



VRIJE  
UNIVERSITEIT  
BRUSSEL



# ANIMAL CLASSIFICATION AI

Machine Learning

Lennert Bontinck

December 11, 2020

Master Computer Science: AI  
Sciences and Bio-Engineering Sciences

# Abstract

This intermediate report documents the development of an animal classification AI using a more "old-school" approach of Visual-Bag-of-Words models. The AI should be capable of differentiating 12 different animals. These models, and thus the AI, are developed in Python-based Jupyter Notebooks accompanied by this document. This animal classification AI was developed as a fulfilment of the Machine Learning course requirements and was used to compete in the organised Kaggle competition (Rosseau, 2020).

Part I of this report discusses the accompanied code in general. Section 1.1 explains which files are accompanied and which are the most important. Section 1.2 describes the ideology used to create the code. To make testing multiple models easier, a form of *pipeline* was created and is discussed in section 1.3.

In part II, the data analysis part of this project is discussed. Section 2.2 talks about the unbalanced data distribution. In the next section, section 2.3, a deeper look is taken into the data and possible preprocessing is discussed. The last 2 sections of this part, 2.4 and 2.5, discuss how the feature extraction is dealt with and what the numerical representation looks like.

The linear baseline model is discussed in part III. This model is a fine-tuned Logistic Regression model from the SciKit Learn library. This model is often used to compare other models with. Only models that perform better than this baseline model should be considered. This part discusses the parameters used and the road to getting those optimal parameters.

Finally, part IV discusses future plans for this project. Section 4.1 lists possible topics that can be explored to create a better animal classification AI. In the last section, section 4.2, some open issues are discussed.

# Contents

<b>I</b>	<b>About the code</b>	<b>1</b>
1.1	Files accompanied by this report . . . . .	2
1.2	Ideology of the developed code . . . . .	2
1.3	A typical model exploration . . . . .	2
1.4	Technical remarks . . . . .	2
<b>II</b>	<b>Data analysis</b>	<b>3</b>
2.1	About this part . . . . .	4
2.2	Data distribution . . . . .	4
2.3	Deeper look at the training data . . . . .	5
2.4	Feature extraction . . . . .	5
2.5	The numerical representation . . . . .	6
<b>III</b>	<b>Linear baseline model</b>	<b>8</b>
3.1	About this part . . . . .	9
3.2	Fine-tuning the input . . . . .	9
3.3	Fine-tuning the validation set . . . . .	9
3.4	Fine-tuning the model parameters . . . . .	9
3.5	The optimal settings for this model . . . . .	10

<i>CONTENTS</i>	iii
<b>IV What's next</b>	<b>11</b>
4.1 Further development . . . . .	12
4.2 Open issues . . . . .	12
<b>More figures</b>	<b>13</b>
Overview of training set . . . . .	13
Overview of features data . . . . .	14
<b>References</b>	<b>15</b>

## Part I

# About the code

## 1.1 Files accompanied by this report

Since this report discusses the development of an AI, a lot of code is discussed as well. This code is not shown inside this document but is available on the GitHub repository (Bontinck, 2020). All code is written in Python-based Jupyter Notebooks.

## 1.2 Ideology of the developed code

The Jupyter Notebooks have many inline comments and markdown blocks to make reading the code easier. If code is extensively discussed in this report a reference to the corresponding section is made inside the code. Some of the gathered results come from time-consuming function calls. These can take multiple hours to receive. To spare some time, these results are saved in a Pickle file so they can be loaded in without having to do the function call. The Notebooks are written in a way that makes testing multiple models easy, paying extra attention to reusability.

## 1.3 A typical model exploration

The testing of a model consists of two main parts. Firstly the input of the model has to be optimized. Afterwards, the (hyper)parameters of the model itself can be optimized. These steps are clearly visible with the discussed models in this report. After optimizing everything in a standalone fashion, it has to be checked that these newly optimized parameters do not influence previously optimized parameters.

## 1.4 Technical remarks

This report was created in L<sup>A</sup>T<sub>E</sub>X by modifying the excellent and well-known VUB themed template from Ruben De Smet (2020). BibLaTeX was used for reference management and natbib was used for more citation control.

Most source files, for this report and the created models, are available on GitHub (Bontinck, 2020). Some files, like the used training images, were not included in this GitHub repository. Details about this can be found on the GitHub page (README file). Rights to this GitHub repository can be asked from the author.

