

RESEARCH

Open Access



Using transfer learning and dimensionality reduction techniques to improve generalisability of machine-learning predictions of mosquito ages from mid-infrared spectra

Emmanuel P. Mwanga^{1,2*}, Doreen J. Siria¹, Joshua Mitton^{2,3}, Issa H. Mshani^{1,2}, Mario González-Jiménez⁴, Prashanth Selvaraj⁵, Klaas Wynne⁴, Francesco Baldini², Fredros O. Okumu^{1,2,6} and Simon A. Babayan²

*Correspondence:
emwanga@ihi.or.tz

¹ Environmental Health and Ecological Sciences Department, Ifakara Health Institute, Morogoro, Tanzania

² School of Biodiversity, One Health, and Veterinary Medicine, University of Glasgow, Glasgow G12 8QQ, UK

³ School of Computing Science, University of Glasgow, Glasgow G12 8QQ, UK

⁴ School of Chemistry, University of Glasgow, Glasgow G12 8QQ, UK

⁵ Institute for Disease Modelling, Bellevue, WA 98005, USA

⁶ School of Public Health, University of Witwatersrand, Johannesburg, South Africa

Abstract

Background: Old mosquitoes are more likely to transmit malaria than young ones. Therefore, accurate prediction of mosquito population age can drastically improve the evaluation of mosquito-targeted interventions. However, standard methods for age-grading mosquitoes are laborious and costly. We have shown that Mid-infrared spectroscopy (MIRS) can be used to detect age-specific patterns in mosquito cuticles and thus can be used to train age-grading machine learning models. However, these models tend to transfer poorly across populations. Here, we investigate whether applying dimensionality reduction and transfer learning to MIRS data can improve the transferability of MIRS-based predictions for mosquito ages.

Methods: We reared adults of the malaria vector *Anopheles arabiensis* in two insectaries. The heads and thoraces of female mosquitoes were scanned using an attenuated total reflection-Fourier transform infrared spectrometer, which were grouped into two different age classes. The dimensionality of the spectra data was reduced using unsupervised principal component analysis or t-distributed stochastic neighbour embedding, and then used to train deep learning and standard machine learning classifiers. Transfer learning was also evaluated to improve transferability of the models when predicting mosquito age classes from new populations.

Results: Model accuracies for predicting the age of mosquitoes from the same population as the training samples reached 99% for deep learning and 92% for standard machine learning. However, these models did not generalise to a different population, achieving only 46% and 48% accuracy for deep learning and standard machine learning, respectively. Dimensionality reduction did not improve model generalizability but reduced computational time. Transfer learning by updating pre-trained models with 2% of mosquitoes from the alternate population improved performance to ~98% accuracy for predicting mosquito age classes in the alternative population.

Conclusion: Combining dimensionality reduction and transfer learning can reduce computational costs and improve the transferability of both deep learning and standard machine learning models for predicting the age of mosquitoes. Future studies



should investigate the optimal quantities and diversity of training data necessary for transfer learning and the implications for broader generalisability to unseen datasets.

Keywords: *Anopheles arabiensis*, Convolutional neural network, Standard machine learning, Generalisability, Dimensionality reduction, Transfer learning

Background

Malaria currently kills approximately one child every minute [1]. In 2020, there were 241 million cases and 627,000 deaths, nearly all in Sub-Saharan Africa [1]. Currently, the most widespread and cost-effective method of malaria prevention is based on controlling the mosquitoes that transmit the disease. Since 2000, insecticide-treated nets (ITNs) and indoor residual spraying (IRS) have so far contributed nearly 80% of all global malaria decline [2]. However, the direct impact of individual control programs on the mosquito populations and on malaria transmission at the sites of intervention remains difficult to measure. To guide further efforts against the disease, evaluating the performance of these and other vector control interventions is crucial for measuring their impact in different settings. The World Health Organization (WHO) now recommends that surveillance be integrated as a core component of malaria control programs [3].

This necessitates scalable, simple-to-implement and low-cost methods for quantifying key biological attributes of mosquitoes, such as age, infection status, and blood meal preferences, which are essential for understanding pathogen transmission dynamics. The age and survivorship of key *Anopheles* vectors are especially important in determining the likelihood that the mosquitoes will live long enough to allow complete parasite development (the extrinsic incubation period), and subsequent transmission to humans [4]. The assessments are essential for monitoring the impacts of interventions such as ITNs and IRS, which primarily kill adult mosquitoes in the field [5].

The current "gold standard" for estimating the age of malaria mosquitoes is to dissect their ovaries to estimate how many times they have laid eggs [5, 6]. Despite their low technical demands, such procedures are time-consuming and labour-intensive. Age-grading dissections can also be imprecise because of gonotrophic discordance, which is common in Afrotropical malaria vectors [7], or of their reliance on the availability of host blood meals, which determines when and how frequently a mosquito blood-feeds.

We and others have demonstrated that spectroscopic analysis of mosquitoes using near infrared (12,500–4000 cm^{-1}) or mid-infrared (MIR) (4000–400 cm^{-1}) frequencies can identify key biochemical signals that vary with age [8, 9]. These methods, when combined with specific machine learning (ML) techniques, allow for rapid estimation of mosquito ages [9, 10].

Despite early successes, these infrared-based applications have limitations such as their portability to mosquitoes from different locations or laboratories [10] and the substantial computational requirements for retraining such models. Indeed, the inherent variability of mosquitoes from different environmental and genetic backgrounds may limit the generalisability of models trained on infrared spectra. The models could also be misled by signals in MIRS that are associated with confounding factors introduced during sampling (e.g., atmospheric contamination with water vapour, temperature variations and high humidity in the laboratory), thus learning features that are not strictly

related to the biochemical trait being investigated. Therefore, machine learning models must be regularly updated with new data from target mosquito populations.

