vertical, Stacked, Grouped bar charts |
| MATLAB | Default | Width of bar(14-16) | 1000 | |
| R | Default,Plotly R Library ggplot2 | space (0-3) | 1144 | |
| JAVA | XChart 3.5.1 jfreechart:1.0.192 javafx.scene | PlotOrientation (vertical or horizontal) with error bars | 1144 | |
| LINE CHARTS | | | | |
| Language | Library | Parameters | Number of plots | Types(Line with) |
| Python | Matplotlib v2.1.2 Plotly v2.5.1 Seaborn v0.8.1 | Linestyle ['-', '-', '-.', ':'] | 1000 | Markers, Multiple Lines |
| MATLAB | Default Plotly | MarkerStyle ['o', '*', '.', '+', 'x', 's'] markersize [8-10] | 1000 | |
| R | Default,Plotly R Library ggplot2 | | 1644 | |
| JAVA | XChart 3.5.1 javafx JFreeChart | MarkerSize (12-16) | 1644 | |
| Box Plots | | | | |
| Language | Library | Parameters | Number of plots | Types(Box with) |
| Python | Matplotlib v2.1.2 Plotly v2.5.1 Seaborn v0.8.1 | | 1000 | Notches, Multiple Boxes |
| MATLAB | Default | | 1000 | |
| R | Default,Plotly R Library ggplot2 | 19 | 1644 | |
| JAVA | XChart 3.5.1 | LegendPosition | 1644 | |

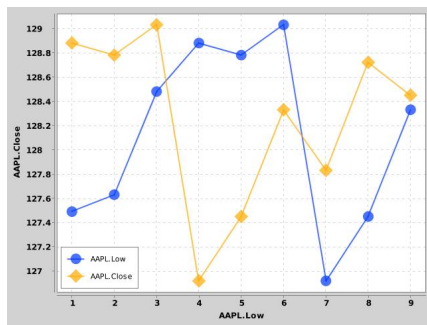


(a) Matlab stacked bar chart (bar width 16)

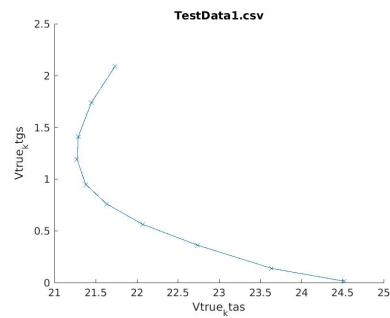


(b) R horizontal bar chart (bar width 16)

Figure 3.2: Example Bar Charts

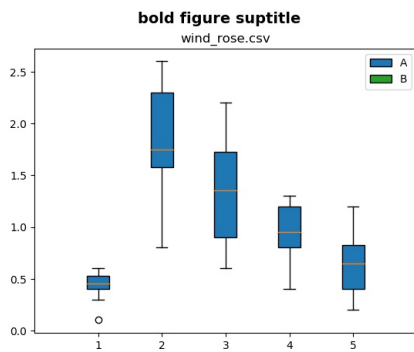


(a) Java line chart with diamond and circular markers

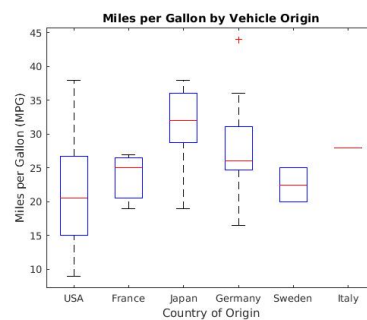


(b) simple Matlab line chart with Asterix marker

Figure 3.3: Example Line Charts



(a) Vertical multiple boxplots in python



(b) Vertical multiple boxplots in Matlab

Figure 3.4: Example Box Plots

3.3 Training, Validation and Test set

Figure ?? summaries how our dataset was split for various activities. In machine learning, a model is an artifact created during the training phase, this model can be likened to a function that maps specific features to their respective labels, the model is trained with a portion of the dataset called the train set. As the model is trained, the validation set is used to decide and pick which metric out of hyper-parameters, early stopping and architecture considerations yields the best performance. These help to adjust and optimize the model [22]. Finally, the test-set in machine learning, is the other portion of the dataset which was not used during the training phase. The test-set checks how well the trained model performs on unseen data by giving an unbiased assessment of the model. For this work, the dataset is divided into train, validation and test set. The validation set is not shown in the figure since it forms part of the train data.

