

OSNABRÜCK UNIVERSITY

STUDY PROJECT

---

# Computer Vision Based Asparagus Classification

---

*Authors:*

Maren Born,  
Michael Gerstenberger,  
Katharina Groß,  
Richard Ruppel,  
Sophia Schulze-Weddige,  
Malin Spaniol,  
Josefine Zerbe

*Supervisors:*

Dr. Ulf Krumnack  
M.Sc. Axel Schaffland

*The study report is submitted in partial fulfillment of the requirements  
for the degree of Master of Science*

*in the*

Research Group Computer Vision  
Institute of Cognitive Science

May 4, 2020

OSNABRÜCK UNIVERSITY

## *Abstract*

School of Human Sciences  
Institute of Cognitive Science

Master of Science

### **Computer Vision Based Asparagus Classification**

by Maren Born, Michael Gerstenberger, Katharina Groß, Richard Ruppel,  
Sophia Schulze-Weddige, Malin Spaniol, Josefine Zerbe

In the agricultural industry, the process of sorting food products into different quality classes is nowadays primarily achieved through sorting machines, using classical, hard-coded programming methods. One food product with many possible quality classes is white asparagus.

In this project, our team tackles the issue of asparagus classification by exploring different classical and state of the art machine learning and computer vision techniques. For this purpose, a range of supervised, semi-supervised, and unsupervised learning approaches are tested and evaluated. The data set is based on the output of the asparagus sorting machine Autoselect ATS II, which uses image data for classification. The result of the project shows that different approaches are suitable for solving the issue. For comparison with classical methods used in the industry and to provide a practical application, further testing of the approaches is still needed.

# Contents

<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 The project . . . . .	3
1.2 Background on computer vision based classification tasks . . . . .	4
1.3 Background on sorting asparagus . . . . .	5
<b>2 Data acquisition and organization</b>	<b>9</b>
2.1 Roadmap of the project . . . . .	9
2.2 Organisation of the study group . . . . .	12
2.2.1 Communication . . . . .	12
2.2.2 Teamwork . . . . .	13
2.3 Data collection . . . . .	14
2.4 Literature research . . . . .	19
<b>3 Preprocessing and data set creation</b>	<b>20</b>
3.1 Preprocessing . . . . .	20
3.2 Automatic feature extraction . . . . .	22
3.2.1 Length . . . . .	23
3.2.2 Width . . . . .	23
3.2.3 Rust . . . . .	24
3.2.4 Violet . . . . .	24
3.2.5 Curvature . . . . .	25
3.2.6 Flower . . . . .	26
3.3 The hand-label app: A GUI for labeling asparagus . . . . .	26
3.4 Manual labeling . . . . .	30
3.4.1 Sorting criteria . . . . .	30
3.4.2 Sorting outcome . . . . .	33
3.4.3 Agreement Measures . . . . .	34
3.4.4 Reliability . . . . .	35
3.5 The asparagus data set . . . . .	37
<b>4 Classification</b>	<b>39</b>
4.1 Supervised learning . . . . .	39
4.1.1 Prediction based on feature engineering . . . . .	40
4.1.2 Single-label classification . . . . .	43
4.1.3 Multi-label classification . . . . .	49
4.1.4 A dedicated network for head-related features . . . . .	55
4.1.5 From features to labels . . . . .	57
4.2 Unsupervised learning . . . . .	61
4.2.1 Principal Component Analysis . . . . .	61
4.2.2 Autoencoder . . . . .	71
4.3 Semi-supervised learning . . . . .	73

<b>5 Summary of results</b>	<b>78</b>
<b>6 Discussion</b>	<b>80</b>
6.1 Classification results . . . . .	80
6.2 Methodology . . . . .	82
6.3 Organization . . . . .	83
<b>7 Conclusion and outlook</b>	<b>85</b>
<b>A Work Overview</b>	<b>87</b>
A.1 Task list . . . . .	87
A.2 Report list . . . . .	87
<b>B Additional figures</b>	<b>89</b>
B.1 Single-label CNN . . . . .	90
<b>Bibliography</b>	<b>98</b>

# List of Figures

1.1 Asparagus Classification Tree . . . . .	7
2.1 Timetable of the Project . . . . .	9
2.2 Planned Roadmap . . . . .	10
2.3 Actual Roadmap . . . . .	11
2.4 The Autoselect ATS II at Gut Holsterfeld . . . . .	15
2.5 Output Trays of the Sorting Machine . . . . .	16
2.6 Example Asparagus Images . . . . .	17
2.7 Example Asparagus Images Querdel's Hof . . . . .	18
3.1 Cropping Based on Center of Mass . . . . .	21
3.2 The Labeling Dialog of the Hand-Label App . . . . .	27
3.3 UML Diagram for the Hand-Label App . . . . .	29
3.4 Example Image Feature Fractured . . . . .	30
3.5 Example Image Feature Hollow . . . . .	31
3.6 Example Image Feature Flower . . . . .	31
3.7 Example Image Feature Rusty Body . . . . .	31
3.8 Example Image Feature Rusty Head . . . . .	32
3.9 Example Image Feature Bent . . . . .	32
3.10 Example Image Feature Violet . . . . .	32
3.11 Example Image Not Classifiable . . . . .	33
3.12 Manual Labeling Output CSV-File . . . . .	34
3.13 Agreement Measure-Wise Comparison of all Features . . . . .	36
3.14 Feature-Wise Comparison of Agreement Measure Scores . . . . .	36
4.1 Feature Engineering Learning Curve for Color-Based Prediction . . . . .	41
4.2 Feature Engineering Learning Curve For Angle Based Prediction . . . . .	42
4.3 Feature Engineering ROC Curve . . . . .	43
4.4 Single-Label CNN ROC Curve . . . . .	47
4.5 Single-Label CNN Example Images Feature Hollow . . . . .	48
4.6 Single-Label CNN Example Images Feature Bent . . . . .	49
4.7 Multi-Label Binary Cross-Entropy Loss . . . . .	52
4.8 Multi-Label Hamming Loss . . . . .	52
4.9 Multi-Label Costum Loss . . . . .	53
4.10 Multi-Label $L_1$ Regularization . . . . .	53
4.11 Multi-Label $L_2$ Regularization . . . . .	54
4.12 Head Features CNN Training Sample . . . . .	56
4.13 Head Features CNN Learning Curve . . . . .	56
4.14 Head Features CNN ROC Curve . . . . .	57
4.15 Screenshot of Streamlit App . . . . .	58
4.16 Random Forest Classifier Confusion Matrices . . . . .	59
4.17 Distribution of Class Labels . . . . .	59
4.18 First Ten Eigenvalues and Eigenasparagus of Feature Length . . . . .	63
4.19 First Ten Eigenvalues and Eigenasparagus of Feature Hollow . . . . .	64
4.20 First Ten Eigenvalues and Eigenasparagus of Feature Bent . . . . .	65

4.21 First Ten Eigenvalues and Eigenasparagus of Feature Violet . . . . .	66
4.22 First Ten Eigenvalues and Eigenasparagus of Feature Width . . . . .	67
4.23 First Ten Eigenvalues and Eigenasparagus of Feature Rusty Body . .	68
4.24 First Ten Eigenvalues and Eigenasparagus of Feature Flower . . . . .	69
4.25 Feature-wise Scatterplots . . . . .	70
4.26 Latent Asparagus Space for MLP-VAE and Reconstruction Manifold	72
4.27 Network Structures for Semi-Supervised Learning . . . . .	74
4.28 Latent Asparagus Space for Semi-Supervised CNN-VAE . . . . .	75
B.1 File Overview Example of Original Image Data . . . . .	89
B.2 Single-Label CNN Example Images Feature Violet . . . . .	90
B.3 Single-Label CNN Example Images Feature Fractured . . . . .	91
B.4 Single-Label CNN Example Images Feature Flower . . . . .	91
B.5 Single-Label CNN Example Images Feature Rusty Head . . . . .	92
B.6 Single-Label CNN Example Images Feature Rusty Body . . . . .	92
B.7 Single-Label CNN Example Images Feature Very Thick . . . . .	93
B.8 Single-Label CNN Example Images Feature Thick . . . . .	93
B.9 Single-Label CNN Example Images Feature Medium Thick . . . . .	94
B.10 Single-Label CNN Example Images Feature Thin . . . . .	94
B.11 Single-Label CNN Example Images Feature Very Thin . . . . .	95
B.12 Single-Label CNN Example Images Feature Not Classifiable . . . . .	95
B.13 Single-Label CNN Accuracy and Loss . . . . .	96

# List of Tables

1.1	List of Asparagus Quality Classes . . . . .	6
3.1	Manual Labeling Feature Representation . . . . .	33
4.1	Feature Engineering Color-Based Prediction . . . . .	41
4.2	Feature Engineering Curvature Prediction . . . . .	42
4.3	Single-Label CNN Classification Results . . . . .	46
4.4	Multi-Label Model Structure . . . . .	51
4.5	Head Features CNN Performance . . . . .	56
4.6	Classification Report of the Random Forest Classifier . . . . .	60
4.7	Performance of the Semi-Supervised Convolutional VAE . . . . .	76
4.8	Performance of the Semi-Supervised Convolutional Autoencoder . . . . .	76

## Chapter 1

# Introduction

An essential step in the commercial asparagus cultivation is the sorting of the harvested asparagus into different quality classes. Depending on size, shape and color, a label is assigned to the individual asparagus which determines the quality class of the asparagus (United Nations Economic Commission for Europe (UNECE), 2017). Nowadays, this task is usually performed automatically, with modern asparagus sorting machines achieving a throughput of up to twelve asparagus spears per second (Ting, 2015). However, the accuracy of such machines can be unsatisfying. Manual resorting is usually necessary, and thereby causes substantial effort.

The aim of this study project was to investigate how techniques from computer vision, both classical and deep learning based, can be applied to improve classification results for asparagus sorting machines. The metric for the sorting corresponds to the quality of an asparagus. The quality in turn is defined based on a number of features that represent differences or even flaws in color, shape, or texture. Hence, improving sorting algorithms for the classification of asparagus requires reliable means of measuring these features. As such, we were interested in the question how the quality-defining features of asparagus can be estimated using state of the art computer vision and machine learning techniques.

The hands-on experience on a long-term project gave the chance to improve our project management skills, learn how to distribute work and how to organize ourselves in a larger team.

The project was supported by a local asparagus farm and a company specializing in asparagus sorting machines. They provided training data and expertise on the classification of the asparagus. The idea to apply new machine learning approaches to a commercial problem of the agricultural industry seemed promising and challenging at the same time. As the data received from the farm was unlabeled, a larger focus on the preprocessing of the recorded image data became inevitable. Further, the variance in quality classes and the subjective view of human sorting behavior toughened the perspective on a practical application into the actual sorting machine. Nevertheless, the original goal of studying different techniques in computer vision and testing their usability in a fixed hardware setting for asparagus classification gave valuable insights into the practical application of previously theoretically assessed knowledge.

Over the course of one year, the team members worked on the implementation of different approaches to tackle the issue of image classification. After intense engagement with the subject, the methodologies and the project outcome were documented and critically reviewed. On the following pages, the different stages of the project together with the results of the applied computer vision approaches are described in detail. The report chapters are mostly in chronological order, each with the focus on a different stage of the study project throughout the year.

In chapter 1 **Introduction**, the idea of the project is presented together with an introduction to the current standards of image classification in machine learning as

well as a general background on the quality classes, the sorting process, and classification of asparagus in the agricultural industry.

Chapter 2 **Data acquisition and organization** comprises the process of data acquisition and organizational background to the project. This includes task management and teamwork, as well as a thorough evaluation of our planning abilities as a group. Further, the process of data collection and a first inspection of the sorting machine are elaborated, with an emphasis on the first issues of unlabeled data. Unfortunately, the research for scientific literature, specifically regarding our topic of agricultural application of deep learning approaches, proved mostly disappointing. Nevertheless, it gave an insight into practical possibilities and a rough guideline for our endeavour. The collected literature that was found to be close to our topic is reported at the end of this chapter.

In chapter 3 **Preprocessing and data set creation**, the preprocessing phase is captured with the building of the data set at its end section. The first classical machine learning techniques for automatic feature extraction on image data are explored, followed by the process of manual feature extraction with the usage of a labeling application. The resulting outcome is reported, evaluated, and the labeled image data is transformed into a data set for the training of the chosen models.

Chapter 4 **Classification** constitutes the different approaches that were chosen to tackle the issue of image classification. All approaches are critically reviewed and discussed in the subchapters. An emphasis was put on the application of deep learning techniques for the classification of the asparagus. As the collected data includes labeled and unlabeled samples, methods from the range of supervised learning to unsupervised learning were tested.

Chapter 5 **Summary of results** contains the overall results of the classification approaches. The outcome of each approach is described and compared to the other methods as well as to the original classification accuracy of the sorting machine at the asparagus farm. Further, the overall results are discussed with an evaluation of their practical application in the food industry in chapter 6 **Discussion**. The overall outcome of the project is judged as a scientific study and as a team management experience.

In chapter 7 **Conclusion and outlook**, the findings are set into a broader perspective. The project and its results are summarized on a scientific basis as well as on the organizational level while future prospects of the project are assessed.

Before analyzing the specific context of software development for the application in classification tasks further, a thorough insight into the idea of the study project is given in the following chapter. Additionally, background on the addressed topics of machine learning and the sorting of asparagus is provided. In 1.1 **The project**, the objective of the study project is introduced. The next section 1.2 **Background on computer vision based classification tasks**, gives an overview on the machine learning techniques that explore how image classification can be implemented and how these approaches can provide a solution to our issue. The first chapter concludes with 1.3 **Background on sorting asparagus**, where the process of asparagus classification in the commercial food industry is illustrated and the different quality classes of asparagus are defined.

## 1.1 The project

The objective of the study project was to find both, conventional and deep learning computer vision based approaches which can be tested for their practical application in the commercial sorting of white asparagus. The different methods were implemented based on the image data that was received from the automatic sorting machine Autoselect ATS II employed by the asparagus farm “Gut Holsterfeld”.<sup>1</sup> The aim was to check whether approaches from the fields of machine learning and computer vision can be applied to improve the classification behavior of the asparagus sorting machine and, thus, if they can be used in a more industrial than scientific setting regarding the specific problem of asparagus classification. The initial intention to directly implement new software into the machine was postponed to allow for intensive research on the different approaches and on fine-tuning the received data to build a practical data set to train neural networks.

The idea for the project was formed by one of the students of the study group. She has a familial relationship to the asparagus farm Gut Holsterfeld and had received note of the unsatisfying classification performance of its sorting machine. The practical application to a real world problem sparked the interest of the group and the general curiosity towards computer vision, in particular neural networks, further inspired to deal with the sorting issue in a deep learning context, as well as with classical methods.

All project members are students in the field of Cognitive Science at the University of Osnabrück. The study project is part of the Master Program in Cognitive Science at the University of Osnabrück. It is supervised by Dr. Ulf Krumnack and M.Sc. Axel Schaffland. The intention of the study project is to confirm the ability of its participants to independently formulate and solve an unknown problem from the scientific context of one subject area using the methods and terms they previously learned (Universität Osnabrück, 2019a; Universität Osnabrück, 2019b). This includes the documentation and presentation of the results, the methodologies as well as the reflection on the work process. Within the scope of the project, for example, the development of software, analysis and interpretation of statistical data material are practiced. A further aspect of the study project is to deepen the communicative and decision-making competence of its participants (Universität Osnabrück, 2019a). The idea is to train independent project work in groups of students under conditions that are common for research projects in science or industry.

As the project took part in the scope of an examination for the Cognitive Science master program at the University of Osnabrück, most of the work was developed at the university, that is, at the **Institute of Cognitive Science (IKW)**. Additional image data was received from the asparagus farm “Querdel’s Hof”.<sup>2</sup> Further cooperation existed with the mechanical engineering company HMF Hermeler Maschinenbau GmbH that developed the sorting machine Autoselect ATS II and provided valuable expertise on the sorting issue.<sup>3</sup> All associated software is stored online, in our Github repository and at the institute internal storing system.<sup>4,5,6</sup>

---

<sup>1</sup>see the website of Gut Holsterfeld at <https://www.gut-holsterfeld.de/>

<sup>2</sup>see the website of Querdel’s Hof at <https://www.querdel.de/>

<sup>3</sup>see the website of HMF Hermeler at [www.hmf-hermeler.de](http://www.hmf-hermeler.de)

<sup>4</sup>The documentation to the project can be found at  
<https://asparagus.readthedocs.io/en/latest/>.

<sup>5</sup>The GitHub repository CogSciUOS/asparagus can be found at  
<https://github.com/CogSciUOS/asparagus>.

<sup>6</sup>The internal storing folder to the project at the University of Osnabrück can be accessed via /net/projects/scratch/winter/valid\_until\_31\_July\_2020/asparagus until 31/07/2020.

## 1.2 Background on computer vision based classification tasks

In this chapter, the current standard of computer-based image classification is described. The main focus will be on **Artificial Neural Networks (ANNs)** and classical computer vision techniques. A broad overview of relevant topics will be given and their importance for this project will be underlined. Computer vision is a field of computer science that aims to automatically extract high-level understanding from image data and provide appropriate output. It is closely linked with machine learning, which describes the ability of a system to learn and improve from experience rather than being specifically programmed. Machine learning is frequently used to solve computer vision tasks.

Image classification is one of the main subfields of computer vision which gained a lot of attention in the scientific as well as the economic world. Besides the classical computer vision techniques that use algorithmic approaches to determine patterns, edges, and other points of interest that can help to classify images, artificial intelligence was introduced to the field in the 1960s. Since then, more and more creative and complex artificial neural networks were introduced to solve numerous classification tasks including letter, face, and street sign recognition among others (Mirończuk and Protasiewicz, 2018; Balaban, 2015; Stallkamp et al., 2011).

Some experts in the field claim that artificial neural networks revolutionized the field of image classification, yielding better results than ever before (He et al., 2016a; Krizhevsky, Sutskever, and Hinton, 2012a). More and more challenges, for example ImageNet, were introduced and computer scientists all over the world implemented creative solutions for the proposed problems. Additionally, influential companies like Google have a deep interest in finding solutions for image-based classification problems and push the research on this and related topics even further. In the era of optimization, computer-based classification becomes indispensable in many industries and also finds its way into agriculture which is the field of interest in this study project. In the following, a short introduction in both classical computer vision techniques and artificial neural networks is given which will span a bridge to our current classification problem.

Neural networks are used for image classification in many domains. In contrast to the algorithmic approaches, it is not determined by the programmer how and what exactly the neural network learns. A large amount of data is provided to the model, which then extracts relevant information to learn the classification task. However, what is relevant to the model is not necessarily relevant or interpretable to a human observer. This is one of the biggest disadvantages compared to the classical computer vision approaches, which are understandable and, therefore, interpretable and adjustable. Another problem that comes along with this is that the bias for those approaches often does not lie in the code itself but in the data, which is by far more difficult to detect and surpass. For this reason, many see artificial neural networks as a black box, which may yield great results but can never be fully understood. Recent advances try to tackle that issue by researching artificial intelligence systematically aiming at making it explainable (Tjoa and Guan, 2019; Gilpin et al., 2018).

In image classification, usually **Convolutional Neural Networks (CNNs)** are used (Géron, 2019; LeCun and Bengio, 1995), which are inspired by the visual cortex of the brain. The idea is that highly specialized components or filters learn a very

specific task, which is similar to the receptive fields of neurons in the visual cortex (Hubel and Wiesel, 1962). These components can then be combined to high-level features, which can in turn be combined to objects that can be used for classification (Géron, 2019; Bishop, 2006; LeCun and Bengio, 1995). In CNNs, this concept is implemented by several successive convolutional layers in which one or more filters are slid over the input generating so-called feature maps. Each unit in one feature map looks for the same feature but in different locations of the input. In recent years, CNNs improved so much that they outperform humans in many classification tasks (Russakovsky et al., 2015; Karpathy, 2014).

Although machine learning exhibits very promising results and a lot of research and literature is available on the topic, many branches of industry still rely on traditional computer vision techniques in their implementation of image classification. This also applies to the asparagus sorting paradigm. To the best of our knowledge, no asparagus sorting machine is currently on the market that uses artificial intelligence for its classification algorithm. This brings us to the second topic of this chapter, namely the classical computer vision algorithms.

Many classical computer vision algorithms aim to detect and describe points of interest in the input images that can be generalized to features. For these features, low-level attributes such as rapid changes in color or luminance can be used. In contrast to the features learned with the help of machine learning, these features are not specific to the training data set and therefore do not depend on it being well-constructed. Further, they are usually created in a way that is interpretable by humans which makes them very flexible and easily adaptable to specific use cases. This is why in some cases, traditional computer vision techniques can solve a problem much more efficiently and in fewer lines of code than machine learning approaches.

In summary, both approaches have interesting implications for computer-based image classification tasks and provide us with promising techniques for our problem of asparagus classification. While we rely mostly on machine learning for the classification task itself, traditional computer vision algorithms are used to detect important features from the images, which not only helps us in the labeling process but also can be used for some of the binary classification tasks of the hand-labeled features.

### 1.3 Background on sorting asparagus

In this section, the background on sorting asparagus is displayed with a focus on the quality classes assessed by the asparagus farm Gut Holsterfeld. The owner of the asparagus farm Gut Holsterfeld, Mr. Silvan Schulze-Weddige and the CEO of the engineering company HMF Hermeler Maschinenbau, Mr. Thomas Hermeler, were consulted on this issue.

While asparagus accounts for a fifth of the area used for vegetable cultivation in Germany, the harvesting season of white asparagus only spans over a period of two months, beginning in April and ending on the 24th of June (Statistisches Bundesamt (Destatis), 2017; Weber and Quinckhardt, 2018). During this period, the asparagus is harvested, classified, and sold in accord with a price range that is defined by the quality class of the asparagus.

Label	Description
I A Anna	Quality class Extra is represented by I A Anna, I A Bona, and I A Clara. I A is defined as asparagus that is perfectly straight and white. There are no large pressure marks, no rust, no violet color, and no curvature. The identification label Anna marks that the width (diameter) of the spear is in a range of 20 - 26 mm.
I A Bona	Quality class Extra asparagus and the spear is of width 18 - 20 mm.
I A Clara	Quality class Extra and asparagus with a width of 16 - 18 mm.
I A Krumme	The asparagus fulfills all criteria for quality class Extra while it shows a slight curvature.
I A Violett	The asparagus is of a violet complexion, while it complies with all guidelines for quality class Extra.
II A	The asparagus is both curved and violet.
II B	The asparagus is curved, rusty, and/or in any other way damaged or classified as defective product.
Rost	The asparagus shows traces of rust or has rusty parts which occur when the roots of the plant were injured in a preceding year. This should not be confused with a fungal disease that can also be referred to as rust.
Dicke	The asparagus exceeds the norm of 26 mm in width.
Hohle	The asparagus is hollow from the inside.
Blume	The head region of the asparagus is about to bloom or it shows clear outlines of a flower.
Suppe	The asparagus has a width of less than 16 mm.
Bruch	Any asparagus that is below a length of 210 mm. However, another distinction in this class label is made between asparagus that has an intact head and a length of at least 100 mm (Kerze), an asparagus without head (Bruch), and an asparagus head alone (Köpfchen).
Keule	The upper end of the asparagus is thicker than the lower part. The shape is similar to a club, hence the name of the class. This class label was of no concern to the project because no image data of this class label could be recorded.

TABLE 1.1: **List of Asparagus Quality Classes** In this table, 14 quality classes for asparagus categorization are listed and described, according to the asparagus farm Gut Holsterfeld. Except for the class label Keule, all 13 quality classes were used as the goal label for the asparagus classification.

In the European Union, there is a uniform system for the sorting of asparagus into quality classes (Europäische Komission, 1999; United Nations Economic Commission for Europe (UNECE), 2017).<sup>7,8</sup> However, supply and demand usually determine the number and accuracy of these classes. One of the first defining features is the color of the asparagus which comprises four categories: white, violet, violet-green, and green (Europäische Komission, 1999). For this project, only the first two colors are of relevance. Further distinction is made between quality classes Extra, class I, and class II. The class Extra defines the product as perfect quality, class I defines it as good quality, and class II includes products that do not qualify for the other classes but satisfy the minimum requirements for commercial distribution (Europäische Komission, 1999). A last distinction is made in the characteristics

<sup>7</sup>see <https://mlr.baden-wuerttemberg.de/de/unser-service/presse-und-oeffentlichkeitsarbeit/pressemitteilung/pid/nationale-handelsklassen-fuer-frisches-obst-und-gemuese-abgeschafft-1/>

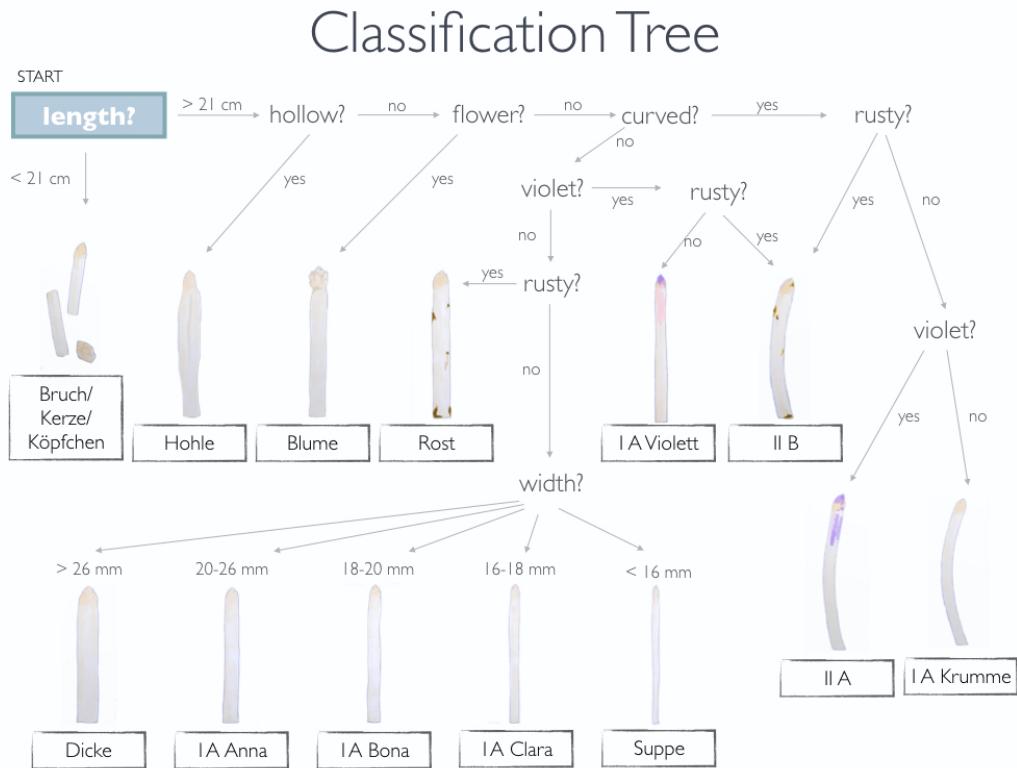
<sup>8</sup>see <https://www.bzfe.de/inhalt/spargel-kennzeichnung-5876.html>

length and width. White and violet asparagus may not exceed 220 mm in length. The minimal length of the asparagus varies but should be above 170 mm for long asparagus. Additionally, there is some level of tolerance accepted for the quality classes. Depending on the quality class, 5% - 15% of wrongly sorted asparagus is tolerated in a package or bundle (Europäische Komission, 1999).

Depending on how carefully the individual farmer further categorizes the asparagus, additional classes can be established. Regional differences to those classes make the challenge for the manufacturers of sorting machines even more complicated. The class labels as stated in this report depend solely on the sorting conduct of the asparagus farm Gut Holsterfeld and do not necessarily apply to any other classification system of the asparagus industry.

In Table 1.1, 14 quality classes are shortly described, of which 13 classes are relevant to this project. The classification tree in figure Figure 1.1 illustrates the decision process that underlies the categorization of the relevant 13 asparagus classes.

One of the challenges of asparagus classification with a machine lies in the human sorting error. The machine can assess exact calculations about, for example, the length or the width of a spear, while the human observer might not be able to detect minor differences in these features. Thus, a threshold has to be defined if a spear is sorted into a certain class. However, the data on which the machine



**FIGURE 1.1: Asparagus Classification Tree** The classification tree shows how each asparagus is attributed with a label of its quality class. It was drawn by the project group and follows the sorting rules of the asparagus farm Gut Holsterfeld. Starting from the upper left corner of the image, mostly binary decisions are made until a label is reached. The width of the asparagus is further subdivided at the end. Any further damaged or defective asparagus automatically belongs in category II B.

calculates its features to characterize an asparagus spear was previously labeled by a human. The subtle differences in color might have escaped the sorter and hence the machine will sort a spear into the class label violet while the spear would be judged as class label I A Anna by the human sorter. Since the human sorting behavior is subjective, the same asparagus can be categorized differently by two independent sorters. Furthermore, an asparagus sorted twice by the same person at different time points might be assigned with two different class labels.

A second problem for the sorting machine poses the interpretation of colors. The color can be perfectly recognized by the program, however, the structure of the color is also important for correct classification. The machine might not be able to distinguish whether the asparagus was, e.g. of brownish color because it is very rusty or simply very dirty in certain cases.

A third factor for classification difficulties is caused by the restricted view of the product. An asparagus can look perfectly shaped from one angle but might be damaged on the side that is not exposed to the camera of the sorting machine.

Another complication poses the question of demand and supply. During some seasons there is more asparagus of a certain class label and less of another one available. The farmer usually distributes the number of spears belonging to a quality class according to the seasonal conditions. For example, during a season with less high quality asparagus of classes I A, the sorting threshold will shift. Spears that are usually sorted into, e.g. a lower class will now belong to a higher class. The challenge for a sorting machine will not only be to sort for the class criteria but also to provide a flexible and individual solution in accord with the farmer's preferences.

These challenges are not impossible to overcome but they make it more difficult to find a solution to the sorting task and compromises might be necessary.

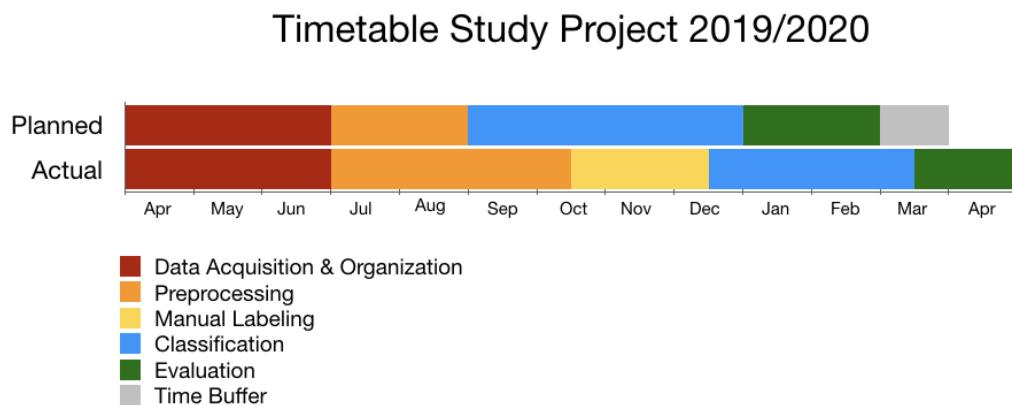
## Chapter 2

# Data acquisition and organization

The first section of the chapter [2.1 Roadmap of the project](#) gives an overview of the time management of the study group. It is followed by the subchapter [2.2 Organisation of the study group](#) in which communication and teamwork are assessed. In [2.3 Data collection](#), the acquisition of the data from the sorting machine Autoselect ATS II is described in more detail. The last section [2.4 Literature research](#) presents the retrieved literature concerning the issue of asparagus classification and agricultural product classification with machine learning based approaches.

### 2.1 Roadmap of the project

At the beginning of the project, a roadmap was created to structure the year into different working stages as well as to have an overview of the tasks and problems that needed to be addressed.



**FIGURE 2.1: Timetable of the Project** The upper timeline shows the estimated time of the study project from April 2019 to April 2020. The lower timeline displays how the time was spent. Both timelines differ in that the year was more optimistically planned than realized. A major factor was the lack of experience of the participants concerning the conduction of a larger project with many co-workers as well as concerning the general implementation of the preprocessing stage for machine learning classification. Another factor influencing the shifted timeline was the working structure of the team.

The timetable in [Figure 2.1](#) gives a broad outline of the major stages of the project. In the upper figure, it was estimated how much time for a specific phase is needed, whereas in the lower figure the time spent for the stage is given. Both figures are

## Planned Roadmap of the Study Project

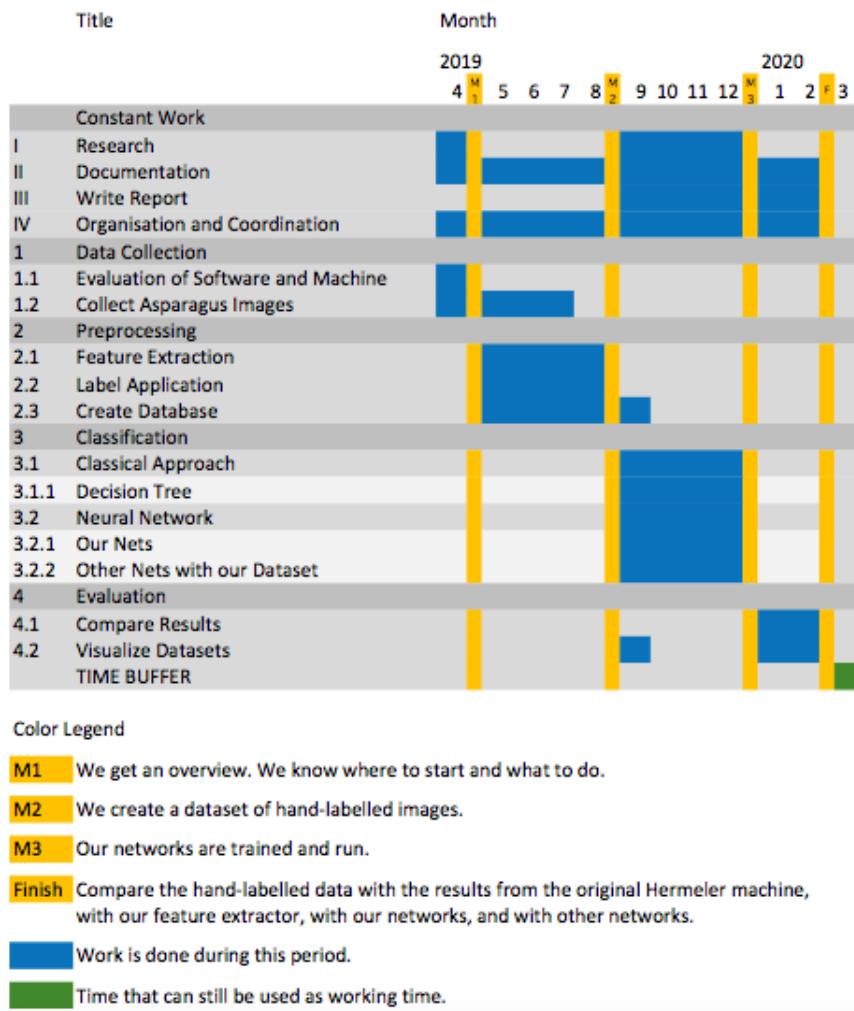


FIGURE 2.2: **Planned Roadmap** The figure shows the planned roadmap of the study project. It reveals how the time needed for each task was estimated in the beginning of the project.

structured to display the year in a line, starting in April 2019 and ending in April/-Mai 2020. The months are represented by the x axis while the colors mark the different working stages.

The project comprises four to five major stages: Data Collection & Organization, Preprocessing, Manual Labeling, Classification, and Evaluation. A detailed representation of the single tasks attributed to each stage can be found in the roadmaps in [Figure 2.2](#) and [Figure 2.3](#). The project started with the data collection and the organization of the study project. During the first stage, the images were recorded with the sorting machine, while the major planning and research for the project took place. In the second phase, most of the preprocessing happened, that is, preparing and labeling the image data. The classification stage includes the time spent on the machine learning approaches which were implemented and trained on the asparagus data. In the last stage, the approaches were evaluated and their results were compared. The different stages overlapped to a certain degree. For the purpose of this figure the start and end time is displayed as a hard boundary.