To increase the generalisability of ML models for a given training dataset, a variety of spectral smoothing and regularisation techniques have been tested, such as penalised regression [11]. These methods are known to be computationally efficient and to improve generalisability [11]. Deep learning (DL) techniques such as convolutional neural networks (CNN) have recently been used on large spectra data [10], improving generalisability through transfer learning (i.e., updating a pre-trained model with a small amount of new data from a different target population). However, when trained on large datasets, such techniques remain computationally expensive and may necessitate repeated sampling of hundreds of mosquitoes from different populations and environments to allow successful generalisability. Alternatively, since standard ML models are less complex than DL, computational time can be kept to a minimum. DL methods are versatile extensions of machine learning that are ideal for complex or large datasets [12]. But are prone to overfitting, such as predicting the training dataset well but failing on previously unseen or new data.

However, unsupervised learning algorithms, which find patterns independent of pre-defined target labels, can aggregate, cluster or eliminate features while retaining dominant statistical information before machine learning training on the spectra data. The resulting dimensionality reduction may improve generalisability, reducing overfitting, increasing the signal-to-noise ratio of the data, as well as lowering computational requirements for training machine learning models. Examples include principal component analysis (PCA) [13–15], which projects a large number of variables into distinct categories that summarise data into a small number of independent principal components, and t-distributed Stochastic Embedding (t-SNE) [16], which clusters datapoints based distances between all their input dimensions.

This study assessed whether the generalisability and computational costs of MIRS-based models for predicting the age classes of female *An. arabiensis* mosquitoes reared in two different insectaries in two locations could be improved by combining dimensionality reduction and transfer learning methods.

Methods

Collection of mosquito spectra data

We analysed mid-infrared spectra from two strains of *An. arabiensis* mosquitoes obtained from two different insectaries, one from University of Glasgow, UK and another from Ifakara Health Institute, Tanzania. The same data had previously been used to demonstrate the capabilities of mid-infrared spectroscopy and CNN for distinguishing between species and determining mosquito age [10]. The insectary conditions under which the mosquitoes were reared (temperature 27 ± 1.0 °C, and relative humidity $80 \pm 5\%$) have been described elsewhere [17].

Mosquitoes were collected from day 1 to day 17 after pupal emergence at both laboratories and divided in two age classes (1–9 day-olds and 10–17 day-olds). Silica gel was used to dry the mosquitoes. For each chronological age in each laboratory, ~120 samples were measured by MIRS on each day. The heads and thoraces of the mosquitoes were then scanned with an attenuated total reflectance Fourier-Transform Infrared (FTIR)

ALPHA II and Bruker Vertex 70 spectrometers both equipped with a diamond ATR accessory (BRUKER-OPTIC GmbH, Ettlingen, Germany). The scanning was performed in the mid-infrared spectral range ($4000\text{--}400\text{ cm}^{-1}$) at a resolution of 2 cm^{-1} , with each sample being scanned 16 times to obtain averaged spectra as previously described [9, 18]. As a result, the spectral dataset contained 1665 spectral features (Fig. 1).

Data pre-processing

The spectral data were cleaned to eliminate bands of low intensity or significant atmospheric intrusion using the custom algorithm [19]. The final datasets from Ifakara and Glasgow contained 1720 and 1635 mosquito spectra, respectively. In these two datasets, the chronological age of *An. arabiensis* was categorised as 1–9 days old (i.e. young mosquitoes representative of those typically unable to transmit malaria) and 10–17 days old (i.e. older mosquitoes representative of those potentially able to transmit malaria) [20].

To improve the accuracy and speed of convergence of subsequent algorithms, data were standardised by centring around the mean and scaling to unit variance [21].

Dimensionality reduction

Principal component analysis (PCA) and t-distributed stochastic neighbour embedding (t-SNE) were used separately to reduce the dimensionality of the data [13–16]. Both PCA and t-SNE were implemented using the scikit-learn library [21].

Separately, t-SNE was used to convert high-dimensional Euclidean distances between spectral points into joint probabilities representing similarities. To cluster the data into three features, the embedded space was set to 3, because the Barnes-hut algorithm in t-SNE is limited to only 4 components. Perplexity was set to 30 as the number of nearest neighbours, which means that for each point, the algorithm took the 30 closest points and preserved the distances between them. For smaller datasets perplexity values ranging from 5 and 50 are thought to be optimal for avoiding local variations and merged clusters caused by small or large perplexity values [16]. The learning rate for t-SNE is generally in the range of 10–1000 [21], thus it was set to 200 scalar.

Machine learning training

Deep learning

DL models were trained and used to classify the *An. arabiensis* mosquitoes into the two age classes (1–9 or 10–17 day-olds). The intensities of *An. arabiensis* mid-infrared

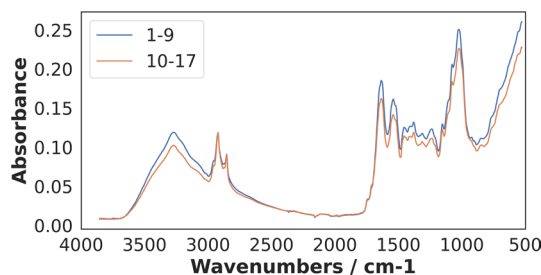


Fig. 1 The Average mid-infrared spectra of dried mosquitoes aged 1–9 days and 10–17 days. The supervised learning was trained on the slight difference between mosquitoes aged 1–9 and 10–17 days

spectra (matrix of features) were used as input data, while the model outputs were the mosquito age classes.

Three different deep learning models were trained; (1) Convolutional neural network (CNN) model without dimensionality reduction, (2) Multi-Layer Perceptron (MLP) with PCA as dimensionality reduction, and (3) MLP with t-SNE as dimensionality reduction. For all models, a SoftMax layer was added to transform the non-normalized outputs of K -units in a fully connected layer into a probability distribution of belonging to either one of two age classes (1–9 or 10–17 days). Moreover, to compute the gradient of the networks, stochastic gradient boosting was used as an optimisation algorithm [22], and categorical cross-entropy loss was used for the classifier’s metric.

