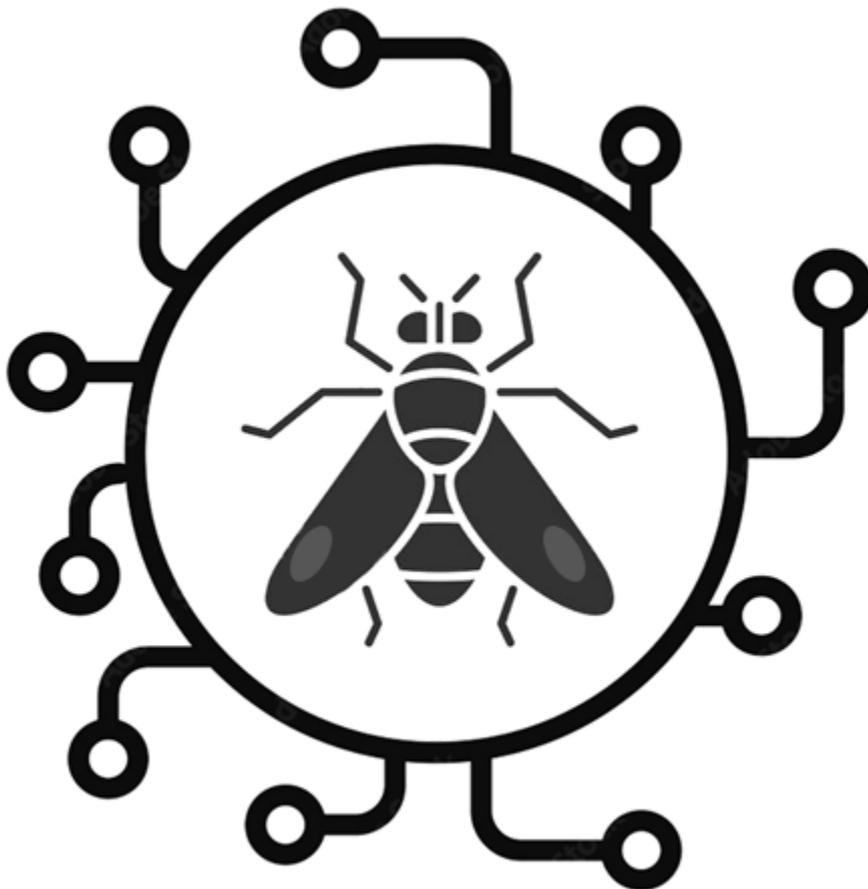


SWD-FREE

A DST for predicting Spotted wing drosophila abundance

Final Report



Peter Skelsey
28/02/2021



The James
Hutton
Institute

Science connecting land and people

Project title: SWD-FREE: A decision support tool for predicting spotted wing drosophila abundance

Project number: SF/TF 145b

Project leader: Peter Skelsey, James Hutton Institute

Report: Final report, 28/02/2022

Previous report: n/a

Key staff: n/a

Location of project: James Hutton Institute, Invergowrie

Industry representative: Scott Raffle, scott.raffle@niab.com

Date project commenced: 1/10/2021

Date project completed: 28/02/2022

TABLE OF CONTENTS

GROWERS SUMMARY	1
SCIENCE SECTION	2
METHODOLOGY	2
DATASETS	2
CLASSIFICATION OF SPRING RISK	2
PREDICTING KEY PERCENTILES OF SWD ACTIVITY	2
MODELLING POPULATION DYNAMICS OVER THE WHOLE SEASON	4
RESULTS	4
CLASSIFICATION OF SPRING RISK	4
PREDICTING KEY PERCENTILES OF SWD ACTIVITY	6
MODELLING POPULATION DYNAMICS OVER THE WHOLE SEASON	13
DISCUSSION	17

Growers Summary

Background

The spotted wing drosophila (*Drosophila suzukii*) is an invasive pest of soft and stone fruit crops and, if left uncontrolled, can result in complete crop loss. Evidence from other countries has shown that early detection and rapid response is crucial to minimising the impact of SWD on soft and stone fruit crops. A successful programme of population monitoring (SF/TF 145a) has confirmed the presence of the pest in Great Britain, and consequently, development of a decision support tool for timely applications of plant protection products is now identified as a top priority for the soft fruit industry.

Main objectives

The main objectives of this project were to develop and test: (i) site-specific population models of abundance (all years of data for each location combined); (ii) a site-nonspecific population model of abundance (all locations and years of data combined); (iii) a range of machine learning (ML) algorithms for predicting both site-specific and site-nonspecific abundance. The aim is for final model to be hosted on a dedicated AHDB webpage as a decision support tool for predicting when key levels of population abundance are reached.