When comparing both timelines, some distinctions can be recognized. The upper

## Actual Roadmap of the Study Project

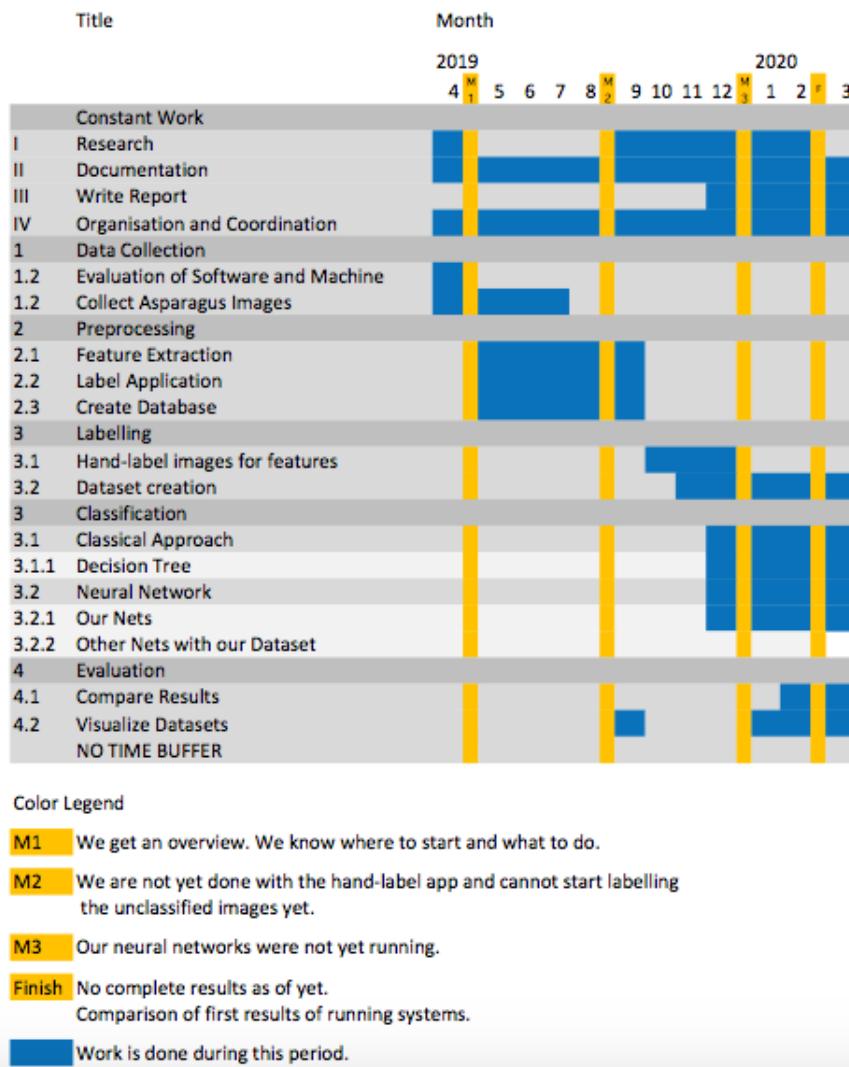


FIGURE 2.3: **Actual Roadmap** A roadmap that shows how the time was spent.

Timeline shows that preprocessing was estimated to be done by September. However, the phase continued until October, as can be observed in the lower timeline. Furthermore, the time for labeling a sufficient amount of images was underestimated, resulting in adjustments of the time attributed to this task. This led to the manual labeling receiving its own phase in the timetable, independent of the preprocessing phase.

The differences are depicted with different color codings. While the main focus of this project was supposed to be the application of different machine learning techniques to classify the data (color-coded in blue and green), the preprocessing phase and the data set creation/manual labeling posed to be most time-consuming (color-coded in orange and yellow).

In Figure 2.2 and Figure 2.3, the stage specific tasks can be seen in more detail. Again, both figures display the estimated time and the actual time, respectively. The headlines serve as a division into the major stages except for the first heading, "Constant Work", which shows the tasks that demanded continuous attention and effort throughout the year. The duration of tasks is represented in blue, while the

yellow lines mark milestones that are explained in the legends.

## 2.2 Organisation of the study group

In this chapter, the management of the work distribution and the communication are taken into focus. For this the tools that were used for communication and organization will be examined as well as the structure of the group work.

### 2.2.1 Communication

The main communication consisted of weekly meetings in which the working process was discussed and new tasks were distributed. In addition to those meetings, different platforms were used, which worked with varying degrees of success. Used platforms were Asana, GitHub and Telegram to facilitate the communication aside from group meetings. The different means of communication will be described and evaluated in the following abstract.

Regular meetings have been the most helpful in organizing our project. During the meetings, which usually lasted from one to two hours, a discussion leader and a protocol writer were picked. The protocols were saved for review in the GitHub project.<sup>1</sup> The meetings were characterized by long discussions about the best approach for the following steps and further possibilities to tackle the next challenge. The project's supervisors were usually present at the meetings, to bring in their expertise and to give the opportunity to ask concrete questions. During the first half of the project, tasks were always distributed at the end of each meeting. In the following meeting the progress of the work was discussed. This procedure was changed in the second half of the project. We started to use a schedule that described the different tasks and deadlines in detail. Additionally, we regularly gathered for co-working. The organizational meetings were continued, during which everyone gave precise and structured reports on their area of responsibility. This helped us to spend less time discussing and debating and more time working on our tasks.

Telegram is a cloud-based instant messaging service for the use on smartphones, tablets and computers. Starting from the first meeting, we had a constant conversation in a group chat on Telegram, in which we informed each other about the status of the project as well as support each other by answering questions. The group chat also created space for mutual motivation when needed. The majority of important information was exchanged via telegram.

Asana is used to distribute tasks and to communicate projects successfully. Many integrations of other applications, such as Slack, can help to achieve this. We used a board view where we listed tasks in different sections. But this function alone did not help us much in our project. The tasks were easier to distribute in direct consultation at physical meetings and demonstrated or discussed after completion. So it happened that Asana was not used enough to be helpful. If we had relied on communication with Slack or other agreed services or applications, it might have made more sense, but Asana alone was proven to be inefficient in our use case.

GitHub is a web-based popular platform using the version control system Git that helps developers to store and manage their code, and track and control changes to their project. During the project, we learned how to use it. Git allowed us to work from anywhere, which encouraged the workflow. Furthermore, we were able to automatically create a documentation via Sphinx. This means that by using the

---

<sup>1</sup>see <https://github.com/CogSciUOS/asparagus>

GitHub project in the right style, the protocols, work schedules, manuals, and code comments were automatically described in our documentation.<sup>2</sup>

### 2.2.2 Teamwork

This subchapter starts by introducing the team members and our previous experiences. This is followed by a description of the practical aspects of teamwork, the work structure, and the distribution of project-relevant tasks.

The project was an initiative of one of the students and, hence, a large part of the project members knew each other from their private environment but had not yet worked together. Further students joined the project after its public announcement to complete the team. Thus, the group consisted of members with varying degrees of knowledge about each other. The team was initially made up by Malin Spaniol, Maren Born, Katharina Groß, Josefine Zerbe, Michael Gerstenberger, Richard Ruppel, Sophia Schulze-Weddige, Luana Vaduva, Thomas Klein and Subir Das. None of the members had yet worked together as such on a project of this scope. During the course of the project, three members left the team for various reasons. Thomas left in July due to a change in his study program. This was an unfortunate occurrence because he posed a valuable source of knowledge in the field of computer vision. Further, Luana and Subir left in October to pursue different study projects.

The members brought a wide variety of backgrounds into the team through different bachelor programs or different majors in the broader field of Cognitive Science. In the beginning of the project, the team members had little to no experience in the application of computer vision or neural networks. The motivation of most students was to pursue new and interesting tasks in these fields. Four students had a theoretical background in computer vision, six students had gained some experience with neural networks through the course “**ANNs with Tensorflow**”, taught at the University of Osnabrück, while some had also taken machine learning classes during their study program. None of the members had prior knowledge about project management or task organization on a broader level. Git was previously only used by three students so far, but none of them were experts on its usage. Further, participants had neither experience with the Grid system of the **IKW**, nor with running jobs on different machines. Thus, the project started with ten members, each of them with a different level of experience in the most relevant fields of machine learning, that is computer vision and neural networks.

As none of the members had any previous experience with team formation, the team lacked some structure and a clear distribution of individual roles in the beginning. One reason for this could have been the harmonious atmosphere between team members due to their friendship. Further tasks such as the trips to the asparagus farm strengthened the team spirit and the social interactions. We therefore started with a very dynamic structuring of tasks distribution by making every decision democratically. Most tasks were done in smaller teams of two to three people. Rather few tasks were done alone. In our meetings, we often formulated possible next tasks, even for the distant future, but without assigning them to specific persons or working groups. This resulted in a lot of unassigned tasks and a discontinuous workflow.

In August, we restructured our own organization. On the one hand, this was due to the fact that the tasks changed and thus a new structure was more appropriate. On the other hand, the strengths of the individual team members were a reason for the

---

<sup>2</sup>see <https://asparagus.readthedocs.io/en/latest/>

restructuring. Some team members had less programming experience than others and therefore had difficulties realizing certain tasks in an equal time period and with the same precision as others. Even if they had good ideas in terms of concept, it was not possible for them to implement these ideas quickly enough so that they could be included into the project. Nevertheless, there was still the possibility to grow skills in new assignments. In addition, the democratic self-organization and difference in programming expertise led to a distortion in the time management of the group and some tasks were not finished in time. To integrate more of the strengths that the single team members brought and to tackle the issue of time management, it was decided to write a work schedule that distributed the work more appropriately, gave an overview of the tasks that still had to be done and how much time was left to do them.

Further, common working hours on campus were introduced. In addition, the responsibilities were shared for both, work tasks and the supervision of those work tasks. The common working hours ensured that questions and decisions that arose could be discussed in person. This was especially helpful when different tasks overlapped and required communication and agreement.

The supervision of the work was divided into manager roles, which means that the work was split into different main fields where each member was responsible for managing their assigned area, distributing tasks and keeping an overview of the relevant work inside their working field. Furthermore, one knew who to consult for questions in need of discussion or feedback. The meetings became more effective due to the new structure, and there was less discussion concerning task distribution. As we distributed roles, we were also more responsive to the strengths and weaknesses of the group members. Therefore one can say that the team structure and distribution of work changed over the course of the project. The strengths of single members were used more efficiently and the supervision of working areas led to a more structured supervision and task distribution.

## 2.3 Data collection

First, the asparagus sorting machine at Gut Holsterfeld is described before reporting about the process of labeled and unlabeled data collection.

The machine Autosort ATS II (2003) is designed for sorting white and green asparagus (HMF Hermeler Maschinenbau GmbH, 2015)(see [Figure 2.4](#)). The asparagus is arranged on a conveyor belt that runs it through the recording section of the machine. Here, a camera takes three pictures per asparagus as seen in [Figure 2.6](#). Small wheels on the conveyor belt rotate the asparagus in the meantime so that it can be photographed from several positions. In the best case, on each image a different side of the asparagus is recorded. The conveyor belt transports the spear further and it is sorted into a tray depending on the chosen class label by the machine. The sorting is based on the parameters width, length, shape, curvature, rust, and color. A total of 30 criteria for classifying an asparagus spear are used to describe these parameters. The calculation of the single features is based on a classical analytical approach. For example, the feature color is composed of eight sub-parameters. Each spear is reviewed at different areas (the head of the asparagus, the area below the head, and the stem) and judged for its hue in percentage. The values are compared and, according to a threshold, the spear is sorted into a color category. For all parameters, there is a minimal threshold and a maximal threshold. As another example, the parameter width is calculated at three points at the asparagus (top, middle, and bottom part). From these three values, an average value is calculated that decides in which category the asparagus is sorted.

If an asparagus exceeds the maximal threshold, it is not recognized and cannot be sorted accordingly. Thus, all parameters have to have an upper and a lower threshold, including parameters like shape, curvature, and color. When evaluating which thresholds to choose, it is recommended to check that most asparagus spears tend to be in between the average value and the maximal threshold, with a larger tendency to accumulate around the average value. All parameters can be freely chosen by the farmer and can in this way be fitted to the needs of the respective asparagus farm.



**FIGURE 2.4: Autoselect ATS II at Gut Holsterfeld** In the figure, the asparagus sorting machine Autoselect ATS II at Gut Holsterfeld can be seen. The conveyor belt transports the asparagus from the right side of the image to the left side. It thereby passes the camera system. The display of the machine gives information on the parameters and the images. The machine is mainly controlled from here.

According to the manual, the number of quality classes is selectable but limited by the number of output trays available to the machine. The machine Autoselect ATS II at Gut Holsterfeld has 16 trays (see Figure 2.5). The user can define the order of parameters. The manufacturer suggests to first sort for length and width, then sort for colors parameters, and analyze the shape parameters last.

Before the first use of the machine, all parameters are selected after a calibrating charge of asparagus has run through the machine. Then, the user can adjust the thresholds accordingly.

The accuracy of the sorting machine is described to be as good as 90% best case by the manufacturer, while the farmer at Gut Holsterfeld reported it to be around 70% at best, with re-sorting being necessary by professional sorters. Especially categories like Blume or Hohle were considered to be inconsistent by both, manufacturer and farmer. Further information could not be given on the software of the machine. A meeting with a representative of the engineering company HMF that manufactured the sorting machine was arranged.<sup>3</sup> Unfortunately, the source code itself was not available to HMF as it was produced by another company.

In the following it is described how the image data was collected. It is possible to save images with the Autoselect ATS II, however, the storage space on the machine is very limited. Further, the selection of images to be saved is restricted to only

---

<sup>3</sup>see [www.hmf-hermeler.de](http://www.hmf-hermeler.de)



**FIGURE 2.5: Output Trays of the Sorting Machine** Depending on the class label, the asparagus is sorted into one of the machine's 16 trays.

1000 images. One workaround to the problem is the installation of the Teamviewer software on the machine and the connection of an external hard drive. After the installation, the process of image collection could be started remotely. This work was very ineffective and time-consuming. The data could not be directly transmitted to another computer because the internet connection at the farm is too slow. An automatic transfer of the images to the external hard drive was not possible until the installation of an automatic file moving service, for which the requirements are described below.

The file moving program needs to transfer the images to a new saving destination and has to run in the background without disturbing the workflow of the sorting machine. After research on background processes and programs, the decision was made to use a service, that is, a system process running independently of any program. As Windows is used as the operating system of the sorting machine, the development was done with the .NET framework in the programming language C#. The package provided is called Topshelf.<sup>4</sup> Topshelf is a service hosting framework for building Windows services using .NET.<sup>5</sup> With the package, it is possible to develop a console application in the development phase, compile it as a service, and install it later via the console. Previously, it was not possible to debug services during the development phase. The function of the service is based on the FileSystemWatcher object from the System.IO namespace.<sup>6</sup> In the main program, a list of files in the source folder is kept. Files that are older than one hour are

<sup>4</sup>For further information on Topshelf, see <https://github.com/Topshelf/Topshelf>.

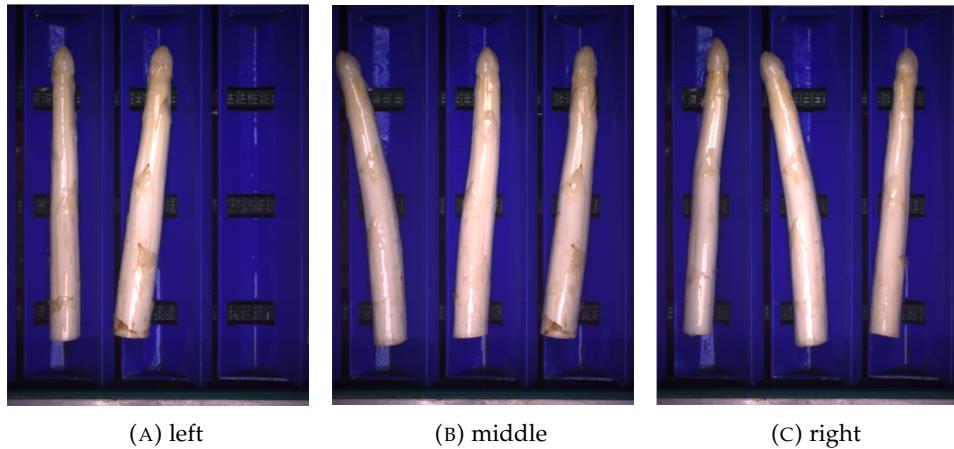
<sup>5</sup>For further information on the .NET framework, see <https://docs.microsoft.com/en-us/dotnet/framework/get-started/overview>.

<sup>6</sup>For further information on the SystemFileWatcher, see <https://docs.microsoft.com/en-us/dotnet/api/system.io.filesystemwatcher?view=netframework-4.8>.

moved to the target folder on the external drive. The selected files are moved by a function that is called, when an event is triggered. The event is triggered by the FileSystemWatcher after subscribing to different flags. Shortly after initialization, the service was adjusted because removing the images from the C disk straight away caused the sorting program to stop. The problem is solved by keeping the most recent 1000 images and moving older images to the external disk.

The project members split in groups of two and exchanged the hard drive two times a week. The collected images were then transferred to storage capacities of the university.

Before the project had started, the farm had not collected any labeled data such that relying on the labels that the machine attributes to each asparagus is not possible. A solution is to send the already sorted asparagus a second time through the machine. Like this, it was assumed that not only labeled data can be acquired but also the parameters calculated by the machine. Unfortunately, it is not possible to extract the parameters from the machine. Another issue poses the fact that a second sorting is not good for the quality of the asparagus. Further, at least one project member has to be involved in the re-sorting and, thus, has to be present at the farm. The sessions of exchanging the external hard drive and collecting labeled image data were then combined.

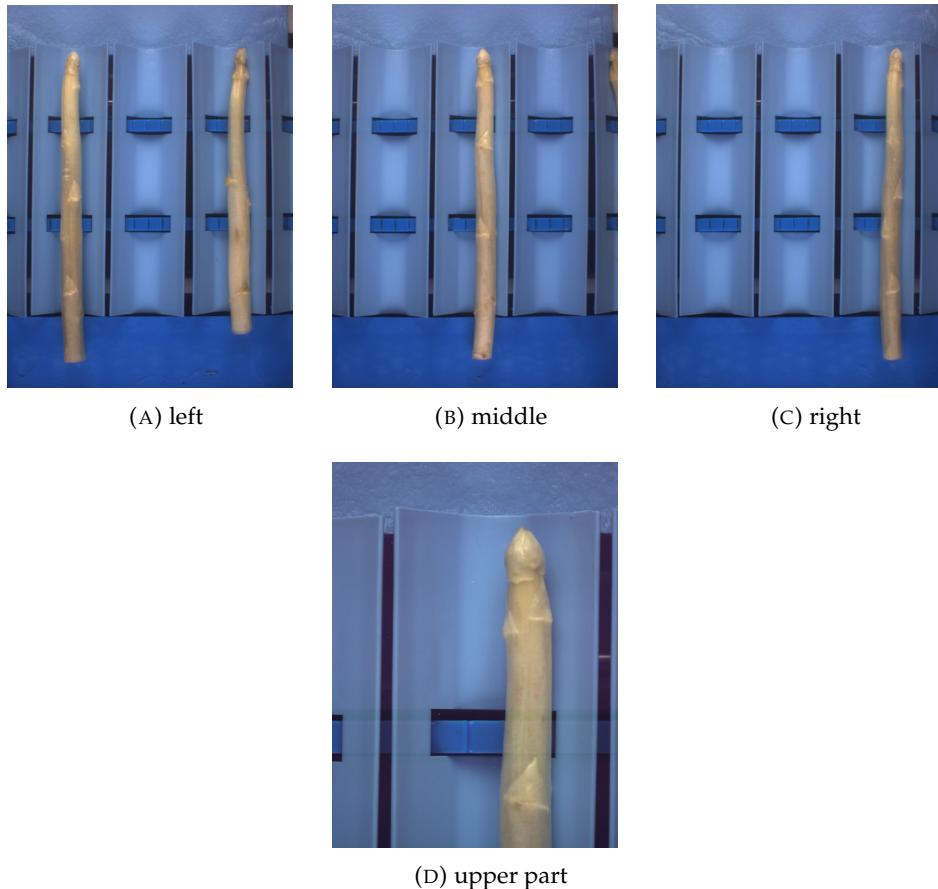


**FIGURE 2.6: Example Asparagus Images** Example pictures of the quality class I A Anna. The asparagus is arranged on a conveyor belt that runs it through the recording section of the machine, where a camera takes three pictures. In picture (A) the asparagus is to the left, in (B) it is in the middle, and in (C) it is to the right. Small wheels on the conveyor belt rotate the asparagus in the meantime so that it can be photographed from several positions. In the best case, on each image a different side of the asparagus is recorded. The conveyor belt transports the spear further and it is sorted into a tray depending on the chosen quality class by the machine.

An example image of the received data can be seen in [Figure 2.6](#). There are three pictures per asparagus. The image resolution is around  $1040 \times 1376$  pixel per image, with an RGB color space.

In total, 612113 images were collected, with each class label being represented with at least 309 images, corresponding to 103 asparagus.

At the asparagus farm Gut Holsterfeld, 591495 labeled and unlabeled images were collected with the Autoselect ATS II. The number of unlabeled data is 578226 images, thus, around 192742 different asparagus spears. Of the labeled data, we collected 1005 images with the label I A Anna, 908 images for I A Bona, 513 images of I A Clara, 936 images of I A Krumme, 1514 images of I A Violett, 1051 images of II A, 1468 images of II B, 1169 images of Rost, 749 images of Dicke, 775 images of Hohle, 1717 images of Blume, 1157 images of Suppe, and at least 309 images labeled as Bruch. The image number does not represent the number of different asparagus, as each asparagus is represented by three distinct images.



**FIGURE 2.7: Example Images Querdel's Hof** Example pictures of the asparagus of the farm Querdel's Hof. In picture (A), the asparagus spear is to the left, in (B), it is in the middle, and in (C) it is to the right. A fourth picture (D) is taken with a second camera with focus on the upper part of the asparagus.

Additionally, a few images could be recorded at another asparagus farm, Querdel's Hof, in Emsbüren.<sup>7</sup> The farm sorts the asparagus with an updated version of the Autoselect ATS II at Gut Holsterfeld, that is, it uses the same software but other hardware. In particular the resolution of the camera was improved and a second camera was installed that focuses on the head region of the asparagus. At Querdel's Hof, 20616 images were collected in total, 76 from the class label "normal", 152 from the class label "violet/flower", and 20388 unlabeled images. At Querdel's Hof, each asparagus is represented by four images: three images show the asparagus from different perspectives and a fourth image depicts solely the head region.

<sup>7</sup>see <https://www.querdel.de/>

An example image can be seen in [Figure 2.7](#). No internet connection could be established to the farm, thus, no further images were collected. Moreover, the data format of the images from Querdel's Hof is different to the data from the farm Gut Holsterfeld, due to the additional head image. Therefore a combination of both data sets was not convenient.

## 2.4 Literature research

In the classification of food products, there are numerous possibilities to apply classical machine learning and [ANN](#) approaches for classification tasks on image data ([Bhargava and Bansal, 2018](#); [Brosnan and Sun, 2002](#)).

None of the found papers were suitable as blueprints for the asparagus classification project. However, some of them helped to get an idea of how to proceed with the project. For example, some papers show how the preprocessing phase could be structured ([Mery, Pedreschi, and Soto, 2013](#)), or they evaluate the machine learning methods that were already used on other food classification tasks ([Bhargava and Bansal, 2018](#)). Further, some of the literature is concerned with the classification of food products but not with differentiating between as many classes as 13 ([Diaz et al., 2004](#); [Kılıç et al., 2007](#)). Often, the variance in the food products used is either too high ([Zhang and Wu, 2012](#)) or too low ([Kılıç et al., 2007](#); [Al Ohali, 2011](#)) in comparison to the variance in our project data. One paper evaluates the sorting of asparagus, however, it only does so on a small data set with three categories of green asparagus ([Donis-González and Guyer, 2016](#)). Further papers on food classification are not detailed enough in their explanations and do not share the information needed for replication ([Pedreschi, Mery, and Marique, 2016](#)). Another paper is mainly about the implementation of a certain toolbox ([Mery, Pedreschi, and Soto, 2013](#)).

As the available data was only sparsely labeled ([Olivier, Bernhard, and Alexander, 2006](#); [Zhu, 2005](#)), further research was done to evaluate the use of a semi-supervised learning approach. Details about the corresponding literature can be found in [section 4.3](#). In regards to deep learning-based approaches, classical neural networks – such as AlexNet, [VGG16/VGG19](#), GoogleNet, Capsule Networks, DenseNet, ResNet or [Network in Network \(NIN\)](#) – were assessed and presented for better understanding of the range of possible pre-trained networks and ideas for network structures ([Krizhevsky, Sutskever, and Hinton, 2012b](#); [Simonyan and Zisserman, 2014](#); [Szegedy et al., 2015](#); [Sabour, Frosst, and Hinton, 2017](#); [Huang et al., 2017](#); [He et al., 2016b](#); [Lin, Chen, and Yan, 2013](#)). Also, classical computer vision approaches were considered, like multiclass [Support Vector Machines \(SVMs\)](#) ([Prakash et al., 2012](#)).

## Chapter 3

# Preprocessing and data set creation

In this chapter, the different preparatory steps for the recorded data are described, including the creation of a data set which is usable for any machine learning or computer vision approach to analyze the image data.

In [3.1 Preprocessing](#) the data is assessed and simplified for any further processing. The second section, [3.2 Automatic feature extraction](#), deals with the creation of feature scripts that were researched and implemented for an automatic recognition of single features. The results were combined in an application which is described in detail in [3.3 The hand-label app: A GUI for labeling asparagus](#). In [3.4 Manual labeling](#), the process of hand-labeling the images for their feature class with the label application is described, followed by a section analyzing the results and comparing the overall agreement of the labelers. The last section [3.5 The asparagus data set](#) concludes with the creation of the final data set, used for the later training of the neuronal networks and other approaches to detect the label of a spear from its three images.

### 3.1 Preprocessing

Before implementing any approach that allows to predict a label to an asparagus spear, the recorded image data has to go through multiple preprocessing steps. In the following, these are elaborated in detail.

As described in [section 2.3](#), each asparagus can be found in three pictures, one in each of the three positions – left, center and right. The image names are used to combine the three relevant images and determine in which position the asparagus is captured. The images are cut into three pieces and renamed in a way that makes it clear which images belong together. Each asparagus gets a unique identification number and the three perspectives are denoted with “a” for left, “b” for center, and “c” for right. For example, the image 42\_b.png is the center image of the asparagus piece with the identification number 42.

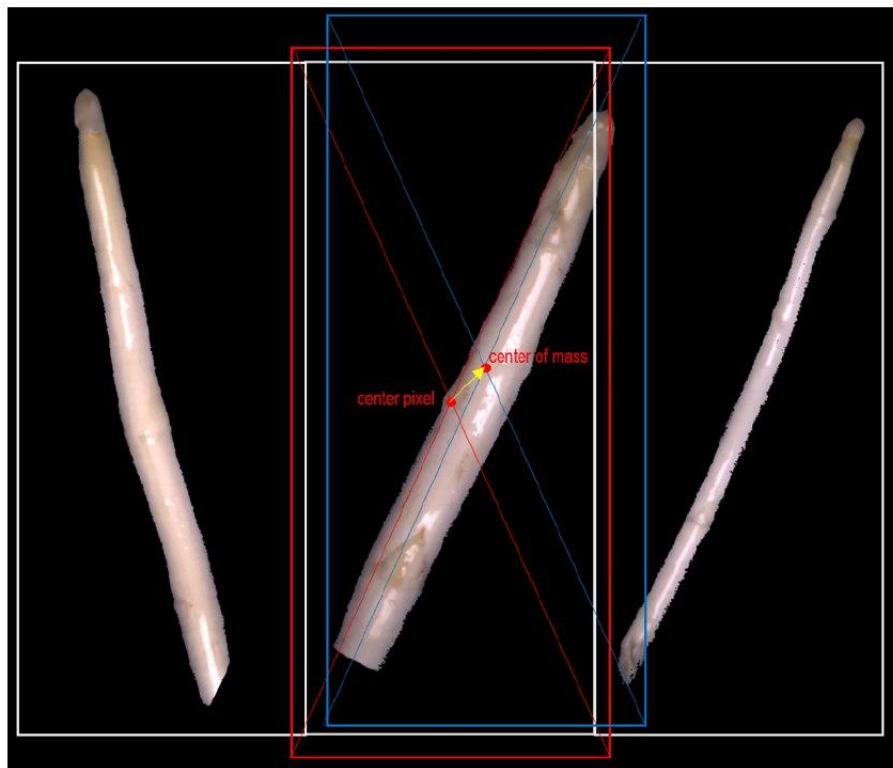
Another step is to remove the background. As the conveyor belt is blue, there is a high contrast to the bright asparagus spears, aiding the background removal. Hence, it is possible to mask the asparagus spear using the hue of the HSV representation of each image. All pixels with a blue hue and very dark regions are marked as background by thresholding the value component. This is particularly important for the automatic feature extraction (see [section 3.2](#)).

To combine the three perspectives of each spear, the cropping regions are set to three patches that are little larger than the compartments of the conveyor belt. They are located such that they cover the piece of interest. As it has been found that some tilted spears span across the borders between the original images, the cropping window in the concatenated images was moved. The new coordinates are determined by shifting the center of the current cropping box to the center of mass of the contained foreground pixels. When repeating this step in a second iteration,

A: Example for source images that relate to one asparagus piece



B: Moving the region of interest to the center of gravity of pixels



**FIGURE 3.1: Cropping Based on Center of Mass** (A) The upper part shows three unprocessed images that depict an asparagus piece from three angles. The respective piece is in the center of the second depiction. (B) The lower part shows the effect of moving the cropping window (red) to a new location (blue) by moving its center to the center of gravity (yellow arrow). The approach is applicable for all three images given that the initial coordinates are set to the respective positions.

the cropping boxes contain each of the three perspectives. However, small parts of the neighboring depiction potentially end up in the region of interest as well. These remainders are removed by masking out all pixels that do not belong to the main blob of connected pixels

Additionally, we trim as much of the unnecessary information as possible. Therefore, the asparagus is rotated upwards to reduce the variance in angle. The rotation angle is achieved by binarizing the asparagus piece image into foreground and background pixels, calculating the centerline as the mean pixel location along the vertical axis and fitting linear regression to the centerline pixels. The estimates for the angles are retrieved using arcus tangens and the images are rotated accordingly. By doing that, any distortions like reflections or single noisy pixels are set to zero and only the asparagus piece itself remains. The only faulty areas that can remain with this approach are areas that are connected to the asparagus. The position of the asparagus within each image part was further improved by centering the fragments of the image around the center of mass using the same algorithm that is used to crop the image patches. This way the asparagus is always in the middle. If the new cutting line falls out of the image boundaries, the missing area is filled with zeros.

Another preprocessing step, which generates an additional collection of images, is converting the images to color palette images. An appropriate palette and the respective mapping of RGB values to palette colors is determined using clustering in color space. First, a set of RGB tuples is collected by adding pixel values for the depiction of 10000 asparagus spears. Second, the resulting list of RGB tuples is converted to an image and third, a palette is determined using a clustering algorithm that determines the position of cluster centers while maximizing the maximal coverage. The resulting cluster centers can be displayed as a list of RGB values and represent the color palette (Python Image Library, 2020; Zarandi, Davari, and Sisakht, 2011). Each image of the downscaled data set is transformed to the palette representation. Visual inspection shows little quality loss such that it can be assumed that the relevant information for image classification is well preserved.

Several additional collections of preprocessed images are computed based on the data without background. This holds for downscaled versions as well as for a version that contains the asparagus heads only. To compute the latter, the images are padded to avoid sampling outside of the valid image boundaries and the uppermost foreground row is detected. Subsequently the center pixel is determined and the image is cropped such that this uppermost central pixel of the asparagus is the center of the uppermost row of the snippet. The resulting partial images of asparagus heads are rotated using the centerline regression approach described above. The approach has proven reliable and the resulting depictions are used to train a dedicated network for head related features(see subsection 4.1.4).

## 3.2 Automatic feature extraction

In this chapter, the different automatic feature extraction methods will be described, alongside with the results that were achieved and future steps that could be taken to improve the results further. For each feature detection method, the images with removed background are used.

According to the producer of the used sorting machine, the software currently on the market uses very similar methods, which gives us reason to believe that the classical methods might not be adequate to yield a better classification than the current status quo.

The results vary greatly between the different feature extraction methods. While the functions to detect the width and length and whether the asparagus is violet or bent are good enough to be integrated into the hand-label app (see section 3.3), other features turned out to be more difficult.

### 3.2.1 Length

The length detection uses a pixel-based approach. It counts the number of rows from the highest to the lowest pixel that is not a black background pixel and therefore not zero. The asparagus is rotated upwards, as described in section 3.1. This is done to improve the results, as the rows between the highest and the lowest pixel are counted and not the pixels themselves. This technique is a simplification, which does not represent curved asparagus very well, because it will have a shorter count than it would have if the pixels were counted along the asparagus piece. But in reality, there are not a lot of asparagus spears close to the decision boundary between fractured and a whole piece. Usually, the asparagus is picked a few centimeters longer than necessary and then cut to the desired length. The only asparagus shorter than that length are the ones that break while in the sorting process. Moreover, if they break they generally break closer to the center of the asparagus rather than at the ends of it. Therefore, the difference in length detection does not matter for our classification.

All in all, the length detection yields good results that are very helpful for the hand-label app. The next step would be to train a decision boundary that determines which number of pixels should be the threshold to differentiate between fractured and not fractured. At first, we tried to calculate this threshold by finding a conversion factor from pixel to millimeter, as we know the cut off in millimeters. But this approach appeared to be more difficult than anticipated, because the conversion factor varies in the different image positions. This problem only became apparent after the asparagus season had ended, for which reason we could not reproduce the camera calibrations in retrospective in order to take well-measured images, for example from a chessboard pattern. Accordingly, the threshold needs to be deduced from the data or learned with a machine learning approach.

### 3.2.2 Width

The width detection uses a very similar approach as the length detection. It takes the pixel count from the left-most to the right-most pixel in a certain row as a width measure. But in contrast to the length, the width was measured at several different image rows from which the mean width was taken. As the width detection functions reliably it is integrated in the hand-label app.

The algorithm operates as follows: Firstly, the images are binarized into foreground and background, which means setting all pixels that are not zero, and therefore not background, to one. After that, the upper-most foreground pixel is detected and the length is calculated with the length detection function as described above. The length of the asparagus is used to divide it into even parts. This is done by determining a start pixel and dividing the remaining rows that contain foreground pixels by the number of positions one wants to measure at. This way several rows are selected in which the number of foreground pixels is counted. One can interpret each row as a cross-section of the asparagus, therefore the number of foreground pixels is a direct measure for the width. Then, the mean of these counts is calculated and used as the final width value. As the head of the asparagus can be of varying form and does not represent the width of the whole asparagus well, it is excluded from the measure. This is done by selecting a start pixel below the head area instead of naively choosing the uppermost pixel. To be precise, the start pixel is chosen 200 pixels below the uppermost pixel in order to bypass the head area with certainty. As described in the length detection, also the width detection might lead to slightly

different outcomes on curved asparagus spears than the true values. Again, this difference is regarded as irrelevant in our case.

### 3.2.3 Rust

The rust detection finds all pixels that fall in the range of RGB values that correspond to the color brown. Those pixels are counted and normalized by the maximal number of possible pixels that could be rusty, namely the number of foreground pixels. But as it is impossible that the whole asparagus is rusty and hence that all the foreground pixels fall into the relevant range of RGB values, this normalization yields small numbers as results. To give an example, an output value of 0.13 is already considered moderately rusty. The lower and upper bound for the RGB values are set to [50,42,30] and [220,220,55], respectively. That means, all pixels that have a red value between 50 and 220, a green value between 42 and 220 and blue value between 30 and 55 are considered as rust.

Visual inspection shows that the rust detection algorithm works well to detect rusty areas and barely misses any rusty parts. Setting a threshold for the number of pixels needed to be classified as rusty still turns out to be difficult. Only clusters of brown pixels are a reliable indicator for rust. Many pixels with a brown color distributed over the whole spear are not supposed to be classified as rust. To make this intuition more clear, it could be the case that two asparagus have the same number of brown pixels, but in one case they are all connected building a cluster and in the other case they are evenly distributed over the whole asparagus spear. That would mean that, although both yield the same output value, only one of them should be considered as rusty. However, it should be mentioned that this is merely an artificial example to display one draw-back of the implementation. In reality, it is very unlikely that an asparagus has a large number of brown pixels that are not in fact rust. Nevertheless, it remains unsolved to set a robust threshold that works well on the whole data set.