To begin, we trained a one-dimensional CNN model with four convolutional layers and one fully connected layer when the dimensionality of the data was not reduced (Fig. 2A), and therefore consisting of 1666 training features from the data. The one-dimensional CNN was used because it is effective at deriving features from fixed-lengths (i.e. the wavelengths of the mid-infrared spectra), and it has been previously been used

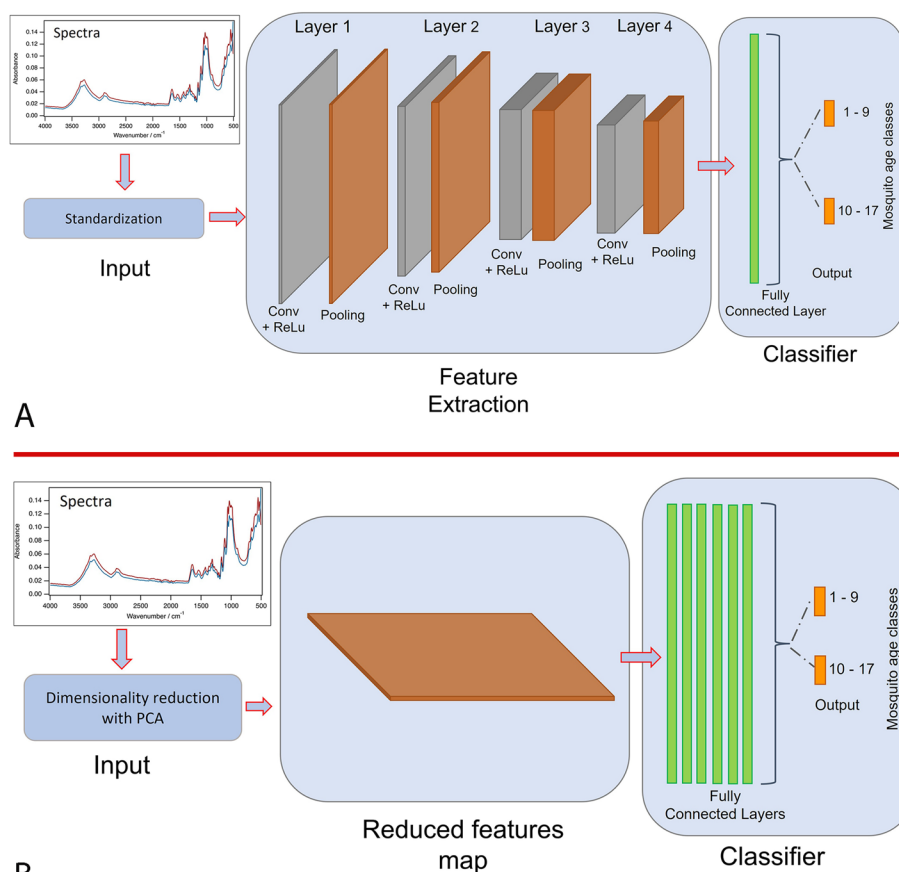


Fig. 2 A schematic representation of a deep learning models that uses mosquito spectra as input to predict mosquito age classes. **A** CNN—no dimensionality reduction is applied: standardised spectral features are fed as input through four different convolutional layers, followed by one fully connected layer, with the predicted age classes shown as the output layer. **B** MLP—dimensionality reduction is used: spectral features that have been reduced in dimension using PCA or t-SNE are fed as input through 6 fully connected layers, with the predicted age classes shown as the output layer

efficiently with spectral data [17]. To extract features from spectral signals, the deep learning architecture used convolutional, max-pooled and fully connected layers. The convolutional operation was carried out with kernel sizes (window) of 8, 4, and 6, and a kernel window shift size (stride) of either 1 or 2. For each kernel size, 16 filters were used to detect and derive features from the input data. Furthermore, given the size of the training data, the fully connected layer consisted of 50 neurons to reduce the model’s complexity.

Moreover, batch normalisation layers were added to both models to improve model stability by keeping mean activation close to 0 and activation standard deviation close to 1. To reduce the likelihood of overfitting, dropout was used during model training to randomly and temporarily remove units from the network at a rate of 0.5 per step. Furthermore, after 50 rounds, early stopping was used to halt training when a validation loss stopped improving.

Dimensionality reduction

We trained two additional deep learning models, in this case Multi-Layer Perceptron (MLP), with PCA or t-SNE transformed input data (Fig. 2B). The models were trained with only fully connected layers (n=6) containing 500 neurons each, given the limited number of training features to ensure performance and stability. To control for overfitting, the procedure was similar to that of the CNN above, except that early stopping was used to halt training when a validation loss stopped improving after 500 rounds.

Transfer learning

The Ifakara dataset was used as the source domain for pre-training the ML models. The Ifakara dataset was divided into training and test sets, and estimator performance was assessed using K-fold cross-validation (k=5) [23], (Fig. 3). We therefore determined what percentage of the new spectra data from the alternate location as target domain was required for ML models to learn the variability between the insectaries. To put transfer learning options to the test, either 82 or 33 spectra were randomly selected from the 1635 of the Glasgow data, accounting for 5% and 2% of the dataset, respectively. The

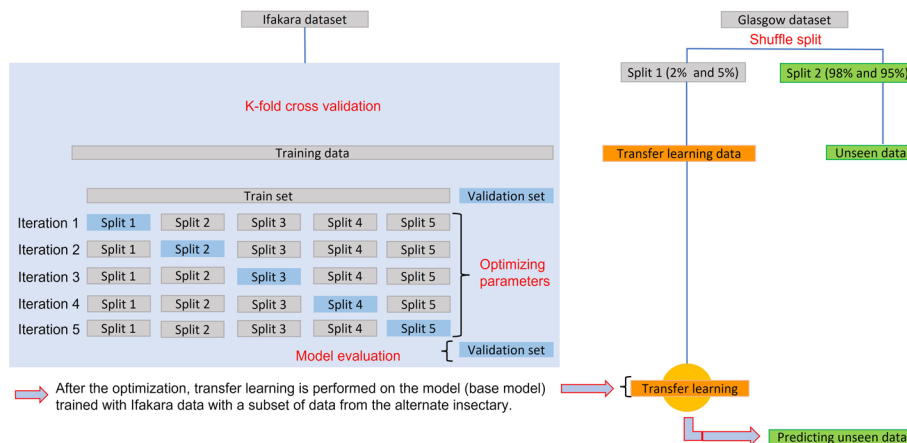


Fig. 3 Schematic illustrating the process of data splitting, model training, cross-validation, and transfer learning

learning process in this case relied on a pre-trained model (trained with Ifakara data), avoiding the need to start training from scratch (Fig. 3). The ML models pre-trained with Ifakara dataset were fine-tuned using 2% or 5% subsets of the Glasgow dataset. The output was compared to that of a model trained solely with Ifakara data (i.e., no transfer learning).

