Machine Learning on Smartphone-Captured Diffraction Data

Ashish Shivajirao Jadhav^{1,2}, Andreas Backhaus² and Udo Seiffert^{1,2}

1 - Otto von Guericke University Magdeburg - Institute for Information and Communication Technology, Magdeburg - Germany

2 - Compolytics GmbH, Barleben - Germany

Abstract. This study presents a novel approach for classifying oily or cream-like substances using diffraction data captured on a smartphone camera, applied specifically to assessing engine oil quality. Utilising the COMPOLYTICS[®] TapCorder approach, optical diffraction patterns were analysed with a tailored feature extraction method. The performance of three machine learning paradigms – Multilayer Perceptrons (MLP), Learning Vector Quantization (LVQ), and Radial Basis Function Networks (RBFN) – was analysed in classifying new and used oil samples. MLP achieved the highest accuracy, while LVQ required the least computation time, highlighting trade-offs relevant for consumer-focused applications. This work clearly demonstrates the feasibility of accessible, low-cost chemical substance analysis via smartphone-based systems.

1 Introduction

The qualitative and quantitative analyses of chemical substances are typically conducted using suitable wet chemical methods. However, these methods are not real-time capable, involve significant investment and operating costs, and are not readily accessible to consumers. Additionally, the process often results in the consumption of the available sample.

The use of smartphones beyond pure communication applications for use as measuring devices / testers makes the determination of ingredients cheaper and generally suitable for consumers. The COMPOLYTICS[®] TapCorder [1] is a novel data acquisition tool for capturing optical characteristics of oily or cream-like substances directly via smartphone. By utilizing geometric optics, it measures diffraction effects on a thin film sample applied to the camera lens. This method, based on optimally designed measurement patterns, enables precise data collection of texture and diffraction features. For some applications, the use of the Siemens star has proven to be the optimal measurement pattern, a circle divided into sectors of alternating contrasting colours (typically black and white), whose radially symmetrical spatial frequencies increase steadily towards the centre. By selecting the number of sectors and thus the distance between them, the value range of the spatial frequencies can be easily parametrised.

Since the relationship between the image properties modulated by the grease film applied on the lens and the qualitative and quantitative properties of the measurement sample defined in the context of the measurement task is typically not given analytically, machine learning methods are used for this modelling. In order to enable data processing within a consumer app as well, the model must not only be sufficiently precise and robust, but also optimised in terms of the required computing time and resources.

This set-up, following the principle of a soft sensor, consisting of an innovatively implemented physical measurement principle and tailored statistical modelling, opens up numerous application perspectives. In the context of an application for assessing the deterioration of engine oil, this paper illustrates the extraction of application-specific optimised features from an extensive image data set and the use of three different machine learning concepts based on an exemplary classification task of new vs. used oil. For illustrative reasons, a third class of 'clean' lens was introduced.

2 Data Acquisition and Feature Engineering

After the application of the available oil samples to the camera lens, the image data was acquired. All images were taken with a regular smartphone, with the Siemens star (36 spokes) positioned in the background (approx. 20 cm distance). Two sets of image data were acquired using two separate kinds of oil, namely new and used oil. Additionally, images with clean lens, also with the Siemens star in the background, were acquired. Ultimately, about 100 images of each type were obtained. To ensure balanced classes, the number of images per class was intentionally kept the same. All images were orientation corrected and cropped to a uniform size of 1850×1850 pixels at 8 bits, and saved as lossless Portable Network Graphics (PNG) files. The complete data set is available for download [2].

Histogram features: Corresponding to the distinct black and white fields of the Siemens star, a bimodal histogram distribution was obtained, modulated by the properties of the oil composition. The initial features were extracted by fitting both a generic Gaussian Mixture Model (GMM), yielding mean and covariance per distribution, and a better shape approximating Weibull Mixture Model (WMM), yielding alpha and beta per distribution.

Texture features: Further relevant features were extracted, such as Shannon entropy and grey-level co-occurrence matrix yielding contrast, dissimilarity, homogeneity, energy, correlation, and angular second moment (ASM).

