

IDENTIFY THE PLANT AND DISEASE BASED ON THE LEAF PROVIDED

Data Mining and Machine Learning II

Pedro Acosta
School of Computing
National College of Ireland
Dublin, Ireland
x21138745@student.ncirl.ie

Brendan O'Dwyer
School of Computing
National College of Ireland
Dublin, Ireland
x21145172@student.ncirl.ie

Femi Adeboye
School of Computing
National College of Ireland
Dublin, Ireland
x21137684@student.ncirl.ie

Antonio Milian y Albacar
School of Computing
National College of Ireland
Dublin, Ireland
x19172125@student.ncirl.ie

Abstract—Plant disease identification, techniques results.
Index Terms—Convolutional Neural Networks

I. INTRODUCTION

Motivation

Emerging technologies and diverse applications make it possible nowadays to prevent or early detect challenging situations where consequences impact not only financially but in the quality of life of people. Nowadays geopolitical conflicts urge all stakeholders in society to reduce errors and waste of resources, especially on food, the current project aims to contribute and make possible the early detection through images of the type of plant and serves to early detect diseases on them. From small farmers that use their mobile phones to industrialized companies taking images with drones, the main objective is to make available the results and techniques from this study for them to be integrated or embedded as part of a simple, complex or wider solution that includes not only indoor analysis but also outdoor data.

Discussion Research questions and objectives

Since one of the core objectives of the project is to provide a reliable model that identifies plants and diseases based on outdoor images, this adds a level of complexity due to the nature of the situation, noise or information factors such as illumination, background, disease signs distributed randomly along with different type of plants around, all contribute to a right application of XXXXXXXXXXXXXXXX

Overview following sections XXXXXXXXXXXXX

Identifying plant diseases at an early stage allows home gardeners, and rural and industrial farmers to take action and reduce crop loss. In the same way, the usage of basic technology such as smart phones as well as more sophisticated and emerging technologies such as drones are leading to consideration of alternatives that could allow stakeholders to identify plant diseases using the technologies already mentioned. This project will allow for the identification of the type of plant between apple, tomato and pepper detecting if they are healthy or carrying some type of disease. This model could

be used as part of a system where images can be the input either by uploading them via PC, or using mobile phones, it could also receive images from drones or pictures taken in laboratories. The results could be used as a part of a more complex system that applies pesticides wherever a crop disease is identified.

In this study we train convolutional neural network models using images of leaves of diseased plants to identify the plant disease. The research question we are attempting to answer in this study is: how do CNN models trained on images of leaves of diseased plants in order to classify plant disease compare when trained using laboratory data, outdoor data and both combined and tested on laboratory and outdoor data? There are two hypotheses in this study:

- The performance of a CNN model trained with both laboratory and augmented outdoor data should perform substantially better than a model trained with both laboratory and outdoor data without augmentation when tested on laboratory and outdoor data.
- The performance of a CNN model trained with a combination of laboratory and augmented outdoor data should perform substantially better than a model trained with laboratory data alone when tested on outdoor data.

This report is structured as follows: the Related Work section provides a review of related literature. The Methodology section describes our methodology. The Results and Evaluation section presents our results and evaluates our method. The Conclusions and Future Work section discusses our conclusions and proposed future work. The References section provides the references in this report.

II. RELATED WORK

1 page with citations into 20 or more works

This should not only summarise related work, but also critically evaluate the strengths and weaknesses of the cited works with respect to the topic under study, i.e. how well/badly does the cited work answer your question(s), what aspects are useful to consider, what are the limitations? You should also discuss any foundational papers that either substantiate

your study design or upon which you are building. ne part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations.

III. METHODOLOGY

The CRoss Industry Standard Process for Data Mining (CRISP-DM) data mining methodology has been followed in this study. The data mining lifecycle is described by the six phases of this methodology. It is a flexible and iterative process. Phases can be revisited as the project progresses. This