## Part II

# Data analysis

## 2.1 About this part

Before rigorously testing different models available, it's important to take a look at the data that is supplied. The supplied data consists of 2 main groups of images, labelled training images and unlabeled test images. As per the requirements of the Kaggle competition, the test images should only be used for evaluating the model on the Kaggle page by submitting a CSV of the prediction results. Thus the test images can't be used for creating the model in any shape or form. This means that only the labelled training images can be used to create and validate the model in development. To avoid altering the model to perform well on the supplied test data and not in general, only the training data will be analysed. All code used for this part is available under the developed code folder on GitHub, in the Jupyter Notebook *data\_analysis.ipynb*. This part describes how data can be analysed with the given code. This code can be easily changed to analyse other variations of the data, e.g. using another descriptor.

## 2.2 Data distribution

The provided labeled training data consists of 12 different classes. There is a total of 4042 labelled training images supplied, the distribution of which is shown in figure 1. As visible in this figure, the distribution between classes is not balanced. This has to be taken into account when fitting a model since some models will show unwanted behaviour when fitted with unbalanced data. Luckily many solutions exist to minify the impact of this unbalance.

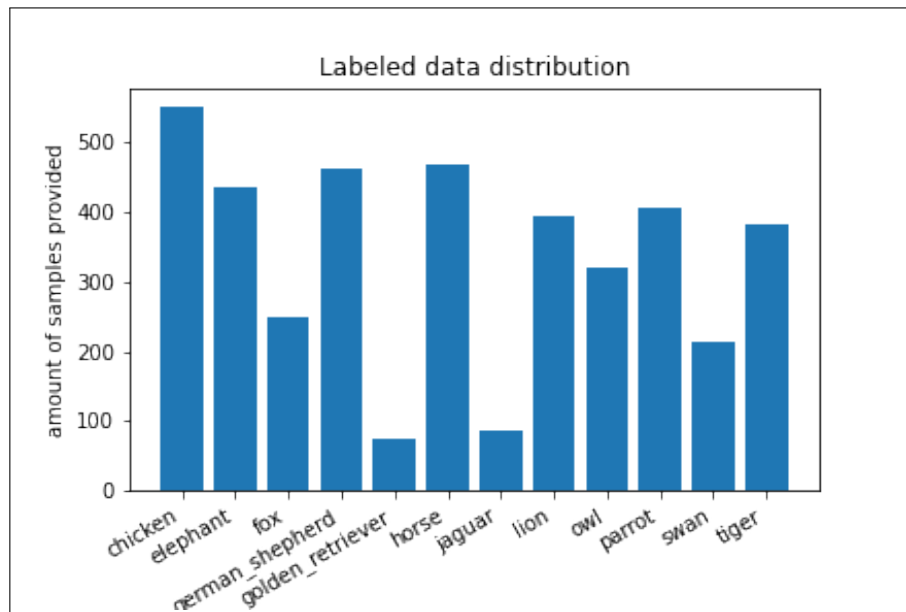


Figure 1: The data distribution of the supplied training set.



## 2.3 Deeper look at the training data

Whilst noting that the available data isn't balanced over all the classes is very important, there are also different aspects of the data that need analysing. An overview of the supplied training data is given in figure 4, available in the figures list at the end of this report. This figure shows the first five images of each class. From this, it becomes apparent that multiple factors of the data aren't *optimal*. This knowledge is important since it can aid in better preprocessing and in finding a better model in general. The most noteworthy findings are listed here:

- Images vary in shapes, some are taken in portrait, others in landscape.
- Images vary in size, some are high resolution whilst others are relatively low resolution.
- The framing of the subject(s) varies a lot. Sometimes the labelled animal is completely visible and centred in the frame. In some images there are multiple animals spread across the image, others show a close-up of the animal.
- Some images have a detailed background that makes up for a lot of the image, in others the background is blurry and its impact is presumably less.
- Some images have very vibrant colours in broad daylight, others are black and white in dimly lit environments.

This diversity in the provided training set is expected since it has been scraped from the web. This also means that *noise* can be expected, another important factor to keep in mind when choosing and optimizing models. Many of the listed things can be minified by doing some clever preprocessing of the images.