Combining both feature domains, the initial dataset had 27 features. Twelve highly correlated features with a correlation of approximately 0.99 were removed from the dataset. A Pearson correlation matrix of the 15 final feature set is shown (see Fig. 1).

Before training the models, an exploratory data analysis was carried out to visualise the general class separability using a Linear Discriminant Analysis (LDA), see Fig. 2 for details. This enabled more effective discrimination between classes and facilitated learning from the features due to class differences.



Fig. 1: Correlation matrix of the selected 15 features presented as heat map. In principle, the lower the correlation between a feature and several other features, the more it contributes to the overall feature set.

3 Machine Learning Models

In selecting neural network paradigms specifically for classification tasks, it is essential to include conceptually distinct and widely used models that are particularly well-suited or even ideally adapted for this purpose. This study therefore focuses on three paradigms, each representing a unique approach to classification: separator-based Multilayer Perceptrons (MLPs) [3, 4], which rely on learned boundaries between classes; Radial Basis Function Networks (RBFNs) [5, 6], which employ radial basis functions to capture localised patterns in data; and prototype-based Learning Vector Quantization (LVQ) [7, 8], which uses representative prototypes to categorise input utilising a specific implementation from [9]. Together, these paradigms provide a comprehensive exploration of classification techniques within neural network architectures.

As a common preprocessing step, the features were standardised using z-score normalisation. Labels were encoded using one-hot encoding to train MLP and RBFN, and integer encoding was used to train LVQ. All models were trained using ten-fold cross-validation, with one fold for validation and the rest for training. To ensure robust results, each cross-validation step was repeated ten times. Mean and standard deviation (SD) of training and validation accuracy over repetitions were calculated. Each paradigm's training set-up started with the



Fig. 2: All samples with their class affiliations, represented in two lowdimensional components using Linear Discriminant Analysis (LDA). The samples indicate the presence of oil: blue for new, red for used, yellow for absent. The no-oil class is clearly separated from the other two, whereas new and used oil are also separable. The within-class variability of the no-oil samples is clearly smaller.

minimum possible design for the architecture type. The complexity of the design was increased until it met our desired validation accuracy target of 95%. The mean computation time per repetition was also recorded.

For the MLP model, the Nguyen-Widrow method was applied to initialise weights, and the Levenberg-Marquardt (LM) optimization technique [10] was used for training. The hidden layer employed a hyperbolic tangent activation function, while the output layer used a linear activation function, with mean squared error (MSE) as the loss metric. In the RBFN model, the least mean square error (LMSE) algorithm provided the optimisation method. The hidden layer utilised radial basis functions as the activation function, while the output layer was linear, with MSE again used as the loss function. For LVQ, class means served as the initial prototypes, and the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm [11] was used for optimisation.

4 Results

The results are summarised in Table 1, which illustrates the variability in performance across network architectures and configurations, highlighting the tradeoffs in accuracy and computational efficiency.

Training: The minimal design and satisfactory design per architecture type are reported here. The MLP with 15-5-3 architecture demonstrated the highest level of validation accuracy. It achieves a validation accuracy of 98.2 % with a relatively short training period of 43 ms. The LVQ prototype, comprising a minimal possible design with one prototype per class, already successfully passes the validation criteria with an accuracy of 95 %. It is unnecessary to increase the complexity of the LVQ; however, a considerably longer training period of 269 ms is required. The RBFN with 16 RBF neurons satisfies the validation criteria but requires the longest training time of 519 ms.

Recall: In terms of a potential user application, the computational time required for a newly acquired image of the test pattern (recall) is definitely more important. The LVQ has the lowest mean validation time per sample, at 1.6 μ s. In comparison, the MLP requires approximately 50 times and the RBF approximately 70 times the validation time of the LVQ.

Apart from the mean value for each result in Table 1, the standard deviation was also calculated to assess the variability in training accuracy due to different random initializations, the consistency of recall results. Additionally, mean computation time of network paradigms considered was accessed when processing new samples. With regard to the calculation of computation time, feature extraction was not considered here.

