

Image Classification of Aquatic Macroinvertebrates

ALEX JUSTESEN KARLSEN

Aarhus University - Department of Engineering
au515295@post.au.dk

May 2, 2019

Abstract

The paper documents my work on the in-class Kaggle Challenge; 'AU-ENG-CVML2019'. The paper is a mandatory submission in the Computer Vision and Machine Learning 2019 course at Aarhus University. The machine learning challenge is to obtain better classification result on the FIN-Benthic dataset of 29 classes of aquatic macroinvertebrates. The paper explains the dataset and a way to augment the dataset to gain better results. Furthermore the paper describes the machine learning methods used and the results.

1. INTRODUCTION

Life below water is a critical ecosystem, that experiences a lot of challenging changes. It is so important a challenge, that it has found its' way to the UN Sustainable Goals [1]. The Finnish Environment Institute [2] SYKE claims, that the loss of aquatic biodiversity and associated Ecosystem Services is one of the most pressing problems on Earth. SYKE launched the DETECT project [3], which aims to identify the various species of aquatic macroinvertebrates by imagery and classification using computer vision and machine learning.

This paper is a part of the Computer Vision and Machine Learning course at Aarhus University. It documents my work on solving the in-class Kaggle Challenge. The aim of the Kaggle challenge is improving the classification accuracy of classifying the 29 classes of aquatic macroinvertebrates. Machine Learning and state-of-the-art Deep Learning approaches utilizing Convolutional Neural Networks, are used to examine the different methods and get hand-on experience with a real classification challenge.

The rest of the paper is structured as follows; section 2 describes the data set in more depth and how the structure specific structure of the dataset have been used to improve the accuracy. Section 3 describes the machine learning and deep learning methods used in the project. Section 4 describes and discusses the results obtained from applying the methods. Section 5 is a discussion of the overall results and lastly in section 6 the project is concluded upon.

2. DATA SET

The FIN-Benthic [4] dataset 2 are used in this paper. The dataset consists of 29 different classes of aquatic macroinvertebrates.

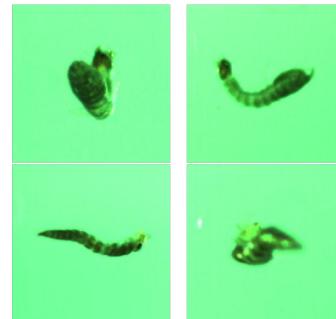


Figure 2.1. Example of aquatic macroinvertebrates

The dataset counts 11,588 images divided into training set (5830 images), validation set (2298 images) and test set (3460 images). The training and validation set are annotated with class labels, whereas the test set are not. This paper is a part of a Kaggle¹ Challenge with the aim of getting the best classification score on the test set, that are evaluated on Kaggle. For the challenge feature vectors produced by an AlexNet[5] trained on the raw images. The features are the output of the last convolutional layer and flattened into a one-dimensional vector with 4096 features. All images and vectors come in pairs of images of the same bug taken from perpendicular angles, an attempt to take advantage of the systematic series the images are all images vertically stacked. If we have a pair of x_1 and x_2 , then x_1 and x_2 are stacked as a combination of the two;

$$x_1 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, x_2 = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}$$

The vectors now contain $2 \times 4096 = 8192$ features. The concatenated version of the dataset are referred to as FIN-Benthic-Concatenated. Furthermore the training and validation sets are merged to form a combined set, to enable Cross-Validation for the different classification algorithms tested in this paper. Cross-Validation and Machine Learning methods used are described in the next section.