## 2.4 Feature extraction

Since the focus of this competition is on developing great models and not necessarily on data preprocessing and feature extraction, some feature extraction has already been provided. More info on the preprocessing and feature extraction provided is available in the provided notebook *creating\_vbowl.ipynb*. In short, images are converted from their typical RGB representation to a numerical representation of interesting points, which can be used as input for our model. How this is done will briefly be discussed here.

Instead of using the whole image as data, only a select few of *interesting points* of the image are taken into consideration. These interesting points of an image are found by using the *Shi-Tomasi corner detector*. The following important parameters for the *features.extractShiTomasiCorners* function call are used for the supplied features:

- number of features = 500
- minimum distance between features = 20

Shown in figure 2 is an example output of interesting points found by the Shi-Tomasi corner detector. It's clear that this is far from optimal, but finding interesting points isn't an easy task and thus the results are better than they might seem on first sight.

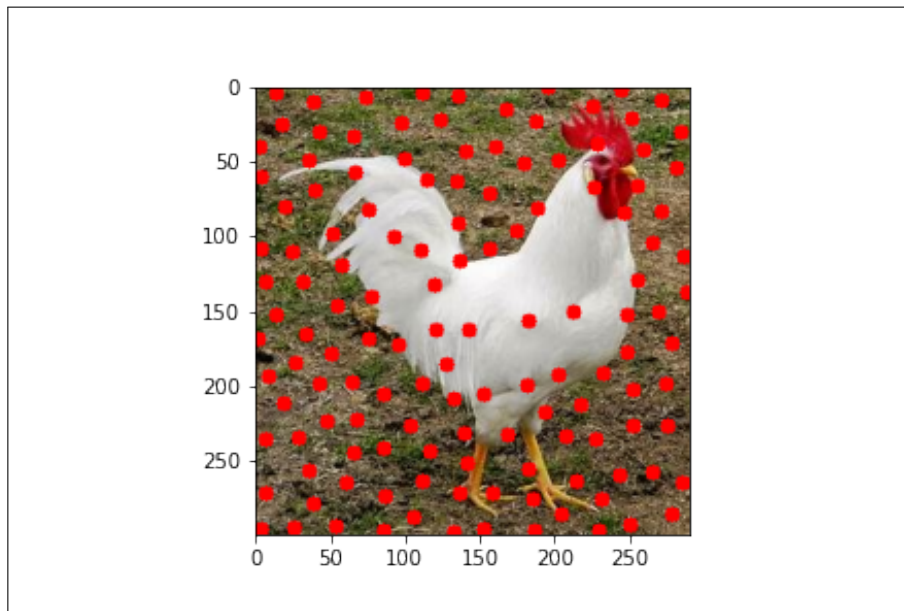


Figure 2: Example of points of interest found by Shi-Tomasi corner detector.

Finding interesting points is only half of the work. These interesting points now need to be represented by numerical values that have actual meaning. Remember from section 2.3 that the provided images differ a lot and thus a descriptor has to be used that minimizes the impact of different lighting, scaling... The following descriptors are used and their outputs are provided: DAISY, ORB, FREAK, LUCID, VGG, BoostDesc, SIFT. Whilst SURF is another great descriptor, it's not provided nor is the license available for this project. SIFT is often referred to as the most famous and successful of these descriptors, but all of them should be explored.

## 2.5 The numerical representation

As discussed in section 2.4, the images are stored as numerical representations of different features using descriptors. To save time, these numerical representations for all the descriptors are stored in a separate *Pickle* file. Since these representations form the input of a model, it's important to get a grip on how these look. The provided *createCodebook* function allows for easily loading in this data. It also allows specifying how many features per image are wanted with the *codebook\_size* parameter. As discussed in section 2.4, at most 500 of these are available per image by default. The function returns 2 lists, one containing the labels for the image represented at a certain index, the other containing the requested amount of features per image.

An overview of the data from such features given by the SIFT descriptor is given in figure 4, available in the figures list at the end of this report. From this, it is visible that the values seem to be normalized by the SIFT descriptor. This would have to be checked for all descriptors used and perhaps some outliers would need to be removed.

In figure 3 the correlation matrix is shown for the first 30 features of the SIFT descriptor. The correlation between these values doesn't seem too dramatic, which is mostly positive for our model building.