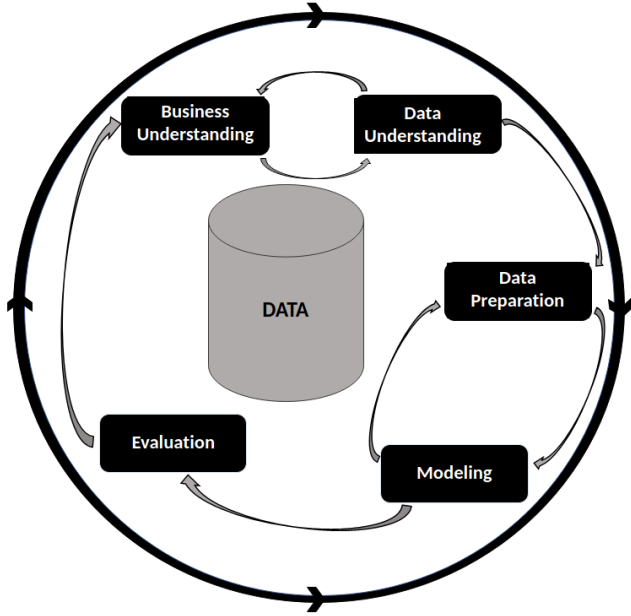


Fig. 1. The phases of the CRISP-DM cycle. Business understanding, data understanding, data preparation, modelling and evaluation in this case there is not deployment phase as the evaluation is the output.

study will involve training different models and as a result phases will be revisited repeatedly. Business understanding is the first phase in CRISP-DM. The data mining goals are determined in this phase and an understanding of the objectives and requirements of the project developed. In the introduction the objectives of the study have been outlined. The business understanding phase is addressed through these objectives as well as through the understanding developed from the related work publications. In the remainder of this report a description of the other CRISP-DM phases is provided

A. Data Understanding

Exploration and collection of data is carried out in the CRISP-DM data understanding phase. An understanding of the data is obtained as well as an understanding of how the objectives can be addressed through the data.

Two datasets are used in this study. The first, PlantVillage, contains images of leaves of diseased and healthy plants with the associated labels. The images are taken indoors in laboratory conditions. The data are made available through the PlantVillage online platform. In total there are over 50,000 images of crop plant leaves that are infected or healthy. There

are 14 different types of plant species with 38 classes in total. Each class corresponds to a plant-disease pair.

The second dataset, Plant Doc, is for detecting plant disease from images. Unlike the first dataset the images are taken outdoors rather than in laboratory conditions. There are 2,598 images in total. There are 13 species of plants in total with 17 disease classes. The data were obtained by scraping images from webpages and annotating them.

The PlantVillage dataset contains coloured, segmented and greyscale versions of each image. The coloured versions of the images were used. In the case of Plant Doc only colour images are provided. For the 13 classes used in this study in total 18,317 images from the PlantVillage dataset were used.

B. Data Preparation

This phase of CRISP-DM involves manipulating data acquired in the data understanding phase in order to prepare it for the modelling phase. The datasets were explored using Python. The feature selection is based on the commonalities between the previously described datasets. Three plant species present in both the PlantVillage and Plant Doc datasets were selected apple, pepper and tomato. In total there are 13 disease classes for these plants common to both datasets. Below the classes and the number of items per class and dataset.

TABLE I
DATA SETS USED

| Classes Objective | Data sets | | | |
|----------------------|-----------|---------|-------------|----------|
| | Indoor | Outdoor | O.Augmented | Combined |
| apple scab | 630 | 87 | 1307 | 1937 |
| apple healthy | 275 | 89 | 1392 | 1667 |
| apple rust | 1645 | 91 | 1325 | 2970 |
| bell pepper healthy | 997 | 71 | 1355 | 2352 |
| bell pepper spot | 1478 | 61 | 1581 | 3059 |
| tomato early blight | 2127 | 107 | 1256 | 3383 |
| tomato septoria | 1591 | 62 | 1392 | 2983 |
| tomato healthy | 1000 | 83 | 1379 | 2379 |
| tomato bacterial | 1771 | 148 | 1272 | 3043 |
| tomato late blight | 1909 | 111 | 1306 | 3215 |
| tomato mosaic | 373 | 54 | 1131 | 1504 |
| tomato yellow virus | 5357 | 75 | 1172 | 6529 |
| tomato mold | 952 | 91 | 1444 | 2396 |

^aOriginal Data set classes and values.

Two test datasets were created. The first uses only images from the dataset of indoor laboratory data. The second consists of images taken outdoors from the Plant Doc dataset. For each of the PlantVillage and Plant Doc datasets the remaining data are assigned randomly to training datasets, validation datasets and also datasets used for tuning hyperparamteters. 72All images were adjusted to have the dimensions 240 pixels by 240 pixels, when using only PlanDoc data set the script complains about the size and resizes it to 224 x224. The PlantVillage dataset has more images for each class than the Plant Doc dataset. As a result of this data augmentation techniques were applied to the Plant Doc dataset. These techniques included rotation, shearing, zooming, flipping horizontally and adjusting the brightness range. The data were preprocessed. The pixel values are rescaled.