3. METHODS

Several experiments on this supervised learning challenge have been conducted using a variety of different techniques for classification are performed on the two encoded datasets; FIN-Benthic and FIN-Benthic-concatenated. The classification methods are; k-Nearest Neighbors (kNN), Linear Discriminant Analysis (LDA), Support Vector Machines (SVM), Perceptron and Neural Networks (NN). Commonly for all these supervised learning algorithm is the need for

¹Kaggle is an online community of data scientists and machine learners. Read more: <https://www.kaggle.com/getting-started/44916>

training. Training an algorithm is an optimization process of a function and a loss criterion, that explains the goodness of the model as training progresses until convergence. Cross Validation (CV) are used for model selection to evaluate the best performing model. All the methods have been applied to the datasets on full dimensions and additionally with dimensionality reduction as a preprocessing step in the classification pipeline. Grid search have been used to search for good hyperparameters. Hyperparameters are user specified parameter to control the algorithms behaviour [6]. Hyper parameter search have been applied for all classifiers. Furthermore the degree of dimensionality reduction have been investigated in the same manner. Lastly existing Convolutional Neural Network architectures have been used on the raw images and using the promise of transfer learning to use a pretrained model to gain faster and better performing models on small datasets. For the transfer learning task the existing train, validation and test split have been used.

k-Nearest Neighbors (kNN) K-Nearest Neighbor is a distance-based classifier and as the name implies, classifies to the most common occurring class of the k-nearest neighbors [7]. There exist a variety of distance metrics, some examples are; euclidean, minikowski and manhattan distance [8].

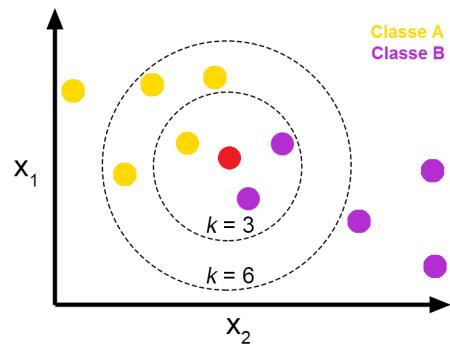


Figure 3.1. k-Nearest Neighbor classifies the red sample to class B for $k = 3$, however for $k = 6$ it classifies to class A. Credits: Italo José

Linear Discriminant Analysis (LDA) Linear Discriminant Analysis is a linear transformation of the data. The linear transformation is an optimization, that tries to maximize the distance between classes and minimize the distance within a class[8]. Figure 3.2 shows how LDA discriminates two classes.

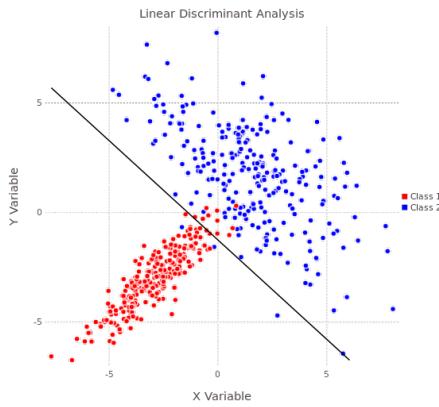


Figure 3.2. LDA creates a linear decision function. The linear decision function is a decision hyperplane, which discriminates classes. Credits: Tim Thatcher

LDA can also be used for dimensionality reduction, this would be further explained later.

Support Vector Machine (SVM) SVM can be used for classification. SVM tries to find the decision hyperplanes, that maximizes the distance between classes [9]. The decision hyperplane are given by a linear function $w^T x + b$, when the function is positive it predicts the class to be present, and intuitively the class is not present, when the function is negative [6]. The Support Vector Machine uses support vectors to define the decision hyperplane. The support vectors are vectors of the training data, that maximizes the margin between two classes. The original SVM proposed by Boser et al. [10] utilize a hard-margin. The hard-margin is the optimal decision hyperplane, however it only works well when the data is linearly separable. Vapnik et al. [11] proposed a soft-margin, which is more robust to noise by introducing some slack by introducing the hinge loss function.

One of the most promising contribution associated with SVMs is the kernel trick. The kernel trick makes not only SVM useful even when data is not linearly separable, but all machine learning algorithms, that can be rewritten exclusively in terms of dot product between examples. The kernel trick is a mapping of the data into a higher dimensional space - even with indefinite dimensions. The kernel trick is equivalent to transforming all data points into higher dimensions, but with a lot less computational demands. In the high dimensional space a linear decision boundary can be obtained. One of the most widely used kernel is the gaussian kernel or radial basis function (RBF), which maps to an indefinite dimensional space [6].