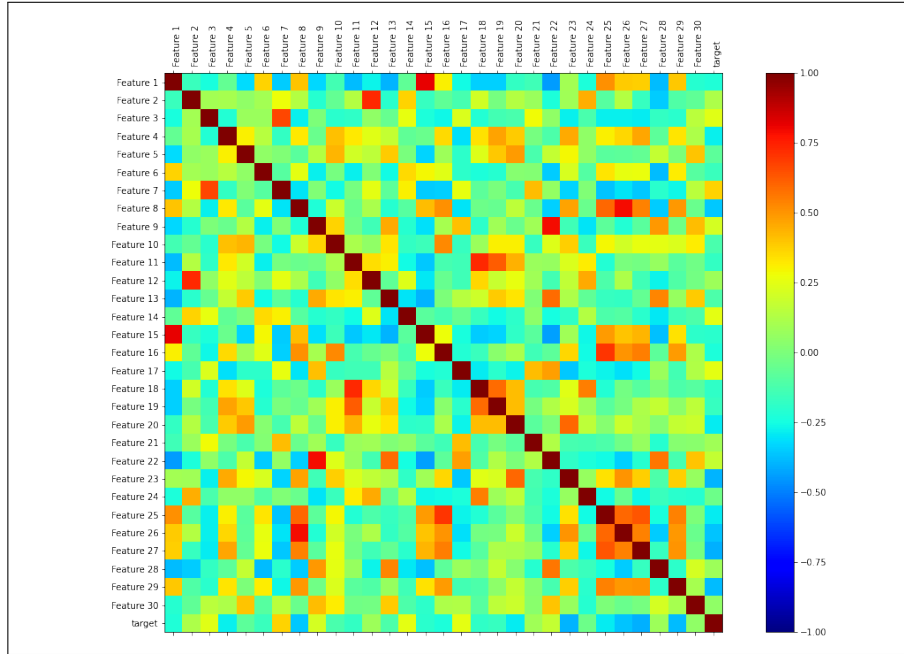


Figure 3: Example of points of interest found by Shi-Tomasi corner detector.

## Part III

# Linear baseline model

### 3.1 About this part

A good linear model available in the SciKit Learn library is the Logistic Regression model. This model is often used as a linear baseline model to compare other models with. Only models scoring better than this linear baseline model should be considered. This part discusses the parameters found to be optimal for this model in this setting and the road to finding these optimal parameters. The Python-based Jupyter Notebook corresponding with this part is *linear\_baseline\_model.ipynb*.

### 3.2 Fine-tuning the input

The first step in finding optimal settings for the model is finding optimal settings for the input of the model. In this case, the parameters that can control the input is the number of features each image has and the descriptor used as discussed in part II. SIFT is often referred to as the most famous and successful of the descriptor, but all of them should be explored. In general, more features often correspond to a better score, however, including many features can lead to overfitting of the model. The following values were tried

- Descriptors: DAISY, ORB, FREAK, LUCID, VGG, BoostDesc, SIFT.
- Feature amounts: 5, 20, 50, 100, 150, 250, 500.

A comparison was done by averaging the multi-class Log Loss score over 5 iterations for each of these feature amounts and descriptors. All of the resulting scores plotted on a graph can be found in the figures list at the end of this report.

### 3.3 Fine-tuning the validation set

Since the training data is further split into a training set and a validation set, this splitting can also be fine-tuned. The parameter that can be fine-tuned in this case is *test\_size* and whether or not to take into account that the data set is unbalanced. The latter is quite obvious as discussed in section 2.2. The splitting should thus be done with the unbalance in mind since using random splitting could lead to test or train sets without instances of specific classes. To find an optimal test size, expressed in a percentage, a similar testing method to the one discussed in section 3.2 is used.