Precision, recall, and F1-scores were calculated from predicted values for each age class to demonstrate the validity of the final models in predicting the unseen Glasgow data. Keras and TensorFlow version 2.0 were used for deep learning process [24, 25].

Standard machine learning

We also compared the prediction accuracy of CNN and MLP to that of a standard machine learning model trained on spectra data transformed by PCA or t-SNE. Different algorithms were compared, including K-Nearest Neighbour, logistic regression, support vector machine classifier, random forest classifier, and a gradient boosting (XGBoost) classifier. The model with the highest accuracy score for predicting mosquito age classes was optimised further by tuning its hyper-parameters with randomised search cross-validation [21]. The cross-validation evaluation used to assess estimator performance in this case was the same as that used in deep learning. The fine-tuned model was used to predict mosquito age classes in previously unseen Glasgow dataset.

Python version 3.8 was used for both the deep learning and standard machine learning training. All computations were done on a computer equipped with 32 Gigabytes of random-access memory (RAM) and an octa-core central processing unit.

Results

DL mosquito age classification with and without dimensionality reduction, did not generalise between the two locations

In the initial analysis, only spectra from the Ifakara insectary were used to train the CNN. During model training, the CNN classifier achieved 99% training accuracy without any dimensionality reduction (Fig. 4A). When given new held-out data from the same Ifakara insectary (test set), the model predicted mosquitoes aged 1–9 days with 100% accuracy and those aged 11–17 days with 99% accuracy (Fig. 4B). However, when the same model was used to predict age classes for Glasgow insectary samples, the

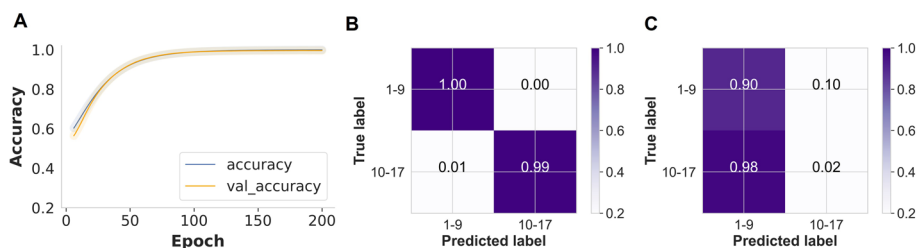


Fig. 4 CNN generalisation and prediction of mosquito age using data from a single insectary (Ifakara) with no dimensionality reduction. **A** Training and validation classification accuracy for mosquito age classes improved from ~60 to 95% as training iterations increased (200 epochs). **B** A normalised confusion matrix displaying the proportions of correct mosquito age class predictions achieved on the held-out Ifakara data (test set) during model training. **C** Proportions of correct mosquito age class predictions based on unseen data from the alternate insectary (Glasgow)

overall accuracy was 46%, and therefore indistinguishable from any random classifications (Fig. 4C).

In addition, a CNN classifier required 200 epochs for training, with a running time of 7.2–7.8 s per epoch when no dimensionality reduction on the input data was used (Table 1).

PCA was used to project the data into lower dimensional space using singular value decomposition [13, 26], with the goal of achieving the best summary using optimal number of principal components (PCs) with up to 98% of variance explained (Fig. 5A). Further, when the impact of PCs on accuracy was assessed, a greater prediction accuracy was found, leading to the selection of 8 PCs. (Fig. 5B).

When PCA was used to reduce the dimensionality of the data, the MLP model trained with only Ifakara spectra predicted the held-out data from the same insectary (Ifakara) with an overall accuracy of 91% but could attain only 58% accuracy for predicting age classes of Glasgow mosquitoes (Table 1). Similarly, when t-SNE was used as the dimensionality reduction technique, the model predicted the held-out Ifakara data (test set) with an accuracy of 85% but failed to accurately predict age classes of Glasgow data (Table 1).

Furthermore, when PCA or t-SNE were used to transform the input data, a MLP classifier needed 5000 epochs to train, with a running time of 0.7–0.8 s per epoch (Table 1).

Transfer learning improves DL accuracy and generalisability

To improve generalisability (i.e., the ability of the models to predict the age classes of samples from other sources), we tuned the pre-trained CNN models with 2% or 5% of the spectra from Glasgow (i.e., 2% or 5% target population samples for transfer learning) and used the updated model to predict the unseen Glasgow dataset. When no dimensionality reduction was used, the pre-trained model predicted the held-out test (Ifakara dataset) with 99% accuracy and transferred well to the Glasgow dataset when 2% and 5% target population samples were used for transfer learning, achieving 100% and 96% accuracies, respectively (Table 1).

However, when PCA or t-SNE were used to reduce the dimensionality of the data, the MLP classifier was trained with only fully connected layers in this case to allow the model to learn the combination of features with the network's learnable weights. Using PCA, the pre-trained model predicted the held-out test (Ifakara dataset) with 91% accuracy, but when 2% transfer learning was applied, the model transferred well to the Glasgow dataset, achieving 97% accuracy when predicting the mosquito age classes, and 96% accuracy with 5% target population samples (Table 1, Fig. 6A–C).

When using t-SNE, the pre-trained predicted the age classes in the held-out data (test set) with 83% accuracy but failed to achieve generalisability for the Glasgow data when either 2% or 5% transfer learning was applied, achieving only 50% and 55% accuracy, respectively (Table 1, Fig. 6D–F).

Transfer learning also reduced training time while improving the performance of both DL and standard machine learning models in predicting samples from the target population. Transfer learning took less than two minutes for both models to produce the desired results (Table 1).