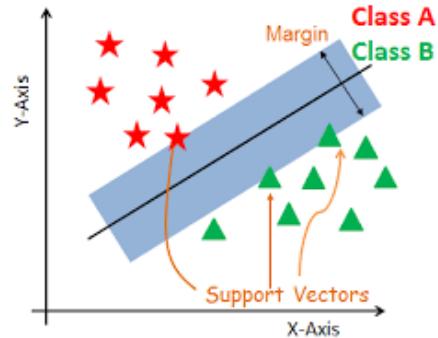


Figure 3.3. A binary SVM constructing the decision hyperplane, that maximizes the margin between the support vectors of the two classes. Credits: Avinash Navlani

Perceptron Perceptrons is a single layer Neural Network. The perceptron was set forth by Frank RosenBlatt [12], as a way of mimicking human perception. Perceptrons produces a linear decision function and only performs well on linearly separable data.

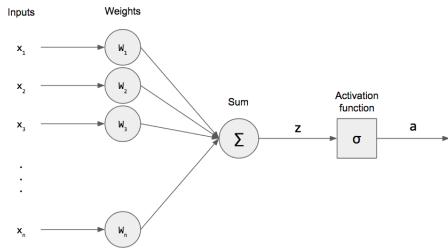


Figure 3.4. The perceptron takes a feature vector as input. All the feature are multiplied by a weight. The weight are then summed to what is called a weighted sum. An activation function are applied to the weighted sum. The output of the activation function is the classification. Credits: Stanley Obumeme Dukor

The perceptron are trained by iteratively updating the weights to produce the correct output for the training data. The loss function are rather simple, if the sample is correctly classified th loss function will output zero and it will output 1 if misclassified. This means the perceptron's weight are only updated when is fails to classify a sample.

Neural Network (NN) Neural Networks extends the perceptron. This ideas was sparked by progress within biology mapping human cognition into a network of perceptrons i.e. neurons. The work working Aritifical Neural Network was completed by A. G. Ivakhnenko and V. G. Lapa [13]. The perceptrons where extended by adding more layers between input and output referred to as hidden layers. The NN produces non-linear decision functions, thus overcomes the short comming of the perceptron.

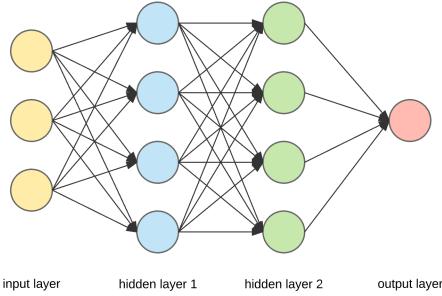


Figure 3.5. Neural Network with two hidden layers and each layer is fomed by a number of neurons. A fully-connected dense neural network receives all output from previous layer as input. The input are multiplied by the neuron's weight and adding bias. At last the nonlinear acitvation function is applied. Credits: Arden Dertat

Training a NN is way more complex than a perceptron. The training are done using feed-forward and back-propagation. Training vectors are given to the network, where all nodes do their operations upon receiving and ninput anf feeding it forward to the next layer, when it reaches the end of the network and error term is computed using the loss function. The error are then backpropagated. throuout the network, where the nodes' contribution to the error is computed. Using gradients the network weight are optimized.

Convolutional Neural Network (CNN)

Dimensionality Reduction Dimensionality reduction as the name implies downscales the feature space into a lower dimnesional space. The techniques used in this paper is the already covered LDA and Principal Component Analysis (PCA).

Principal Component Analysis (PCA) PCA is an unsupervised leraning technique. It learns an othogonal linear transformation onto a smaller feature space. The number of dimensions i.e. principal components are a hyperparameter. PCA tries to preserve as much as the information i.e. variance as possible. The first component i.e. dimension is a projection a long

the axis of the data with the variance, the second component orthogonal to the first axis. The second component describes the axis with second most variance. The third the third most variance and so on, as the figure 3.6 shows.