TODO XXX

### 3.4 Fine-tuning the model parameters

Now that all of the parameters available for the input are fine-tuned, the parameters of the model itself can be optimized. As found in the documentation from the *LogisticRegression*

function of the SciKit Learn library there are multiple parameters available (Pedregosa et al., 2011). The ones that are the most interesting are:

- *solver*
  - Specifies which solver should be used for the optimization problem in the model.
  - *lbfgs* is used as default and whilst a little slow, this parameter doesn't require further fine-tuning.
- *penalty*
  - Since the *lbfgs* solver is used, the default *l2* penalization norm is the only one that can be used.
- *C*
  - The regularisation hyperparameter *C* defaults to 1. Fine-tuning might be required.
  - TODO XXXX
- *max\_iter*
  - Since there's convergence with the optimal settings, fine-tuning of this parameter is not required.
- *fit\_intercept*
  - Boolean that specifies if a constant (a.k.a. bias or intercept) should be added to the decision function.
  - The results with this parameter set to true and false were checked and TODO XXX

### 3.5 The optimal settings for this model

After all the fine-tuning discussed in the previous sections, the following settings were found to be optimal:

- Descriptor used: TODO XXX
- Feature amounts: TODO XXX
- Sample size: TODO XXX
- *C*: TODO XXX
- Fit intercept: TODO XXX

## Part IV

# What's next

## 4.1 Further development

Due to limited time available, this intermediate report and the current state of the project isn't overwhelming in any stretch of the imagination. Because of this, a special thanks is given to the teacher and teaching assistants who've softened the requirements for the intermediate report.

In the time available until the final deadline many new possibilities for creating a better model will be explored, this might include but is not limited to:

- Preprocessing can be explored.
- More available models need to be tested and compared to the linear baseline model.
- Nesting multiple good performing models might be an option, using weighted probabilities.

Some aspects of the created models and pipeline might also be modified further once the open issues discussed in section 4.2 are resolved.

## 4.2 Open issues

While developing the first models many questions arose. Most are already answered by googling or discussion with colleagues. Open questions that are not yet answered are listed below, for which guidance from the teaching assistants is asked.

- The data is not available as a pickle file for the SURF descriptor.



# More figures

Some figures are referred to in the text but not placed directly under the text. These are included in this list. All figures are high resolution thus zooming in the PDF should be viable to get a clearer view.

## Overview of training set

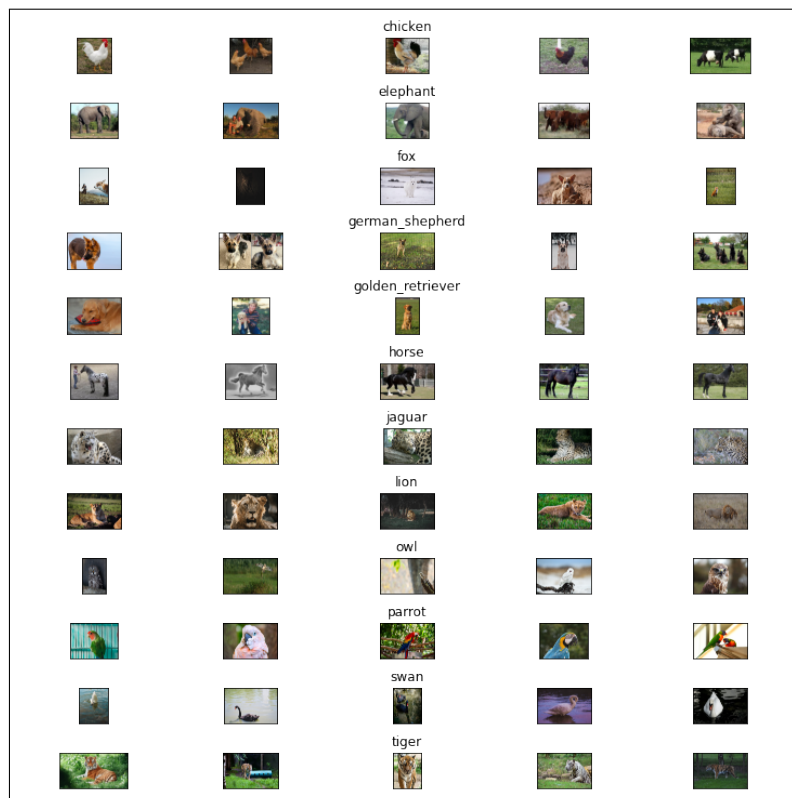


Figure 4: An overview of the supplied data per class.