Table 1 The performance of deep learning and standard machine learning models for predicting mosquito age classes from the same or alternate insectaries, with and without dimensionality reduction (DR) and transfer learning

Models	Dimensionality reduction (DR) technique	Training data sources	Transfer learning	Base Model runtime	Transfer learning runtime	Predictions for age of mosquitoes from same insectary (Ifakara) -Test accuracy (%)	Predictions for age of mosquitoes from alternate insectary (Glasgow)—unseen data accuracy (%)
CNN-1	No DR	Ifakara	No TL	7.2 s/iteration	N/A	99	46
CNN-2	No DR	Ifakara	2% (33 of 1635)	7.2 s/iteration	1 min	99	100
CNN-3	No DR	Ifakara	5% (82 of 1635)	7.8 s/iteration	2 min	99	96
MLP-1	PCA	Ifakara	No TL	6.5 s/iteration	N/A	91	58
MLP-2	t-SNE	Ifakara	No TL	1 s/iteration	N/A	84	58
MLP-3	PCA	Ifakara	2% (33 of 1635)	0.8 s/iteration	35 s	91	97
MLP-4	PCA	Ifakara	5% (82 of 1635)	0.7 s/iteration	51 s	91	96
MLP-5	t-SNE	Ifakara	2% (33 of 1635)	0.7 s/iteration	47 s	83	50
MLP-6	t-SNE	Ifakara	5% (82 of 1635)	0.7 s/iteration	49 s	83	55
XGB-1	No DR	Ifakara	No TL	645 s/iteration	N/A	92	48
XGB-2	No DR	Ifakara	2% (33 of 1635)	975 s/iteration	1 s	92	98
XGB-3	No DR	Ifakara	5% (82 of 1635)	861 s/iteration	1 s	92	98
XGB-4	PCA	Ifakara	No TL	60 s/iteration	N/A	90	48
XGB-5	t-SNE	Ifakara	No TL	66 s/iteration	N/A	68	55
XGB-6	PCA	Ifakara	2% (33 of 1635)	54 s/iteration	1 s	90	98
XGB-7	PCA	Ifakara	5% (82 of 1635)	54 s/iteration	2 s	90	97
XGB-8	t-SNE	Ifakara	2% (33 of 1635)	60 s/iteration	1 s	81	43
XGB-9	t-SNE	Ifakara	5% (33 of 1635)	60 s/iteration	1 s	82	49

*CNN—1 to 3: Different versions of convolutional neural network, MLP—1 to 6: Different versions of Multi-Layer Perceptron, XGB-1 to 9: Different versions of XGBoost classifier (standard machine learning), *No DR*: No dimensionality reduction, *PCA*: Principal component analysis, *t-SNE*: t-distributed stochastic neighbour embedding, *No TL*: No Transfer learning, *N/A*: Not applicable. The highest prediction accuracy as a result of transfer learning with less computational time is shown in the bold

Comparison between deep learning and standard machine learning models in achieving generalisability

The XGBoost classifier (Fig. 7A), when trained with Ifakara data only, failed to predict age classes of mosquitoes from the Glasgow insectary, with or without dimensionality

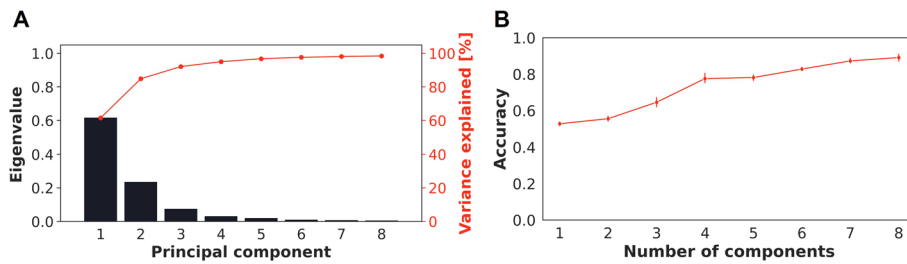


Fig. 5 **A** cumulative explained variance and eigenvalues as the function of principal components. **B** Number of principal components included in the XGB classifier (i.e. from 1:8 PCs)

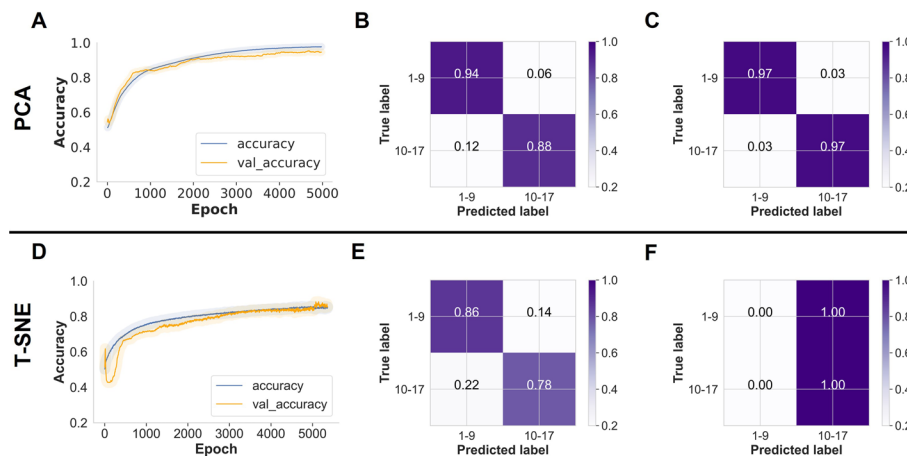


Fig. 6 MLP trained on PCA-transformed Ifakara dataset plus 2% new target population samples: **A** As training time increased (5000 epochs), training and validation classification accuracy for mosquito age classes increased from 50 to 91%, **B** A normalised confusion matrix displaying the proportions of correct mosquito age class predictions achieved on the held-out Ifakara test set during model training, **C** Proportions of correct mosquito age class predictions achieved on unseen Glasgow dataset. MLP trained on t-SNE-transformed Ifakara dataset plus 2% new target population samples: **D** As training time increased (5000 epochs), training and validation classification accuracy for mosquito age classes increased from 60 to 83%, **E** A normalised confusion matrix displaying the proportions of correct mosquito age class predictions achieved on the held-out Ifakara test set during model training, **F** Proportions of correct mosquito age class predictions achieved on unseen Glasgow dataset