		Accuracy				Computation Time	
		Mean	SD	Mean	SD	Mean	
Type	Design	Train		Validation		Train	Validation
MLP	1-3	0.867	0.010	0.659	0.035	$43 \mathrm{ms}$	$57.3 \ \mu s$
MLP	5-3	0.994	0.002	0.982	0.015	$43 \mathrm{ms}$	$51.5 \ \mu s$
LVQ	1-1-1	0.971	0.006	0.950	0.042	$269 \mathrm{~ms}$	$1.6 \ \mu s$
RBFN	1-3	0.676	10^{-16}	0.676	10^{-16}	23 ms	$70.1 \ \mu s$
RBFN	16-3	0.960	0.005	0.950	0.039	$519 \mathrm{~ms}$	$70.4 \ \mu s$

Table 1: Summary of accuracy and computation time results for the considered neural network paradigms and configurations (MLP, LVQ, and RBFN). For each model, mean and standard deviation (SD) values are reported for both training and validation accuracy, as well as the mean computation time required for training and validation phases. For MLP and RBFN, design means [the number of neurons in hidden layer - number of neurons in output layer], and for LVQ the number of prototypes per class.

5 Conclusion and Future Work

The purpose of this study is to evaluate the technical feasibility of the underlying use case in accordance with the patent. It uses a simple dataset to determine whether classification is possible on such data. It also highlights applicationspecific feature engineering required in this scenario.

The comparative analysis of MLP, LVQ, and RBFN models demonstrates distinct strengths and limitations across classification accuracy and computational efficiency, with MLP achieving high validation accuracy and showing minimal variability at the same time compared to other methods achieving the desired accuracy. RBFN model shows minimal variability. These findings suggest that the choice of a suitable neural network paradigm in a subsequent consumer application should consider not only accuracy requirements but also the specific computational constraints of the application, particularly in real-time or resource-limited settings of a smartphone app. Moreover, in terms of interpretability and potential model size, prototype-based neural networks, such as LVQ, generally offer interesting perspectives [12].

In the future, the complexity of the data will need to be increased by increasing the number of classes (i.e., oil types and conditions), potentially the number of features, and also the number of images.

References

- Udo Seiffert, Andreas Backhaus, and Andreas Herzog. Method for examining substances, computer program for use in said method, and mobile electronic device. *European Patent*, EP 3 999 838 B1, 2024.
- [2] Ashish S. Jadhav and Udo Seiffert. The TapCorder Data Set. available at: https://doi.org/10.24352/ub.ovgu-2024-096, 2024.
- [3] Marvin Minsky and Seymour Papert. Perceptrons: An introduction to computational geometry. *MIT Press*, 1969.
- [4] David Rumelhart, Geoffrey Hinton, and Ronald Williams. Learning representations by back propagating errors. *Nature*, 323:533–536, 10 1986.
- [5] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximations neural networks. *Neural Networks*, 2:359–366, 1989.
- [6] John Moody and Christian J. Darken. Networks for pattern recognition: Alternative to Backpropagation. In Proceedings of the IEEE International Conference on Neural Networks (ICNN), pages 881–885, 1989.
- [7] Teuvo Kohonen. Learning vector quantization. Neural Networks, 1(3):303–319, 1988.
- [8] Michael Biehl, Barbara Hammer, and Thomas Villmann. Prototype-based models in machine learning. Wiley Interdisciplinary Reviews. Cognitive Science, 7(2):92–111, 2016.
- Petra Schneider, Kerstin Bunte, Han Stiekema, Barbara Hammer, Thomas Villmann, and Michael Biehl. Regularization in matrix relevance learning. *IEEE Transactions on Neural Networks*, 21(5):831–840, 2010.
- [10] Martin Hagan and Mohammad Menhaj. Training feedforward networks with the Marquardt algorithm. *IEEE Transactions on Neural Networks*, 5:989–93, 1994.
- [11] Roger Fletcher. Practical Methods of Optimization. Wiley, 1987.
- [12] Andreas Backhaus and Udo Seiffert. Classification in high-dimensional spectral data: Accuracy vs. interpretability vs. model size. *Neurocomputing*, 131:15–22, 2014.