One problem that cannot be solved algorithmically, is dirt in the sorting machine. If the machine is not cleaned thoroughly and regularly, dirt can be falsely classified as rust as it often falls in the same color range. Another problem can be a change of lighting when taking the images. Both issues can be controlled for, but have to be communicated well to the farmers.

### 3.2.4 Violet

Especially when exposed to light asparagus spears tend to change in color. An initial shift from the desired white appearance to a slightly rose or violet tone can be observed first. This change in color, that is considered a flaw according to german quality standards, typically affects the head of asparagus spears. However in some cases a non-uniform distribution of more or less violet areas can be observed. Moreover, it has been shown in practice that color impression is highly subjective (Luo, 2000). This manifested in discussions of edge cases during labeling. Effects of meta contrast arguably occurred when many similar spears were labeled in succession which makes minor color differences more visible (Reeves, 1981). Hence, effects of meta contrast potentially affect the color impression even on an individual level. The lack of a formal definition for violet asparagus spears poses challenges to approaches of measuring this feature.

In a simple approach to measure whether an asparagus spear is violet or not, color hues are evaluated. More precisely, this strategy is based on evaluating histograms

of color hues that are calculated for foreground pixels of the asparagus images after converting them to the HSV color space. Pale pixels are removed from the selection by thresholding based on the value component of the HSV representation. Finding the optimal threshold was proven difficult and corresponds to the abovementioned question of what it means for an asparagus spear to be violet. A value of 0.3 was considered a good compromise. All three perspectives were taken into account to compute one histogram per asparagus spear. A score was calculated by summing up the number of pixels that lie within a violet color range. A second threshold was used as the decision boundary for violet detection. The direct and intuitive feedback in the hand-label app showed the relation between varying thresholds and the prediction. It has been shown that lowering thresholds also means the feature extractor becomes more sensitive at the price of a reduced specificity. Best overall matches (accuracies) with the subjective perception were found for very low thresholds. In many cases, however, measurements based on this definition of violet did not match the attributed label.

Hence, another sparse descriptor was derived from the input images. However, instead of thresholding pale values and calculating the histograms of color hues this approach relies directly on the colors that are present in the depiction of an asparagus spear. As the 24 bit representations contain a vast amount of color information in relation to the number of pixels it is, however, unfeasible to use these as an input. Instead, the color palette images can be used. Histograms of palette images can serve as the basis to define the feature violet in a way that captures more of the initial color information while being simple and understandable enough to allow for customizations by users of sorting algorithms or machines. As a consensus regarding such an explicit definition is hard to achieve and somewhat arbitrary, the descriptor was used to learn implicit definitions of the feature by examples (see subsection 4.1.1).

It has been shown that explicit ways of directly measuring, in how far an asparagus spear is violet can be implemented but heavily depend on the definition of this feature. As color perception is highly subjective across and even within subjects (Reeves, 1981), machine learning approaches that are trained on human labeled data appear to be promising. Using them can help generalizing a definition for the degree to which an asparagus piece is violet.

### 3.2.5 Curvature

Multiple curvature scores can easily be computed based on regression fits to the centerline of an asparagus spear. For example, the parameters of linear or polynomial regression can be interpreted as a description of the curvature of an asparagus spear.

Deriving sparse descriptions is based on a two stage approach. In the first stage, the centerline of an asparagus spear is computed because it is considered to be a good description of the curvature of asparagus spears. Each asparagus spear is approximately oriented vertically. This means that also for curved spears the head relies within the top center of the image (see section 3.1). The centerline is computed by binarizing the image into foreground and background and computing the mean of pixel locations along the vertical axis (i.e. for each row). The resulting binary representation shows a single pixel line. It serves as the input to the second stage of curvature estimation.

In the second stage, curves are fit to the pixel locations of the centerline. For a simple score, linear regression is employed and the sum of squared errors is thresholded and interpreted as a curvature score. This score is small for perfectly straight

asparagus spears and increasingly large for bent ones. As an S-shaped piece is arguably perceived as bent even when the overall deviations from the center line are small, a second descriptor was computed as the ratio between the error of a linear fit and polynomial regression of degree three. Thresholding values and employing a voting scheme for the results for all three perspectives yields a rule to measure curvature (e.g. at least one of the three perspectives indicates curvature). However, it has again proven difficult to set thresholds appropriately to reliably capture the visual impression. Hence, another sparse representation was calculated by dividing the spears into several segments and fitting linear regression to each segment. A Multilayer Perceptron (MLP)) was trained on the resulting 18 angles per piece (see autoref{subsec:FeatureEngineering}).

Calculating a score for curvature is fast and efficient. While the respective approach is suitable to define curvature it does not necessarily meet up with the subjective perception of asparagus curvature. Just like histograms of palette images, curvature scores are the results of feature engineering: The use of extensive domain knowledge to filter relevant features (Zheng and Casari, 2018). They can serve as an input to a machine learning approach that maps this sparse representation to the target categories (see subsection 4.1.1).

### 3.2.6 Flower

The implementation of the flower detection function turned out to be difficult to realize. Several approaches have been tested, but none of them generated sufficiently good results. Two main notions were tried. The first approach uses the shape of the head as an indicator for a flower. The idea is that asparagus spears with a flowery head exhibit a less closed head shape. In other words, the head looks less round and has no smooth outline, but shows fringes. The second approach focuses on the structure within the head. Supposedly, asparagus with flowery heads exhibit more edges and lines in the head area. In both cases, it is challenging to find a way to discriminate between asparagus with and without flowery heads. One reason for that is the poor resolution of the camera that is installed in the sorting machine. With a pixel to millimeter ratio of around one to four, it is even difficult to detect flowers with the human eye. Likewise, the current software in the machine struggles greatly with the classification of this feature as well.

## 3.3 The hand-label app: A GUI for labeling asparagus

Providing a sufficiently large data set that contains information about target categories is one of the major non-algorithmic challenges in the application of machine learning for classification tasks (Al-Rawi and Karatzas, 2018). In many cases, however, the respective labels are missing. This is especially problematic if traditional supervised learning methods such as feed forward CNNs are employed: If the variance in the input images is high, a large number of samples is required (Russakovsky et al., 2015).

The options to reduce the variance and hence the need to attribute labels to a large number of samples are limited. We employed preprocessing to reduce the variance. In addition, strategies on the algorithmic domain such as the use of pretrained networks for transfer learning (Gupta, 2018), semi-supervised learning (see section 4.3) or manual feature engineering (see subsection 4.1.1) promise to reduce the labeling effort. This does not mean, however, that labeled samples are completely irrelevant. A sufficiently large number of labeled samples is still required for training

and more importantly performance evaluation. Without labels quality metrics such as accuracies, sensitivity and specificity cannot be calculated. Hence labels had to be manually attributed.



**FIGURE 3.2: The Labeling Dialog of the Hand-Label App** The depiction shows the main dialog used for labeling. In the left part you can see all three available images (perspectives) for the asparagus with the ID 40. The current question that targets at one of the features of interest is displayed in the area below the images. They are phrased such that the user can answer them with yes or no using the respective buttons or keyboard controls. On the right side you can see the results of the automatic feature extraction. The upper right panel shows the histogram of color hues.

Annotating labels manually requires plenty of effort. “Data set annotation and/or labeling is a difficult, confusing and time consuming task” (Al-Rawi and Karatzas, 2018). Human performance is often acknowledged as the baseline or “gold standard” that image classifiers are evaluated by. Hence, in many scenarios data is labeled by human coders such that machine learning algorithms can be applied. This holds especially for image classification.<sup>1</sup> In the present case some features were reliably measurable by means of computer vision (e.g. the width or the length). For features such as a flowering asparagus head or the evaluation whether or not a spear is affected by rust, this has proven to be difficult (see subsection 3.2.3). Considering the amount of data that could potentially be labeled, a custom interface is required that allows for time efficient attribution of labels.

For this reason, a custom application was built, compromising two user interfaces. A startup window that allows for a preview of asparagus images and the attributed labels (represented by `Ui_Asparator`) as well as the main labeling interface (`Ui_LabelDialog`). Using the labeling interface is possible only after the user selects the source folder of images and specifies or loads a file that contains the attributed labels. This ensures that the input and output file path are set. A dictionary that maps indices to images can be parsed from filenames and the minimum index is determined. As such, the label dialog and the respective controller class

<sup>1</sup>For example, the performance of GoogLeNet is compared to human level performance using the ImageNet data set (Russakovsky et al., 2015).

always resides in a valid initial state. For labeling, the user answers questions that are displayed alongside the images that depict each asparagus spear from three perspectives. This can be done using the respective buttons or the arrow keys. The result is saved and the next question will automatically appear upon answering. Automatic feature detection can be selected as an alternative for manual labeling for specific features. The result is displayed and saved to a file. This flexible approach was chosen as it was initially unclear and disputed in how far automatic feature extraction yields results that meet up with the individual perception. It also allowed to improve automatic feature extraction methods and to develop a direct intuition for the relation with the data. On top of that, it has proven to be useful for debugging automatic feature extraction methods that failed for some images.

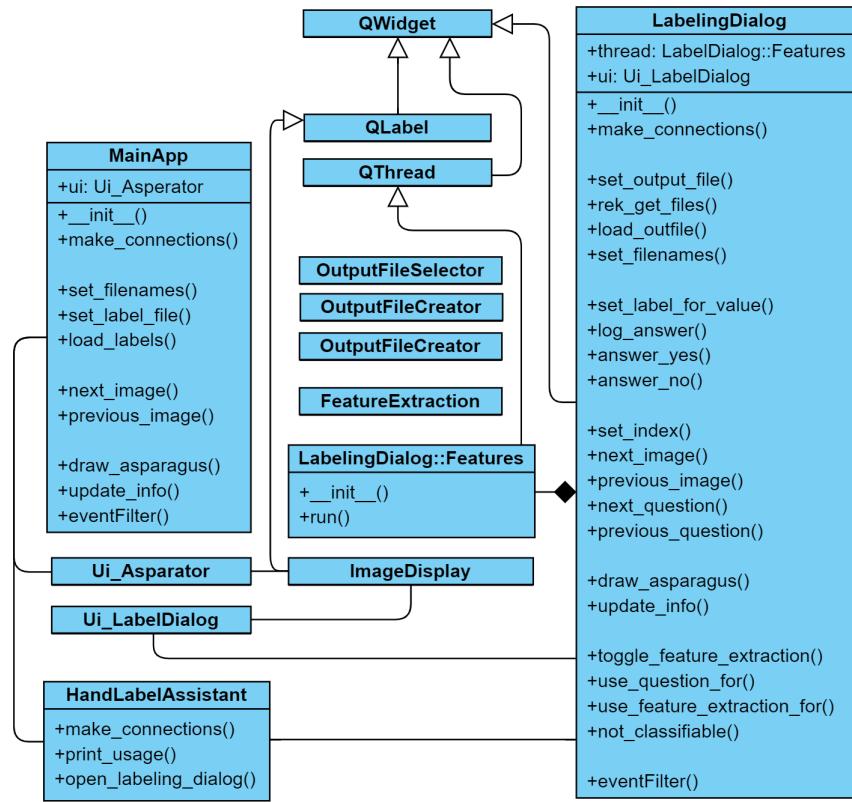
The development of the app was accompanied by three major challenges. First, handling a large data set of several hundred gigabytes that is accessible in a network drive. Second, changing requirements that resulted from group decision processes with respect to automatic feature extraction as well as from unforeseen necessities in (parallel) preprocessing. This made substantial changes of the initial architecture and the reimplementation of parts of the code necessary. The third challenge was the related question of the handling of internal states of the app. The latter may be further explained.

Internal states of the app are handled such that it is possible for the user to navigate into invalid states for which no images are available. The reason for this choice may be explained shortly. Note that preprocessing was done such that each asparagus spear has a unique identification number and a specifier for perspectives a, b and c in its filename. While generally the identification numbers are in a continuous range from zero to n, some indices are missing. As preprocessing jobs were scheduled to run in parallel and preprocessing failed for few corrupted files, it has proven almost inevitable to end up with few missing indices although a dictionary of input filename and output filename was passed to each grid job. In addition, the large amount of data did not allow to save all files in a single directory. In summary, this means that one could not simply iterate over asparagus identification numbers (represented by the state of a spin box in the user interface), determine the file path and display the related images. Instead, parsing filenames from a slow network drive is necessary which requires limiting the number of selected pieces. As GUI elements such as spin boxes and keyboard controls allow for setting an integer, and it was a requirement that this integer relates to the asparagus identification number, one ends up with the following situation. Either one prevents that the asparagus identification number is set or incremented freely to a value that does not exist, or one allows to navigate into an invalid state. It was decided that the last way would be the easiest solution.<sup>2</sup> Hence, all cascades of methods of the app including preprocessing functions that require the respective images as an input were adjusted such that they can handle this case.

The app is implemented using the PyQt5 framework while coarsely following the model, view controller principle. No separate database and the respective model class is used. The data model is implicitly parsed from the filenames and administered in attributes of the view. The labels are administered as a Pandas DataFrame and serialized as a csv-file. Upon state change (i.e. index increment) images are

---

<sup>2</sup>The earlier approach showed to have several implementation specific drawbacks. Note for example that upon entering multiple digits in an input field an event is triggered multiple times. Upon entering the value 10 for the asparagus identification number, one ends up with the value being set to 1 before being set to 10 where 1 relates to a potentially missing identification number. This means the user could not freely enter IDs when setting them to certain values is impossible.



**FIGURE 3.3: UML Diagram for the Hand-Label App** The depiction shows the class diagram for the hand-label app in UML (unified modeling language).

loaded from the network drive that was mounted on the level of the operating system. The views are designed using QtDesigner. Four manually coded classes are essential for the architecture of the app: The class HandLabelAssistant in which the PyQt5 app is instantiated, the controllers MainApp and LabelingDialog as well as Features which is member class of the latter, of type QThread and represents the API to the automatic feature extraction methods in FeatureExtraction. Ui\_Asparator, UiLabelDialog and classes for file dialogs represent the views. A class ImageDisplay is required to display images with the correct aspect ratio. The Figure 3.3 shows the UML class diagram alongside methods and attributes that are considered relevant to understand the architecture of the app.

Developing a custom app for the labeling process required substantial time resources. However, it was found that existing solutions did not meet the specific requirements. Our custom hand-label app allowed us to attribute labels for more than 10000 labels in a manageable amount of time. Details of the labeling process are described in the next section.

## 3.4 Manual labeling

In this section, the process of manually labeling the data with the help of the hand-label app is laid out. This includes the sorting criteria which allocate each spear to a single asparagus class, the practice of manually labeling the preprocessed data for its features, the outcome of the labeling process, and one approach to measure the agreement of the manual labeling.

The images were labeled by all members of the group. As none of the group participants were experts in asparagus labeling, a general guideline for the labeling had to be established. The guideline was written in accordance with the owner of the asparagus farm Gut Holsterfeld, Mr. Silvan Schulze-Weddige. He was consulted in all questions regarding the sorting of the asparagus.

General challenges in the manual sorting in front of a computer screen, including the respective image quality, and the variance in the agreement of the project members were expected from the start. As the task relies on the subjective view of individual humans, opinions about the affiliation to one asparagus class label can diverge. By consulting Mr. Schulze-Weddige on difficult decision-making for the team, it became clear that some examples are difficult to classify even for experts.

To tackle the issue and to have an overview of the general agreement of the sorting between group members, a measurement unit was researched and applied, namely the Kappa Agreement. The Kappa Agreement was used to assess the degree of accordance in sorting between the single members and monitor how the sorting agreement developed during the manual labeling process.

### 3.4.1 Sorting criteria

The class label of an asparagus spear is decided by several factors, ranging from its shape to its color. Put together, these single features give the class label for the spear. As it was decided to label the images for their features rather than finite class labels, the features were checked by the labelers. Even though the features chosen closely resemble the class labels defined by Gut Holsterfeld, they are different and should not be confused with one another. The images are displayed with the hand-label app and the features are divided as follows: fractured, hollow, flower, rusty body, rusty head, bent, violet, very thick, medium thick, thick, thin, very thin.

Further, images that could not be sorted thoroughly were sorted as not classifiable (e.g. when the spear is cut off the image). In the following text, the different criteria for manually labeling the data images for their corresponding features are described.

#### Fractured

An asparagus spear includes the feature fractured if it is broken, or if it does in any other way not fulfill the required length of 210 mm (see [Figure 3.4](#)).

For this feature, the parameter length was used, which is automatically calculated within the hand-label app.



FIGURE 3.4: Feature Fractured

### Hollow

The feature hollow indicates if the spear has a cavity inside. This might be expressed by a bulgy center and a line running vertically along the spear's body. Another, more distinct indicator is when the asparagus looks like two spears fused together, forming a single asparagus (see [Figure 3.5](#)). A hollow asparagus can be confused with a very thick asparagus.

The feature can be easily checked when you have physical access to the asparagus. If the asparagus is actually hollow, it will have a hole at its bottom noticeable when turning the spear around. Unfortunately, this cannot be done when only looking at the spears from the side.

The feature hollow sometimes occurs without showing a clear line or obvious bulge at its center. Therefore, there is a high risk of wrong classification.



**FIGURE 3.5: Feature Hollow**

### Flower

The asparagus is graded as flower when the bud shape of the head is more developed (for more information, see [subsection 3.2.6](#)).

When a bud is in full bloom, it is clearly visible (see [Figure 3.6](#)). However, it can be quite difficult to distinguish between a spear with clearly cut but closed petals and a spear that has just begun to develop a flower. It has been decided not to label an asparagus without a fairly clear flower, as it does not have the characteristic flower. With this decision, there is less correspondence error between hand sorters.



**FIGURE 3.6: Feature Flower**

### Rusty body

If a spear has rust on its body, it is visible as a dark brown color.

It often starts at the tips of becoming leaves or at the bottom part (see [Figure 3.7](#)). The color is not to be confused with bruises, pressure marks, or a slightly yellow complexion, which can occur in a ripe asparagus. The latter coloring is neglected. Therefore, rust is set to be present even when only the tip of a leaf shows a dark spot. Other brownish bruises are not classified as rust.

We split the feature rust into the sub-features rusty body and rusty head. In the case of rust being only on the body, it is removable by peeling.



**FIGURE 3.7: Feature Rusty Body**

### Rusty head

If there is a dark spot on or close to the head region recognizable as rust, it is captured by this feature (see [Figure 3.8](#)). The head part is usually distinct in shape and color from the body part of the asparagus.

In principle, the same guidelines which apply to rusty body can be transferred here, with an explicit focus on the top part (until around 1 cm below the collar) of the asparagus. As rust at the top part of the spear cannot be removed without damaging the head, it is more decisive for the later categorization into a price class, if the head has rust. Therefore, this separation is of importance concerning the quality of a spear.



**FIGURE 3.8: Feature Rusty Head**

### Bent

An asparagus is categorized as bent, if the shape of the asparagus is curved and not straight (see [Figure 3.9](#)). If it is only slightly curved but can otherwise be thought of as straight, that means fitting next to other straight spears without standing out, it is sorted as straight. If the spear looks close to the same on all three pictures regarding its shape, it might indicate that it is heavily bent and therefore cannot be turned on the machine's conveyor belt. The feature bent has a broader range of shapes where the asparagus is not obviously deformed but also not completely straight. An exception holds for S-shaped spears, which always count as bent.



**FIGURE 3.9: Feature Bent**

### Violet

The feature violet indicates whether an asparagus is of violet color.

The shift of color from white to violet occurs most often around the head region - either at the tip of the head or just below the collar of the head region (see [Figure 3.10](#)). Here, the spear will darken further after the sorting. Thus, even a slightly pink spear is sorted as being violet.



**FIGURE 3.10: Feature Violet**

### Thickness and length

The thickness and the length are features hardly recognizable by view. Therefore, both measures are extracted from the computer-vision based feature extractions (see subsection 3.2.2 and subsection 3.2.1). The division into different categories of thickness can be inferred by the overall thickness of the spear. The feature ‘very thick’ is attributed to asparagus that is more than 26 mm in width. Feature ‘thick’ corresponds to 20 - 26 mm, feature ‘medium thick’ to 18 - 20 mm, and feature ‘thin’ to 16 - 18 mm. Every asparagus with less than 16 mm in width is described with the feature ‘very thin’.



FIGURE 3.11: Feature Not Classifiable

### Not classifiable

Whenever the spear was unrecognizable, the head part of the spear was severed, two spears were present in one picture, the spear was cut off by the image, or other unusual circumstances occurred, it fell into the category of being unclassifiable (see Figure 3.11).

### 3.4.2 Sorting outcome

In this section, the process and the results of the manual sorting are presented. The labels of all manually sorted images were saved in a csv-file. As shown in Figure 3.12, the features are noted in the first column. The first entry is the image identification number. Every feature can be of value 0, value 1 or empty. Whenever a feature is present in an image, the value is set to 1. If the feature is absent, it is set to 0. For the images labeled as not classifiable, all feature values remain empty. The image path to every of the three images for one spear is saved in the label file in a separate column. After the labeling process, the individual csv-files with the labels are merged into one large combined\_new.csv file. The content of this file is later used for the classification of the data with the different approaches. It can be found in the study project’s Github repository.<sup>3</sup>

Feature	%	Feature	%
hollow	3.3%	fractured	3.5%
flower	12.9%	very thick	4.0%
rusty body	45.5%	thick	29.1%
rusty head	14.7%	medium thick	18.7%
bent	40.0%	thin	17.9%
violet	7.9%	very thin	30.3%
		not classifiable	2.1%

TABLE 3.1: Feature Representation in the Data Set

In this table, the representation of each feature in the manually labeled 13319 asparagus samples is reported in %.

<sup>3</sup>see <https://github.com/CogSciUOS/asparagus>

1	id;is_brunch;is_hollow;has_brunche;has_rust_head;has_rust_body;is_bended;is_violet;very_thick;thick;medium_thick;thin;very_thin;unclassified;au
2	0;0;0;1;0;0;0;1;0;1;0;0;;;;;1;0;;221;;14;0;241.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
3	1;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;238;;12;0;186.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
4	2;0;0;0;0;0;0;1;0;0;0;;;;;1;0;;219;;22;0;34.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
5	3;0;0;0;0;0;1;0;1;0;0;0;;;;;1;0;;237;;23;0;5.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
6	4;0;0;0;0;0;1;0;1;0;0;0;;;;;0;0;;236;;17;0;68.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
7	5;0;0;0;0;0;0;0;0;0;0;0;;;;;0;0;;238;;19;0;33.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
8	6;1;0;0;0;0;1;0;1;0;0;;;;;0;0;;240;;21;0;801.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
9	7;0;0;0;0;0;0;0;1;0;0;;;;;0;0;;241;;18;0;63.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
10	8;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;237;;21;0;11.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
11	9;0;0;0;0;0;0;0;0;0;0;;;;;0;0;;238;;17;0;4.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
12	10;0;0;0;0;0;0;0;0;0;0;;;;;0;0;;239;;14;0;7.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
13	11;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;238;;15;0;35.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
14	12;0;0;0;0;1;0;1;0;0;0;;;;;0;0;;239;;21;0;40.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
15	13;0;0;0;1;0;1;0;0;0;;;;;0;0;;241;;18;0;7.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
16	14;0;0;0;1;0;1;0;0;0;;;;;2;0;;236;;19;0;12.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
17	15;0;0;0;1;0;1;0;1;0;0;;;;;1;0;;234;;23;0;28.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
18	16;0;0;0;1;0;1;0;0;0;;;;;0;0;;238;;19;0;8.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
19	17;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;237;;11;0;49.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
20	18;0;0;0;1;0;1;0;0;0;;;;;0;0;;237;;16;0;25.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
21	19;0;0;0;0;0;0;0;0;0;0;;;;;0;0;;237;;14;0;22.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
22	20;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;239;;14;0;12.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
23	21;0;0;0;0;0;0;0;0;0;0;;;;;0;0;;235;;18;0;10.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
24	22;0;0;0;0;1;0;1;0;1;0;0;;;;;0;0;;236;;18;0;49.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
25	23;0;0;0;0;1;0;1;0;1;0;0;;;;;0;0;;237;;17;0;44.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
26	24;0;0;0;0;1;0;0;0;0;0;;;;;0;0;;237;;19;0;3.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
27	25;0;0;0;1;0;1;0;0;0;0;;;;;1;0;;238;;26;0;23.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
28	26;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;238;;16;0;24.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
29	27;0;0;0;0;0;0;1;0;0;0;;;;;0;0;;235;;14;0;171.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v
30	28;0;0;0;0;0;1;0;1;0;0;0;;;;;0;0;;237;;12;0;241.0;['/Users/Josefine/Documents/Cognitive Science MSc/new_gate/net/projects/scratch/summer/v

FIGURE 3.12: **Label File** The feature labels extracted by the manual sorting process were saved in a csv file. This image shows the beginning of combined\_new.csv in which all label files were later combined to one file.

The manual sorting lasted over the period of November 2019 to January 2020. A session usually consisted of 500 images, with an asparagus spear being viewed in three positions from the same perspective. A session of sorting 500 images took between two and four hours. One minor factor sometimes influencing the time spent for sorting were difficulties with the external ssl connection to the **IKW** storage, as all images are stored on the university servers. The average time spent for sorting one asparagus is around 27 - 48 seconds.

All in all, 13319 triples of images were labeled for their features. There is a large variance in the presence of the features in the data as can be seen in [Table 3.1](#). Of the acquired 13319 images, the feature most present in the data is rusty body with 45.5%, followed by the feature bent with 40%. Features whose representation is below 10% include the features violet (7.9%), hollow (3.3%), fractured (3.5%), and very thick (4%). The feature not classifiable shows least presence with 2.1%.

It emerges that many features are only sparsely present in the data. This poses an imbalance that is relevant for later classification tasks and the usage of the data set. Further, every member of the group participated in the sorting but not everybody sorted the same number of images. Due to this circumstance, the sorting bias of certain members is more present in the data than of others.

The manual sorting was stopped when the amount of classified data exceeded 13000 samples. Reason for this is that the less present features (hollow, violet, and flower) were in advance assumed to have a presence of around 8%. In order to have at least 1000 images with such a feature present, a threshold of 13000 labeled samples was set.

### 3.4.3 Agreement Measures

When different annotators label data, it is indispensable to verify the degree of agreement among raters. In order to judge how consistently the data is labeled, several statistical methods (inter-rater-reliability) can be applied.

For the current purpose, different agreement measures, all implemented by scikit learn, are used. The first is Cohen's Kappa. It is seen as a more robust measure than a simple agreement percentage, such as a measure of true positive and true negatives, which was traditionally used for those cases (Cohen, 1960). Cohen's Kappa is more robust, as the rate of agreement occurring by chance is included in the calculation. This method is applicable to compare the rating of two raters on a classification problem. The degree of agreement is always within -1.0 and 1.0 inclusive. The higher the Kappa value, the higher the agreement. Values around zero indicate no agreement and negative values indicate negative agreement, which can be interpreted as systematic disagreement. Values between 0.41 - 0.6 are seen as moderate agreement, 0.61 - 0.8 as substantial agreement, and everything above as almost perfect agreement. All scores below 0.4 are interpreted as unacceptable (McHugh, 2012).

Another statistical method used to measure agreement is the F1 score. The F1 score is used for the evaluation of binary classification. It relies on both, precision, as well as recall of a test's accuracy. A F1 score value lies between 0.0 and 1.0, the higher the F1 score, the higher the agreement.

Lastly, calculated the accuracy measure. For a normalized accuracy score, the values lie between 0.0 and 1.0, and the best performance is 1.0. This measure returns the fraction of correctly classified samples. It is a less robust measure than Cohen's Kappa score (McHugh, 2012).

### 3.4.4 Reliability

In order to evaluate the degree of agreement of our data, we made agreement measures at two points in time.<sup>4</sup> The first time, six different annotators labeled images to 13 class labels out of each asparagus group. We ensured that always two different annotators labeled the same set of images. The Kappa scores varied strongly between groups and features from -0.03 to 0.76, while the accuracy scores ranged from 0.49 to 1. It seems untypical to us, that the agreement scores are low, even though the raters gave the same label to many of the asparagus spears. This is an acknowledged problem (Powers, 2012; Sim and Wright, 2005; Feinstein and Cicchetti, 1990; Flight, 2018).

One reason for our results could be that we compared the agreement class-wise. The occurrence of 1s and 0s per class is therefore very unbalanced. For Kappa scores, if the distribution of 0s and 1s is not balanced, disagreement of the underrepresented value is punished more heavily.<sup>5</sup> Therefore, we decided to repeat our agreement measure feature-wise on non-labeled images, so that the annotators can not anticipate a specific group label. In order to better understand the reliability of our data, we additionally decided to look at the accuracy score and the F1 score. Beforehand, we labeled another 50 images all together, clarified classification boundaries again and discussed unclear images.

The second time, 50 images were labeled by four annotators. The agreements were measured annotator-pair-wise, and then averaged. The results in the Cohen's Kappa score vary between and within features, and between annotator pairs as well. The highest aggregated kappa score over all annotator pairs is reached for

<sup>4</sup>for our API documentation see

[https://asparagus.readthedocs.io/en/latest/api/measure\\_agreement.html](https://asparagus.readthedocs.io/en/latest/api/measure_agreement.html)

<sup>5</sup>see also

<https://stats.stackexchange.com/questions/47973/strange-values-of-cohens-kappa>

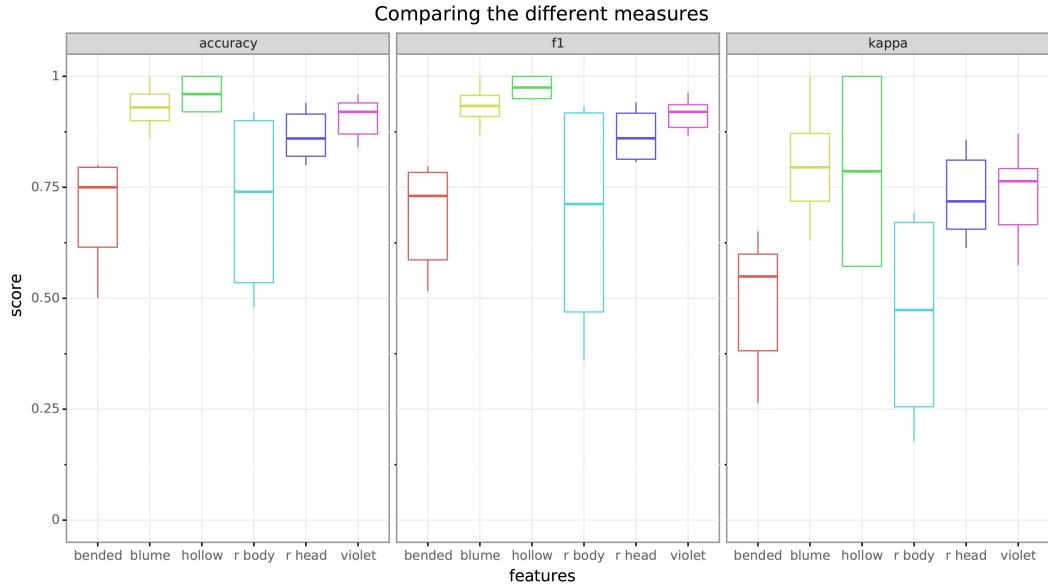


FIGURE 3.13: **Comparing Measures** The figure shows the agreement measures accuracy, F1 and Cohen’s Kappa, separately for each manually labelled feature. Shown are the box-plots, so the middle line indicates the median, the box indicates the IQR. All scores are aggregated scores over all annotator pairs.

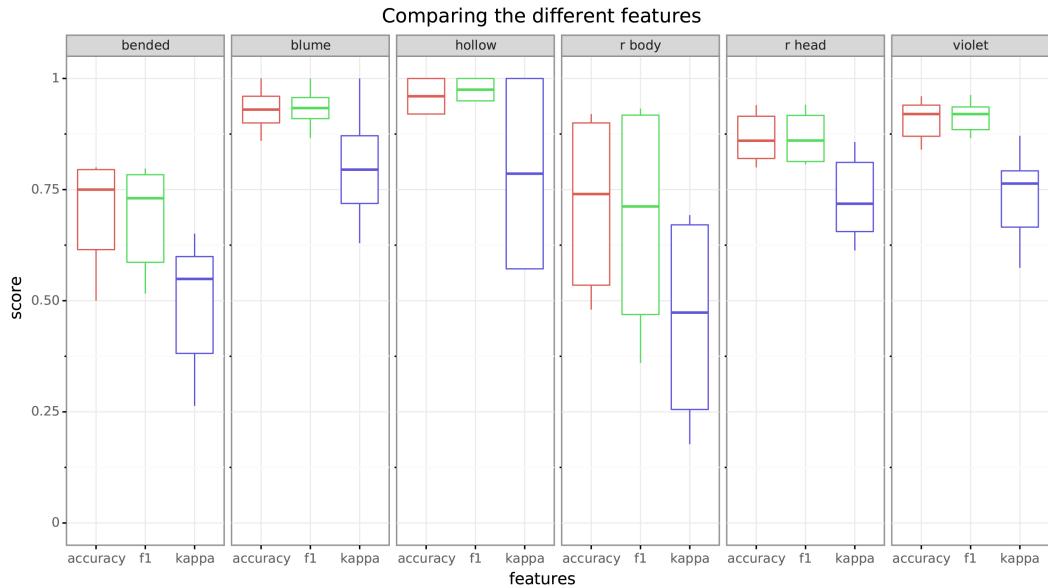


FIGURE 3.14: **Comparing Features** The figure shows each feature separately. For each feature, the corresponding accuracy, F1 and Cohen’s Kappa score is given. Shown are the box-plots, so the middle line indicates the median, the box indicates the IQR. All scores are aggregated scores over all annotator pairs.

the features flower (0.79) and hollow (0.79), then violet (0.76), rusty head (0.72), bent (0.55) and lastly for rusty body (0.47). For the features flower, rusty head and violet, the interquartile range (IQR) is quite small, whereas the IQR for hollow, bent and rusty body is much larger (see Figure 3.13 and Figure 3.14).

The agreement scores accuracy and F1 yield very similar results. Results are slightly better than the Kappa scores, in total and for each feature. The highest accuracy score is reached for the feature hollow (0.96), then flower (0.93), then violet (0.92), then rusty head (0.86), then bent (0.75) and then rusty body (0.74). The order is the same for the F1 scores. The median F1 scores lie between (0.71 and 0.97).

All in all, it can be said that the agreement scores indicate moderate up to substantial agreement. Regardless of the measurement, agreement is highest for the features flower, violet and rusty head.

### 3.5 The asparagus data set

In this subchapter, it is described how the images are processed to generate varying data sets. After giving a theoretical overview of the different possibilities of creating pipeble data sets for the tensorflow framework, a description of the in practice used methods is given. It is described which common methods are used to prepare and transform the data. We introduce the methods of simply using raw data, creating a numpy file, saving a Tenors or a TFRecord file and the tf.data API and the tf.data.data set API. These methods are compared and followed by a recommendation for further work with the data set.

One question arising in regards to modeling a data set was: What is the best way to provide a model with the image input and corresponding information? The most straight-forward variant is to build the model first, then read in the data, and load the samples one after the other. APIs such as OpenCV, or Imageio help with this in the project, the possibility is given by using Python. However, the runtime is slow and further processes, which are described below, are not efficient. Additionally, there may be an overflow and large amounts of data does not fit into the memory.

One of the used data sets is constructed with all the labeled images in a single numpy array, which can be stored and loaded. As the three perspectives of each asparagus spear are concatenated horizontally, they appear to lay next to each other in the resulting image. These concatenated images are then combined to the final file. The first out of four dimensions in this file depicts the number of labeled asparagus spears. The second and third dimension represent the height and the width of the images, respectively. Further, the fourth dimension represents the three RGB values. In addition, the images are downscaled to facilitate the training process and reduce memory consumption. Initially, each image is downsampled by a factor of six, that means every 6th pixel is used in the reduced image. This factor can be easily changed and a new data set can be created.

In the following, Tensorflow's own binary storage format TFRecord is introduced. This approach facilitates the mix and match of data sets and network architectures. The large amount of data that we have has a significant impact on our import pipeline and, therefore, on the total training time. The file format is optimized for images and text data. These are stored in tuples which always consist of file and label. The difference in reading time in our case is especially significant, because the data is stored in the network and not on a SSD hard disk on the PC. The serialized file format allows the data to be streamed through the network efficiently. Another advantage is that the file is transportable over several systems, regardless of the model one wants to train.

Two data sets are created in this file format. First, one with all files we preprocessed, the other one with the labelled data. The binary file includes images with background in png format and has a size of 225 GigaByte (GB). The TFRecord file of

the same size does not only take up less memory capacity, but can also be read more efficiently. The second data set with all the labels included is reduced in memory space as well.