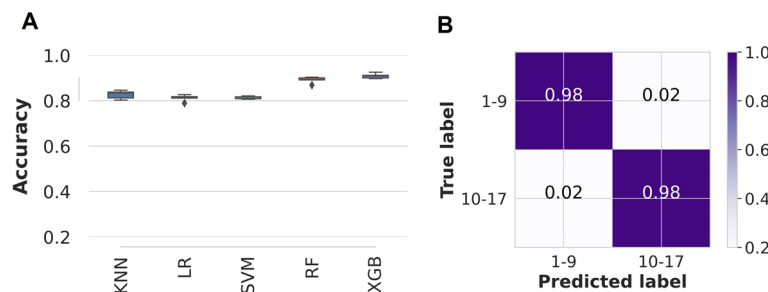


Fig. 7 Standard machine learning models' predictive accuracies and generalisability when trained with PCA-transformed Ifakara data plus 2% new target population. **A** Comparison of standard machine learning models for mosquito age classification; *KNN*: K-nearest neighbours, *LR*: Logistic regression, *SVM*: Support vector machine classifier, *RF*: Random Forest classifier, and *XGB*: XGBoost. **B** proportions of correct mosquito age class predictions achieved on unseen Glasgow dataset

reduction (Table 1). However, when the classifier was updated with 2% target population samples, the model correctly classified individual mosquito age classes with 98% for both 1–9 days old and 10–17 days old mosquitoes (Fig. 7B). Increasing the samples for transfer learning to 5% of the training set had no effect on the accuracies (Table 1). However, when t-SNE was used for dimensionality reduction, transfer learning with either 2% or 5% Glasgow samples did not improve the generalisability of the XGBoost classifier (Table 1).

Table 2 shows how the performance of deep learning and standard machine learning was evaluated using other metrics such as precision, recall, and f1-scores. When it comes to mosquito age classification, the XGBoost classifier matches the deep learning model in both specificity (precision) and sensitivity (recall).

Further to that, standard machine learning models were trained with 10 iterations, and still the computing runtimes were generally shorter than those for CNN models when PCA and t-SNE were used to transform the input data, in some cases by up to 5 times (Table 1).

Discussion

This study demonstrates that transfer learning approaches can substantially improve the generalisability of both deep learning and standard machine learning in predicting the age class of mosquitoes reared in two different insectaries. We evaluated 1635 mosquito spectra from Glasgow-reared mosquitoes and show that using transfer learning and dimensionality reduction techniques could improve machine learning models to predict mosquito age classes from alternate insectaries. Furthermore, reducing the dimensionality of the spectral data reduced computational costs (i.e. computing time) when training the machine learning models.

The current study adds to the growing evidence of the utility of infrared spectroscopy and machine learning in estimating mosquito age and survival [8, 27–29]. In the past, most applications of infrared spectroscopy in estimating mosquito vector survival relied on near-infrared frequencies (12,500–4000 cm^{-1}). A recent study used mid-infrared spectra (from 4000 to 400 cm^{-1} frequencies) and standard machine learning to distinguish mosquito species with up to 82% accuracy, but found lower age prediction accuracy in several alternate settings [9]. González et al., suggested

Table 2 Precision, recall, and f1-score of the best deep learning model for classifying mosquito age classes from alternate sources compared to the best standard machine learning algorithm (i.e. XGBoost classifier)

Model name	Age class (Days)	Precision	Recall	f-1 score	No. of samples per age class
MLP-3	1–9	0.98	0.97	0.98	895
	10–17	0.97	0.97	0.97	707
XGB-6	1–9	0.98	0.99	0.98	895
	10–17	0.98	0.98	0.98	707

*MLP-3: Multi-Layer Perceptron trained with PCA as a dimensionality reduction technique and 2% transfer learning, XGB-6: XGBoost classifier trained with PCA as a dimensionality reduction technique and 2% target population samples used for transfer learning

that machine learning underprediction may be explained by the small training dataset and ecological variability between the training and validation sets [9].

In our study, despite categorising mosquito chorological age into two classes (young: 1–9-day olds and old: 10–17-day olds), deep learning and standard machine learning approaches both remained unable to generalise, even after reducing the dimensionality of the spectra data. This result is consistent with Siria et al. [10], where CNN underperformed as a result of the difference in data distribution between the training and evaluation data driven by non-genetic factors such as ecological variation. When near-infrared spectroscopy was used to predict the age of *Anopheles* mosquitoes reared from wild populations, a similar limitation was reported [8, 27].

Nonetheless, Siria et al. [10] also observed that using transfer learning to correct the difference data distribution between training and evaluation data improved deep learning generalisation, achieving 94% accuracy in predicting both species and mosquito age classes. Furthermore, in the latter study, the performance of the classifier was improved by incorporating a subset ($n = 1200$ – 1300 spectra) of the evaluation data into the training data.

The present study shows performing transfer learning using 2% of the spectra from the target domain (33 of 1635) as well as dimensionality reduction resulted in the improved generalisability of both deep learning and standard machine learning models achieving overall accuracy of $\sim 98\%$. In this case, we expected that all models to which transfer learning was applied would outperform the baseline models as previously demonstrated [10, 30]. However, as the proportion of data from the target domain in the training increased, the performance slightly dropped for the deep learning. The reason for the deterioration in performance after turning the pre-trained/base model with 5% transfer learning could be that the model overfitted random noise during training, which negatively impacted the performance of these models on unseen data. Other studies have proposed alternative transfer learning approaches, such as adaptive regularisation to address cross-domain (i.e. source domain and target domain) learning problems [31], transferring knowledge gained in the source domain during training to the target domain [32], and integrating dimensionality reduction to transform features of the source to ensure data distribution in different domains is minimised [33], such as transfer learning with multi-target regression approach to exploit orthologous genes to capture similarities in metabolic responses in mice and humans [34, 35].