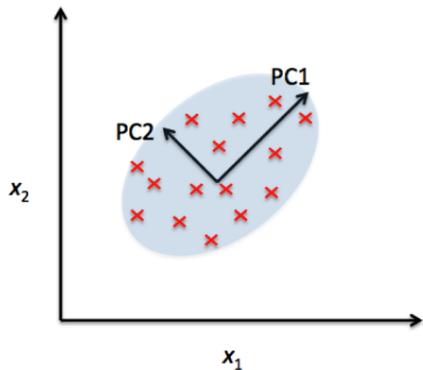


Figure 3.6. Principal Component Analysis in a 2 dimensional space. PCA rotates the data to be perpendicular to the principal components axes. Credits: Sebastian Raschka

Linear Discriminant Analysis (LDA) As mentioned LDA can be used for dimensionality reduction as well. One must specify the n number of components i.e. dimensions to transform the data onto, if the n are less than the original data's dimensions D, the the number of dimensions will be reduced.

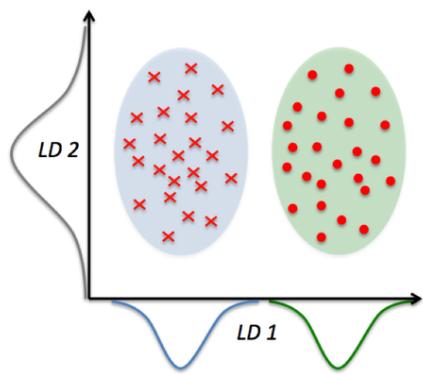


Figure 3.7. Linear Discriminant Analysis project the data onto an axis, that maximizes the distance between the classes and minimizes the distance within the respective classes. Credits: Sebastian Raschka

Model Selection Model selection are procedures to train the best model and avoid overfitting the model to the training data. Overfitting is the pitfall of having a model, that describes the training data too well, thus suffers to not generalising the true underlying structure of the data. Although other methods exist the two used in the paper are Validation Holdout and Cross Validation.

Validation Holdout Validation Holdout takes a fraction of your available labelled data and hold it out of the training procedure.

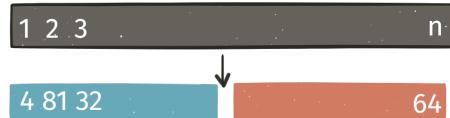


Figure 3.8. The validation subset is then used to validate the model's performance on data it has not been training.

The assumption is, that the holdout set resembles the true data, since it has been drawn from the same distribution [6]. If the model performs well on both training and validation it is assumed to be a good model. The split is typically made 70/30 for training data and validation data respectively.

Cross Validation (CV) Cross Validation extends on the holdout validation. k -fold Cross Validation creates k experiments compared to one. The entire data set are split in k folds, for every k step one subset is for validation and the remaining $k - 1$ fold are used for training [6]. Figure 3.9 illustrates the idea.



Figure 3.9. Cross Validation ensures that the model have been trained on all the data and validated as well. The selection folds are typically either $k = 5$ or $k = 10$.

Hyperparameter Search There are mainly two approaches to search for hyperparameters; Random Search and Grid Search other than heuristic guessing. In this paper Grid Search are used.

Grid Search Grid Search is an approach to seek better hyperparameters for your model. You specify a set of options for all hyperparameters and then search all possible combination in the grid.

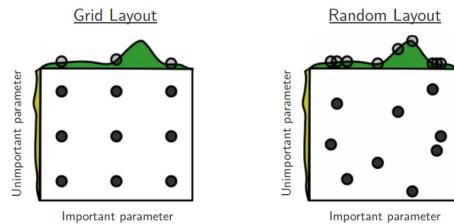


Figure 3.10. Grid search is a systematic approach to explore a broader space of hyperparameters although at the cost of additional time. The down-side of this approach is, that it is still based on intuition of good hyperparameters. The random approach are also widely used and might be preferred, as it may find a sweet spot, that goes beyond human intuition, however that is no guarantee, and it is still time costly [14]. Credits: Bergstra and Bengio

4. RESULTS

Two main experiments have been conducted; the first experiment involves applying the aforementioned machine learning methods on the feature extracted vector. The second experiment

involves applying deep learning practices on the raw image data.