Working with these “like liked lazy list” files it simplifies the next steps of transformations. With the tf.data API complex input pipelines from simple, reusable components are created, even for large data sets. The preferred pipeline for our asparagus project can apply complex transformations to the images and combine them into stacks for training, testing and validating in arbitrary ratios. A data set can be changed, e.g. by using different labels, just by transformations like mapping, repeating, batching, and many others. These dozens of transformations can be combined. A frequently used order is described in the following. The principle of these “like liked lazy list” can be found in most mainstream languages like C#’s LINQ or Java 8’s streams.

To summarize, our order for the different transformations is:

1. create the data set
2. shuffle (with big enough buffer size) 3, repeat
3. map with the actual work (preprocessing, augmentation...) using multiple parallel calls
4. batch
5. prefetch

Besides the described functional transformations of the input pipeline under tf.data, furthermore an iterator gives sequential access to the elements in the data set. The iterator stays at the current position and gives the possibility to call the next element as a tuple of tensors. Initializable iterators go through the data set several times. In addition, different parameters are passed to start the call. This is especially handy when searching for parameters.

In summary there are two advantages. On the one hand, it is possible to build a data set with different data sources. On the other hand, there are many functional transformations and an iterator with sequential access to the data.

Derek Murray recapitulates: “If you compare a tf.data pipeline to the equivalent queue-based pipeline, we use a similar structure of queues and threads, but importantly there are no Python queue runner threads on the critical path, and so the tf.data pipeline is not constrained by the Global Interpreter Lock and it scales to much higher throughputs.” (Murray, 2019).

Further we tried to add our data set to the TFDS.<sup>6</sup> TFDS enables all users to access the data set directly using the tensorflow api. For this, we need to publish the data. This is possible in terms of data integrity. However, the documentation was not elaborate at that time, so understanding it was too time consuming. Especially the large amount of data was a problem. Questions like: How do we deal with the fact that only a part of the images has labels? How should we pass the labels: each as a feature, in a list, or as several features? It would have been better faster and more helpful for the development of the networks, if we had continued to search and to integrate the TFRecord files with the tf.data api in our pipeline early.

---

<sup>6</sup>see [https://www.tensorflow.org/datasets/beam\\_datasets](https://www.tensorflow.org/datasets/beam_datasets)

## Chapter 4

# Classification

Given the structure of our data set, namely image data with a sub data set with corresponding class labels, and a sub data set with corresponding feature labels, different machine learning and computer vision methods were chosen, to tackle the problem of image classification.

Image classification refers to the method of identifying to which category an image belongs to, according to its visual information. Classification problems can be divided into three different types: binary, multiclass and multi-label. Moreover, there are different methods on how to approach image classification. Those can be divided into three main groups: supervised learning, semi-supervised learning and unsupervised learning.

Image classification refers to the method of identifying to which category an image belongs to, according to its visual information. Classification problems can be divided into three different types: binary, multiclass and multi-label. Whereas a binary classification only distinguishes between two different classes and therefore classifies an image into one of the two classes, a multiclass classification distinguishes between multiple exclusive classes. A multi-label classification also works for multiple classes. However in the latter, a single image can belong to none, one, several or all of the classes (Har-Peled, Roth, and Zimak, 2003).

Binary classification can be used to decide whether a certain feature is present at one certain asparagus piece, or not. Multiclass classification solves this problem as well, but is also applicable for data with more than two classes. As the class label results from a combination of the presence of certain features, and the absence of others, it is also reasonable to go for a multi-label classification approach.

During our group work, algorithms of the different classification types as well as of the different learning types were applied. In the long run, an integrated model was aimed at predicting all features of a single asparagus piece, from which the final class label can be inferred. However, as intermediate steps towards that goal, the focus was to optimize models on identifying the presence of the features described in subsection 3.4.1.

This chapter gives a general background of the different approaches chosen for our image classification problem, as well as a detailed overview of the concrete implementations of the models and the mechanisms of their hyperparameters.

## 4.1 Supervised learning

In machine learning, there are different approaches for an application to be trained on a set of data (Géron, 2019; Bishop, 2006). Depending on the level of supervision that the system receives during the training phase, the learning process is grouped into one of four major categories (Géron, 2019). One of these categories is supervised learning. For supervised learning approaches, the training data includes not

only the input but also its corresponding target labels. The objective is to find a mapping between object ( $x$ ) and label ( $t$ ) when a set of both ( $x, t$ ) is provided as training data to the application (Olivier, Bernhard, and Alexander, 2006). An advantage of supervised learning is that the problem is well defined and the model can be evaluated in respect to its performance on labeled data (Daumé III, 2012; Olivier, Bernhard, and Alexander, 2006). In other words, the labels are used as a direct basis for the model optimizing function during training.

Supervised learning spans over a large set of different methods, from decision trees and random forests, to **Support Vector Machine (SVM)** and **Artificial Neural Networks** (Caruana and Niculescu-Mizil, 2006; Géron, 2019).

A classical task for supervised learning systems is the classification of received data and mapping it to one of a finite number of categories (Bishop, 2006). The disadvantage of supervised learning is the effort of receiving enough labeled data. It can be challenging to obtain fully labeled data because labeling experts are needed to classify the data which is usually a time consuming and expensive task (Zhu, 2005; Figueroa et al., 2012).

In the following subchapters, different supervised learning methods were chosen to solve the classification task using the data that was manually labeled for features as described in chapter 3. In subsection 4.1.1, an approach using an **MLP** is described for feature classification. In the second subsection 4.1.4, a **CNN** is used to train solely on the head image data for the features flower and rusty head. The subsection 4.1.2 is concerned with labeling the input images with a **CNN** in a binary setup for their designated features. In subsection 4.1.3, a neural network is trained on the data to label it not only for one feature but all features at the same time. Finally, in subsection 4.1.5, a random forest approach is described to map the features of the image data to their class label.

### 4.1.1 Prediction based on feature engineering

Besides approaches that directly use images as an input one may use high level feature engineering to retrieve sparse representations that contain relevant information in a condensed form and apply classical machine learning classifiers such as **MLPs** to predict labels (Zheng and Casari, 2018). These classifiers are comparatively simple, fast to train and only few network hyperparameters have to be defined. One may argue that this is one of the major benefits as compared to networks of higher complexity (e.g. deep **CNNs**). As a consequence, finding suitable parameters for **MLPs** (i.e. the number of hidden layers and neurons per layer) is comparatively easy.<sup>1</sup>

An extensive search in hyperparameter space is practicable because of its small size and because training on sparse representations limits the number of neurons in the networks which results in fast training that allows for many experiments. In **CNNs** for deep learning suitable means are required to avoid the vanishing gradient problem. (Wang, 2019) Further, kernel sizes, strides, the number of kernels and other parameters must be defined. Therefore, designing shallow **MLPs** appears to

---

<sup>1</sup>The challenge of finding appropriate network parameters is well known in the deep learning community: “Designing and training a network using backpropagation requires making many seemingly arbitrary choices [...]. These choices can be critical, yet there is no foolproof recipe for deciding them because they are largely problem and data dependent” (LeCun et al., 2012). The requirement to specify hyperparameters is a disadvantage of neural networks (including **MLPs**) as compared to parameter free methods (Scikit-Learn, 2019). Due to the combinatorial explosion, the above mentioned challenge of finding suitable parameter settings is harder for more complex networks as more options must be considered.

be easier: There is simply less one could possibly do wrong. It deserves no further explanation that underfitting can potentially be due to an unsuitable network design or because predictions are impossible due to incongruencies or missing information in the sparse data set (LeCun et al., 2012). If the learning task is simple enough to be accomplished by **MLPs** (e.g. finding combinations of partial angles that correspond to the impression of curvature), one may speculate that underfitting can rather be explained by incongruencies or missing information in the labels than a result of issues in design- and training hyperparameters of the network. This classical machine learning approach that relies on feature engineering was applied to predict features based on color and partial angles of asparagus spears.

### Violet and rust prediction based on color histograms

The initial approach of measuring the feature violet that is based on distribution of sufficiently intense color hues in the violet range faces at least two drawbacks: First, it requires to define two thresholds. Second, the impression of a violet asparagus spear could potentially be affected by the combination of colors that are potentially outside the violet range or are too pale to be considered (see subsection 3.2.4). The same holds for the features rusty body and rusty head. Hence, in a second approach histograms were computed for foreground pixels after transforming the images to palette images with 256 color hues (see section 3.1). The resulting representation is a sparse descriptor that allows to predict color features using explicitly defined rules or trainable machine learning models.

	False positive	False negative	True positive	True negative	Sensitivity	Specificity
is_violet	0.04	0.03	0.05	0.88	0.62	0.96
has_rust_body	0.19	0.14	0.33	0.34	0.71	0.65

TABLE 4.1: **Color-Based Prediction** Performance of color histogram based predictions with a **MLP**.

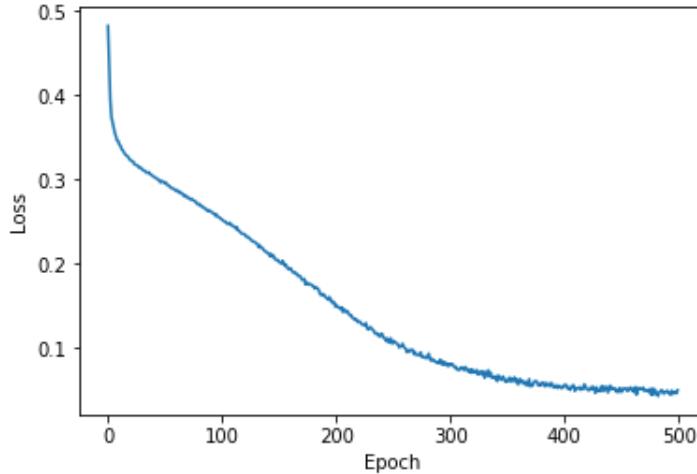


FIGURE 4.1: **Learning Curve for Palette-Color Histogram based Prediction** The depiction shows the loss per training episode for the **MLP** trained on histograms of palette images.

A simple **MLP** with four hidden layers and 128 neurons in each of them was trained on the resulting normalized histograms of palette colors (ReLU activation / sigmoid activation in the final layer). Hyperparameters were optimized and the network was trained for a total of 500 epochs as the learning curve indicated convergence at this point.

### Curvature prediction based on partial angles

Although the accuracies are far from perfect the results appear to be promising. The hit rate for curvature detection is high: 82% of all bent spears were identified as such. In comparison, this holds for 71% of the spears affected by rust and 62% of the violet spears. In contrast, almost all spears that are identified not to be violet are labeled accordingly (96% specificity), whereas the specificity for rust (65%) and curvature (67%) is lower.

	False positive	False negative	True positive	True negative	Sensitivity	Specificity
is_bended	0.19	0.07	0.34	0.4	0.82	0.67

TABLE 4.2: **Curvature Prediction** Performance of curvature prediction based on angles

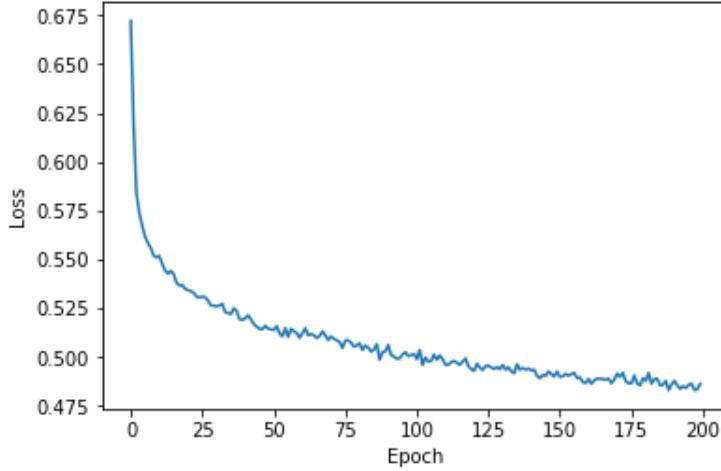


FIGURE 4.2: **Learning Curve For Angle Based Prediction** The depiction shows the loss per training episode for the **MLP** trained on partial angles of the centerline of asparagus spears.

The receiver operating characteristic reveals that the prediction is of better quality for violet and curvature prediction as compared to rust prediction which is reflected in a smaller area under the curve. Possibly this reflects that rather small brown spots were considered rust by some raters but not by others.

Considering the low agreement in labeling, higher values for the specificity and sensitivity of the classifier were not expected. A likely explanation is that the model generalizes deviating understandings or perceptual color-concepts we had when attributing labels and this affected the reliability of the data. Arguably only little information was discarded by computing the sparse representations that served as an input for training. Information about irregularities in the outline are not reflected in the sparse representation but potentially contribute to the perception of curvature. The same holds for the spatial distribution of colored pixels which might

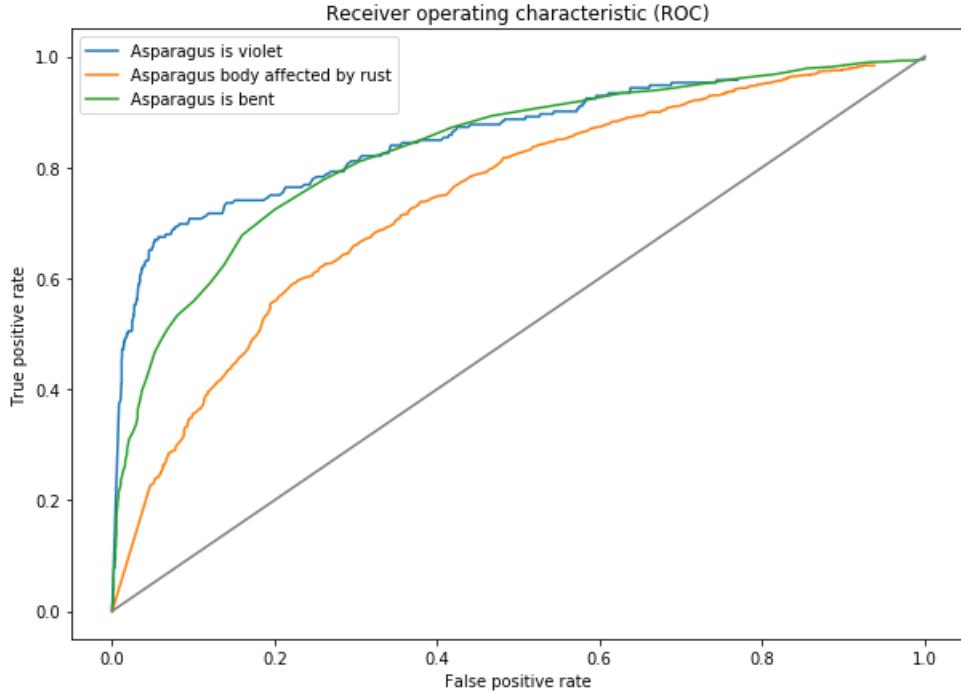


FIGURE 4.3: **ROC for MLPs Trained on High Level Features** The depiction shows the **Receiver Operating Characteristic (ROC)** for the classifiers that were trained on the features retrieved via feature engineering. This allows to compare the performance. A larger area under the **ROC** curve indicates better performance while a curve close to the diagonal line indicates poor results.

contain additional information regarding rust and violet-detection. However, the major criteria are captured. As **MLPs** have few hyperparameters, they are suitable for non-linear mappings and as the task of mapping high level features to human estimates appears to be rather simple, one could argue that there is little potential to improve the predictive quality using other techniques. By introducing a bias, the sensitivity of the classifier can be adjusted at the cost of more false positives. Here, introducing a bias means that the threshold that is used to convert the floating point outputs of a neural network to booleans that indicate whether a feature is present or not is set to values other than 0.5. The possibility of making the classifier more or less sensitive appears to be a good option to be implemented as a feature for customization by the user in asparagus sorting machines.

#### 4.1.2 Single-label classification

In the following chapter, a **Convolutional Neural Network** is described, which is used for single-label classification on features. The approach was tested on the 13319 hand-labeled data images. 13 models were created, each predicting one feature.<sup>2</sup>

A general model structure was needed as inspiration for the **CNN**. For example, the **Visual Geometry Group (VGG)** networks with varying depth seem to be a good choice for image classification as their **VGG16** won the ImageNet challenge of 2014 and is often implemented for image classification tasks (ul Hassan, 2018c; Simonyan and Zisserman, 2014). However, there are two major drawbacks on

<sup>2</sup>The 13 features to predict are fractured, hollow, flower, rusty head, rusty body, bent, violet, very thick, thick, medium thick, thin, very thin, and not classifiable.

using them, namely their depth and the amount of fully-connected nodes which makes them slow to train and in need of a lot of memory storage (ul Hassan, 2018c; Zhang et al., 2015). Part of these problems also arise with even deeper networks like ResNet (He et al., 2016b; ul Hassan, 2018b). Thus, AlexNet is chosen as a blueprint for the CNN because it is small in relation to other networks while still performing comparatively good (ul Hassan, 2018a; Krizhevsky, Sutskever, and Hinton, 2012b; Géron, 2019). As the variance in the data images is relatively small it is assumed that not as many layers are needed as employed in deeper networks like VGG (Géron, 2019).

The model architecture is roughly based on AlexNet but it is strongly simplified to the level of variability and complexity of the underlying data. The network comprises four hidden layers: a convolutional layer, followed by a pooling layer, a second convolutional layer, and a dense layer. The input is an array of multiple horizontally stacked images with no background removed and reduced by a factor of six. This input is trained on a set of binary labels containing information on whether the respective feature label applies to the current image. The output of the network gives a prediction on each entered image gated by a sigmoid function on a range between zero and one. The rounded integer values of this output give a prediction of the apparent feature label. For the training phase of the model, the Adam optimizer is used because of its general acceptance as the state of the art optimizer for backpropagation (Bushaev, 2018; Kingma and Ba, 2014). As a loss function, binary cross entropy is used as it promises good results for binary single-label classification tasks (Géron, 2019; Godoy, 2018; Dertat, 2017).

When training Artificial Neural Networks, it can be difficult to find clear guidelines on how to implement an architecture such that an optimal training performance is given (Heaton, 2015; Géron, 2019; Bettilyon, 2018). Hence, the idea was to start with the simplest form of a CNN and then gradually increase the complexity of the network. While AlexNet provides a good baseline for an image classification network, its architecture is still assumed to be unnecessarily complex for the given task. First, the architecture was reduced to the minimum number of layers and parameters needed for a CNN. Over the period of training optimization, various processing steps and hyperparameters were implemented and compared according to their performance. During this process, the data was split between 12000 samples for training data and 1319 samples of validation data in order to have a reasonable overview on the possible test performance and to directly check for overfitting. The data used as test data was randomly chosen from the whole data set and newly created for every trained model.

In the following, the course of the hyperparameters is explained.

The batch size was initiated comparatively low with 64 samples per batch but soon adjusted to a value of 512 samples. The larger batch size was implemented in order to guarantee the convergence of the training data, since smaller batch sizes result in jumping gradients, which are not able to converge into any minimum (Bengio, 2012).

Various learning rates were tested during the optimization process. The initial learning rate of 0.003 (which is the standard learning rate for Adam optimizers in Tensorflow (Kingma and Ba, 2014; Géron, 2019)) was soon found to be too large to guarantee convergence of the algorithm. Thus, the learning rate was decreased and found to be most effective in the range of  $1e^{-5}$  to  $1e^{-8}$ . A gradually decreasing learning rate was tested in order to make the training more effective (Bengio,

2012).<sup>3</sup> It did not give better results and thus a constant learning rate of  $1e^{-5}$  is implemented.

Starting with the smallest layer size possible, in the course of time more layers were added. For example, in order to detect the feature fractured one layer for the edge detection should be sufficient. For other, higher-level features, such as bent, more convolutional layers were expected to be more helpful (Géron, 2019). During the training process, it was settled to a model with the architecture described above.

A small number of kernels is used compared to AlexNet since for single-label classification fewer kernels are expected to be needed. The kernel size is picked to be eight  $5 \times 5$  kernels for the first convolutional layer and 16  $3 \times 3$  kernels for the second convolutional layer. The hidden dense layer consists of 32 neurons. These sizes are assumed to be sufficient, since increasing the number of kernels does not lead to any better results but to overfitting on the training data.

Both convolutional layers are built with batch normalization. The gradients were inspected visually and the results give no reason for assuming exploding or vanishing gradients (Pascanu, Mikolov, and Bengio, 2012).

A next step was to weight the loss function because of the largely unbalanced data (He and Garcia, 2009; Batista, Prati, and Monard, 2004). However, this did not lead to any changes. The model still tended to make an unbalanced prediction by classifying all values as negative samples. Another idea is to reduce the data set to make the number of images with the regarding feature present even to the number of images where it is absent. This would mean to throw away valuable data, which can otherwise provide information about negative cases to support the model in its training (Batista, Prati, and Monard, 2004). It was decided to keep all data images and instead balance the data by multiplying the minority of samples to match the number of contrary samples. As there was no feature positively exceeding a presence of 50% in the data, solely positive labels were oversampled. The balancing was only performed on the training data, while the test data was not changed.

To prevent overfitting of the negative data samples,  $L_2$  regularization was applied. To improve training performance, some kinds of data augmentation, like horizontal flipping or small changes in the angle (up to 5°), were tested but are not used in the final version (Brownlee, 2019).

Around 2700 to 5400 training steps were performed for the training, translating roughly to 120 epochs (with an exception for the feature fractured, which was trained for 60 epochs). Due to the balancing of the data, the number of training data, and therefore the ratio of training steps to epoch, varied when training the CNN on the different features.

The CNN was trained on 13 features separately, resulting in 13 trained models. For some features, there were no labels in the csv-file but they could be calculated from the parameters length and width of an asparagus. That is, the features very thick, thick, medium, thin, and very thin were all calculated from the thickness measured by the automatic feature extraction algorithm for width, according to the boundaries for each class label (for reference, see Table 1.1 and Figure 1.1 in section 1.3, and further subsection 3.4.1). For the feature fractured, the length was set to a threshold of 210 mm, with all asparagus of smaller length labeled as fractured. Additionally, all asparagus labeled as not classifiable was included into the feature fractured. The reason is that fractured asparagus without a head part was previously labeled as not classifiable (see the feature descriptions in subsection 3.4.1). The feature not classifiable was also trained on separately. Further, for all other

---

<sup>3</sup>In theory, the Adam optimizer already manages learning rate decay (Kingma and Ba, 2014)

	Sensitivity	Specificity	Validation Accuracy	Balanced Accuracy	Training Steps
fractured	0.8846	0.9976	99.32%	94.11%	2690
hollow	0.7308 (0.7692)	0.9821 (0.9831)	97.58% (97.77%)	85.65% (87.62%)	5390 (4270)
flower	0.5603 (0.6983)	0.8975 (0.8288)	85.96% (81.41%)	72.89% (76.36%)	4790 (2900)
rusty head	0.4311 (0.5210)	0.8695 (0.8106)	79.86% (76.38%)	65.03% (66.58%)	4670 (3320)
rusty body	0.6420 (0.6400)	0.7767 (0.8049)	71.15% (72.51%)	70.94% (72.25%)	2990 (2270)
bent	0.6541 (0.6753)	0.7533 (0.7500)	71.25% (71.93%)	70.37% (71.27%)	3230 (2470)
violet	0.4643 (0.5119)	0.9652 (0.9452)	92.45% (91.00%)	71.48 (72.86%)	5150 (3550)
very thick	0.9778	0.9939	99.32%	98.59%	5390
thick	0.9525	0.9688	96.41%	96.07%	3950
medium thick	0.8917	0.9249	91.87%	90.83%	4550
thin	0.9399	0.9280	93.13%	93.40%	4190
very thin	0.9733	0.9579	96.13%	96.56%	4310
not classifiable	0.5217	0.9961	98.79%	75.89%	5390

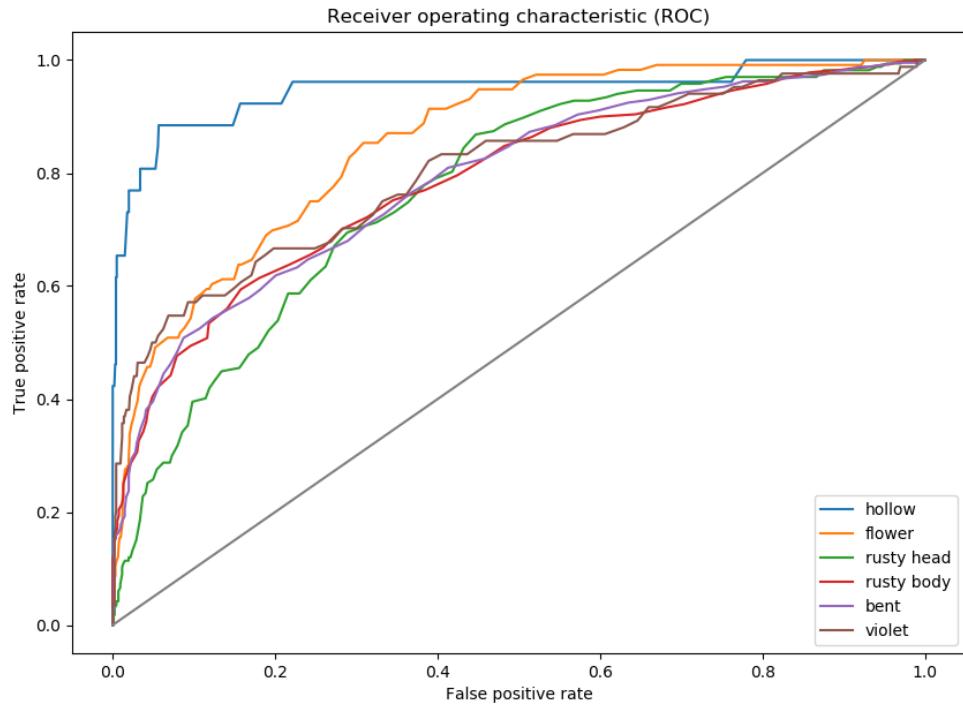
**TABLE 4.3: Single-Feature Label Classification Results** In this table, the sensitivity, specificity, validation accuracy, balanced accuracy, and the number of training steps are given for each feature model after training on 12000 original hand-labeled data images. The numbers in brackets indicate the best result for the model in relation to its balanced accuracy. For most features, training lasted for 120 epochs. However, the number of data samples (and thus training steps) varies between features because of the data balancing.

features, not classifiable samples were removed before training to prevent a bias in the occurrence of false positives.<sup>4</sup>

In [Table 4.3](#), the 13 features are listed on which the [CNN](#) architecture was trained 13 times. It further shows the results for the sensitivity, specificity, validation accuracy, and balanced accuracy of each feature after a certain number of training steps, indicated in the last column of the table. Sensitivity is a measure to assess the performance of the model in labelling positive samples correctly, while the specificity shows how correctly the model predicts negative samples. The validation accuracy is the accuracy of the validation set which is a representative subset of the entire data set. The balanced accuracy of a feature label is the mean of the sum of its sensitivity and specificity. It represents the accuracy of the feature label if positive and negative samples in the data set were evenly balanced. Additionally, the performance of six features is calculated in a [ROC](#) curve in [4.4](#). The features which indicate the thickness and length, as well as the feature not classifiable are excluded. In the following discussion, it is referred to the best result of a model (depicted in brackets in [Table 4.3](#)) and not the last result.

The results reveal that for every feature, the sum of sensitivity and specificity exceeds 1. This corresponds to a balanced accuracy over 50% for each feature, which is better than chance level. For all features, the balanced accuracy is above 65%. Best results are achieved for features that indicate the thickness of an asparagus. The feature very thick reaches the best results with 98% sensitivity, 99% specificity,

<sup>4</sup>Not classifiable asparagus was not sorted for the presence of any other features (see [subsection 3.4.1](#)). If a sample is sorted by a model, e.g. that is detecting the presence of feature bent, and the sample shows the sorting criteria, i.e. it is bent, the model will classify it as positive but its feature label will be negative. This will disturb the model and, thus, it was decided to exclude not classifiable samples with an exception for training on feature fractured and for feature not classifiable.

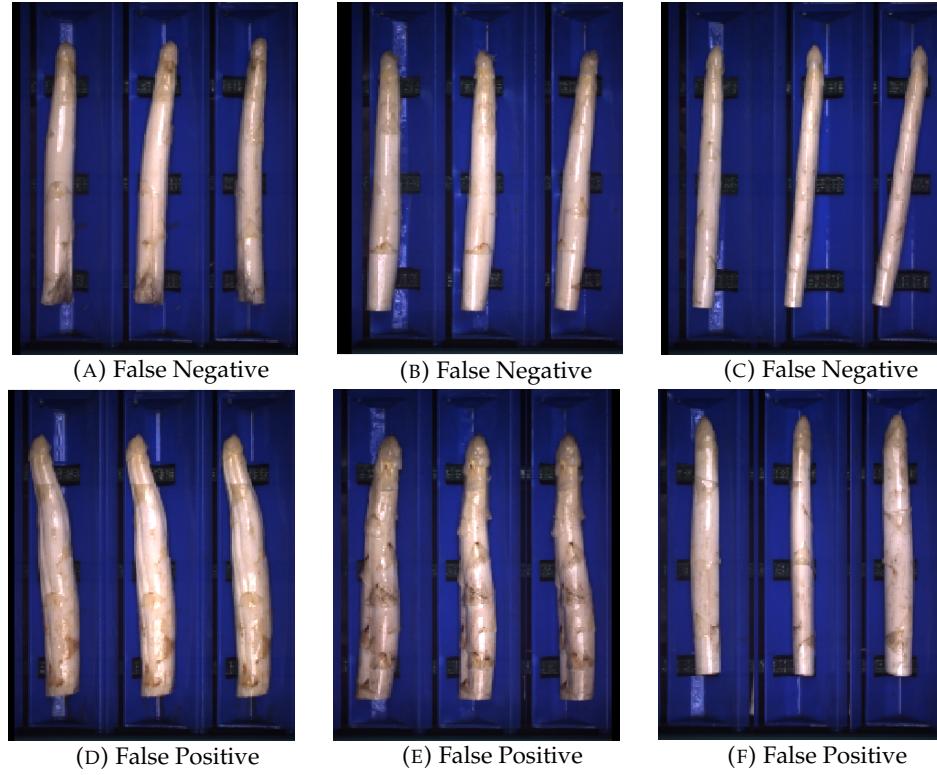


**FIGURE 4.4: Feature ROC Curve** The figure shows the ROC curves for the six features hollow, flower, rusty body, rusty head, bent, and violet. For the features rusty body and bent, the curve is comparatively smooth as their representation is more balanced. The curves are calculated from 1319 data samples.

and a balanced accuracy of 98.5%. Besides features for thickness, the best prediction is observed for the feature fractured, which relies on the parameter length. It reaches a sensitivity of 88%, specificity of 99.8%, and balanced accuracy of 94%. On average, for the features hollow, flower, rusty head, rusty body, bent, and violet a balanced accuracy above 72% is reached. Feature rusty head performed worst with 52% sensitivity, 81% specificity, and 67% balanced accuracy. In general, the specificity of all features is relatively high. Most features reach a specificity above 90%, except for the features flower (83% specificity), rusty head (81% specificity), rusty body (80% specificity), and bent (73% specificity). For all features, the sensitivity is above 50%. For most features, no form of overfitting was detected. Further, example images of wrong classification are shown for feature hollow in Figure 4.5 and for feature bent in Figure 4.6. Additional example images of false negative and false positive classification for the other feature classes can be found in the appendix in section B.1. For the feature fractured, training and test loss as well as training and test accuracy are also plotted as an exemplar in the appendix in Figure B.13.<sup>5</sup>

The results indicate, that the CNN architecture is able to learn every feature. Especially features relying on the parameters of length and width achieve a good performance, with both validation accuracy and balanced accuracy above 90%. The

<sup>5</sup>The log files for accuracy and loss of all features can be found at <https://github.com/CogSciUOS/asparagus>.



**FIGURE 4.5: Feature Hollow** Randomly chosen example images of false negatives and false positives of the feature hollow after 5400 training steps. In image (B), the presence of feature hollow is not obvious. Image (A) and (C) might have even been sorted wrongly. Of the false positives, image (E) might have also been sorted wrongly, as the vertical line can be seen on the lower part of the asparagus (seesubsection 3.4.1). The thickness of the asparagus might have been an indicator to the model that the feature is present. This can be seen in image (D), which shows a thick asparagus, and image (F), where the asparagus' thickness varies depending on the position.

prediction of features like hollow or flower was expected to be more difficult. However, the balanced accuracy for both is above 75%, with feature hollow even reaching a balanced accuracy of 88%. This shows that the model predicts them relatively well. Features that depend on color (like rusty body, rusty head, and violet) do not reach equivalent results. It should be further tested whether an increase in model depth might lead to better results for features depending on color or for features depending on a more complex shape (like the feature bent).

An inspection of false positive and false negative images at the end of the training process suggests that training performance might be influenced by mislabeled data to a certain extent.

The random images for feature hollow indicate that the model might use the thickness of an asparagus as an indicator. Some of the samples look like hollow asparagus that might have been labeled incorrectly (see [Figure 4.5](#)). For the feature bent, the randomly selected examples in [Figure 4.6](#) could reveal a labeling bias by the human labelers. When the feature is only slightly present, it seems to become a random choice whether the sample was labeled as positive or negative by the human labelers. If true, this can make it difficult for the machine to perform above a certain level of accuracy. It might also be that the difference between the samples is not prominent enough for the model. Again, these images are just randomly chosen examples and can mean that the models are simply not able to label them correctly.



**FIGURE 4.6: Feature Bent** Randomly chosen example images of false negatives and false positives of the feature bent after 3240 training steps. In all cases the presence of the feature bent is disputed. This becomes evident when comparing the false negative examples with the false positives. In the false positive examples, the asparagus is only slightly bent making the difference to the false negative samples not evident (see the criteria for feature bent in [subsection 3.4.1](#)). It gives the impression that the human labelers did not sort with a clear threshold regarding the feature bent.

The results suggest that correct labeling might sometimes be difficult for the model because of an inconsistency in the labeling behavior of the human labelers.

On a general note, the architecture of the model is very flexible. It can be applied to many tasks, i.e. predicting different features, without much preprocessing of the image data beforehand. Further, the model is quite small, which makes it fast and robust for practical applications. However, instead of having the same architecture for all features, more precise adjustment of each model to its feature is needed.

In conclusion, the results of the single-label **CNN** seem promising. Due to a very long debugging phase, there was not enough time to further test and improve the performance of the model. In turn, this means there is still a lot of fine-tuning and possible improvement. A restriction poses the labeling bias of the manually labeled images.

#### 4.1.3 Multi-label classification

Building on the standard single-label classification we were further interested in how well a model that predicts several feature labels at the same time performs. A multi-label classification model hereby gets an image as the input and learns to predict the presence or absence of the feature labels. For this model, we use a small **CNN** as described below and the features that we labeled by hand. Each of the six features (hollow, flower, rusty head, rusty body, bent and violet) is encoded by a

binary output in the target vector, indicating whether the asparagus exhibits the feature in question or not.

Multi-label classification is a useful tool for classification problems in which several classes are assigned to a single input. In contrast to a multiclass classification, where the model is supposed to predict the most likely class for an input, the multi-label classification makes a prediction for each class separately, determining whether the class is present in the image or not. While the different classes are mutually-exclusive in the multiclass classification, they can be related in the multi-label classification. Further, there is no limit on how many classes can be depicted in one image. It is possible that all or none of the classes are present.

Multi-label classification tasks can be thought of as consisting of different sub tasks. Therefore, the problem can be transformed to multiple binary classification tasks. In this transformation, a new model for each feature is trained, which are then combined to give a single output. That means that all features are independent of one another, because they are learned separately. This can be seen as one of the major drawbacks as it is not always clear whether features are related, but in many cases they are. Therefore, we decided to not only use single-class classification as described in Chapter 4.1.2 but to explore the possibilities of multi-label classification. A second approach to transform a multi-label classification task is to interpret each possible combination of features as one class. Hereby, the problem is redefined as a multiclass classification task. For a classification problem with six features that means there are  $2^6 = 64$  classes to be learned. The problems with this approach are, on the one hand, the exponentially increasing number of classes, and on the other hand, the sparsity of samples per class. In many cases, some of the classes are highly underrepresented or even empty. For that reason, we decided not to elaborate this approach further and implement a model for multi-label classification without transforming the task to a multiclass problem.