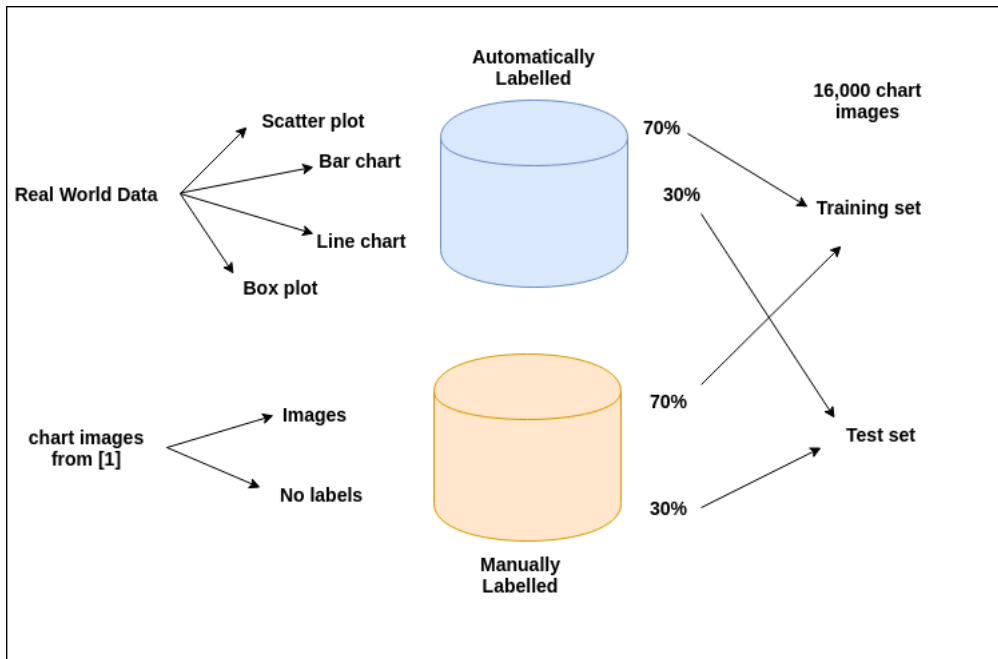


Figure 3.5: Steps involved in building the Classification model

3.4 The Architecture

3.5 Image Preprocessing

Image preprocessing is a very important step in any image based application. This process involves taking the image and improving it in a way that enables easier machine understanding. This ultimately makes it easier for the machine to extract important features for other operations.

4 Evaluation

5 Summary

6 Bibliography

- [1] V Karthikeyani and S Nagarajan. Machine learning classification algorithms to recognize chart types in portable document format (pdf) files. *International Journal of Computer Applications*, 39(2), 2012.
- [2] Xiao Liu, Binbin Tang, Zhenyang Wang, Xianghua Xu, Shiliang Pu, Dapeng Tao, and Mingli Song. Chart classification by combining deep convolutional networks and deep belief networks. In *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, pages 801–805. IEEE, 2015.
- [3] Manolis Savva, Nicholas Kong, Arti Chhajta, Li Fei-Fei, Maneesh Agrawala, and Jeffrey Heer. Revision: Automated classification, analysis and redesign of chart images. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 393–402. ACM, 2011.
- [4] V Shiv Naga Prasad, Behjat Siddiquie, Jennifer Golbeck, and Larry S Davis. Classifying computer generated charts. In *Content-Based Multimedia Indexing, 2007. CBMI'07. International Workshop on*, pages 85–92. IEEE, 2007.
- [5] Jihen Amara, Pawandeep Kaur, Michael Owonibi, and Bassem Bouaziz. Convolutional neural network based chart image classification. 2017.
- [6] Richard Boyle, Bahram Parvin, Darko Koracin, Nikos Paragios, and Syeda-Mahmood Tanveer. *Advances in visual computing*. Springer, 2007.
- [7] Yan Ping Zhou and Chew Lim Tan. Bar charts recognition using hough based syntactic segmentation. In *International Conference on Theory and Application of Diagrams*, pages 494–497. Springer, 2000.
- [8] Mingyan Shao and Robert P Futrelle. Recognition and classification of figures in pdf documents. In *International Workshop on Graphics Recognition*, pages 231–242. Springer, 2005.
- [9] Akihiro Inokuchi, Takashi Washio, and Hiroshi Motoda. An apriori-based algorithm for mining frequent substructures from graph data. In *European*