Experiment 1: Machine Learning

The machine learning experiments are implemented in Python 3 using the SciKit Learn library, except the Neural Network that are implemented in Keras using Tensorflow backend.

All classifiers have been run on the FIN-Benthic dataset and on the FIN-Benthic concatenated dataset using Grid Search and CV to find the best hyperparameters, except the neural networks. As grid search has been used to find the best set of parameters for each of the classifiers and on the different datasets, it has not given the same results for best hyperparameter settings. Only the best performing setting of hyperparameters for each classifier are presented in this paper.

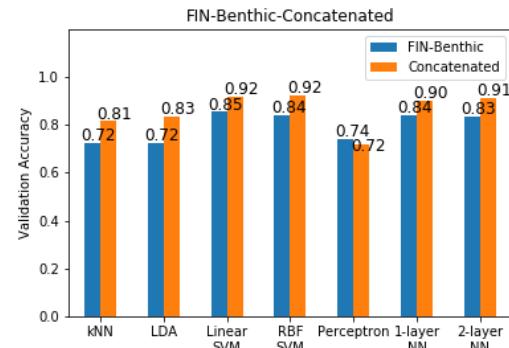


Figure 4.1. The best setting for each of the classifiers in terms of accuracy. Generally the classifiers perform best on the concatenated dataset. The assumption about concatenating the pair sample of the same bug holds true.

I added PCA to the classification pipeline and ran the same grid search on all the classifiers on both datasets. For all runs the dimensionality have been reduced to $n = \{128, 256, 512, 1024\}$ principal components. PCA add additional time to the pipeline, however the time to train and predict on fewer dimensions significantly improved. Altough the dimensionality have been

reduced with a large margin it does not have huge impact on the classification accuracy.

Bring some results using dimensionality reduction

Experiment 2: Deep Learning

5. DISCUSSION

6. CONCLUSION

7. REFERENCES

- [1] “Transforming our world : the 2030 agenda for sustainable development,” *UN General Assembly*, accessed: 8 April 2019.
- [2] K. Meissner, “Detect.” [Online]. Available: https://www.syke.fi/en-US/Research_Development/Research_and_development_projects/Projects/DETECT
- [3] E. Riabchenko, K. Meissner, I. Ahmad, A. Iosifidis, V. Tirronen, M. Gabbouj, and S. Kiranyaz, “Learned vs. engineered features for fine-grained classification of aquatic macroinvertebrates,” in *2016 23rd International Conference on Pattern Recognition (ICPR)*, Dec 2016, pp. 2276–2281.
- [4] K. Meissner, “Fin-benthic.” [Online]. Available: <https://etsin.avointiede.fi/dataset/urn-nbn-fi-csc-kata20170615175247247938>
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, ser. NIPS’12. USA: Curran Associates Inc., 2012, pp. 1097–1105. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2999134.2999257>
- [6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [7] A. Karpathy, “Image classification.” [Online]. Available: <http://cs231n.github.io/classification/#knn>
- [8] A. Iosifidis, “Course notes: Introduction to machine learning,” Blackboard, 2017.
- [9] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press, 2008.
- [10] B. E. Boser, I. M. Guyon, and V. N. Vapnik, “A training algorithm for optimal margin classifiers,” in *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*. ACM Press, 1992, pp. 144–152.
- [11] C. Cortes and V. Vapnik, “Support-vector networks,” *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, Sep. 1995. [Online]. Available: <https://doi.org/10.1023/A:1022627411411>
- [12] F. Rosenblatt, “The perceptron: A probabilistic model for information storage and organization in the brain,” *Psychological Review*, pp. 65–386, 1958.
- [13] A. G. Ivakhnenko and V. G. Lapa, “Cybernetic predicting devices,” p. 250, 04 1966.
- [14] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 281–305, Feb. 2012. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2503308.2188395>