Inspiration for the model gave a blogpost (Franky, 2018) which aims to classify images of the MNIST fashion data set in the context of multi-label, rather than multiclass classification. The author altered the data set in such a way that each input image contains four randomly selected items from the MNIST fashion data set. The model then learns to predict which classes are present in the image. The target vector has ten values, one for each class, which are either 0 or 1 depending on whether that class can be found in the input image or not.

This model was chosen as inspiration for two main reasons. Firstly, it tackles a similar problem as ours and the number of classes is similar.

Secondly, the model uses a data set of similar size. Despite the rather small data set in comparison to many other image classification problems, good results with an accuracy of 95-96% were reached (Franky, 2018). This leads us to think, it might be a model with a good complexity for our problem too, as it is complex enough to model the underlying distribution, but not too complex for the medium-sized data set.

A classical **CNN** was chosen for the multi-label classification task. It consists of five blocks of convolution layers with max pooling layers followed by a global average pooling layer and a dense layer.

In contrast to multiclass classification models, where usually a softmax activation function is used in the last layer together with a categorical cross entropy loss, the multi-label classification model uses a sigmoid activation function and a binary cross entropy loss.

As the input of the model a concatenated image of the three perspectives of each spear is used in order to maximize the information the model gets. This yields input

Layer (type)	Output Shape	Param #
<hr/>		
conv2d_1 (Conv2D)	(None, 224, 182, 32)	896
<hr/>		
max_pooling2d_1 (MaxPooling2D)	(None, 112, 91, 32)	0
<hr/>		
conv2d_2 (Conv2D)	(None, 112, 91, 32)	9248
<hr/>		
max_pooling2d_2 (MaxPooling2D)	(None, 56, 45, 32)	0
<hr/>		
conv2d_3 (Conv2D)	(None, 56, 45, 32)	9248
<hr/>		
max_pooling2d_3 (MaxPooling2D)	(None, 28, 22, 32)	0
<hr/>		
conv2d_4 (Conv2D)	(None, 28, 22, 32)	9248
<hr/>		
max_pooling2d_4 (MaxPooling2D)	(None, 14, 11, 32)	0
<hr/>		
conv2d_5 (Conv2D)	(None, 14, 11, 32)	9248
<hr/>		
max_pooling2d_5 (MaxPooling2D)	(None, 7, 5, 32)	0
<hr/>		
global_average_pooling2d_1 (GlobalAveragePooling2D)	(None, 32)	0
<hr/>		
dense_1 (Dense)	(None, 6)	198
<hr/>		
Total params: 38,086		
Trainable params: 38,086		
Non-trainable params: 0		

---

TABLE 4.4: **Multi-Label Model Structure** The structure of the multi-label classification model is shown. This summary describes which layers are implemented, how the output changes in each layer and how many parameters are trained in each layer and in total.

images that look like the three asparagus spears are laying side by side. Further, the images are downscaled by a factor of six to facilitate training (see section 3.5). The output of the model is a vector of length six in which each position encodes one of the six hand-labeled features (hollow, flower, rusty head, rusty body, bent and violet). Each feature can either be present in the input or not, which leads to a 1 or 0 in the target vector, respectively.

Three loss functions are tested to improve the models performance. The first two losses are in-built functions from keras, namely binary cross entropy loss and hamming loss, which uses the fraction of the wrong labels to the total number of labels. Additionally, a custom loss function was implemented, that penalizes falsely classifying a feature as present more than falsely classifying one as absent. The motivation for this custom loss was the fact that the two labels 0 and 1 are highly unbalanced. As previously stated, there are noticeably more 0s than 1s in many classes. To be more precise, the model can reach an accuracy of 77% by labeling all features as 0. By penalizing this error more, we intended to counteract the unbalanced data set. But at the end, the binary cross entropy loss remains the one with the best results.

Further, it was tested whether regularization would improve the performance of the model on the validation data by preventing overfitting. For this, the model was trained by adding  $L_1$  or  $L_2$  regularization, respectively, to all five of the convolutional layers. Hereby, a kernel regularization was implemented with a value of 0.01.

$L_1$  and  $L_2$  regularization can both be interpreted as constraints to the optimization that have to be considered when minimizing the loss term. The main difference between the two is that  $L_1$  regularization reduces the coefficient of irrelevant features to zero, which means they are removed completely. Hence,  $L_1$  regularization allows for sparse models and can be seen as a selection mechanism for features. The inputs to this model are images that consist of a large number of pixels, and additionally a large portion of those pixels are black, because the background was removed. Therefore, it appears to be a good idea to reduce the number of features taken into account in the early layers.  $L_2$  regularization, on the contrary, does not set coefficients to zero, but punishes large coefficients more than smaller ones. This way the error is better distributed over the whole vector.

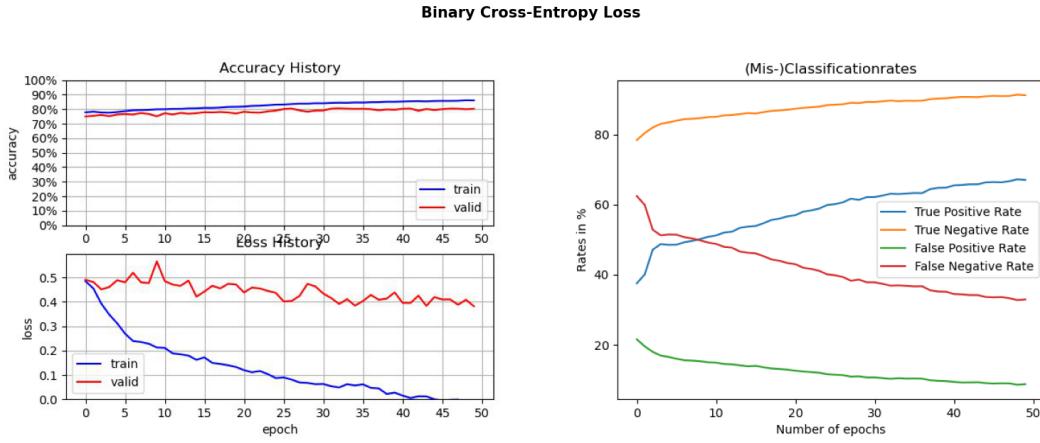


FIGURE 4.7: **Binary Cross-Entropy Loss** These graphs show the evaluation of the training with binary cross-entropy loss. The model was trained over 50 epochs and the accuracy and loss was measured. Further, the false/true positive/negative rates were determined.

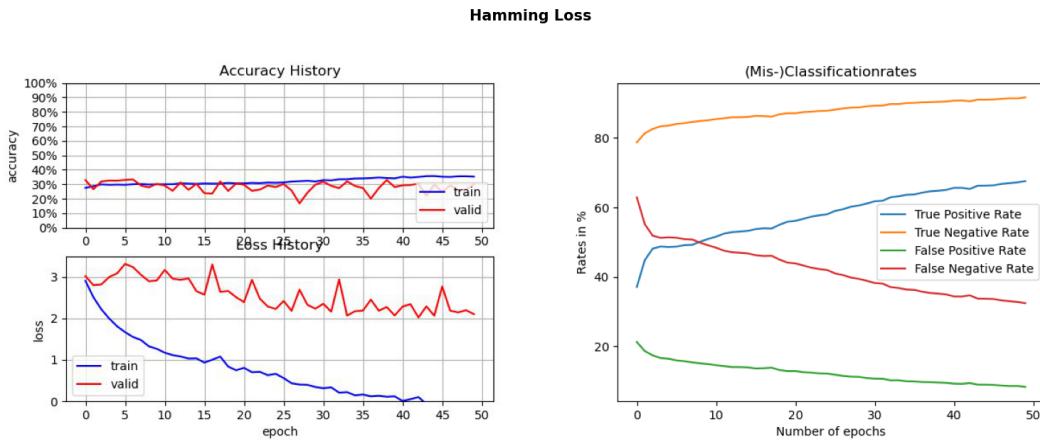


FIGURE 4.8: **Hamming Loss** These graphs show the evaluation of the training with hamming loss. The model was trained over 50 epochs and the accuracy and loss was measured. Further, the false/true positive/negative rates were determined.

As shown in Figure 4.7, Figure 4.8 and Figure 4.9, all the different approaches explained above show a similar behavior in accuracy and loss values. The training

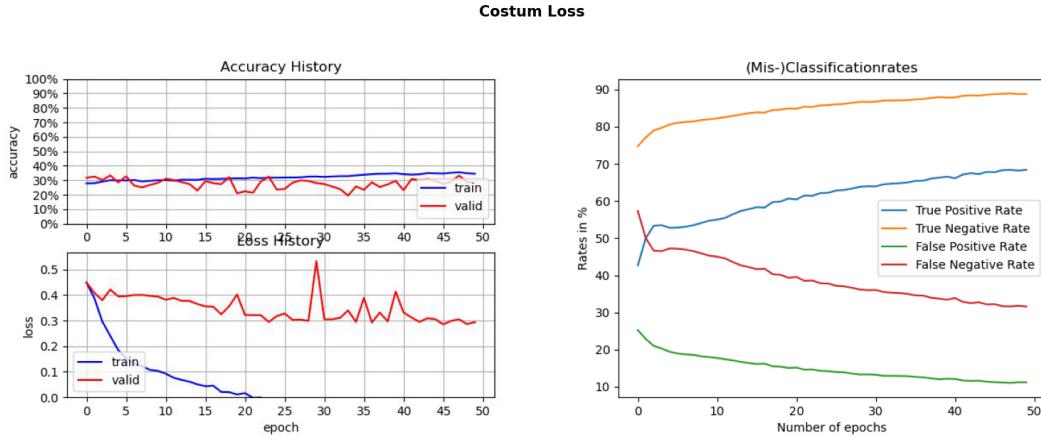


FIGURE 4.9: **Costum Loss** These graphs show the evaluation of the training with costum loss that punishes falsely classified ones more than falsely classified 0s. The model was trained over 50 epochs and the accuracy and loss was measured. Further, the false/true positive/negative rates were determined.

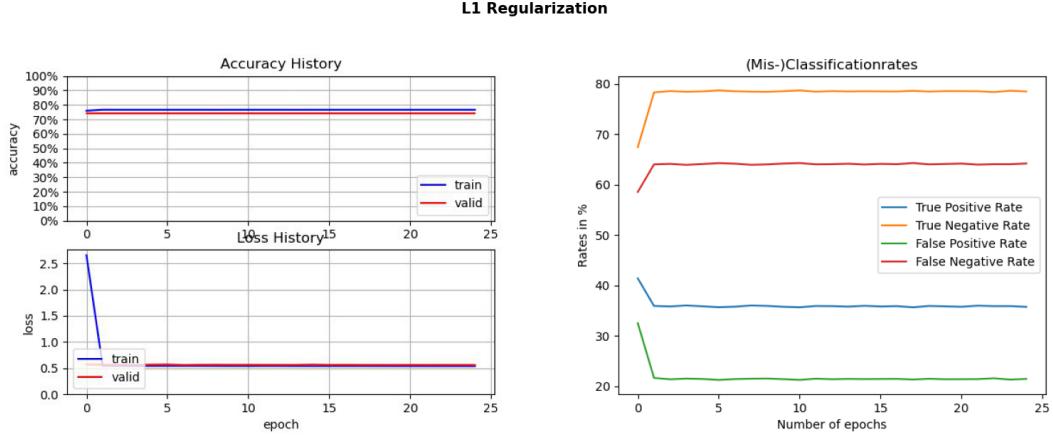


FIGURE 4.10:  **$L_1$  Regularization** These graphs show the evaluation of the training with  $L_1$  regularization. As the loss function the binary cross-entropy loss was used. The model was trained over 25 epochs and the accuracy and loss was measured. Further, the false/true positive/negative rates were determined.

and validation accuracy increase slowly but steadily with the training accuracy always being a little higher than the validation accuracy. The training loss decreases rapidly, while the validation loss only decreases very little and shows random fluctuations. This can be an indicator for overfitting. Usually,  $L_1$  and  $L_2$  regularization are used to prevent overfitting, but in our case it did not improve the results, as shown in Figure 4.10 and Figure 4.11.

When looking at the sensitivity and specificity, it becomes apparent that both increase during the training process, while the false negative and false positive rates decrease with the same slope. The false positive and false negative rates are mirror images to the true positive and true negative rates with the mirroring axis at the 50% mark. It can be observed that the rates change rapidly in the first two to four epochs, after which the change progresses slowly in the same direction with no greater disturbances. The model trained with the  $L_2$  loss is an exception as it does

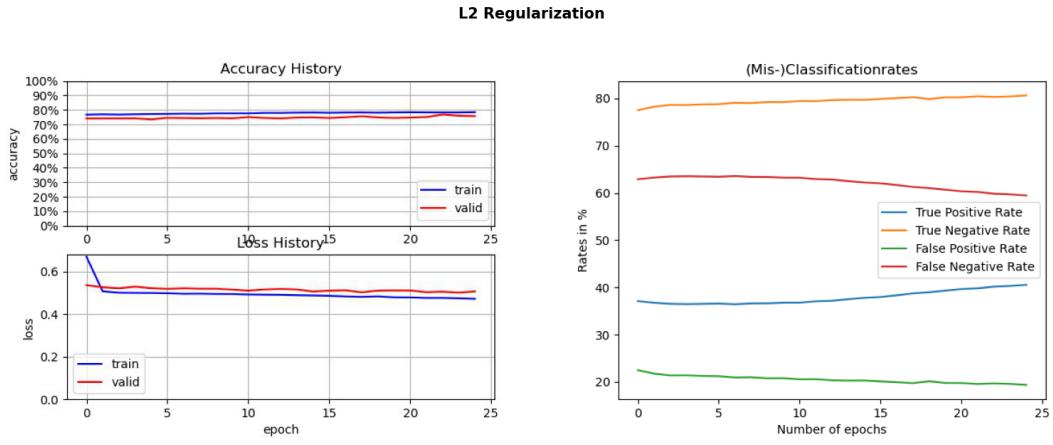


FIGURE 4.11:  **$L_2$  Regularization** These graphs show the evaluation of the training with  $L_2$  regularization. As the loss function the binary cross-entropy loss was used. The model was trained over 25 epochs and the accuracy and loss was measured. Further, the false/true positive/negative rates were determined.

not show these large changes in either of the rates.

When comparing the three different loss functions, it is noticeable that the binary cross entropy loss has significantly larger accuracy values than the hamming loss and the costum loss. The values of the binary cross entropy loss start at 75%, while they start at around 30% for the other two loss functions. The behavior of the curves and the (mis-)classification rates, however, are very similar in all three approaches. The specificities start off very high with values around 78.46% for the binary cross entropy loss, 78.73% for the hamming loss and 74.74% for the costum loss, and increase further during the training. The highest values are reached with the binary cross entropy loss (91.41%), closely followed by the hamming loss (91.37%) and the costum loss (88.75%). The sensitivity values start off lower, at around 37% to 42%, and increase rapidly in the first few epochs, after which the rates proceed to increase but with a narrower slope. They reach values of up to 67.27% with the binary cross entropy loss, 68.17% with the hamming loss and 68.17% with the costum loss. As stated above, the false negative and false positive rates show the same slope but in the opposite direction.

The accuracy values of the models that are trained with  $L_1$  or  $L_2$  regularization, respectively, do not change over the epochs. The same holds for the validation loss. The training loss decreases in the first few epochs and remains stable thereafter. While the (mis-)classification rates of the model trained with  $L_1$  regularization behave similarly to the ones trained with no regularization, the rates of the model trained with  $L_2$  regularization show a smaller increase and lack the fast change in the first epochs.

The slopes of all curves indicate that the model is learning, because they are increasing in the case of the accuracy, sensitivity and specificity and decreasing in the case of the loss, false positive and false negative rates until the end of training. Hence, one might think that a longer training period will lead to better results. But the training loss decreases very rapidly while the validation loss does not. This suggests overfitting of the model, a problem which gets worse when increasing the training steps. Therefore, a longer training period most likely will not increase performance unless overfitting is prevented. As shown in the result section neither  $L_1$  nor  $L_2$  regularization alone were able to prevent overfitting.

Another common practice that can be tested is the drop-out, in which a certain amount of nodes are left out in different backpropagation steps. This way the model learns to not rely on a small number of nodes but distribute the information between all nodes available. Hence, the coefficients remain smaller. Another method to prevent overfitting is to reduce the model's complexity. A model with fewer parameters to train, is less prone to overfitting. A fitting degree of complexity should be found to model the data sufficiently good without losing the possibility of generalization.

Accuracy alone might not be a good indicator to evaluate a multi-label model (Gibaja and Ventura, 2015). As it highly depends on the loss function, it may have misleading results. This can be seen in the comparison between the three different loss functions. Although the sensitivity and specificity show similar values, the accuracy values suggest that the binary cross entropy loss outperforms the other two loss functions by far. The accuracy of the model trained with the binary cross entropy loss has an accuracy more than twice as high. But when looking at the slope of the curve, it appears that the model with the binary cross entropy loss does not perform better than the other two models. All three have an increase of accuracy of roughly 10% and a similar sensitivity and specificity. This indicates that the slope of the accuracy function can be considered to evaluate the training process of the model, but the real values should be interpreted with caution.

One thing that comes to attention when looking at the (mis-)classification rates is that the sensitivities are a lot lower than the specificities. A reason for that might be that the positive values are more difficult to learn because they are only sparsely present. The model might have learned that, if the allocation to an output class is not clear, a 0 is the more likely guess.

The  $L_1$  and  $L_2$  regularization both seem to prevent the model from learning all together instead of only preventing overfitting. A reason for that might be that the regularization factor is too high. More experiments should be conducted with varying values to test this hypothesis.

In summary, the model improves its sensitivity and specificity, but it seems like it does so by overfitting the training data. Therefore, the next step should be to prevent the model from overfitting. Afterwards, it should be tested whether additional changes can improve the performance of the model further.

#### 4.1.4 A dedicated network for head-related features

As mentioned before, some features relate to the asparagus heads only. Hence, it was assumed that classification is easiest when training a CNN on depictions of the respective region of the asparagus. Therefore, a data set that consists only of images of the head area is used. Images of all three perspectives are appended horizontally such that each sample contains the information from all available viewpoints. This is especially important as rust affected spots are sometimes only visible from some angles. The Figure 4.12 shows one sample of rust affected asparagus heads.

A simple feedforward CNN was trained on the images. The features flower and rusty head are chosen as target categories. Hence, the model is another example for multi-label classification. The network comprises the input layer, three convolutional layers with kernel size two, a fully connected layer with 128 neurons as well as the output layer. For the final layer, the sigmoid activation function is applied while the hidden layers have ReLU activations. A dropout layer was added to avoid overfitting. The network was trained using mean squared error MSE as an error function. The development of loss in the learning curve indicates convergence after 40 epochs.



FIGURE 4.12: **Training Sample** The depiction shows a sample for the preprocessed and vertically aligned asparagus heads. Mind that images from all perspectives were stacked horizontally and used as a single input for the **CNN**.

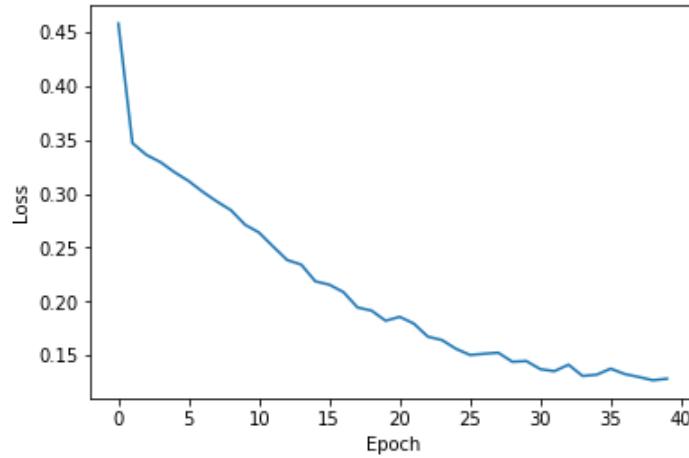


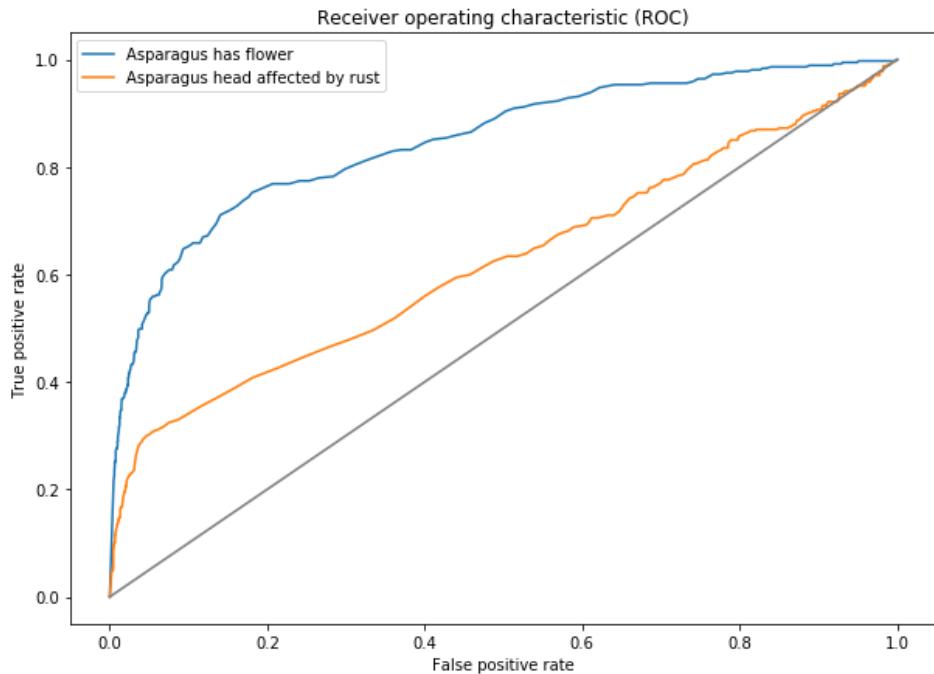
FIGURE 4.13: **Learning Curve for the Head Features CNN** The depiction shows the loss per training episode for the **CNN** trained on asparagus heads.

The results for both features show to be highly specific. In contrast, the sensitivity is rather low. Only 55% of the asparagus spears labeled as flower are identified as such whereas the true positive rate is only 19% for rusty head. Given the low labeling agreement for these criteria (see subsection 3.4.4), these mediocre results are not surprising.

The **ROC** curve indicates how the classifiers respond to the introduction of a bias and shows the overall prediction quality. The area under the curve is small for the feature rusty head. Beside incongruencies in the labels this is possibly due to the choice of the size of the head region. It might be the case that brown spots in regions other than the cropped head were considered as an indicator for a rusty head when attributing labels. Improvements by increasing the cropped head region appear to be possible.

	False positive	False negative	True positive	True negative	Sensitivity	Specificity
has_blueme	0.04	0.06	0.08	0.82	0.55	0.95
has_rost_head	0.02	0.13	0.03	0.83	0.19	0.98

TABLE 4.5: **Performance** Performance of the **CNN** trained on asparagus heads.



**FIGURE 4.14: ROC for the Head Features CNN** The depiction shows the **Receiver Operating Characteristic (ROC)** for the predictions of the **CNN** trained on asparagus heads. It allows to compare the performance. A larger area under the **ROC** curve indicates better performance while a curve close to the diagonal line indicates poor results.

#### 4.1.5 From features to labels

Approximately 200 asparagus spears per class label are pre-sorted and served as ground truth mappings between input images and output class labels. These images were manually annotated with features (see [section 3.4](#)). This allows to divide the classification process into two steps: In the first step we predict feature values from images and in a second step we predict class labels from those feature values.

This chapter deals with the second step: Using supervised learning to predict 13 class labels based on manually labeled features. These 13 class labels refer to the classes at the asparagus farm Gut Holsterfeld. We built a unified interface to load different models, train them and analyze their predictions. It provides compatibility for scikit-learn as well as for keras models. To explore the data and visualize the predicted results, we built a streamlit app.<sup>6</sup> This has the advantage that the user can easily load and train a model, select an asparagus spear, see the corresponding pictures and the predicted class label (see [Figure 4.15](#)).

The user can also inspect the distribution of the selected data ([Figure 4.17](#)) and the absolute and relative number of correctly and incorrectly classified asparagus spears in a confusion matrix ([Figure 4.16](#)).

<sup>6</sup>inside code/pipeline in the GitHub repository

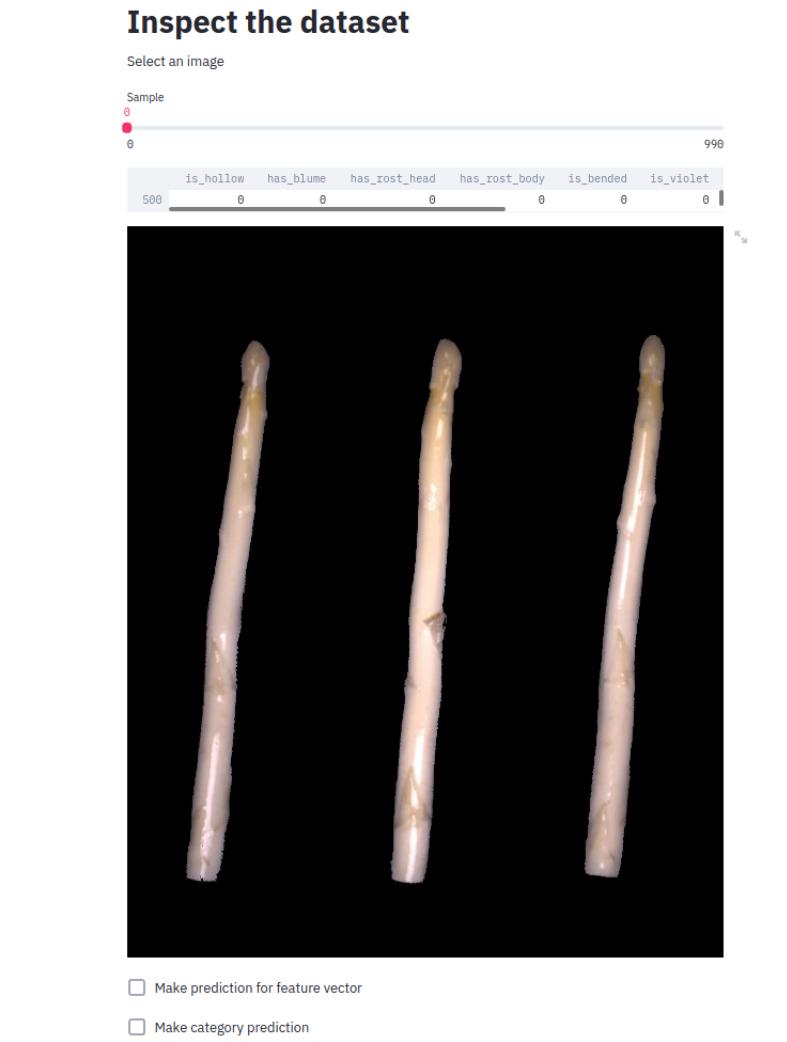


FIGURE 4.15: **Screenshot of Streamlit App** Screenshot of the streamlit app showing the three images of one asparagus with the corresponding labeled feature.

The precision and recall are measures on how “useful and complete” the results are (Wikipedia contributors, 2020). Precision is the ratio between true positives and the set of all true and false positives, recall the ratio between the true positives and the set of all true positives and false negatives. The confusion matrix (Figure 4.16) gives us insight about which kind of errors the models make. We can observe that classes such as 1A Anna, Dicke, Hohle, Rost, Köpfe, and Suppe can be recalled well (relative recall  $\geq 0.8$ ), while other classes, such as 2A and 2B have much lower recall ratings (relative recall  $\leq 0.6$ ). That means that 2A and 2B are the most commonly mislabeled classes (see Table 4.6).

Two exemplary models were implemented and tested to predict the classes from the features. The first one is a random forest (Breiman, 2001) with 100 trees. Random forests are a popular machine learning approach, because they reach good results in both regression and classification tasks. Further, they are based on decision trees which offer an intuitive interpretation. Although decision trees themselves are powerful tools, they are prone to overfitting. Random forests aim to avoid this problem, while still using the advantages of decision trees, namely that they are flexible, easy to use and fast to train. They do so by training several decision trees

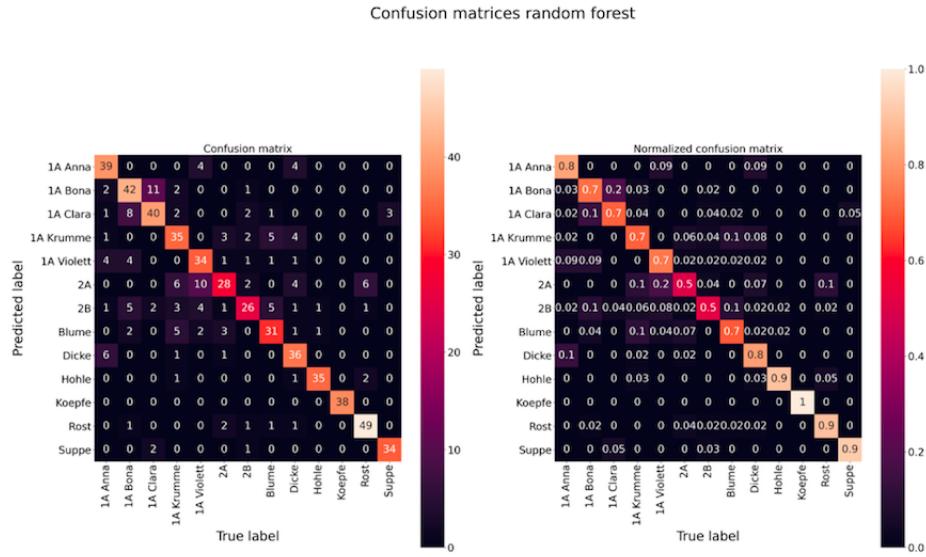


FIGURE 4.16: **Random Forest Classifier Confusion Matrices** Confusion matrices showing the absolute and relative number of true positives of the random forest model.

on different random parts of the data, after which a majority vote decides on the final output or class. An additional trick is to use a random selection of features for each tree. This reduces the influence of highly correlated features and thereby makes the random forest more robust.

The second model to predict classes from features is an **MLP** with six fully connected layers. It is only implemented to show how to integrate other models, but achieves a similar score as the random forest classifier when trained for 500 epochs (score of 0.76). However, it takes longer to train than the random forest classifier.

The score of 0.76 on the validation set is the accuracy of the random forest classifier, that means that 76% of the data is predicted correctly — random guessing would yield an accuracy of 0.08 for uniformly distributed classes, thus the result is satisfying.

Although the results point into the right direction, it would be interesting to see

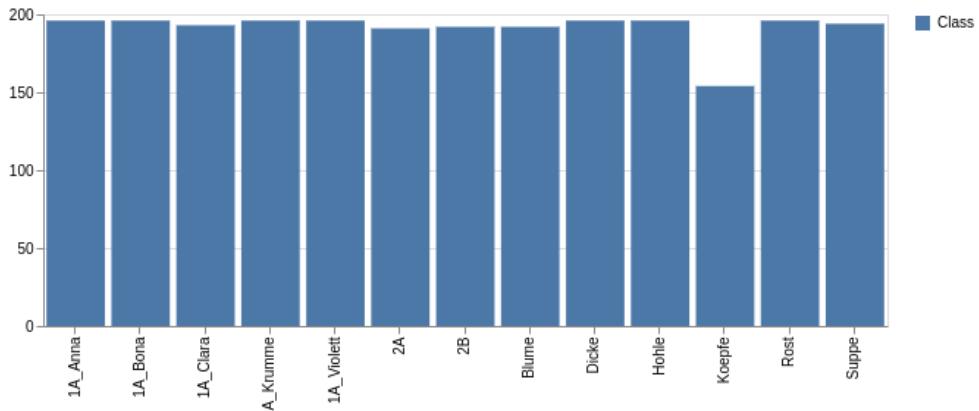


FIGURE 4.17: **Distribution of Class Labels** Absolute number of the asparagus spears for which a ground truth class label is available.

	precision	recall	f1-score	support
1A Anna	0.7360	0.8936	0.8077	47
1A Bona	0.7000	0.7241	0.7119	58
1A Clara	0.7368	0.7368	0.7368	57
1A Krumme	0.6667	0.6400	0.6531	50
1A Violet	0.6481	0.7609	0.7000	46
2A	0.6842	0.4643	0.5532	56
2B	0.7179	0.5600	0.6292	50
Blume	0.6939	0.7556	0.7234	45
Dicke	0.7000	0.7955	0.7447	44
Hohle	0.9459	0.8974	0.9211	39
Koepfe	1	1	1	38
Rost	0.8305	0.8909	0.8596	55
Suppe	0.9444	0.9189	0.9315	37
accuracy	0.7588	0.7588	0.7588	0.7588
avg	0.7696	0.7722	0.7671	622
weighted avg	0.7591	0.7588	0.7549	622

TABLE 4.6: **Classification Report of the Random Forest Classifier** The classification report shows key metrics for the trained random forest classifier.

how the selection of features and the agreement on the values of the labeled features influence the performance of the described classifiers. Further work is required to implement the decision process described by the local farmer (see [section 1.3](#)). One could test if that decision tree based on expert knowledge can outperform the random forest trained on the training samples.

## 4.2 Unsupervised learning

Unsupervised learning are all kinds of machine learning algorithms which deal with unlabeled data. More specifically, they work without a known goal, a reward system or prior training, and are usually used to find structure within data. Dimension reduction algorithms and clustering algorithms have been identified as the two main classes of unsupervised machine learning algorithms which are used in image categorization (Olaode, Naghdy, and Todd, 2014).

Multivariate data sets are generally high dimensional. However, it is common that some parts of that variable space are more filled with data points than others. A large part of the high dimensional variable space is not used. In order to recognize a structure or pattern in the data, it is necessary to reduce the number of dimensions. For this, both linear- as well as non-linear approaches can be applied. Linear unsupervised learning methods for which also descriptive statistics can be acquired are e.g. **Principal Component Analysis (PCA)**, Non-negative matrix factorization, and Independent component analysis (Olaode, Naghdy, and Todd, 2014). Some examples for non-linear approaches are Kernel **PCA** (Olivier, Bernhard, and Alexander, 2006), Isometric Feature Mapping, Local Linear Embedding, and Local Multi-Dimensional Scaling. We employed the linear dimension reduction algorithm **PCA** as well as autoencoders for nonlinear dimension reduction.

### 4.2.1 Principal Component Analysis

**Principal Component Analysis** was chosen, because it is one of the standard unsupervised machine learning methods. Moreover, it is a linear, non-parametric method and a widespread application to extract relevant information from high dimensional data sets. The goal is to reduce the complexity of the data, by only a minor loss of information (Shlens, 2014). Besides being a dimension reduction algorithm, **PCA** can also be useful to visualize or compress the data, filter noise or extract features.

Our initial aim was to reduce the dimension of our data set for further models. As **PCA** was applied at the beginning of the data inspection, we also had the aim to visualize our data in a three-dimensional space, in order to get a better understanding of the data distribution. The images that comprise our original data set have a high quality, and instead of only reducing the pixel size, we aimed for reducing the information contained in named depictions by analyzing the principal components in a first step and projecting all relevant images into the lower dimensional space. This information could serve as the input to supervised machine learning algorithms or a simple lookup scheme to retrieve the label of the example with the most similar low dimensional representation.

**PCA** relies on linear algebra. A general assumption, especially with respect to images, is that large variances are accompanied by important structure. The covariance matrix of the data reveals information about the overall structure and orientation of the data points in the multivariate space. The axis with the largest variance is set as the first principal component. An additional assumption is that the principal components are orthogonal to each other. Therefore, the second principal component is the highest variability of all directions which are orthogonal to the first one (Bohling, 2006). The covariance between each pair of principal components is zero, as they are uncorrelated. It generally holds that the higher the eigenvalues are, the more useful they are for the analysis. As eigenvalues specify the variance of the data, higher eigenvalues indicate more relevant features.

The result of a **PCA** is a representation of the data in a new, smaller coordinate system, which depends on the axes of the largest variance. The dimensionality of this lower coordinate system should depend on the magnitude of the eigenvalues. When plotting the data along the axes of the principal components, it is often easier to understand and interpret the data, than in the original variable space.

It is widely agreed upon that **PCA** can be a good method to apply dimension reduction on images (Turk and Pentland, 1991; Lata et al., 2009). However, there are several different ways on how to do so. First of all, we performed the **PCA** on black and white images. When working on black-and-white images, only one data point per pixel is given. Therefore, performing a **PCA** is less computationally expensive, and finding structure is easier. However, as we need to be able to recognize violet as well as rust, and therefore be able to differentiate between color nuances, it was decided to work on colored images.