Furthermore, dimensionality reduction was used in conjunction with transfer learning to reduce noise, redundant features, and computational time. Based on our findings, dimensionality reduction alone cannot achieve generalisability of machine learning models. The PCA improved model stability because the eigenvectors of the correlation matrix in PCA provide new axes of variation to project new data while preserving the original distance between the points in the data. The model with t-SNE as a dimensionality reduction technique failed to achieve generalisability on the new data, the reason for poor performance could be t-SNE is a probabilistic technique with a non-convex cost function [16], causing the output to differ from multiple runs, and may not preserve the original distances between the points in the data. In this study, PCA is considered a better choice than other dimensionality reduction

technique for training machine learning models from spectra data because it is simple to implement, computationally efficient, and produces good results.

Furthermore, incorporating dimensionality reduction substantially reduces model training time and thus, computational requirements. When compared to models trained without dimensionality reduction, the computing runtimes for models trained with dimensionality reduction were less than five-fold. Moreover, transfer learning in general was fast, tuning the pre-trained models in under two minutes on our machine (standard laptop). This makes the technique applicable and reproducible even to users with low computing power and capacity providing they have access to pre-trained models.

This study only included *An. arabiensis* reared in the laboratory from two insectaries. Future research should put the techniques to the test with samples from more laboratories, field settings, and mosquito species, as these factors can affect the model's predictive capacity. The optimal ratio of transfer learning data required to achieve best generalisability in predicting mosquito age class has yet to be determined, so future studies could investigate this gap. Furthermore, because dimensionality reduction reduced the computational requirements in this study, we suggest that clustering spectra with algorithms such as PCA can be a beneficial strategy for models trained on MIRS.

Conclusion

This study found that using transfer learning and dimensionality reduction with principal component analysis (PCA) improved the generalizability of machine learning models for predicting mosquito age classes from 56 to $\geq 97\%$. This suggests that these techniques could be scaled up and further evaluated to determine the age of mosquitoes from different populations. In addition, when dimensionality reduction and transfer learning are used, simpler machine learning algorithms, such as the XGBoost classifier, can reduce computational time while still achieving performance close or equal to deep learning. This could help entomologists reduce the amount of time and work required to dissect large numbers of mosquitoes. Overall, these approaches have the potential to improve model-based surveillance programs, such as assessing the impact of malaria vector control tools, by monitoring the age structures of local vector populations.

For future research, our goal is to create a large database of spectra data and use transfer learning to build a pipeline that can predict the age of wild malaria mosquitoes across different populations in order to support vector surveillance in malaria-endemic areas. Here we have presented a new technique that uses transfer learning and dimensionality reduction to improve the generalizability of machine learning predictions. However, the optimal proportion of new data from target populations required for generalizability is still unknown, and warrants further optimisation.

Abbreviations

CNN	Convolutional neural network
ITNs	Insecticide treated nets
PCA	Principal component analysis
t-SNE	T-distributed stochastic neighbour embedding

Acknowledgements

The authors express their gratitude to for the research and administrative support teams at Ifakara Health Institute and the University of Glasgow.

Author contributions

EPM, SAB, DJS, FB, MGJ and FOO designed the study. DJS supported in data collection semi-field experiments. EPM performed data analysis. JM provided technical support to EPM during data analysis. EPM wrote and revised the manuscript. EPM, SAB, IHM, PS, FOO, and FB reviewed and revised the manuscript. All authors read and approved the final manuscript.

Funding

This research was supported by the Medical Research Council (MRC) grant (Grant No. MR/P025501/1). EPM and DJS were also supported by the Wellcome Trust International Masters Fellowships in Tropical Medicine and Hygiene, Grant Nos. WT214643/Z/18/Z and WT 214644/Z/18/Z respectively.

Availability of data and materials

The mid-infrared spectral data generated and/or analysed during the current study are deposited and available in the Enlighten database at <https://doi.org/10.5525/gla.researchdata.1235>.

Declarations**Ethical approval and consent to participate**

All methods in this study were performed in accordance with the relevant guidelines and regulations from IHI, National Institutes of Medical Research (NIMR), and UofG. At IHI, ethical approval for the study was obtained from the IHI Institutional Review Board (Ref. IHI/IRB/EXT/No: 005-2018), and NIMR Ref: NIMR/HQ/R.8c/Vol.II/880. At the UofG, Ethical approval for the supply and use of human blood for mosquito feeding was obtained from the Scottish National Blood Transfusion Service committee for governance of blood and tissue samples for non-therapeutic use (submission Reference No 1815).

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 26 July 2022 Accepted: 26 December 2022