C. Modelling

The model typology to be used to analyse the data is CNN, to take advantage of the properties of the models in image recognition. Modelling techniques are used to build and assess models in the CRISP-DM modelling phase. A batch size of 32 was used. As convolutional layers two types are taken into consideration, Inception ResNet, in this case the pretrained convolutional neural network Inception-ResNet-v2 was used for transfer learning. Weights trained on the ImageNet dataset, which has one thousand classes and 1.4 million images, were used. On top of this model a new classifier is added. The layers of the base convolutional model are frozen. The second type is MobileNet. In addition to Inception-ResNet-v2 the pretrained convolutional neural network MobileNet V2 was also used for transfer learning. The data were again preprocessed in preparation for use with the MobileNet V2 model. Weights trained on the ImageNet dataset were used. As before a new classifier was added on top of this model and the layers of the base convolutional model were frozen.

1) *Hyperparameter Optimisation:* Adding to the convolutional layers described the optimizers Adam and Adagrad were also tested. Combining the different values of convolutional layers and optimizers several models are proposed as potential candidates. The models are passed to a build model function in which several parameters are proposed for the rate of drop on the and the learning rate. A few values are provided for both variables. The optimization is centered in finding the best values for drop rate and from the learning rate. The model with the Convolution layer is frozen and a range of values proposed for both variables. The values for learning rates ranging from 0.01 to 0.00001, and the drop rates used range from 0.1 to 0.2. From the algorithm point of view, the library keras tuner has been used to ascertain the optimal values. From keras tuner, the Hyperband tuning algorithm is chosen. The algorithm has the championship bracket approach, running the models and choosing only the half best performers to be carried onto the next epoch. This approach is faster and guarantees better results than a simple Grid optimization approach. The keras tuner implementation has been run for 6 epochs in each trial, obtaining the optimal learning rate for each model as well as the optimal number of epochs. Accuracy over the validation is the metric used to find the best parameter values. The goal is to decide which model's architecture render the best validation accuracy and continue with that model through the fine-tuning process. Same methodology has been used for the three different implementations; the values gathered have subsequently used in the tested models. The three models were trained using the three different data sets, Indoors, Outdoor and Combined. The first model is using the Outdoor data set, with the augmentation images for training. Applying the Hyperparameter optimization over the datasets the below results are obtained, "

2) *Indoor data model:* Applying the Hyperparameter tuning to the indoor data set, the model that offers the best results has a convolutional layer based on MobileNet v2,

which gives not only better accuracy but is also 4 times faster than the ResNet v2 alternative. The optimizer that performs better is The images are transformed into 224 x 224 pixels, grouped into batches of 32, the convolutional method used is MobileNet at 7 x 7, it renders a total of x parameters, it is passed to a drop layer of 0.02. The dense layer has a $(13 \times 1280) + 13$ trainable parameters. The performance of models trained with different optimisers was compared. For the PlantVillage indoor laboratory dataset a model was trained using transfer learning with Inception-ResNet-v2 and the Ada-grad optimiser. The model was trained for ten epochs before fine tuning. The test accuracy obtained was 97.75%. The same optimiser and loss function were used as for the model using Inception-ResNet-v2. In total our model contained 2,306,662 parameters of which 48,678 were trainable. The base model contains 154 layers. The layers from 100 upwards were used in fine-tuning the model. For this model the total number of parameters was again 2,306,662, but in this case 1,910,118 of these were trainable. The models were again evaluated on the test dataset for the indoor laboratory data and also on the test dataset for the outdoor data. The pretrained model's weights are not updated. The top layers of the model are unfrozen when fine-tuning and trained along with the new classifier layers that have been added. The softmax activation function is used in the prediction layer. The Adam optimiser is used to train the model. Hyperparameter tuning is used to determine the optimal learning rate. The categorical crossentropy loss function was used. In total our model contained 54,356,717 parameters of which 19,981 were trainable. Models were fit using the training and validation datasets. For training the models the number of epochs which were optimal were determined. The training accuracy and loss and validation accuracy and loss were determined. Metrics for precision, recall and auc were calculated. Our base model contains 780 layers. The layers from 700 upwards were used in fine-tuning the model. For this model the total number of parameters was again 54,356,717, but in this case 12,989,645 of these were trainable. The models were evaluated on the test dataset for the indoor laboratory data and also on the test dataset for the outdoor data.