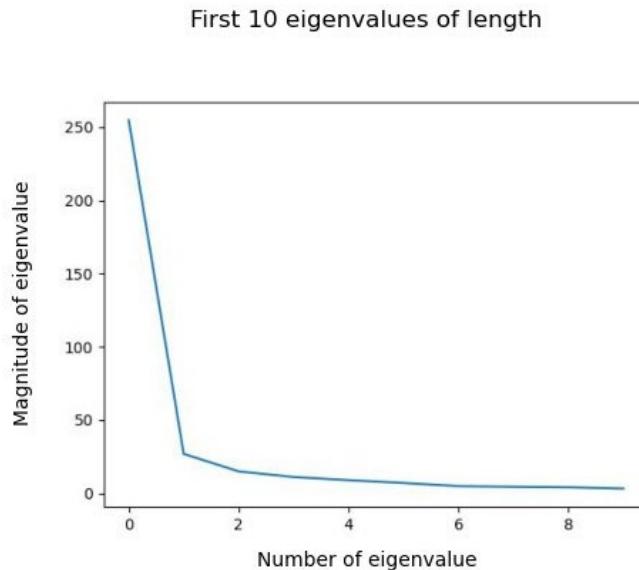
There are different possibilities regarding our project which can be considered as useful ways on how to perform a **PCA**. First, it can be performed on images of different classes at the same time – similar as to capture several images of several people in one database, on which one **PCA** is performed. In our case this would mean that a **PCA** is applied to a data set of several input images of all 13 class labels. The second way would be to perform a **PCA** separately on each group. This way, an “Eigenasparagus” in each group would be calculated, and distances between the Eigenasparaguses of different groups could be measured. Thirdly, **PCA** can be employed feature-wise. In this case, the data set would consist of a collection of images with a certain feature present vs a collection of the same size with the feature absent.

After trying several different approaches, we decided to perform our final **PCA** on sliced RGB images with background, that are labeled with features by us with the hand-label app, as this seems to yield the best results. An amount of 400 pictures per feature was used to perform a binary **PCA** for each feature (either the feature is absent or present).

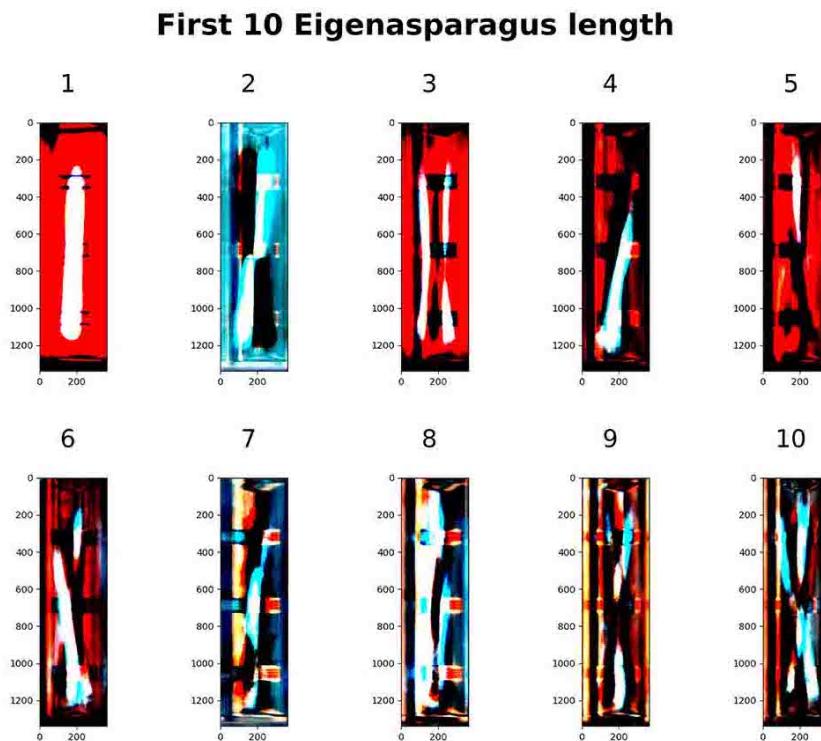
The 200 pictures where a certain feature is present as well as the 200 pictures where a certain feature is absent are extracted loops over the csv-file, where all hand-labeled information is stored as well as the path to the labeled pictures. For each feature a matrix is created, storing 200 pictures with the present feature and 200 pictures without the feature. E.g., `m_hollow` is the matrix created for the feature hollow (`shape = 400, img_shape[0] × img_shape[1] × img_shape[2]`). The first 200 entries in the matrix are pictures of hollow asparagus, the last 200 pictures show asparagus, which is not hollow. These matrices were calculated for the features hollow, flower, rusty head, rusty body, bent, violet, length and width. The data points of those 400 images in 2D space can be seen in [Figure 4.25](#).

For all these features a **PCA** is calculated by first standardizing the matrix pixel-wise (dimensions:  $1340 \times 346 \times 3$ ), calculating the covariance matrix and then extracting the ordered eigenvalues. The principal components are calculated multiplying these eigenvectors with the standardized matrix. The feature space, the principal components and the standardized matrices are saved to later perform a classification function. The highest ten eigenvalues were plotted to visually decide where to set the threshold of how many principal components will be further included. The first ten eigenvalues and the first ten Eigenasparaguses for each feature can be seen in [Figure 4.18](#) to [Figure 4.24](#).

The classification function is a control function, which performs on unseen images and tries to predict if a feature is absent or present. It reads in a new picture of one



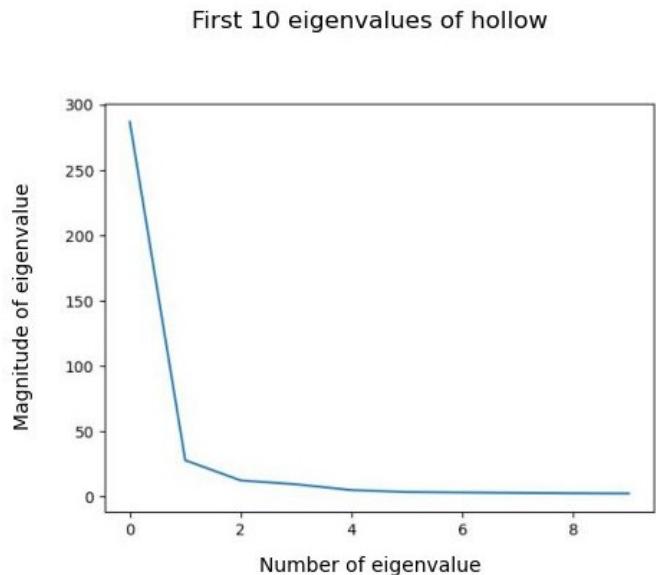
(A) This plots shows the magnitude of the first ten eigenvalues for the feature length.



(B) These images show the first ten Eigenasparagus of the feature length.

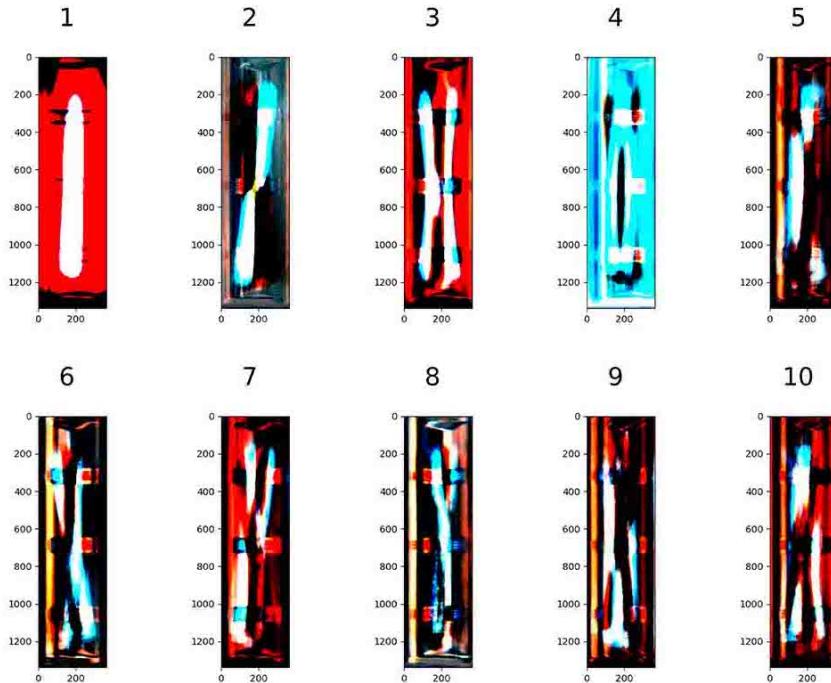
**FIGURE 4.18: Eigenvalues and Eigenasparagus of Feature Length** (A) This plots shows the magnitude of the first ten eigenvalues for the feature length. (B) These images show the first ten Eigenasparagus of the feature length.

asparagus, which is not part of the asparagus database, meaning not within the 400 images that were used to calculate the **PCA**. Then, it searches for the asparagus that is most similar to it in the feature matrices (which are 200 pictures of asparagus carrying the feature, 200 pictures of asparagus without the feature). This comparison



(A) This plots shows the magnitude of the first ten eigenvalues for the feature hollow.

**First 10 Eigenasparagus hollow**

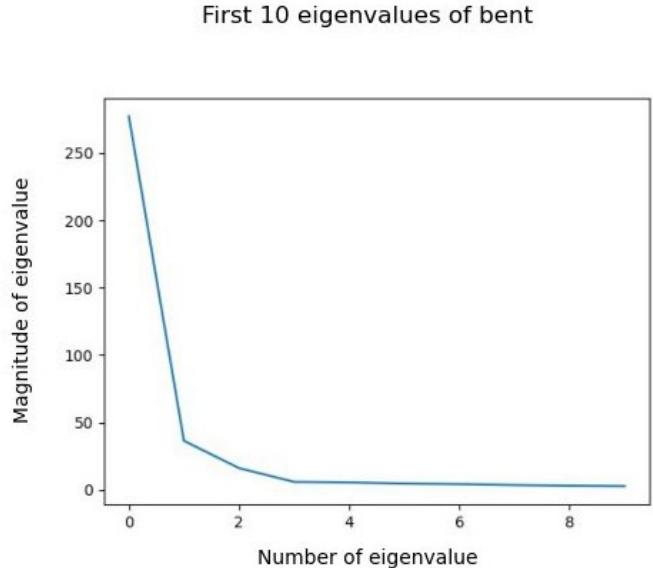


(B) These images show the first ten Eigenasparagus of the feature hollow.

**FIGURE 4.19: Eigenvalues and Eigenasparagus of Feature Hollow**

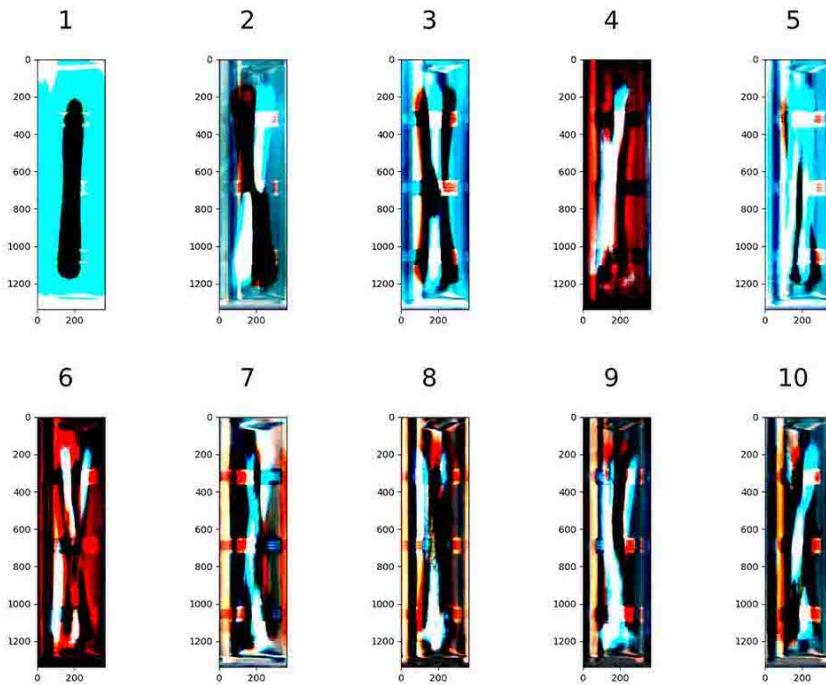
is made with the reduced representation of the original pictures.

In greater detail, the input picture is first centered by subtracting the mean asparagus and then the picture is projected into the corresponding feature-space. That means the picture is translated into the lower dimensional space, in order to compare it to the known 400 pictures. The comparison is made through calculating the distance between the single centered Eigenasparagus and the 400 pictures in the



(A) This plots shows the magnitude of the first ten eigenvalues for the feature bent.

### First 10 Eigenasparagus bent

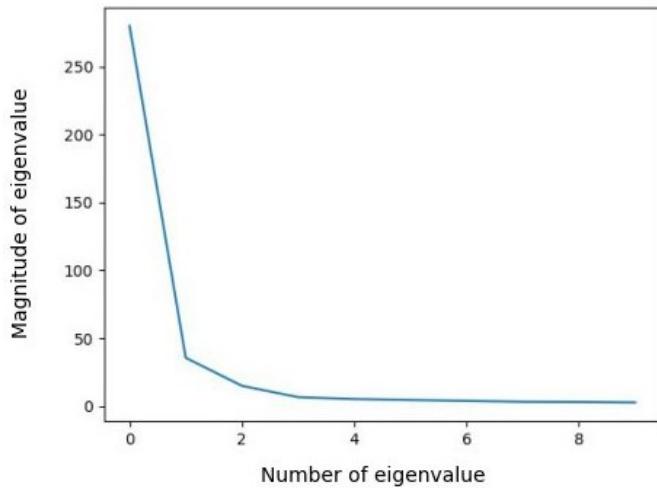


(B) These images show the first ten Eigenasparagus of the feature bent.

FIGURE 4.20: Eigenvalues and Eigenasparagus of Feature Bent

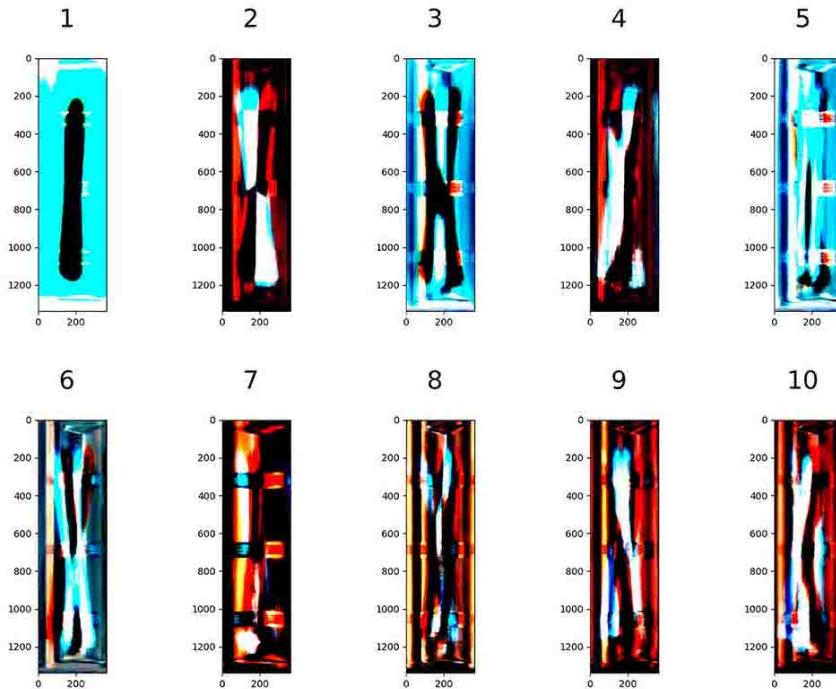
feature space, by using the `cdist` function of the SciPy package. The smallest distance is considered as the most similar asparagus. If the index of the most similar asparagus is smaller than 200 we know that the feature is present, if it is above 200 the feature is absent. By comparing this to the information of the single asparagus spear, we know if the new asparagus has the same feature as its closest asparagus in the feature space, or not. By doing this for several images, we can already presume if the two features are likely to be easily separable or not. By evaluating

First 10 eigenvalues of violet



(A) This plots shows the magnitude of the first ten eigenvalues for the feature violet.

First 10 Eigenasparagus violet

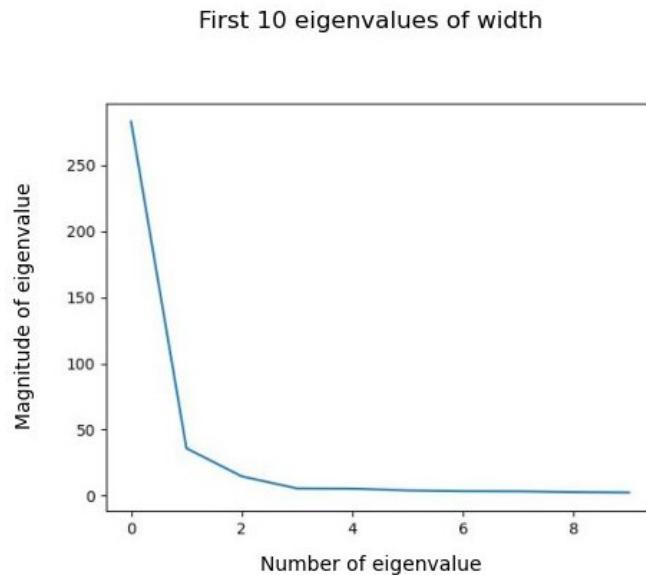


(B) These images show the first ten Eigenasparagus of the feature violet.

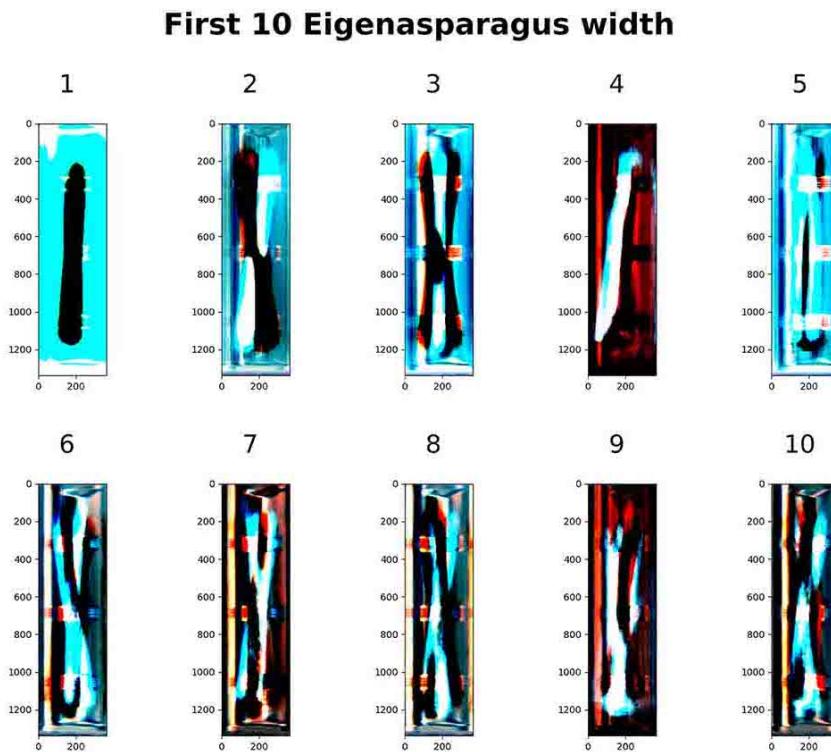
FIGURE 4.21: Eigenvalues and Eigenasparagus of Feature Violet

this, we have a measure of how well our used principal components capture the distinguishing information of each feature.

The features on which results are given are: Hollow, flower, violet, rusty body, bent, width and length. All features are binary. For the first five features, the manually hand-labeled information was obtained. For the last two features (length and width), the information, which was automatically extracted by the algorithms used



(A) This plots shows the magnitude of the first ten eigenvalues for the feature width.



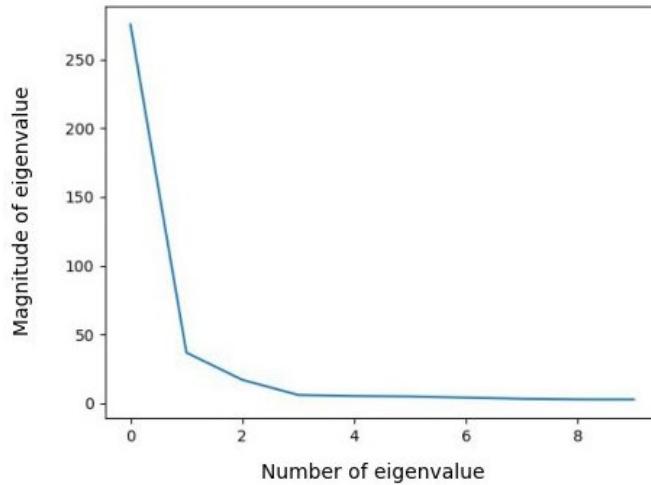
(B) These images show the first ten Eigenasparagus of the feature width.

**FIGURE 4.22: Eigenvalues and Eigenasparagus of Feature Width**

in the hand-label app were taken. The decision boundary for the first five therefore depends on our predefined criteria on labeling (see [subsection 3.4.1](#)). The decision boundary for width was narrower or equal to 20 mm or wider than 20 mm, for length shorter or equal to 210 mm or longer than 210 mm.

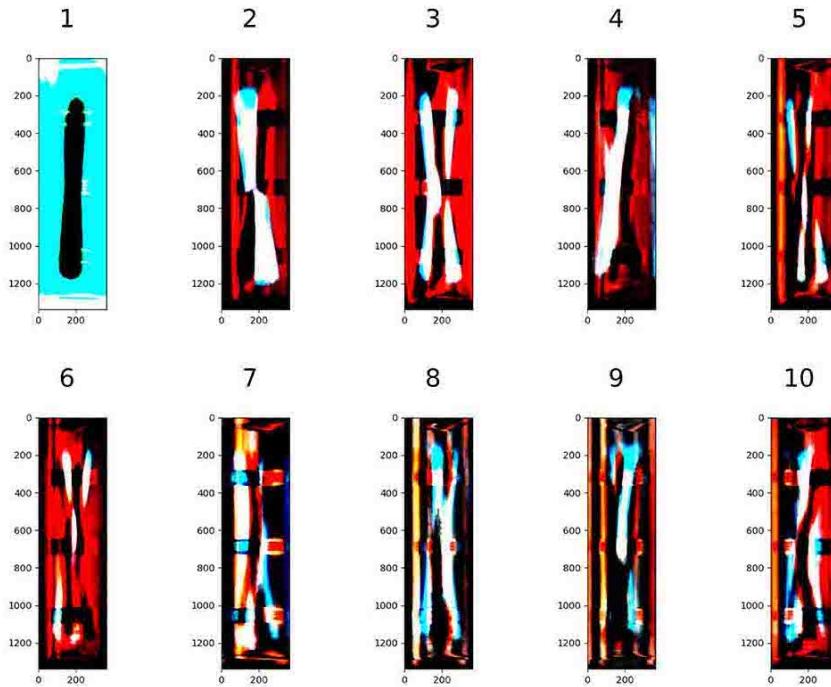
There are no results for the feature fractured, even though it is one of the initial main features, as there were only five labeled pictures for this category. Therefore, it was excluded for further analysis. For the feature rusty head, a calculation problem

First 10 eigenvalues of rusty body



(A) This plots shows the magnitude of the first ten eigenvalues for the feature rusty body.

First 10 Eigenasparagus rusty body

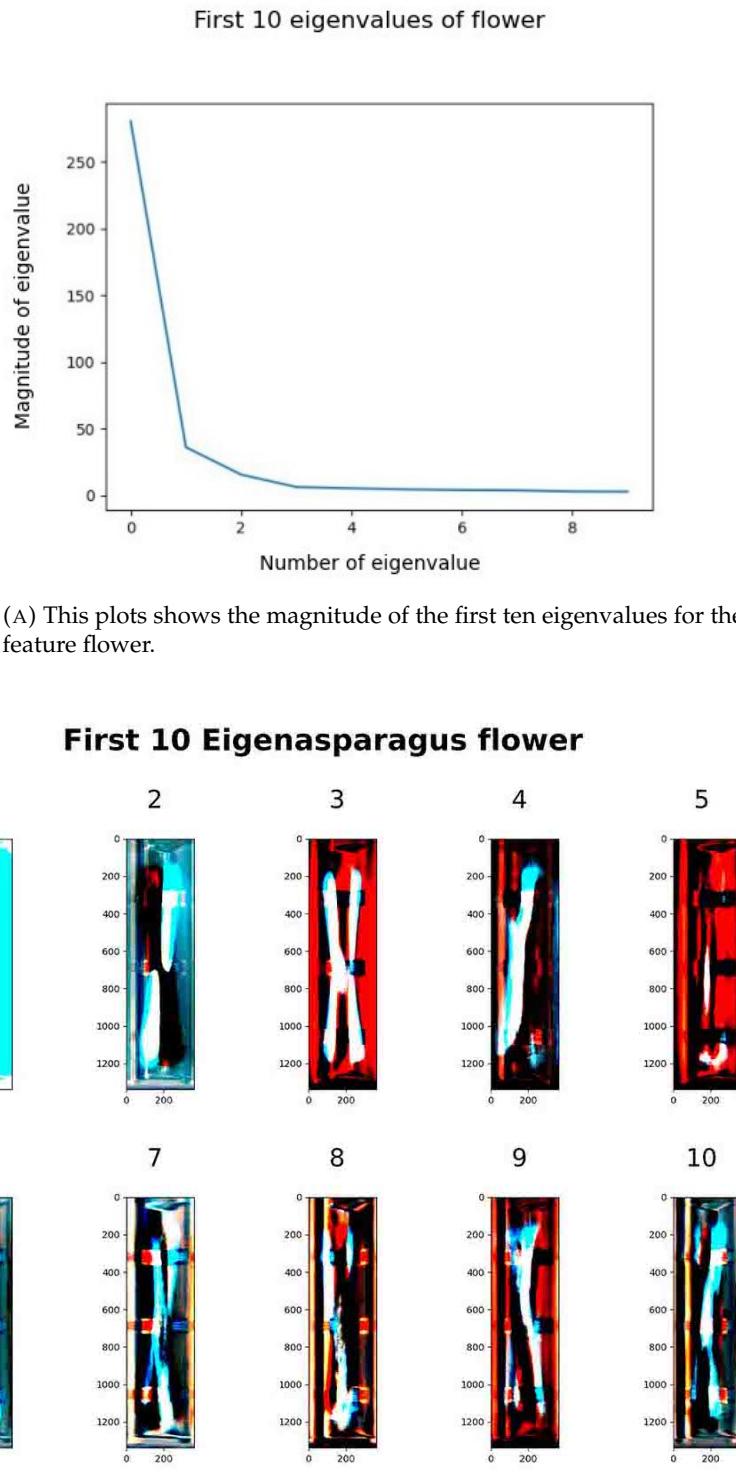


(B) These images show the first ten Eigenasparagus of the feature rusty body.

FIGURE 4.23: Eigenvalues and Eigenasparagus of Feature Rusty Body

emerged during the process. For an unexplainable reason, there were complex values in the calculation of the principal components. Due to time constraints, this problem was not solved, yet. Therefore, plotting of the feature rusty head is not possible.

By applying the approach for each feature, the eigenvectors, eigenvalues, and principal components for each feature were computed and stored. In all cases, the first



(B) These images show the first ten Eigenasparagus of the feature flower.

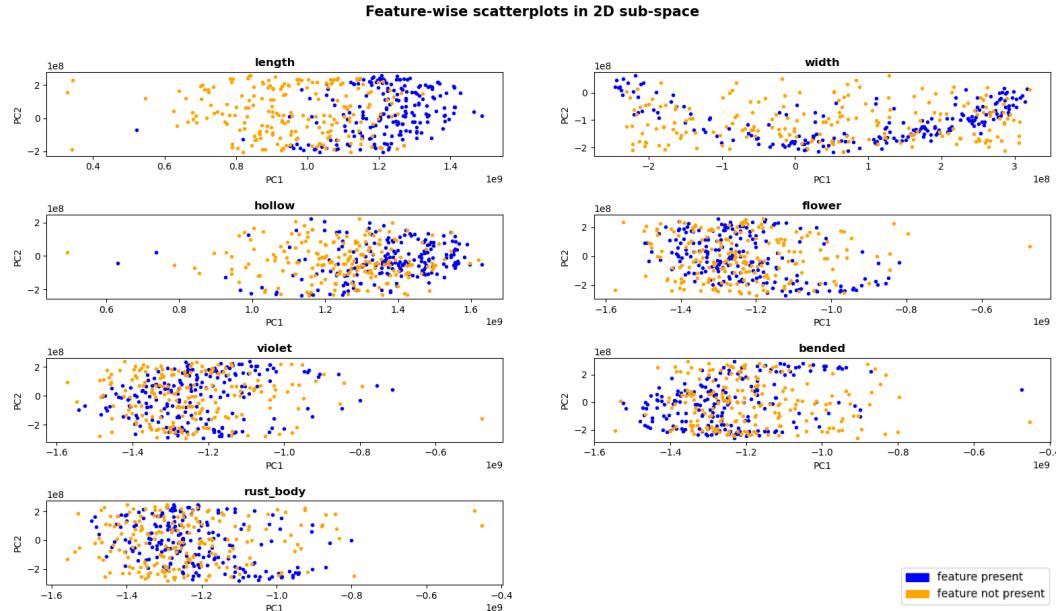
FIGURE 4.24: Eigenvalues and Eigenasparagus of Feature Flower

principal component is quite high (between 286,94 and 254,78), and drops down rapidly afterwards. The values of the principal components after the fourth principal component are very low ( $\leq 9$ ).

After inspection of the graph of the first ten principal components, only the first four principal components are used for the analysis, as there is either a sharp drop

in the slope after the first four or to maintain continuity between the different features. Following, the corresponding feature space is the space spanned by the first four principal components.

The scatterplots in [Figure 4.25](#) show the data of each feature lined up along the axes of the first two principal components of each feature.



**FIGURE 4.25: Feature-wise Scatterplots** These scatterplots show the 200 data points where the feature is present in blue and the 200 data points where the feature is absent in orange. The data is projected into a 2D subspace, which is spanned by the first two principal components of each feature.

For the classification function, we used the same ten images, to test each feature. The classification works best for the features length and hollow (10/10 classified correctly) and then width (8/10 classified correctly). It performs around chance-level for flower (6/10 classified correctly), violet and rusty body (5/10 classified correctly) and extremely poor for bent (2/10 classified correctly).

From the images of the first ten principal components, we can visually assume that there is information about the length and shape stored in the first principal component, as a clear asparagus spearpiece can be seen. The following images leave a lot of room for interpretation, about what information is contained there.

We performed the [PCA](#) on each feature separately, to extract the principal components of each feature. It is interesting to see that the pictures of the different features are all very similar. One reason for this might be that many of the 400 pictures for each feature are overlapping between the remaining features. Another reason might be, that even though the images vary between features, the general information of all asparagus images is very similar.

From the results of the classification function, one can see that there are large differences between the features in how well our [PCA](#) performs (20% - 100%). One reason for this could be that certain features are simply more difficult to distinguish than others. Another reason for this large variation can be that certain features are also more difficult to label consistently (see [subsection 3.4.3](#)), and that the results are due to inconsistencies within the data. One indicator that this is a considerable reason is that the performance of the width and length features, which

is information that is not hand-labeled, is very high. Moreover, the poorest results can be observed for the features bent and rusty body. Those are the features, for which the agreement measures show the largest discrepancies between annotators (see subsection 3.4.4).

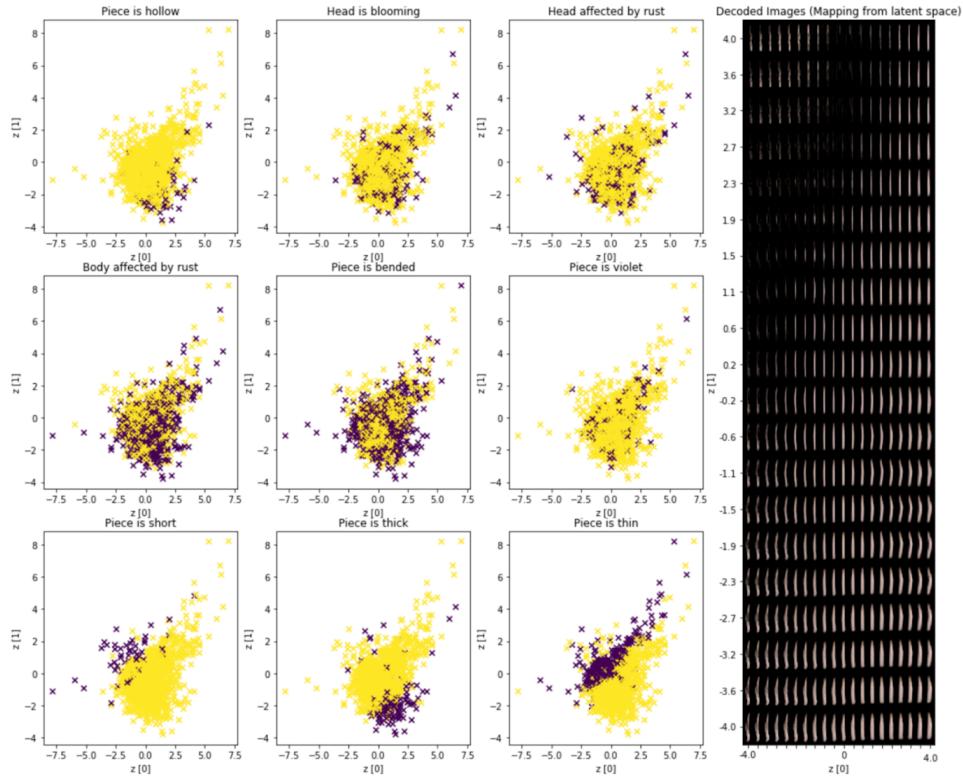
Another reason why the results are partly only moderate, is that RGB image data possesses complicated structures, and by representing it in a linear low dimensional feature space, it might be that simply too much information is lost. Even though there are papers reporting good results on **PCA** on image data (Turk and Pentland, 1991; Lata et al., 2009), there are other papers claiming that nonlinear dimension reduction algorithms are needed for this kind of image data (Olaode, Naghdy, and Todd, 2014).

### 4.2.2 Autoencoder

Beside **PCA**, there are further techniques for dimension reduction. An alternative that can be employed to deduce sparse representations and automatically extract features by learning from examples are autoencoders. Simple autoencoders, in which the decoder and encoder consist of **MLPs**, were already proposed as an alternative to **PCA** in early days of artificial neural networks when computational resources were still comparatively limited (Kramer, 1991). Today one can choose from a multitude of network architectures and designs that all have one property in common: A bottleneck layer. For image classification it is common practice to use convolutional autoencoders. There are numerous papers that employed convolutional autoencoders in various domains. Examples range from medical to aerial radar images (Chen et al., 2017). These include not only shallow networks but more recently the benefits of deep autoencoders were demonstrated (Geng et al., 2015). In addition, more complex architectures combining autoencoders with generative adversarial models were proposed recently (Bao et al., 2017). In many cases the purpose of autoencoders is dimension reduction and feature extraction or in other words the learned latent space of the bottleneck neurons. In other cases, autoencoders are used to map images from one domain to another, for example camera recordings to label-images such that after training a labeled image can be retrieved from the decoder's output layer (Iglovikov and Shvets, 2018). In short, there are many possible ways to apply autoencoders and many architectures to realize them.

This motivates the question of how autoencoders work. As mentioned, all autoencoders have a bottleneck layer. If applied for dimension reduction autoencoders are usually used to predict the input, in this case the image, with the input itself. The sparse representation corresponds to the activation of the latent layer for a given input. Autoencoders consist of an encoder that contains the initial layers as well as the bottleneck layer and the decoder that maps the respective latent space back to the image. The desired mapping function of the input to a sparse representation is generated as a side product of the optimization in end to end training, as weights of the decoder are trained such that meaningful features are extracted. The main difference to **PCA** is that nonlinear functions can be approximated. Feedforward neural networks such as the encoder are non-linear function approximators. Networks with multiple layers are especially well known for establishing named nonlinear correlation. Hence, autoencoders allow for non-linear mappings to the latent space (Kramer, 1991). This means that in the latter multiple features may be represented in a two dimensional space. It shows that compared to **PCA** where

one dimension typically corresponds to one feature more information can be represented in fewer dimensions. Different properties of the input are mapped to different areas of the latent space. We used **Variational Autoencoders (VAEs)**. In VAEs, a special loss function is used that ensures features in latent space are mapped to a compact cluster of values. This allows for interpolation between samples by moving on a straight line because regions between points in the latent space lie within the data and hence reconstructions of the decoder are more realistic (Scikit-Learn, 2020). Other than that, **Variational Autoencoders** share most properties with regular autoencoders. The location of a point in latent space refers to a compressed representation of the input. These can be interpreted as features of the input.