## Overview of features data

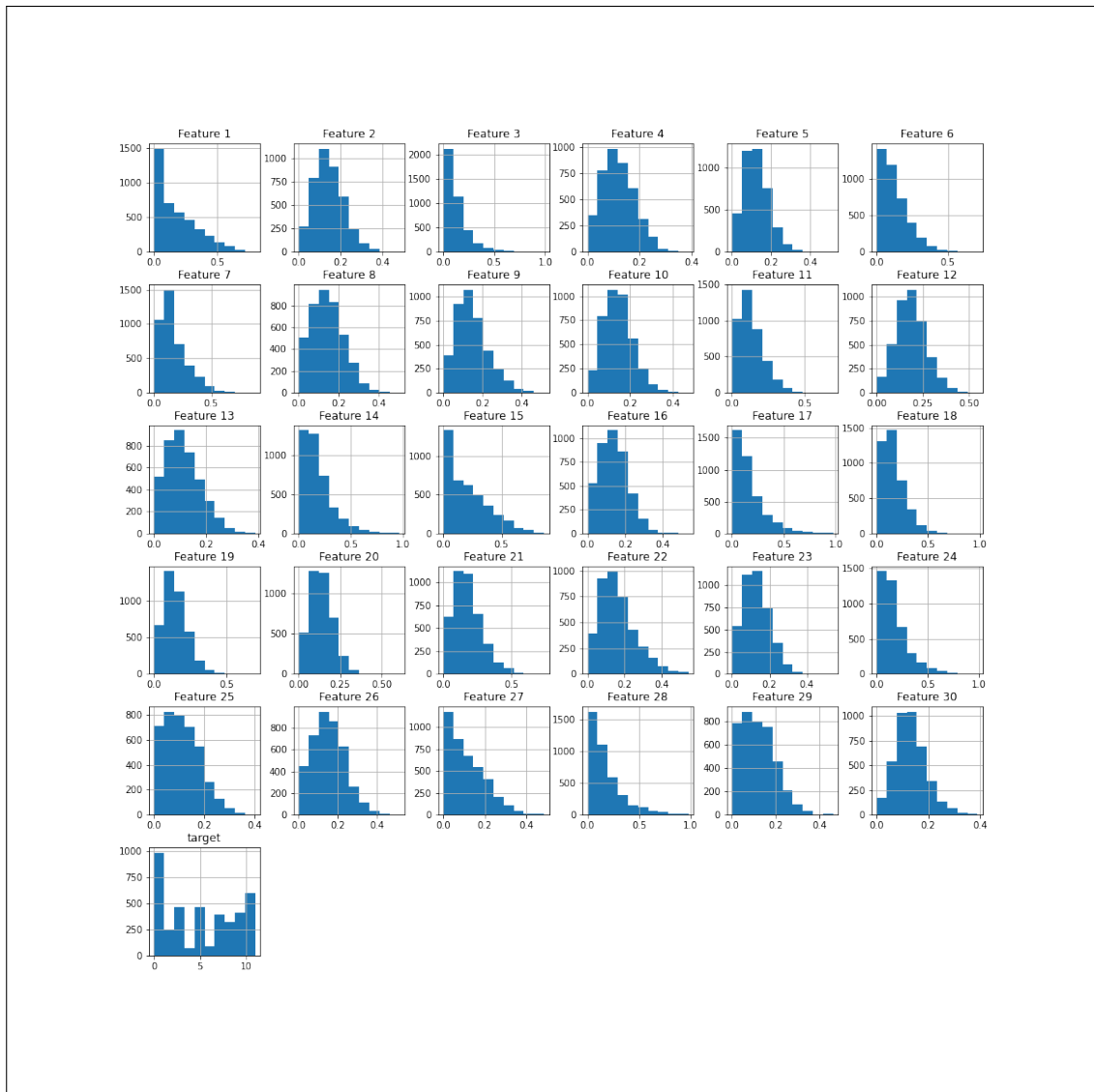


Figure 5: An overview of the first 30 features data from a SIFT descriptor.

# References

- Bontinck, L. (2020). *Machine learning project* [GitHub commit: Todo]. Retrieved December 11, 2020, from <https://github.com/VUB-CGT/ml-project-2020-pikawika>
- De Smet, R. (2020). *Vub latex huisstijl* [GitHub commit: d91f55799abd390a7dac92492f894b9b5fea2f47]. Retrieved November 2, 2020, from <https://gitlab.com/rubdos/texlive-vub>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- Rosseau, A. (2020). *Vub: Animal classification*. Retrieved November 15, 2020, from <https://www.kaggle.com/c/vub-animal-classification-20/>