Research undertaken

SWD abundance data from 16 locations spanning 2014-2020 were integrated with UK Met. Office temperature data and modelled using three different approaches: population modelling (curve-fitting), statistical modelling, and ML. Model fitting was conducted over a range of geographic groupings to determine the most informative / pragmatic spatial scale for development and deployment of a decision support tool: site-, region-, country- and national-scale.

Key findings

- A national-scale ML algorithm was developed for forecasting risk of SWD activity in spring from weather data. It achieved an overall predictive accuracy of 80% and provides support for decisions on the need to start crop protection measures at the beginning of the season (March 1).
- Regression models were derived for predicting key levels of population abundance that precipitate crop protection measures. The day-of-season on which 5% of population abundance was reached was well predicted at the national-scale with an average error of 3.15 days, increasing to 5.96 days for the day-of-season on which 50% of population abundance was reached.
- The SWD data showed large variation between capture sites and years for the timing of the exponential phase of SWD activity (i.e., 0-5% of SWD captures). It was therefore not possible to develop robust predictive tools to forecast this important phase of activity.
- Percentile capture charts were plotted over a range of spatial scales to serve as decision support tools. They show observed values (and uncertainty) for accumulated degree-days that mark key

levels of population abundance and can be used to identify the optimum moments to apply control measures to coincide with these dates.

- Population models were fit to the data to provide growers with tools that can be used to evaluate management decisions across the whole growing season. Model fits were excellent, with average R^2 values of 0.95, 0.91, 0.92 and 0.87 at the site-, region-, country-, and national-scales.

Science Section

Methodology

Datasets

SWD capture data (male + female) from crops in 2013-2020 were available from a total of 16 sites in Scotland and England. Data from all sites for 2013 were excluded from analysis as there were no observations in the first 8 months of the year. Data from 2021 were excluded as there were no corresponding weather data to construct predictor variables. This left a total of 88 datasets for analysis.

The coordinates of capture locations were used to match each site to its nearest weather datapoint in the UK Met. Office Best Data database (a total of 3600 locations across the UK, providing hourly weather variables spanning 2012-2020).

Classification of spring risk

This modelling task aimed to provide growers with a tool that can forecast if SWD activity is expected in spring (March-April-May, MAM), and thus if there is a need to start applying crop protection measures from the beginning of the growing season (March 1). Total MAM captures were calculated for each dataset and discretized to binary risk variables, where 0 = no SWD captured (no risk, no action required) and 1 = SWD captured (SWD risk, action required). Hourly UKMO weather data for January and February were used to provide the following predictor variables: minimum temperature, maximum temperature, temperature sum, hours of relative humidity $\geq 90\%$, average wind speed, average wind gust, total cloudiness, total sunshine, and total precipitation. The latitude and longitude of each site provided two further predictor variables. A Decision Tree algorithm was built to predict risk of MAM activity from the predictor variables, using the MATLAB procedure `fitctree` implemented within a nested k -fold cross-validation procedure with Bayesian optimization for model tuning and selection. The hyperparameter `MaxNumSplits` was set to 10 to avoid deep (complex) trees.

Predicting key percentiles of SWD activity

The second modelling task investigated if the timing of key percentiles of SWD captures could be predicted using accumulated degree-days. Percentiles are values that divide a set of observations into 100 equal parts, e.g., the 5-percentile point is the value below which 5% of data falls. Datasets with very few / low captures were removed from this analysis as there were not enough datapoints to provide robust results for this task. Any site with less than a total of 50 *Drosophila* captured in a year was excluded. Data for site 10a for 2014 and 2015 were also removed as there were no recorded observations for the second halves of those years. Data for site 1300 in 2020 could not be

used as there were no weather data for the 2021 component of the season. After data cleaning, there were a total of 59 unique site-year datasets for the calculation of percentile capture points.

The day of the season of the 5- and 50-percentile of captures was calculated by interpolation using Piecewise Cubic Hermite Interpolating Polynomials. Hourly screen temperature values were used to calculate accumulated degree-days for each day of capture in the data. To do so it was necessary to define a 'starting point' for the season, i.e., a biofix date that signals the start of degree-day accumulations. This was problematic as numerous datasets contained low levels of sporadic captures throughout the year. Each dataset was analysed visually and mathematically to determine local minima in trap captures that signalled the start and end of peak activity periods. A rule was developed whereby the first observed capture after April 1 that was $\geq 1\%$ of the maximum capture value was defined as the start of seasonal SWD activity. This ensured that the start of activity coincided with increase in abundance and was not influenced by occasional sporadic captures between cropping seasons. The median first day of capture for all site-year combinations was day 195, which reflected the seasonality in trap catches (peak abundance in Autumn-Winter and decline over Spring). The median first day of capture was used to create a vector of potential biofix dates for examination in the modelling analyses: 18/06, 02/07, 16/07, 30/07, 13/08, and a season was defined as spanning a 365-day period (366 for seasons that spanned a leap year) from the biofix date.