**FIGURE 4.26: Latent Asparagus Space for MLP -VAE and Reconstruction Manifold** The depiction illustrates the mapping by the encoder and decoder of the VAE: On the left you can see scatterplots illustrating the activation of latent layer neurons for the test data set (the mapping by the encoder). Image-features are mapped to the respective latent asparagus space. There is a scatterplot for each feature of interest where colors indicate positive and negative samples. On the right a manifold of decoded images is shown. The axes relate to the points sampled in latent asparagus space that correspond to the reconstructions (mapping by the decoder).

Different **Variational Autoencoders** were tested in the realm of the project. First, a simple **Variational Autoencoder** with a **MLP** as decoder was implemented. The second approach was a comparatively shallow convolutional **Variational Autoencoder**. The third approach relates to an autoencoder with a deeper encoder that was later used to design the networks for semi-supervised learning. As the third approach is more complex but did not improve the mapping of the properties of asparagus spears to a two dimensional latent asparagus space, the results for the

second of named networks are reported here. It was derived from a standard example for **Variational Autoencoders** (Scikit-Learn, 2020). The presented results are for a network that comprises two hidden convolutional layers in the encoder and two deconvolutional layers in the decoder.

Similar to the presented way of applying **PCA**, batches of images that contain only one perspective were used as input to the network. The downsampled data set was used. Images had to be padded as the implementation does not work with inputs of arbitrary shape. The deconvolutional layers of the decoder can only increase dimensionality by an integer factor: The filters that are used for the deconvolution in the given net double the tensor dimensionality. For padded images, this means an increase for the vertical dimension from 34 to 68 and finally to the desired 136 pixels is achieved in the last three layers of the network (impossible for the original height of 134 pixels). The input shape which also defines the shape of the output layer must be divisible by two without remainder twice. Depiction [Figure 4.26](#) shows the results.

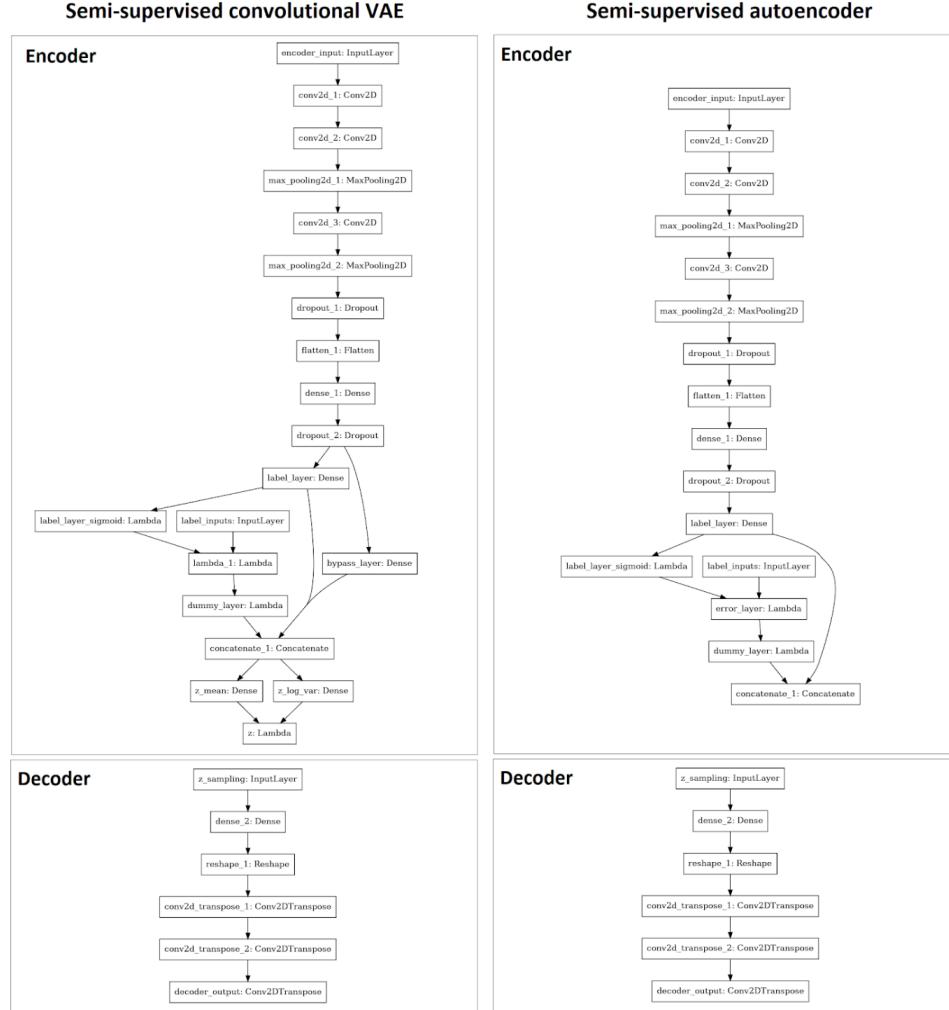
It demonstrates that the features short, thick and thin are mapped to separable clusters. As a tendency curvature correlates with a region in the lower periphery as indicated especially by the deconstruction. The other features (violet, rust and flower) are not mapped adequately and they are not visible in the reconstruction. This shows that only some features of interest were mapped to the latent space and used to decode images. Reconstructions of autoencoders are known to miss many details. Better results could potentially have been achieved using larger input images. The possibility to generate, for example, more or less curved asparagus spears may help to define a clear cut decision boundary and classify images accordingly. As a potential feature for asparagus sorting machines this would allow the user to customize the definition of curvature to the own taste. For this approach to be viable for all features, however, the network performance appears to be insufficient. Some features are poorly separated by the network that was employed on downscaled images.

### 4.3 Semi-supervised learning

We collected more than 100000 samples. Considering the uniform appearance of the images this represents a substantial amount. However, labels had to be manually generated which was done for only around 10000 samples. As a consequence, there is only a small subset of data with attributed labels. Smaller amounts of labeled data mean that predictions can be successful only if the variance in the source values is limited. Hence, for high dimensional data such as images, sparse representations are desirable. Extracting features automatically instead of relying on manual feature engineering is a strategy that is especially appealing if large amounts of unlabeled data are available. In semi-supervised learning features are extracted from the main input (here images) using unsupervised learning if target labels are unavailable (Keng, 2020). Hence, semi-supervised learning promises better results by using not only labeled samples but also unlabeled data points of partially labeled data sets.

In the previous chapter, methods and results for unsupervised learning are presented. One example is a convolutional autoencoder. In this section, it is shown how convolutional autoencoders with additional soft constraints in the loss function can be used for semi-supervised learning. Instead of using unsupervised methods to compute another data set of sparse representations and use the latter to predict labels, in semi-supervised learning sparse representations are retrieved and

mapped to latent features at the same time (Keng, 2020). Bottleneck layer activations represent automatically extracted features. For semi-supervised learning one tries to enforce that latent layer activations of autoencoders correlate with the target categories.

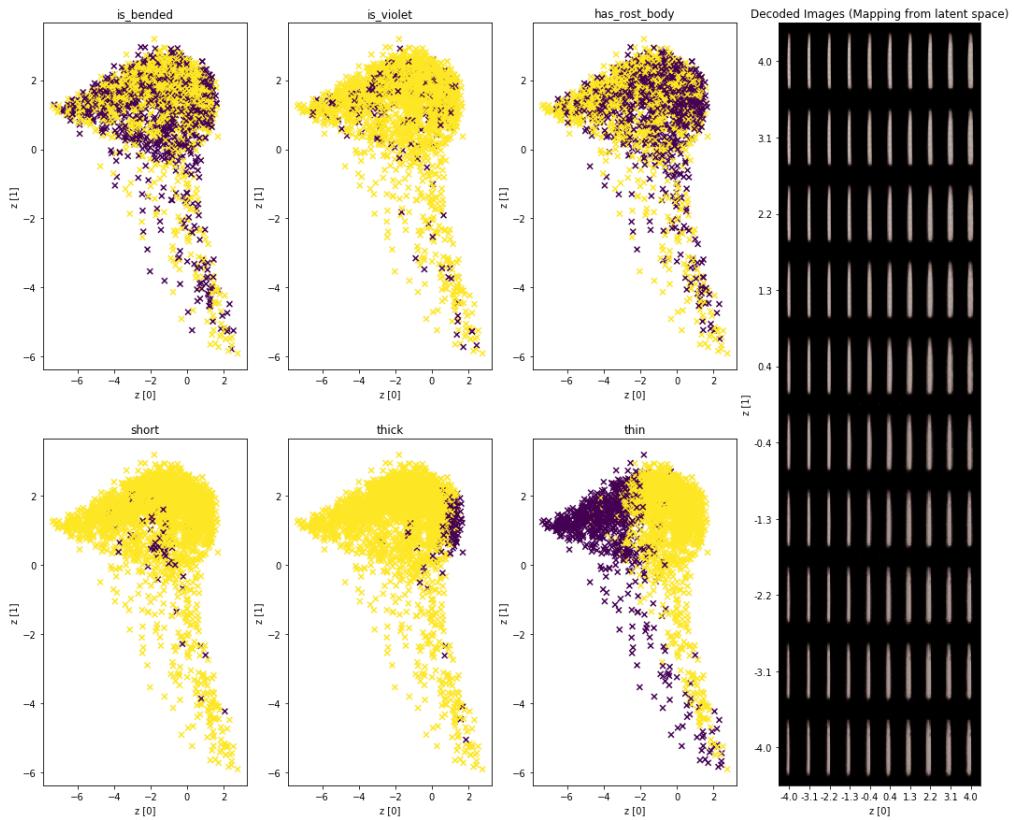


**FIGURE 4.27: Network Structures for Semi-Supervised Learning** The depiction illustrates the network structure for the autoencoders used for semi-supervised learning. Left: The structure for the semi-supervised convolutional Variational Autoencoder. Right: The structure for the semi-supervised convolutional autoencoder. Further explanations can be found in the text.

The network structure is derived from the convolutional autoencoder used for unsupervised learning. The feedforward CNN is replaced by a network that has proven to be suitable to detect at least some features with sufficient adequacy when trained on heads (see subsection 4.1.4). It comprises three convolutional layers with a kernel size of two in the first and three in subsequent layers as well as 32 kernels each. A max pooling layer with stride two is added mainly to reduce the number of total neurons while maintaining a high number of kernels. Additionally, a dropout layer is added to avoid overfitting. In contrast to other implementations for semi-supervised learning (Keng, 2020) the same network is used for the prediction of labels (when they exist for the current batch) as for the encoder that retrieves the sparse representation for reconstructing images. For the decoder we chose the

same one as in the **Variational Autoencoder** presented in the previous chapter. The effects of a bypass layer that contains neurons not being subject to the label layer loss (see [Figure 4.27](#)) were tested. As accuracy did not improve, it was later missed out. Two variations of the network were tested: A convolutional **Variational Autoencoder** and a **Variational Autoencoder** for semi-supervised learning.

A challenge results from training with multiple inputs. As deep learning frameworks usually require a connected graph that links inputs to outputs, a trick is used in order to handle that images as well as labels are given as an input. A dummy layer was introduced where all information which is derived from the labels was multiplied by zero. The output vector is concatenated with the bottleneck of the encoder. As it contains no variance, training and more importantly validation accuracies remain unaffected even though information about the categories that are predicted is added on the input side. Nevertheless, the labels are part of the network graph and could hence be used in the loss function.



**FIGURE 4.28: Latent Asparagus Space for Semi-Supervised CNN-VAE** The depiction illustrates the mapping by the encoder and decoder of the **Variational Autoencoder** with additional constraints in the loss function for semi-supervised learning: On the left one can see scatterplots illustrating the activation of latent layer neurons for the test data set (the mapping by the encoder). Image-features are mapped to the respective latent asparagus space. Mind that there is a scatterplot for each feature of interest where colors indicate positive and negative samples. On the right, a manifold of decoded images can be found. The axis relates to the points sampled in latent asparagus space that correspond to the reconstructions (mapping by the decoder).

A custom conditional loss function is used. If labels are present for the current batch the costum loss is equal to a combined loss that comprises the reconstruction

loss and the label loss. Here, reconstruction loss refers to the pixel-wise loss that is used for the main task of the network - namely mapping of input images back to the same images (fed into the network as target values to the output layer). The label loss is used with the goal of mapping label layer activations to the actual labels. It is low if activations in the sigmoid transformed label layer match the target values i.e. the sum of the error layer is low. The loss that is due to labels is multiplied by a custom factor ( $k$ ). In addition it is defined such that it scales with the current pixel-wise reconstruction loss and converges to a constant ( $c$ ). These values were chosen aiming for an increase of the contribution of the label loss to the combined loss especially in late stages of training.

	False positive	False negative	True positive	True negative	Sensitivity	Specificity
is_bended	0.04	0.37	0.04	0.55	0.10	0.94
is_violet	0.00	0.08	0.00	0.92	0.00	1.00
has_rust_body	0.14	0.27	0.20	0.39	0.42	0.73
short	0.00	0.02	0.00	0.98	0.00	1.00
thick	0.00	0.07	0.00	0.93	0.00	1.00
thin	0.00	0.14	0.16	0.70	0.54	0.99

TABLE 4.7: Performance of the semi-supervised convolutional VAE.

	False positive	False negative	True positive	True negative	Sensitivity	Specificity
is_bended	0.02	0.34	0.07	0.57	0.18	0.96
is_violet	0.00	0.08	0.00	0.92	0.00	1.00
has_rust_body	0.23	0.19	0.28	0.30	0.59	0.57
short	0.00	0.01	0.01	0.97	0.57	1.00
thick	0.00	0.06	0.00	0.93	0.04	1.00
thin	0.02	0.07	0.23	0.68	0.76	0.97

TABLE 4.8: Performance of the Semi-Supervised Convolutional Autoencoder

The results for the semi-supervised **Variational Autoencoder** are illustrated in [Table 4.7](#) and visualized by [Figure 4.28](#). As one can immediately see, the feature thin was adequately mapped to a decisive region in latent space. For the other features, no such clear cut clustering is visible. Reconstructions indicate that the main purpose of the network of predicting the input images was accomplished successfully although the reconstructions look rather uniform. Values for accuracy and sensitivity indicate only poor performance. Only the features rusty body( $se=.42$ ), thin( $se=.54$ ) and bent( $se=.1$ ) sensitivities are above zero.

Compared to the variational convolutional autoencoder for semi-supervised learning, the convolutional autoencoder for semi-supervised learning performs better. However, there are substantial potentials for improvements. [Table 4.8](#) shows a summary of results. Violet detection is not successful at all as indicated by zero sensitivity ( $se = 0$ ). For the other features the network was trained on, mediocre results are achieved. Thickness detection shows little sensitivity ( $se = .04$ ) however a high specificity ( $sp = 1.0$ ). Better results exist for curvature (( $se = .18$ ) rusty body ( $se = .42$ ) and short spears ( $se = .6$ ). The specificity for rusty body is lower ( $sp = .6$ ) as compared to named other features ( $sp = 1$  and  $sp = .96$ ). Thin spears are detected in

76% of all cases and few false positives characterize the detection of named feature ( $sp = .97$ ).

The approach for semi-supervised learning presented here faces two challenges. First, the networks were trained only on one perspective although labels were attributed per asparagus spear. Information that is only visible on one or two out of three perspectives can not be mapped to the desired target category. Image-wise labels would be desirable to improve the approach.

Second, reconstructions using convolutional autoencoders contain little detail. However, small details in the image, such as small brown spots that indicate rust, play an important role for classification. These features are not sufficiently reconstructed by the **Variational Autoencoder**. Arguably, they are not reflected in the sparse representation that corresponds to latent layer activations. One may speculate that a substantial increase of the network size would help to reconstruct more details and hence extract more features. As generative adversarial networks are known to generate more detailed images (Bao et al., 2017) they could possibly be adopted for semi-supervised learning with greater success. However, this is a question that must be answered empirically.

In summary, one may conclude that autoencoders appear as an alternative to manual feature engineering if a large data set is available and only a subset contains labeled samples. However, more research is necessary to find best suitable network structures for asparagus classification.

## Chapter 5

# Summary of results

The current work was conducted to investigate the state of the art on asparagus classification aiming at the development of approaches that could help to improve the current sorting algorithm implemented at the Autoselect ATS II at the asparagus farm Gut Holsterfeld. Data was collected, preprocessed and then analyzed with seven different approaches. Out of our 591495 images corresponding to 197165 different asparagus spears, 13271 were collected by re-sorting in the machine and 13319 images manually labeled by us. This labeled data was considered for the supervised approaches. The semi-supervised approach was in addition based on approximately equally many unlabeled images. The results illustrate that classifying asparagus is not a trivial problem. However, the results also show that it is possible to extract relevant features and to improve current sorting approaches.

For supervised learning, we employed **MLPs** and **CNNs**. Whereas the earlier were trained on sparse descriptions retrieved by high level feature engineering, the latter were directly trained on preprocessed images. They include a dedicated network for head-related features, as well as networks for single-label classification and multi-label classification. All approaches aim to solve the same image classification problem, using supervised learning.

The feature engineering **MLP** for curvature prediction and the single-label classification **CNN** perform binary classification (see subsection 4.1.1 and subsection 4.1.2), whereas the **CNN** for head-related features as well as the multi-label approach perform multi-label classification (see subsection 4.1.4 and subsection 4.1.3).

With the random forest approach, a third option namely multiclass prediction is added that works completely differently. Each approach has drawbacks and benefits. The complexity and the requirement to specify many parameters has proven to be a disadvantage of relatively deep, but also rather shallow **CNNs** as compared to e.g. **MLPs** or random forests. In contrast, stronger preprocessing or even feature engineering is required to successfully employ the latter. After all, however, the most important criterion to evaluate an approach is its predictive performance.

The heterogeneity of approaches with respect to the number of target categories and the variety of performance measures pose challenges for a direct comparison using the overall accuracies. Therefore, feature-wise evaluation appears most promising. As the distribution of some features (e.g. violet) has proven to be very unbalanced in our dataset, even high accuracies might relate to poor predictions (e.g. when the feature is never detected). Hence, feature-wise accuracies are only a coarse indicator of the model's performance that may nonetheless give insights where difficulties lie and what features are more problematic than others. However, for some promising approaches we computed the sensitivity and specificity per feature to reveal a more fine-grained picture of the predictive performance.

The multiclass approach has an overall accuracy of 75%. The performance of the **CNN** for head-related features is indicated by sensitivity and specificity values.

Flower detection reaches 55% sensitivity and 95% specificity while rusty head detection attains only 19% sensitivity at 98% specificity. The overall accuracy of the multi-label **CNN** approach reaches up to 87%. For this model, accuracies are not calculated per feature.

In contrast, feature-wise accuracies for binary classification can be reported. The same holds for feature-wise performance measures that were calculated for some of the other approaches. The feature engineering based approaches show good results on curvature (82% sensitivity, 67% specificity) and similarly for violet detection (62% sensitivity and 96% specificity) as well as for rusty body (71% sensitivity, 65% specificity).

In the single-label **CNN**, best results are achieved for features relying on the thickness and length of the asparagus. All of these features achieve a balanced accuracy above 90%, with best results for the feature very thick (98% sensitivity and 99% specificity). Of the solely hand-labeled features, feature hollow shows the best performance (77% sensitivity and 98% specificity). The feature rusty head has the least performance (52% sensitivity and 81% specificity).

The unsupervised learning approaches, namely **PCA** and the convolutional autoencoder, both deal with dimension reduction. Both were trained on the sample set for which labels are available. While the classification method based on **PCA** targets at binary feature prediction (absence or presence of a feature), the unsupervised autoencoder does not predict labels. The accuracy of **PCA** is promising for length and hollow (100%) but extremely poor for curvature detection (20%). It has to be mentioned, however, that only very few samples were used for training and evaluation of the named approach. As such it is yet to be proven whether or not these results generalize.

In the last approach, a semi-supervised learning method was based on a partially labeled data set. It performs multi-label classification. Unfortunately, the predicted power was rather poor, best for curvature (18% sensitivity, 96% specificity) and worst for violet as it does not detect any violet spears (0% sensitivity).

## Chapter 6

# Discussion

In our study project we pursued three main objectives. The main goal was to develop and improve algorithms for asparagus classification. The second objective is closely linked to this and relates to best practices in relation to applied data science and big data. This included storage of data on remote servers and computationally expensive procedures that are required for training in the computational grid of Osnabrück University: The methodological aspect of our study project. As our work also served as a sample project to learn more about possibilities to effectively organize collaborative work we also targeted at the organizational aspect that is closely linked to project management. In the following, the core results with respect to named objectives are shortly named and discussed.

### 6.1 Classification results

Asparagus spears have several features that we aimed to extract. Some features such as the length and width of asparagus spears were undoubtedly measurable using a pure computer vision approach that does not rely on machine learning. For others, direct filtering was not easily possible because no clear cut definition of features, such as curvature or violet, exists, that is precise enough to be implemented directly. Although relevant information can easily be extracted, the rules to infer the desired binary features are inaccessible: On one side, decision boundaries for binary classifiers have to be found. On the other side, the perception of the features color or curvature has been shown to be very subjective. Moreover, filtering features such as flower has been proven difficult. Named attribute relates to details in a few pixels, comes in different forms, and highly depends on the perspective. Because of these reasons machine learning must be employed to successfully classify asparagus.

We designed several neural networks and applied them in different ways to analyze and classify a large-scale asparagus image data set. Some approaches worked better than others, as can be concluded from [chapter 5](#).

Training **MLP** classifiers on histograms of palette images has proven a promising approach to predict color features. Named histograms contain information about the fraction of foreground pixels that correspond to violet or rusty colors. As **MLPs** have few parameters, the design is rather trivial and the training process quick. The results are considered as a good baseline performance for rusty body and violet. The application on palette images that show asparagus heads only might also be suitable to predict the feature rusty head. This is yet to be tested. **MLPs** have also shown to be suitable when applied to predict the binary curvature feature based on partial angles. The fact that predictions are far from perfect might be due to inconsistencies in the training data. One may assume, however, that the models generalized well and represent rules that relate to average opinions or definitions of highly subjective features such as color or curvature.

Feedforward **CNNs** were applied to predict individual features, for multilabel prediction and predictions based on snippets that depict asparagus heads. In addition, effects of a custom loss function were tested. Promising in the multilabel prediction is that not only individual features, but also the relation between features can be considered in the learning process. Our multilabel **CNN** reaches an accuracy up to 87%, which seems high. However, when looking at the accuracy and loss values over time, one can see that the model does not improve much. While sensitivity and specificity improve, and therefore indicate learning, the validation loss remains high, indicating overfitting. This model seems to be especially sensitive to the imbalance between 0 and 1 in the label vectors. Concerning this, there is still room to play around with our parameters to further improve the architecture.

Applying **PCA** on individual features and projecting the image information into a smaller dimensional subspace showed promising results. It revealed that the first principal components managed to capture most of the information. However, differences between most features seem to be too small to be adequately represented in the low dimensional space. In this approach, the features width, length and hollow seem to be classifiable with high performance, and the features bent and rusty body seem to be most difficult. Width, length and hollow (as hollow asparagus is likely to be thick) are features that can be related to the shape and spatial appearance of the spear in the picture. This walks together with the findings that the first principle components (so the most important ones) only refer to the appearance of the asparagus in the picture and shows the same picture for all asparagus. This leads to the assumption that the spatial appearance is counted as the most important feature, rather than taking the surface of the spears into account. This problem could be improved by generating pictures, where the possible asparagus positions are equally distributed over the pictures. Another reason for the similarities between the first principal component pictures is that one asparagus can have many features. Therefore the same pictures can be used for several feature pcas. As asparagus with only one present feature is rather difficult to find, another solution to attract this problem needs to be found.

Similarly, **Variational Autoencoders** were used to derive a low dimensional representation using unsupervised learning. While some features such as the width and length are mapped to clearly differentiable regions in latent asparagus space, this is not the case for many others. Only as a tendency, spears labeled as bent are for example mapped to regions in the lower periphery. Autoencoders are known for blurry reconstructions. This is a possible explanation for the lack of clusters in latent space for features that relate to details that are not sufficiently reconstructed.

Convolutional autoencoders were used for semi-supervised learning. However the results for this approach can be described as merely mediocre. One problem is arguably the mentioned insufficiency in reconstructing details. As details such as brown spots define target classes (e.g. rusty head) and they are not present in latent space. It is hard to establish a correlation of the respective latent layer activation and the target labels. Larger input image sizes or different network architectures that are suitable to reconstruct higher detail images could potentially help to improve performance of semi-supervised learning.

Detecting the feature rusty head has proven rather difficult even though a dedicated network was trained on snippets that show asparagus heads in rather high resolution. This is potentially the case because details that are hardly visible even to the human eye have to be considered that occur in different locations. Although

better results are achieved for the feature flower, the same most likely holds for this category as well. In contrast, better results are achieved for features that relate to the overall shape of asparagus spears instead of fine details. This holds for the category hollow as well as for curvature. Color features are detected especially well based on histograms of palette images while **CNNs** have proven suitable to detect shape related features.

In summary, we successfully measured the width and height of asparagus spears and were able to develop detectors for the other features that performed surprisingly good, given the moderate inter-coder reliability that was partially due to unclear definitions of binary features such as bent or violet.

## 6.2 Methodology

Looking back, there are several methodological issues, which we would process differently now. We started our study project at the beginning of April. This is as well the beginning of the asparagus harvesting season. Therefore, we were able to start collecting data straight away. On the down-side, we had to start collecting data without a detailed plan in advance. Planning the data collection ahead could have made the data acquisition more efficient, and structured. Afterwards we could not change relevant parameters to answer a lot of organizational and methodological questions such as: How much data do we need? What format do we need it in? Is autonomous calibration possible? How exactly do we store the images effectively and efficiently? Is the setup as we need it? And if not, how can we improve it? What kind of measurements or changes do we want to perform with the camera? Is the illumination as we want it? Could stereo cameras or other 3D view techniques such as depth cameras or laser grid be used? Would we integrate an additional camera taking pictures from the bottom of the spear or from the head region separately? What should our pipeline look like? How can we get labeled data right away?

As already mentioned, there was a misunderstanding between the group and the supporting asparagus farm about the data type necessary. The already existing images were too few, and unlabeled. Therefore, we spent the entire asparagus season with data acquisition instead of starting with preprocessing as planned originally. The number of labeled images that were collected by running re-sorted classes though the machine is arguably insufficient to learn classes using the chosen deep learning approaches (Russakovsky et al., 2013; Russakovsky and Fei-Fei, 2010; Warden, Pete, 2017). Therefore, a lot of time was spent on preprocessing and labeling the data manually.

Another discussion point concerns our data. The image quality in terms of pixel size of our images is really high. Due to limited memory capacity and long run-times of some of the tested networks, images needed to be down sampled. Additionally, images of different file types were used (e.g. png, jpg), which also includes reduction in disc space. We should further investigate to what extent images can be down sampled without losing critical information, in order to optimize the named difficulties.

Even though a variety three images of every asparagus spear are given, they are all taken from the same perspective – from above. In the ideal case, the asparagus spear rotates over the conveyor belt, such that each spear is depicted in the pictures from a different viewpoint. The better the asparagus rotates, the more reliable is a later judgement of the spear in terms of class labels or features. As the rotation

often lacks, because the spear is too bent, an additional viewpoint could improve the rating. Concrete ideas how to improve the setup are given in the outlook.

As previously mentioned, our labels of asparagus features are partly achieved by computer-vision algorithms, partly based on human perception. As previously outlined, human performance is commonly acknowledged as the baseline performance in classification tasks. While the performance of our automatic feature-extraction for length and width is really high, we decided that for the features violet, rusty head, rusty body, bent, hollow and flower a human perception would be more accurate. Even though this is commonly used as the “gold standard”, it can of course also bring more variation, and maybe even inconsistency between raters, than an algorithm.

As explained in the section [Preprocessing](#), during preprocessing for the labeling procedure, the three images of one asparagus spear are joined together and then labeled. Therefore labeling was faster, because three perspectives are labeled at once. However, it could also be tried to use the single perspectives and conclude the classes from a combination of the recognized features.

We kept the features binary, as this is easier to label, and easier to use for our supervised classification approaches. The down side of a binary label is, however, that a clear boundary is set, where in real life there is a smooth transition. Even for our supervising farmer, it is sometimes difficult to decide on a boundary. This is arguably due to the small differences between groups, and vague borders between positive and negative examples. While the binary representation makes certain analyses and classification much easier, it also brings restrictions.

Moreover, we observed difficulties concerning the labels in the communication between the group and the farmer. The communicated need is that the sorting algorithm works “better”. But what does that technically mean? And what is technically possible? For the farmer, the sorting would already be “better”, if the sorting mistakes would be more systematic. This would not necessarily mean that the overall accuracy of correctly sorted asparagus into one out of 13 class labels needs to improve, but that the overall impression of all spears sorted into one tray is more homogeneous.

### 6.3 Organization

This section summarizes and briefly discusses what each member has personally learned on an organizational level during the year.

The team did not only achieve new scientific skills and techniques of data acquisition, preparation, and analysis but also gained valuable new insights into the organization of a large project. It was learned how to structure the team more successfully and purposefully.

First and foremost, this needs excellent communication. Not everyone has to discuss or listen to each specific detail in every working area. Instead, communication should be balanced. Often, it suffices when all team members have a broad overview.

Secondly, it turns out to be helpful when one or two members exchange some task-related work in favor of more management-related work. Democratic decision making does not necessarily exclude the role of a team leader or of a manager who has an overview of the tasks or of a certain task area. The whole team agrees to structure the next project in the same way as done during the second term, including manager roles and a stricter working plan. By this, better team dynamics are gained and time management can be improved.

Further, the strengths of the single members have to be evaluated before the project starts so they can be used efficiently. Although, not everyone has to do the task he or she is best at. One should also have the opportunity to work on tasks that are new, challenging and interesting. It allows each member to broaden their skills and it avoids discouragement.

Finally, the team agrees to have more focus on the overall goal than only think of what directly lies ahead. For this, concrete goals have to be formulated well. Milestones or intermediate goals should be defined and better checked. More time has to be taken into consideration when planning ahead as well as for including adjustments.

In conclusion, the experience of having two different working structures gave us the ability to compare and judge what is essential to successful teamwork. It also helped to understand how each member can contribute to the team regarding personal skills and interests, and what each member wants to improve for future teamwork.

As the main intention of the study project was to enhance our knowledge, we sought out a task that we were highly motivated to do. By that, we did not only practice theoretical knowledge but gained new experiences and sustainably improved our team skills for future work.

## Chapter 7

# Conclusion and outlook

In this project, we successfully adopted various modern computer vision techniques with the aim of classifying asparagus spears according to their descriptive features or class labels. Alongside the development of the respective approaches, we examined the practicability in commercial applications. The specifications and goals of the project that we defined in the beginning helped structuring the agile planning sessions. As a result we implemented several approaches that allow us to classify asparagus spears in dependence of their features or class label. After a time-consuming data collection, laborious labeling and classical pre-processing steps, different trainable models were implemented from the fields of supervised learning, unsupervised learning and semi-supervised learning.

Whether we have succeeded in improving the currently running sorting algorithm could not yet be clarified systematically because of the lack of time and resources for a suitable evaluation. In cooperation with the local asparagus farm Gut Holsterfeld and the manufacturer of the asparagus sorting machine Autoselect ATS II, a method can now be developed. However, we could prove that computer vision and machine learning based techniques are able to solve the asparagus classification problem.

The intention of the study project regarding the training and evaluation of project management (Universität Osnabrück, 2019b; Universität Osnabrück, 2019a) has been followed up. Good self-organisation and team building became essential for the progress of the project. We have learned to organize, work as a team, put theoretical knowledge into practice, develop software, analyze data, and to practice the analysis and interpretation of statistical data under the conditions of a research project.

Due to the direct start of the harvesting season at the beginning of the project, the sorting machine and the hardware setup had to be used as already implemented. However, improvements were discussed with the manufacturer of the machine, Mr. Hermeler. A second camera to capture the head of the asparagus is needed in particular. This is reflected in our results. Especially for features like flower or rust, an additional head camera helps greatly with the classification. Another camera taking an image from the bottom of the spear could improve the detection of the feature hollow. These cameras are already implemented in some new asparagus sorting hardware.<sup>1</sup> Overall, it is an advantage to have even more perspectives on the asparagus, as some features might only be visible from a certain side. Namely, the lighting and reflection are slightly different, which can alter the color values, and distort the images. Hence, the asparagus might be stretched depending on the perspective. This brings in more usable information for classical computer vision algorithms to determine certain features accurately. It is also conceivable that

---

<sup>1</sup>see

<https://www.neubauer-automation.de/uk/asparagus-sorting-machine-espaso-technicaldata.php>

other sensor systems, such as laser technology, could help to assign properties to asparagus. As mentioned in the **Discussion**, the most crucial step for archiving the goal to improve the asparagus sorting is the setup. We therefore have to adapt the preprocessing pipeline and fine-tune our approaches accordingly.

In order to put the results of this project into practice, a control system for the machine has to be implemented. For this purpose, a cooperation with the manufacturing company has to be arranged. After we developed results in the project, this can be the next step in the near future.

If our algorithm controls the sorting machine, it will be possible to create an evaluation. During our meetings we discussed existing difficulties of evaluation. One method is to measure and compare the sorting of the harvesters. Further methods need to be considered. If necessary, the setup for automated procedures can be extended.

Furthermore, it is possible to maintain data collection locally and use it for further improvement. An automated implementation, running locally on the computer of the machine is viable. Moreover locally generated networks could be collected and compared. Furthermore, merging then the different networks of different farms would be conceivable.

Another claim that we started to address in the **Methodology** is that the asparagus classification varies between farmers, as well as within one farm throughout the harvesting season. For example, upon request, asparagus should be placed in a higher category than normal after an unfavourable weather situation. In order to meet this requirement, manual adjustment of the screws and a smooth transition for several features is needed, which is impossible with most neural networks in their current form. Since we introduce a temporal, sequential component, it may be worthwhile to consider **LSTMs** in combination with **CNNs**. According to the literature there are still many wide possibilities, for example applying bayesian networks learning.

In conclusion, we can confirm that modern approaches from computer vision and machine learning bear potentials for the improvement of asparagus classification. Effective means of agile development allow for efficient collaboration in the production of the respective implementations. We demonstrated that the algorithms we selected could not only be used for scientific purposes but also in industrial applications. Modern machine learning approaches can help to improve solutions for the asparagus sorting problem.

## Appendix A

# Work Overview

### A.1 Task list

Below the name of each project member, the tasks that were contributed by the member are listed.

#### Maren Born

- ...

#### Michael Gerstenberger

- ...

#### Katharina Groß

- ...

#### Richard Ruppel

- ...

#### Sophia Schulze-Weddige

- ...

#### Malin Spaniol

- ...

#### Josefine Zerbe

- ...

### A.2 Report list

Here, an overview of the report is given in chapters and subchapters. The name behind each chapter indicates who wrote the respective section.