- Conference on Principles of Data Mining and Knowledge Discovery*, pages 13–23. Springer, 2000.
- [10] ujjwalkarn. An intuitive explanation of convolutional neural networks – the data science blog. <https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/>, August 2016. (Accessed on 05/16/2018).
 - [11] Ifu Aniemeka. A friendly introduction to convolutional neural networks | hashrocket. <https://hashrocket.com/blog/posts/a-friendly-introduction-to-convolutional-neural-networks#relu-layer>, August 2017. (Accessed on 05/17/2018).
 - [12] Adit Deshpande. A beginner’s guide to understanding convolutional neural networks – adit deshpande – cs undergrad at ucla (’19). <https://adeshpande3.github.io/A-Beginner%27s-Guide-To-Understanding-Convolutional-Neural-Networks/>, July 2016. (Accessed on 05/17/2018).
 - [13] Paulo Roberto Silva Chagas Junior, Alexandre Abreu De Freitas, Rafael Daisuke Akiyama, Brunelli Pinto Miranda, Tiago Davi Oliveira De Araújo, Carlos Gustavo Resque Dos Santos, Bianchi Serique Meiguins, and Jefferson Magalhães De Moraes. Architecture proposal for data extraction of chart images using convolutional neural network. In *Information Visualisation (IV), 2017 21st International Conference*, pages 318–323. IEEE, 2017.
 - [14] Rodriguez, J. Convolutional neural networks for the rest of us part iii: Benefits and motivation. <https://medium.com/@jrodthoughts/convolutional-neural-networks-for-the-rest-of-us-part-iii-benefits-and-motiv> 2011. [Online; accessed 5-Jun-2018].
 - [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
 - [16] Editor. what-to-look-for-in-training-dataset. <https://medium.com/@jrodthoughts/convolutional-neural-networks-for-the-rest-of-us-part-iii-bene> 2018. [Online; accessed 5-Jun-2018].
 - [17] James Eagan. Project datasets. <https://perso.telecom-paristech.fr/eagan/class/igr204/datasets>. [Online; accessed 20-April-2018].
 - [18] Michael Vogt. Datplot. <https://vincentarelbundock.github.io/Rdatasets/datasets.html>, 2011. [Online; accessed 20-April-2018].

- [19] plotly/datasets. Latex — Wikipedia, the free encyclopedia. <https://github.com/plotly/datasets>, 2011. [Online; accessed 20-April-2018].
- [20] vincentarel bundock. Rdatasets. <https://vincentarelbundock.github.io/Rdatasets/datasets.html>, 2011. [Online; accessed 20-April-2018].
- [21] Po-shen Lee, Jevin D West, and Bill Howe. Viziometrics: Analyzing visual information in the scientific literature. *IEEE Transactions on Big Data*, 4(1):117–129, 2018.
- [22] Akshay Balwally,. What is the difference between validation set and test set? <https://www.quora.com/What-is-the-difference-between-validation-set-and-test-set>, 2016. [Online; accessed 7-Jun-2018].