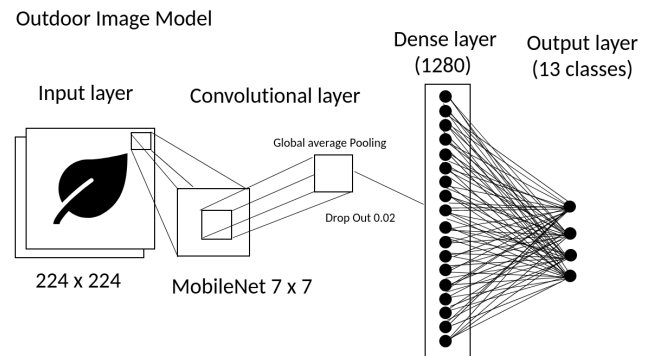


Fig. 2. MobileNet Indoor model schema.

3) *Outdoor data model:* A convolutional neural network model was trained on the PlantVillage indoor laboratory

images and tested on both the test dataset of indoor images and the test dataset of outdoor images. Another CNN model was trained on the Plant Doc outdoor images and tested on both the indoor and outdoor test datasets. The indoor laboratory images and outdoor images were merged into a combined training dataset and a model was trained on this dataset and tested on both the indoor and outdoor test datasets. An additional model was trained on the augmented Plant Doc outdoor images dataset and tested on both test datasets. The augmented outdoor images dataset was merged with the indoor laboratory images dataset and a model trained on the combined dataset and tested on both test datasets.

4) *Combined data model:*

IV. AUTHORS AND AFFILIATIONS

The class file is designed for, but not limited to, six authors. A minimum of one author is required for all conference articles. Author names should be listed starting from left to right and then moving down to the next line. This is the author sequence that will be used in future citations and by indexing services. Names should not be listed in columns nor group by affiliation. Please keep your affiliations as succinct as possible (for example, do not differentiate among departments of the same organization).

V. RESULTS AND EVALUATION

Here you should discuss how you used your methodology to answer the question. How can you determine if your approach is good? If you have to parameterise part of your approach, how have you done that, and why were these choices made? What impacts can different parameterisations have on your results? You should also discuss the results in detail and discuss their impact and implications. What do they show or not show? ". How effective the models are in achieving the business objectives is assessed in the CRISP-DM evaluation phase. The models to evaluate should be Indoors images, Outdoor images with augmented train data and tested on the test data, Indoor and outdoor data combined, model to be tested on both random test subset and on the test outdoor data. Compared results will be then highlighted in the conclusion phase.

Also in this step a decision is made to either proceed to deployment or iterate further.

VI. CONCLUSION AND FUTURE WORK

The CRISP-DM final phase is deployment where the results of the data mining study are communicated to end users. This purpose is served by this report. The objectives of this study were met through the results. Knowledge discovery is the aim of CRISP-DM. This study accomplished this aim in a variety of ways. The lack of availability of images of leaves of diseased plants obtained outdoors in real world environments limited this study. Larger datasets of outdoor images would be of benefit for leaf disease classification. Also in this study only three plant species were considered. In future work classification of plant disease from images could

be extended to include additional plant species. This study could be extended in a number of ways including through transfer learning using other pretrained CNN models to carry out image classification. In addition to the data augmentation techniques applied in this study other techniques could also be used including creating synthetic images of plant leaves using conditional generative adversarial networks. The Plant Doc outdoor images of leaves could be segmented and models, trained using these segmented images, compared against those trained using coloured versions of those images. We could extend this study from classifying the disease from the image of the leaf to grading how diseased the leaf is. In addition we could compare the classification performance of the CNN models with that of other machine learning techniques. Our model could also be included as part of a pesticide prescription system for plant diseases. We could also compare the performance of a CNN model trained from scratch for classifying plant disease from leaf images against that of a CNN model trained by transfer learning. Our model could be used as part of a system to classify plant disease from images obtained using unmanned aerial vehicles. We could develop a lightweight model for devices with limited resources that still provides accurate performance in classifying plant disease from leaf images. The effect of varying the number of images used to train the models on the accuracy of classification could be investigated in future work.