#### 1. Introduction – *Josefine*

- 1.1. The project – *Josefine*
- 1.2. Background on computer vision based classification tasks – *Sophia*
- 1.3. Background on sorting asparagus – *Josefine*

#### 2. Data acquisition and organization – *Josefine*

- 2.1. **Roadmap of the project** – *Josefine*
- 2.2. **Organisation of the study group** – *Richard*
  - 2.2.1. **Communication** – *Richard*
  - 2.2.2. **Teamwork** – *Richard*
- 2.3. **Data collection** – *Josefine, Richard*
- 2.4. **Literature research** – *Josefine*
3. **Preprocessing and data set creation** – *Josefine*
  - 3.1. **Preprocessing** – *Sophia*
  - 3.2. **Automatic feature extraction** – *Sophia*
    - 3.2.1. **Length** – *Sophia*
    - 3.2.2. **Width** – *Sophia*
    - 3.2.3. **Rust** – *Sophia*
    - 3.2.4. **Violet** – *Michael*
    - 3.2.5. **Curvature** – *Michael*
    - 3.2.6. **Flower** – *Sophia*
  - 3.3. **The hand-label app: A GUI for labeling asparagus** – *Michael*
  - 3.4. **Manual labeling** – *Josefine*
    - 3.4.1. **Sorting criteria** – *Josefine*
    - 3.4.2. **Sorting outcome** – *Josefine*
    - 3.4.3. **Agreement Measures** – *Malin*
    - 3.4.4. **Reliability** – *Malin*
  - 3.5. **The asparagus data set** – *Richard, Sophia*
4. **Classification** – *Malin*
  - 4.1. **Supervised learning** – *Josefine*
    - 4.1.1. **Prediction based on feature engineering** – *Michael*
    - 4.1.2. **Single-label classification** – *Josefine*
    - 4.1.3. **Multi-label classification** – *Sophia*
    - 4.1.4. **A dedicated network for head-related features** – *Michael*
    - 4.1.5. **From features to labels** – *Katharina*
  - 4.2. **Unsupervised learning** – *Malin*
    - 4.2.1. **Principal Component Analysis** – *Malin*
    - 4.2.2. **Autoencoder** – *Michael*
  - 4.3. **Semi-supervised learning** – *Michael*
5. **Summary of results** – *Maren*
6. **Discussion** – *Malin, Richard*
  - 6.1. **Classification results** – *Malin, Michael*
  - 6.2. **Methodology** – *Malin*
  - 6.3. **Organization** – *Josefine, Richard*
7. **Conclusion and outlook** – *Richard*

## Appendix B

# Additional figures



**FIGURE B.1: File Overview Example of Original Image Data** View of the data files after data collection. The original file names can be seen. The images belonging to one asparagus are marked with F00, F01, and F02 respectively.

## B.1 Single-label CNN



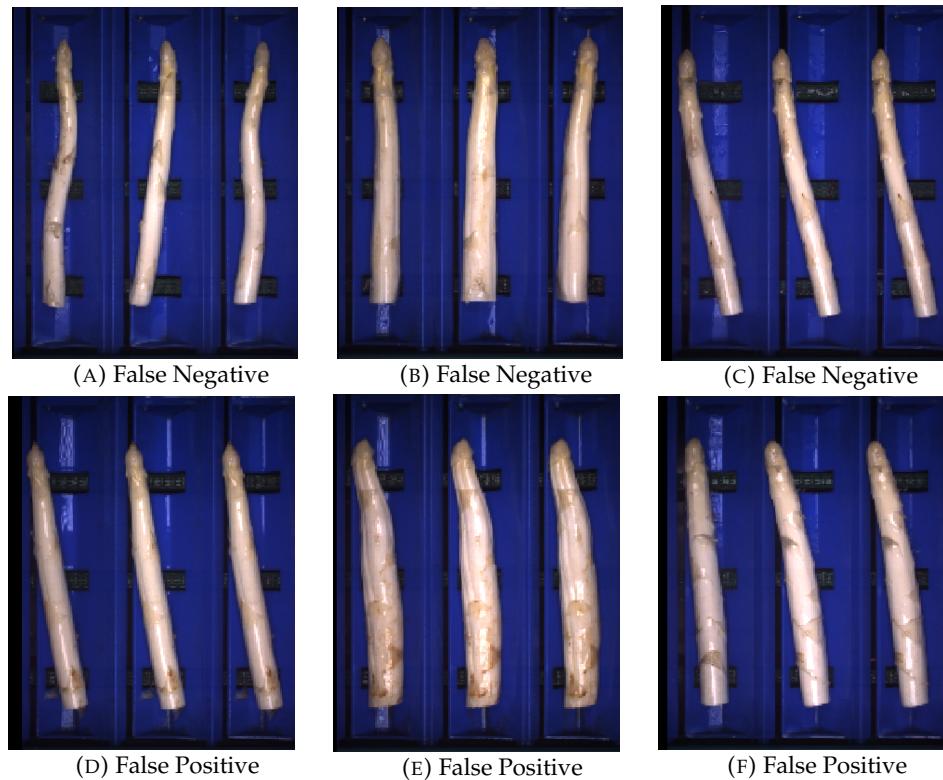
**FIGURE B.2: Feature Violet** Randomly chosen example images of false negatives and false positives of the feature violet after 5160 training steps. In the following, suggestions for the wrong classification are given which apply only to the depicted images. The false negative examples are all cases where the presence of feature violet is only slightly remarkable, with most prominence in sample (A). Looking closer at the false positives, the asparagus in image (D) and (E) might also be slightly rose around the upper part. In image (F), the dark color at the asparagus head is probably rust. The presence and absence of feature violet might have been difficult in cases where the feature was not clearly recognizable as violet but more of a pale rose color.



**FIGURE B.3: Single-Label CNN Example Images Feature Fractured** Example images of false negatives and false positives of the feature fractured after 2700 training steps.



**FIGURE B.4: Single-Label CNN Example Images Feature Flower** Example images of false negatives and false positives of the feature flower after 4800 training steps.



**FIGURE B.5: Single-Label CNN Example Images Feature Rusty Head** Example images of false negatives and false positives of the feature rusty head after 4680 training steps.



**FIGURE B.6: Single-Label CNN Example Images Feature Rusty Body** Example images of false negatives and false positives of the feature rusty body after 3000 training steps.

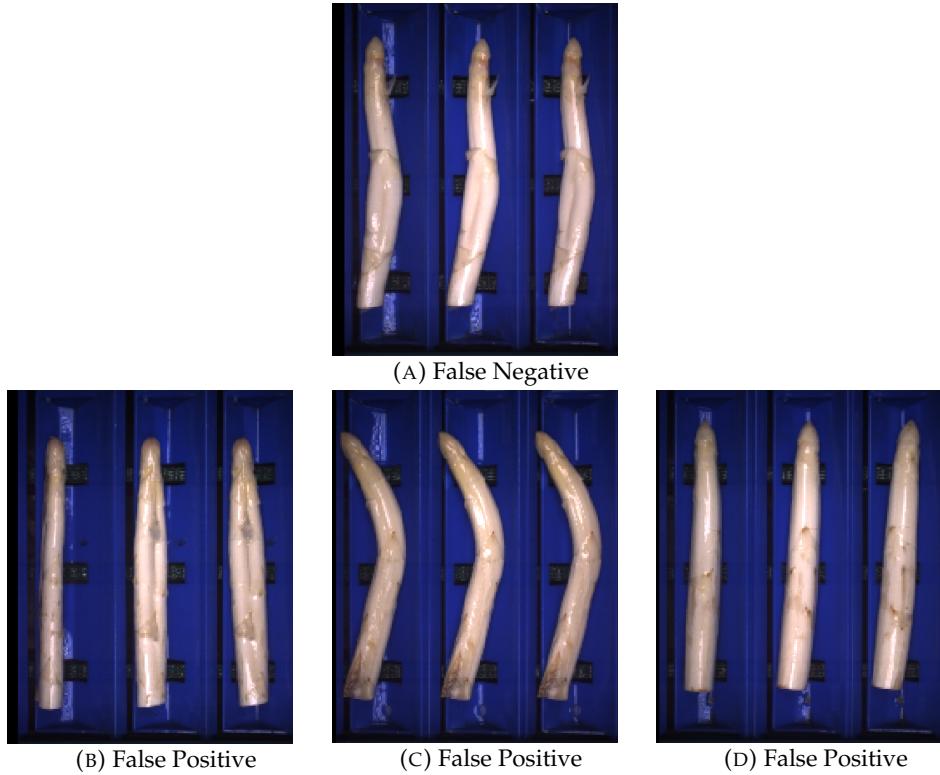


FIGURE B.7: **Single-Label CNN Example Images Feature Very Thick** Example images of false negatives and false positives of the feature very thick after 5400 training steps.

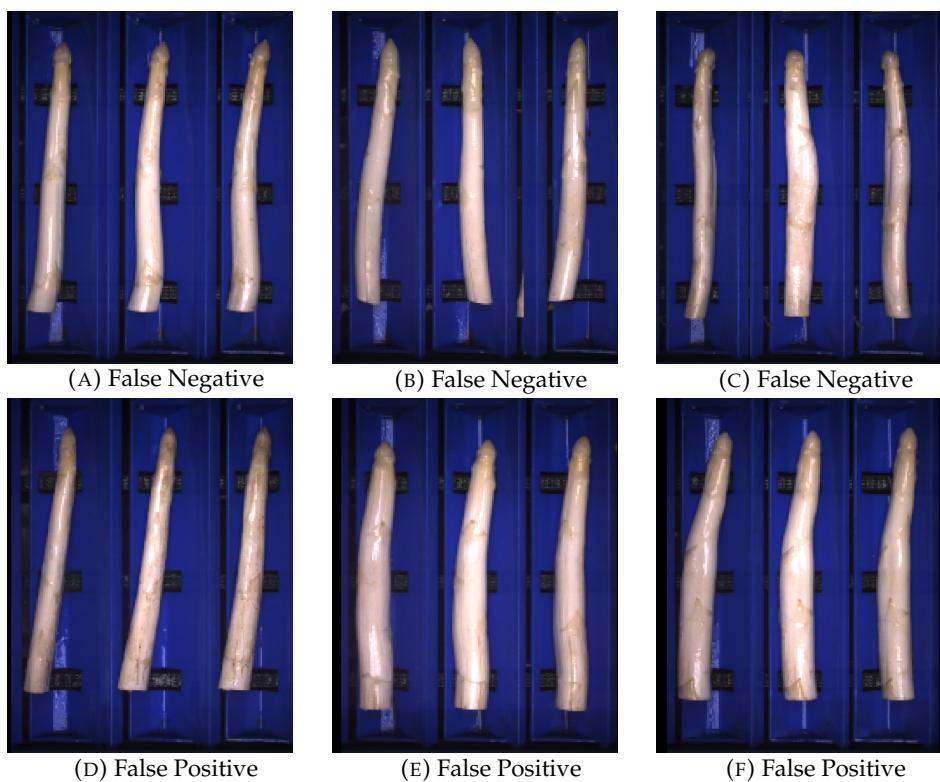
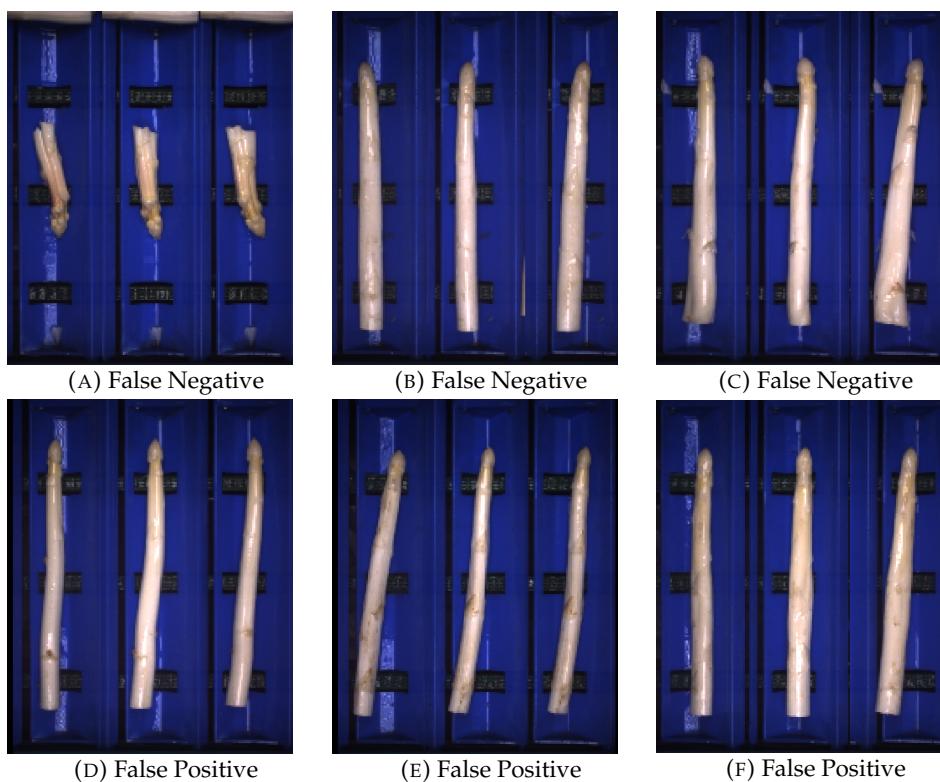


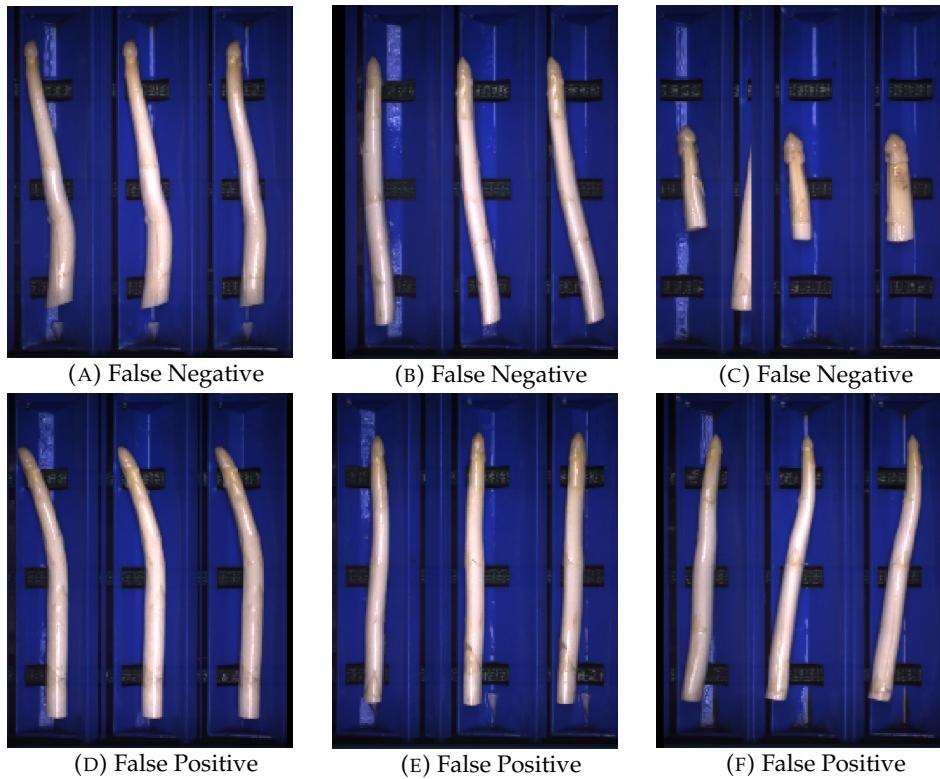
FIGURE B.8: **Single-Label CNN Example Images Feature Thick** Example images of false negatives and false positives of the feature thick after 3960 training steps.



**FIGURE B.9: Single-Label CNN Example Images Feature Medium Thick** Example images of false negatives and false positives of the feature medium thick after 4560 training steps.



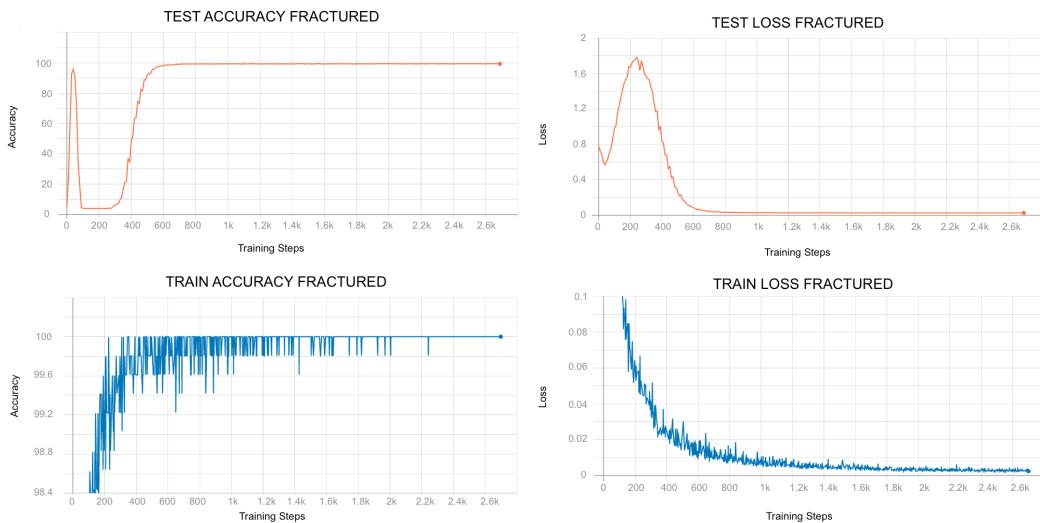
**FIGURE B.10: Single-Label CNN Example Images Feature Thin** Example images of false negatives and false positives of the feature thin after 4200 training steps.



**FIGURE B.11: Single-Label CNN Example Images Feature Very Thin** Example images of false negatives and false positives of the feature very thin after 4320 training steps.



**FIGURE B.12: Single-Label CNN Example Images Feature Not Classifiable** Example images of false negatives and false positives of the feature not classifiable after 5400 training steps.



**FIGURE B.13: Accuracy and Loss of Feature Fractured** Test accuracy and test loss as well as training accuracy and training loss for the feature fractured are depicted above. It was randomly chosen to show an exemplary course of accuracy and loss during training.



# Bibliography

- Al Ohali, Yousef (2011). "Computer vision based date fruit grading system: Design and implementation". In: *Journal of King Saud University-Computer and Information Sciences* 23.1, pp. 29–36.
- Al-Rawi, Mohammed and Dimosthenis Karatzas (2018). "On the Labeling Correctness in Computer Vision Datasets." In: *IAL@ PKDD/ECML*, pp. 1–23.
- Balaban, Stephen (2015). "Deep learning and face recognition: The state of the art". In: *Biometric and Surveillance Technology for Human and Activity Identification XII*. Vol. 9457. International Society for Optics and Photonics, 94570B.
- Bao, Jianmin et al. (2017). "CVAE-GAN: fine-grained image generation through asymmetric training". In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2745–2754.
- Batista, Gustavo EAPA, Ronaldo C Prati, and Maria Carolina Monard (2004). "A study of the behavior of several methods for balancing machine learning training data". In: *ACM SIGKDD explorations newsletter* 6.1, pp. 20–29.
- Bengio, Yoshua (2012). "Practical recommendations for gradient-based training of deep architectures". In: *Neural networks: Tricks of the trade*. Springer, pp. 437–478.
- Bettilyon, Tyler Elliot (2018). *How to classify MNIST digits with different neural network architectures*. URL: <https://medium.com/tebs-lab/how-to-classify-mnist-digits-with-different-neural-network-architectures-39c75a0f03e3> (visited on 04/24/2020).
- Bhargava, Anuja and Atul Bansal (2018). "Fruits and vegetables quality evaluation using computer vision: A review". In: *Journal of King Saud University-Computer and Information Sciences*.
- Bishop, Christopher M (2006). *Pattern recognition and machine learning*. Springer.
- Bohling, Geoff (2006). "Dimension reduction and cluster analysis". In: *Scientist, no. March*, pp. 1–22.
- Breiman, Leo (2001). "Random forests". In: *Machine learning* 45.1, pp. 5–32.
- Brosnan, Tadhg and Da-Wen Sun (2002). "Inspection and grading of agricultural and food products by computer vision systems—a review". In: *Computers and electronics in agriculture* 36.2-3, pp. 193–213.
- Brownlee, Jason (2019). *How to Configure Image Data Augmentation in Keras*. URL: <https://machinelearningmastery.com/how-to-configure-image-data-augmentation-when-training-deep-learning-neural-networks/> (visited on 04/24/2020).
- Bushaev, Vitaly (2018). *Adam — latest trends in deep learning optimization*. URL: <https://towardsdatascience.com/adam-latest-trends-in-deep-learning-optimization-6be9a291375c> (visited on 04/24/2020).

- Caruana, Rich and Alexandru Niculescu-Mizil (2006). "An Empirical Comparison of Supervised Learning Algorithms". In: *Proceedings of the 23rd International Conference on Machine Learning*. ICML '06. Pittsburgh, Pennsylvania, USA: Association for Computing Machinery, 161–168. ISBN: 1595933832. DOI: [10.1145/1143844.1143865](https://doi.org/10.1145/1143844.1143865). URL: <https://doi.org/10.1145/1143844.1143865>.
- Chen, Min et al. (2017). "Deep features learning for medical image analysis with convolutional autoencoder neural network". In: *IEEE Transactions on Big Data*.
- Cohen, Jacob (1960). "A coefficient of agreement for nominal scales". In: *Educational and psychological measurement* 20.1, pp. 37–46.
- Daumé III, Hal (2012). "A course in machine learning". In: *Publisher, ciml.info* 5, p. 69.
- Dertat, Arden (2017). *Applied Deep Learning - Part 1: Artificial Neural Networks*. URL: <https://towardsdatascience.com/applied-deep-learning-part-1-artificial-neural-networks-d7834f67a4f6> (visited on 04/24/2020).
- Diaz, R et al. (2004). "Comparison of three algorithms in the classification of table olives by means of computer vision". In: *Journal of Food Engineering* 61.1, pp. 101–107.
- Donis-González, Irwin R and Daniel E Guyer (2016). "Classification of processing asparagus sections using color images". In: *Computers and Electronics in Agriculture* 127, pp. 236–241.
- Europäische Komission (1999). "Verordnung (EG) Nr. 2377/1999 der Kommission vom 9. November 1999 zur Festsetzung der Vermarktungsnorm für Spargel". In: *Amtsblatt der Europäischen Gemeinschaften*, OJ L 287, 10.11.1999, p. 6–11. URL: <https://eur-lex.europa.eu/legal-content/DE/TXT/PDF/?uri=CELEX:31999R2377&from=DE> (visited on 03/25/2020).
- Feinstein, Alvan R and Domenic V Cicchetti (1990). "High agreement but low kappa: I. The problems of two paradoxes". In: *Journal of clinical epidemiology* 43.6, pp. 543–549.
- Figueroa, Rosa L et al. (2012). "Predicting sample size required for classification performance". In: *BMC medical informatics and decision making* 12.1, p. 8.
- Flight, Laura (2018). *The Disagreeable Behaviour of the Kappa Statistic*. URL: [https://www.sheffield.ac.uk/polopoly\\_fs/1.404095!/file/RSS\\_Poster\\_Laura\\_Flight\\_Final.pdf](https://www.sheffield.ac.uk/polopoly_fs/1.404095!/file/RSS_Poster_Laura_Flight_Final.pdf) (visited on 04/21/2020).
- Franky, Frank (2018). *Multi-Label Classification and Class Activation Map on Fashion-MNIST*. URL: <https://towardsdatascience.com/multi-label-classification-and-class-activation-map-on-fashion-mnist-1454f09f5925> (visited on 04/21/2020).
- Geng, Jie et al. (2015). "High-resolution SAR image classification via deep convolutional autoencoders". In: *IEEE Geoscience and Remote Sensing Letters* 12.11, pp. 2351–2355.
- Géron, Aurélien (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly Media.
- Gibaja, Eva and Sebastian Ventura (Apr. 2015). "A Tutorial on Multi-Label Learning". In: *ACM Computing Surveys* 47. DOI: [10.1145/2716262](https://doi.org/10.1145/2716262).

- Gilpin, Leilani H et al. (2018). "Explaining explanations: An overview of interpretability of machine learning". In: *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*. IEEE, pp. 80–89.
- Godoy, Daniel (2018). *Understanding binary cross-entropy / log loss: a visual explanation*. URL: <https://towardsdatascience.com/understanding-binary-cross-entropy-log-loss-a-visual-explanation-a3ac6025181a> (visited on 04/24/2020).
- Gupta, Vikas (2018). *Keras Tutorial Transfer Learning using pre-trained models*. URL: <https://www.learnopencv.com/keras-tutorial-transfer-learning-using-pre-trained-models> (visited on 01/01/2020).
- Har-Peled, Sariel, Dan Roth, and Dav Zimak (2003). "Constraint classification for multiclass classification and ranking". In: *Advances in neural information processing systems*, pp. 809–816.
- He, Haibo and Edwardo A Garcia (2009). "Learning from imbalanced data". In: *IEEE Transactions on knowledge and data engineering* 21.9, pp. 1263–1284.
- He, Kaiming et al. (2016a). "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- (2016b). "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- Heaton, Jeff (2015). *AIFH, Volume 3: Deep Learning and Neural Networks*.
- HMF Hermeler Maschinenbau GmbH (2015). "Bedienungsanleitung Automatische Spargelsortiermaschine Autoselect". In: HMF Hermeler Maschinenbau GmbH.
- Huang, Gao et al. (2017). "Densely connected convolutional networks". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708.
- Hubel, David H and Torsten N Wiesel (1962). "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex". In: *The Journal of physiology* 160.1, pp. 106–154.
- Iglovikov, Vladimir and Alexey Shvets (2018). "Ternausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation". In: *arXiv preprint arXiv:1801.05746*.
- Karpathy, Andrej (2014). *What I learned from competing against a ConvNet on ImageNet*. URL: <http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/> (visited on 04/21/2020).
- Keng, Brian (2020). *Semi-supervised Learning with Variational Autoencoders*. URL: <http://bjlkeng.github.io/posts/semi-supervised-learning-with-variational-autoencoders/> (visited on 03/01/2020).
- Kılıç, Kivanç et al. (2007). "A classification system for beans using computer vision system and artificial neural networks". In: *Journal of Food Engineering* 78.3, pp. 897–904.
- Kingma, Diederik P and Jimmy Ba (2014). "Adam: A method for stochastic optimization". In: *arXiv preprint arXiv:1412.6980*.
- Kramer, Mark A (1991). "Nonlinear principal component analysis using autoassociative neural networks". In: *AIChE journal* 37.2, pp. 233–243.

- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton (2012a). "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*, pp. 1097–1105.
- (2012b). "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*, pp. 1097–1105.
- Lata, Prof et al. (Apr. 2009). "Facial recognition using eigenfaces by PCA". In: *SHORT PAPER International Journal of Recent Trends in Engineering* 1.
- LeCun, Yann, Yoshua Bengio, et al. (1995). "Convolutional networks for images, speech, and time series". In: *The handbook of brain theory and neural networks* 3361.10, p. 1995.
- LeCun, Yann A et al. (2012). "Efficient backprop". In: *Neural networks: Tricks of the trade*. Springer, pp. 9–48.
- Lin, Min, Qiang Chen, and Shuicheng Yan (2013). "Network in network". In: *arXiv preprint arXiv:1312.4400*.
- Luo, M Ronnier et al. (2000). "A review of chromatic adaptation transforms". In: *Review of Progress in Coloration and Related Topics* 30, pp. 77–92.
- McHugh, Mary (2012). "Interrater reliability: The kappa statistic". In: *Biochimia medica : časopis Hrvatskoga društva medicinskih biokemičara / HDMB* 22, pp. 276–82. DOI: [10.11613/BM.2012.031](https://doi.org/10.11613/BM.2012.031).
- Mery, Domingo, Franco Pedreschi, and Alvaro Soto (2013). "Automated design of a computer vision system for visual food quality evaluation". In: *Food and Bio-process Technology* 6.8, pp. 2093–2108.
- Mirończuk, Marcin Michał and Jarosław Protasiewicz (2018). "A recent overview of the state-of-the-art elements of text classification". In: *Expert Systems with Applications* 106, pp. 36–54.
- Murray, Derek (2019). *Introduction to td.data*. URL: [https://docs.google.com/presentation/d/16kHNtQslt-yuJ3w8GIx-eEH6t\\_AvFeQ0chqGRFpAD7U/edit#slide=id.g254d08e080\\_0\\_67](https://docs.google.com/presentation/d/16kHNtQslt-yuJ3w8GIx-eEH6t_AvFeQ0chqGRFpAD7U/edit#slide=id.g254d08e080_0_67) (visited on 04/21/2020).
- Olaode, Abass, Golshah Naghd, and Catherine Todd (Sept. 2014). "Unsupervised Classification of Images: A Review". In: *International Journal of Image Processing* 8, pp. 2014–325.
- Olivier, Chapelle, Scholkopf Bernhard, and Zien Alexander (2006). "Semi-supervised learning". In: *IEEE Transactions on Neural Networks*. Vol. 20. 3, pp. 542–542.
- Pascanu, Razvan, Tomas Mikolov, and Yoshua Bengio (2012). "Understanding the exploding gradient problem". In: *CoRR, abs/1211.5063*, 2, p. 417.
- Pedreschi, Franco, Domingo Mery, and Thierry Marique (2016). "Grading of potatoes". In: *Computer vision technology for food quality evaluation*. Elsevier, pp. 369–382.
- Powers, David MW (2012). "The problem with kappa". In: *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, pp. 345–355.
- Prakash, J Suriya et al. (2012). "Multi class Support Vector Machines classifier for machine vision application". In: *2012 International Conference on Machine Vision and Image Processing (MVIP)*. IEEE, pp. 197–199.

- Python Image Library (2020). *Image quantization used by the Python Image Library (PIL)*. URL: <https://github.com/python-pillow/Pillow/blob/55e5b7de6c41b0386660b0bee7784ac04f412f4b/src/libImaging/Quant.c> (visited on 03/25/2020).
- Reeves, Adam (1981). "Metacontrast in hue substitution". In: *Vision Research* 21.6, pp. 907–912.
- Russakovsky, Olga and Li Fei-Fei (2010). "Attribute learning in large-scale datasets". In: *European Conference on Computer Vision*. Springer, pp. 1–14.
- Russakovsky, Olga et al. (2013). "Detecting avocados to zucchinis: what have we done, and where are we going?" In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2064–2071.
- Russakovsky, Olga et al. (2015). "Imagenet large scale visual recognition challenge". In: *International journal of computer vision* 115.3, pp. 211–252.
- Sabour, Sara, Nicholas Frosst, and Geoffrey E Hinton (2017). "Dynamic routing between capsules". In: *Advances in neural information processing systems*, pp. 3856–3866.
- Scikit-Learn (2019). *Neural network models (supervised)*. URL: [https://scikit-learn.org/stable/modules/neural\\_networks\\_supervised.html](https://scikit-learn.org/stable/modules/neural_networks_supervised.html) (visited on 03/01/2020).
- (2020). *Example of VAE on MNIST dataset using CNN*. URL: [https://keras.io/examples/variational\\_autoencoder\\_deconv/](https://keras.io/examples/variational_autoencoder_deconv/) (visited on 03/01/2020).
- Shlens, Jonathon (2014). "A tutorial on principal component analysis". In: *arXiv preprint arXiv:1404.1100*.
- Sim, Julius and Chris C Wright (2005). "The kappa statistic in reliability studies: use, interpretation, and sample size requirements". In: *Physical therapy* 85.3, pp. 257–268.
- Simonyan, Karen and Andrew Zisserman (2014). "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556*.
- Stallkamp, Johannes et al. (2011). "The German traffic sign recognition benchmark: a multi-class classification competition". In: *The 2011 international joint conference on neural networks*. IEEE, pp. 1453–1460.
- Statistisches Bundesamt (Destatis) (2017). *Gemüseerhebung 2016 - Anbau und Ernte von Gemüse und Erdbeeren*. URL: [https://www.destatis.de/GPStatistik/receive/DEHeft\\_heft\\_00070882](https://www.destatis.de/GPStatistik/receive/DEHeft_heft_00070882) (visited on 03/25/2020).
- Szegedy, Christian et al. (2015). "Going deeper with convolutions". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9.
- Ting, Marie (2015). "Eine Region im neuen Zeitalter". In: *Frankfurter Allgemeine Zeitung - Verlagsspezial*. URL: [https://www.hering-international.com/fileadmin/media/archive1/downloads/presseberichte/2015-09-23\\_FAZ\\_-\\_Verlagsspezial\\_Industrierregion\\_S%C3%BCdwestfalen\\_-\\_Textilbeton\\_f%C3%A4r\\_BCr\\_Dior.pdf](https://www.hering-international.com/fileadmin/media/archive1/downloads/presseberichte/2015-09-23_FAZ_-_Verlagsspezial_Industrierregion_S%C3%BCdwestfalen_-_Textilbeton_f%C3%A4r_BCr_Dior.pdf) (visited on 03/16/2020).
- Tjoa, Erico and Cuntai Guan (2019). "A survey on explainable artificial intelligence (xai): Towards medical xai". In: *arXiv preprint arXiv:1907.07374*.

- Turk, Matthew and Alex Pentland (1991). "Face recognition using eigenfaces". In: *Proceedings. 1991 IEEE computer society conference on computer vision and pattern recognition*, pp. 586–587.
- ul Hassan, Muneeb (2018a). *AlexNet – ImageNet Classification with Deep Convolutional Neural Networks*. URL: <https://neurohive.io/en/popular-networks/resnet/> (visited on 04/24/2020).
- (2018b). *ResNet (34, 50, 101): Residual CNNs for Image Classification Tasks*. URL: <https://neurohive.io/en/popular-networks/alexnet-imagenet-classification-with-deep-convolutional-neural-networks/> (visited on 04/24/2020).
- (2018c). *VGG16 – Convolutional Network for Classification and Detection*. URL: <https://neurohive.io/en/popular-networks/vgg16/> (visited on 04/24/2020).
- United Nations Economic Commission for Europe (UNECE) (2017). UNECE STANDARD FFV-04 concerning the marketing and commercial quality control of asparagus. URL: <https://www.unece.org/trade/agr/standard/fresh/ffv-standardse.html> (visited on 03/25/2020).
- Universität Osnabrück, Fachbereich Humanwissenschaften (2019a). *Modulbeschreibungen Cognitive Science*. URL: [https://www.uni-osnabrueck.de/studium/im\\_studium/zugangs\\_zulassungs\\_und\\_pruefungsordnungen/fach\\_master/cognitive\\_science\\_msc.html](https://www.uni-osnabrueck.de/studium/im_studium/zugangs_zulassungs_und_pruefungsordnungen/fach_master/cognitive_science_msc.html) (visited on 03/21/2020).
- (2019b). *Prüfungsordnung für den Master Studiengang Cognitive Science*. URL: [https://www.uni-osnabrueck.de/studium/im\\_studium/zugangs\\_zulassungs\\_und\\_pruefungsordnungen/fach\\_master/cognitive\\_science\\_msc.html](https://www.uni-osnabrueck.de/studium/im_studium/zugangs_zulassungs_und_pruefungsordnungen/fach_master/cognitive_science_msc.html) (visited on 03/21/2020).
- Wang, Chi-Feng (2019). *The Vanishing Gradient Problem The Problem, Its Causes, Its Significance, and Its Solutions*. URL: <https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484> (visited on 03/01/2020).
- Warden, Pete (2017). *How many images do you need to train a neural network?* URL: <https://petewarden.com/2017/12/14/how-many-images-do-you-need-to-train-a-neural-network/> (visited on 03/25/2020).
- Weber, Karola and Katrin Quinckhardt (2018). "Heimvorteil Spargel - Selbst angebaut, selbst zubereitet! Anbautipps und Rezepte für Spargel". In: Landwirtschaftskammer Nordrhein-Westfalen. URL: <https://www.landwirtschaftskammer.de/verbraucher/rezepte/spargelrezepte.pdf> (visited on 03/25/2020).
- Wikipedia contributors (2020). *Precision and recall — Wikipedia, The Free Encyclopedia*. [Online; accessed 21-April-2020]. URL: [https://en.wikipedia.org/w/index.php?title=Precision\\_and\\_recall&oldid=945184172](https://en.wikipedia.org/w/index.php?title=Precision_and_recall&oldid=945184172).
- Zarandi, MH Fazel, Soheil Davari, and SA Haddad Sisakht (2011). "The large scale maximal covering location problem". In: *Scientia Iranica* 18.6, pp. 1564–1570.
- Zhang, Xiangyu et al. (2015). "Accelerating very deep convolutional networks for classification and detection". In: *IEEE transactions on pattern analysis and machine intelligence* 38.10, pp. 1943–1955.
- Zhang, Yudong and Lenan Wu (2012). "Classification of fruits using computer vision and a multiclass support vector machine". In: *sensors* 12.9, pp. 12489–12505.
- Zheng, Alice and Amanda Casari (2018). *Feature engineering for machine learning: principles and techniques for data scientists*. " O'Reilly Media, Inc."

- Zhu, Xiaojin (2005). *Semi-Supervised Learning Literature Survey*. Tech. rep. 1530. Computer Sciences, University of Wisconsin-Madison. URL: [http://pages.cs.wisc.edu/~jerryzhu/pub/ssl\\_survey.pdf](http://pages.cs.wisc.edu/~jerryzhu/pub/ssl_survey.pdf) (visited on 02/16/2020).