Published online: 09 January 2023

References

1. WHO. World Malaria report 2021. 2021.
2. Bhatt S, Weiss DJ, Cameron E, Bisanzio D, Mappin B, Dalrymple U, et al. The effect of malaria control on *Plasmodium falciparum* in Africa between 2000 and 2015. *Nature*. 2015;526:207–11.
3. World Health Organization. Global Malaria Programme. Global technical strategy for malaria 2016–2030 [Internet]. World Heal. Organ. 2015. Available from: http://apps.who.int/iris/bitstream/10665/176712/1/9789241564991_eng.pdf?ua=1.
4. Charlwood JD, Smith T, Billingsley PF, Takken W, Lyimo EOL, Meuwissen JHET. Survival and infection probabilities of anthropophagic anophelines from an area of high prevalence of *Plasmodium falciparum* in humans. *Bull Entomol Res*. 1997;87:445–53.
5. Detinova TS. Age-grouping methods in Diptera of medical importance with special reference to some vectors of malaria. In: Monogr Ser World Health Organ. Geneva: World Health Organization; 1962.
6. Polovodova PV. The determination of the physiological age of female *Anopheles* by the number of gonotrophic cycles completed. *Medskaya Parazit*. 1949;18:352–5.
7. Rao V. On gonotrophic discordance among certain Indian anopheles. *Indian J Malariol*. 1947;1:43–50.
8. Mayagaya VS, Michel K, Benedict MQ, Killeen GF, Wirtz RA, Ferguson HM, et al. Non-destructive determination of age and species of *Anopheles gambiae* sl using near-infrared spectroscopy. *Am J Trop Med Hyg*. 2009;81:622.
9. Gonzalez-Jimenez M, Babayan SA, Khazaeli P, Doyle M, Walton F, Reedy E, et al. Prediction of malaria mosquito species and population age structure using mid-infrared spectroscopy and supervised machine learning. *Wellcome Open Res*. 2019;4:76.
10. Siria DJ, Sanou R, Mitton J, Mwanga EP, Niang A, Sare I, et al. Rapid age-grading and species identification of natural mosquitoes for malaria surveillance. *Nat Commun*. 2022;13:1501. <https://doi.org/10.1038/s41467-022-28980-8>.
11. Esperança PM, Da DF, Lambert B, Dabiré RK, Churcher TS. Functional data analysis techniques to improve the generalizability of near-infrared spectral data for monitoring mosquito populations. *bioRxiv*. 2020;2020.04.28.058495. Available from: <http://biorxiv.org/content/early/2020/04/29/2020.04.28.058495.abstract>.
12. Géron A. Hands-on machine learning with scikit-learn and TensorFlow. First Edit. Tache N, editor. Boston: O'Reilly Media, Inc.; 2017.
13. Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemom Intell Lab Syst*. 1987;2:37–52.
14. Lever J, Krzywinski M, Altman N. Principal component analysis. *Nat Methods*. 2017;14:641–2. <https://doi.org/10.1038/nmeth.4346>.
15. Schölkopf B, Smola A, Müller K-R. Kernel principal component analysis. In: International conference on artificial neural networks. Berlin, Heidelberg: Springer; 1997. p. 583–8.
16. Van Der Maaten LJP, Hinton GE. Visualizing high-dimensional data using t-sne. *J Mach Learn Res*. 2008;9:9.

17. Siria DJ, Sanou R, Mitton J, Mwanga EP, Niang A, Sare I, et al. Rapid ageing and species identification of natural mosquitoes for malaria surveillance. *bioRxiv*. 2020;2020.06.11.144253. Available from: <http://biorxiv.org/content/early/2020/06/12/2020.06.11.144253.abstract>.
18. Mwanga EPP, Mapua SAA, Siria DJJ, Ngowo HSS, Nangacha F, Mgando J, et al. Using mid-infrared spectroscopy and supervised machine-learning to identify vertebrate blood meals in the malaria vector *Anopheles arabiensis*. *Malar J*. 2019;18:187. <https://doi.org/10.1186/s12936-019-2822-y>.
19. Jiménez MG. A custom program that imports the IR spectra, cleans and screens them to eliminate the badly measured ones, and extracts the most interesting data from them!. *Wellcome Open Res*. 2019. p. 4:76. Available from: https://github.com/SimonAB/Gonzalez-Jimenez_MIRS/blob/v1.0/Locomosquito.ipynb.
20. Ohm JR, Baldini F, Barreaux P, Lefevre T, Lynch PA, Suh E, et al. Rethinking the extrinsic incubation period of malaria parasites. *Parasit Vectors*. 2018;11:178. <https://doi.org/10.1186/s13071-018-2761-4>.
21. Pedregosa F, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, et al. Scikit-learn: machine learning in Python. *J Mach Learn Res*. 2011;12:2825–30.
22. Sutskever I, Martens J, Dahl G, Hinton G. On the importance of initialization and momentum in deep learning. In: 30th international conference on machine learning (ICML 2013), 2013.
23. Kohavi R. A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection. In: International joint conference on artificial intelligence. 1995.
24. Chollet F. Keras: The Python deep learning library. *Keraslo*. 2015.
25. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, et al. TensorFlow: a system for large-scale machine learning. In: Proceedings of the 12th USENIX symposium on operating systems design and implementation, OSDI 2016. 2016.
26. Halko N, Martinsson PG, Tropp JA. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev*. 2011;53:217–88. <https://doi.org/10.1137/090771806>.
27. Lambert B, Sikulu-Lord MT, Mayagaya VS, Devine G, Dowell F, Churcher TS. Monitoring the age of mosquito populations using near-infrared spectroscopy. *Sci Rep*. 2018;8:5274.
28. Sikulu-Lord MT, Devine GJ, Hugo LE, Dowell FE. First report on the application of near-infrared spectroscopy to predict the age of *Aedes albopictus* Skuse. *Sci Rep*. 2018;8:1–7.
29. Ntamatungiro AJ, Mayagaya VS, Rieben S, Moore SJ, Dowell FE, Maia MF. The influence of physiological status on age prediction of *Anopheles arabiensis* using near infra-red spectroscopy. *Parasites Vectors*. 2013;6:298. <https://doi.org/10.1186/1756-3305-6-298>.
30. Hanczar B, Bourgeois V, Zehraoui F. Assessment of deep learning and transfer learning for cancer prediction based on gene expression data. *BMC Bioinform*. 2022;23:262. <https://doi.org/10.1186/s12859-022-04807-7>.
31. Long M, Wang J, Ding G, Pan SJ, Yu PS. Adaptation regularization: a general framework for transfer learning. *IEEE Trans Knowl Data Eng*. 2014;26:1076–89.
32. Si S, Tao D, Geng B. Bregman divergence-based regularization for transfer subspace learning. *IEEE Trans Knowl Data Eng*. 2010;22:929–42.
33. Pan SJ, Tsang IW, Kwok JT, Yang Q. Domain adaptation via transfer component analysis. *IEEE Trans Neural Netw*. 2011;22:199–210.
34. Pio G, Mignone P, Magazzù G, Zampieri G, Ceci M, Angione C. Integrating genome-scale metabolic modelling and transfer learning for human gene regulatory network reconstruction. *Bioinformatics*. 2022;38:487–93. <https://doi.org/10.1093/bioinformatics/btab647>.
35. Mignone P, Pio G, Džeroski S, Ceci M. Multi-task learning for the simultaneous reconstruction of the human and mouse gene regulatory networks. *Sci Rep*. 2020;10:22295. <https://doi.org/10.1038/s41598-020-78033-7>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