Degree-days were calculated using the continuous integration method on hourly temperature values. This is considered more accurate than the standard approximation methods as it accounts for temperature variations within each day. Degree-days require a baseline temperature below which it is assumed no development occurs. These vary widely for SWD in the literature therefore a heuristic approach was adopted matching that of the biofix date, where a range of potential baseline temperatures (11 values ranging from 0 to 10 degrees) was used to create 11 different degree-day predictor variables for testing. Note that by convention, any temperature below the baseline was set equal to the baseline, but no upper temperature threshold (for development) was set. Model fitting / training proceeded by iterating over each biofix-baseline temperature combination to find the optimum setting for peak predictive performance.

Degree-days accumulated above the vector of 11 baseline values starting from the five biofix dates to the day-of-season of 5- and 50-percentile captures were calculated. Note that 'day-of-season' (starting at the biofix date) and not 'day-of-year' is used, as the period of peak SWD activity sometimes spanned consecutive years. Day-of-season of percentile points were then regressed on accumulated degree-days. To estimate the generalization power of the optimal model, i.e., the ability to predict on new, unseen data, a hold-out cross-validation procedure was used. Data for 2018 (approximately 20% of the full dataset) were kept aside for model testing, with the remainder used for fitting. The optimal settings (biofix-baseline combination) to predict the percentile captures of the training data were determined by calculating the mean absolute error (MAE) between all predicted and observed days-of-season for percentile captures and choosing the values that gave the smallest MAE. Results for the 'optimal model' on the hold-out test set are provided. The optimal model was then refit using all the data and goodness of fit results for each group member are provided. Note that the region- and national-scales were used in this analysis.

A total of 24 machine learning algorithms were also used to predict percentile capture points from accumulated degree days to investigate if improvements in predictive performance over simple regression could be attained. A hold-out cross-validation procedure (75% train, 25% test, partitioned at random) was performed to test the generalization ability of the algorithms.

'Percentile capture charts' were plotted to serve as decision support tools. Mean values for accumulated degree-days (using the optimal biofix-baseline settings) up to the day-of-season of the 5- and 50-percentile captures were used to create line graphs that identify the optimum moments to apply control measures to coincide with these dates.

Modelling population dynamics over the whole season

The third modelling task aims to provide growers with a model of population dynamics that can be used to evaluate management decisions across the whole growing season. Trap captures from the 59 site-year combinations with sufficient data (described above) were accumulated at each site (from each biofix date) and computationally normalized (division by total capture) to produce proportional captures ranging from 0 to 1. Proportional captures were analysed relative to accumulated degree days using six classic models from population ecology: the monomolecular, logistic, log-logistic, Weibull, Gompertz and Richards equations. This made a total of 330 fits for each dataset (5 biofix-11 degree-day baselines-6 models). In addition, the models were fit using four scales of analysis, where observations were grouped by site (16), region (East Scotland, Yorkshire & Humber, West Midlands, East Midlands, East of England, and South East England), country (Scotland, England), and nation (not grouped). This made a total of 5280, 1980, 660, and 330 model fits, respectively. The optimal biofix-baseline-model combination to predict proportional captures at each scale of analysis (best overall for site-, region-, country- and national-scales) was assessed using the coefficient of determination. Separate fit results for the optimal combination are provided for each group within each scale.

A total of 24 machine learning algorithms were also used to predict the proportional capture data to investigate if improvements in predictive performance over population models could be attained. A hold-out cross-validation procedure (75% train, 25% test, partitioned at random) was performed to test the generalization ability of the algorithms.

Results

Classification of spring risk

Approximately 40% of all datasets had a total spring (MAM) capture of 0, providing a balanced dataset for learning. The Decision Tree ML algorithm was able to predict whether SWD would be captured in spring (spring risk) with an overall predictive accuracy of 80% on the hold-out test folds of the nested k -fold procedure (Fig. 1). The optimal hyperparameters determined by the Bayesian optimization method were `MinLeafSize = 4`, `SplitCriterion = deviance`, and `NumVariablesToSample = 4`. The variables longitude, average temperature, minimum temperature, and average wind speed were used in the final model, and their relative importance as predictors of spring SWD risk is given in Fig. 2. Fig. 3 provides a graphic description of the final model.

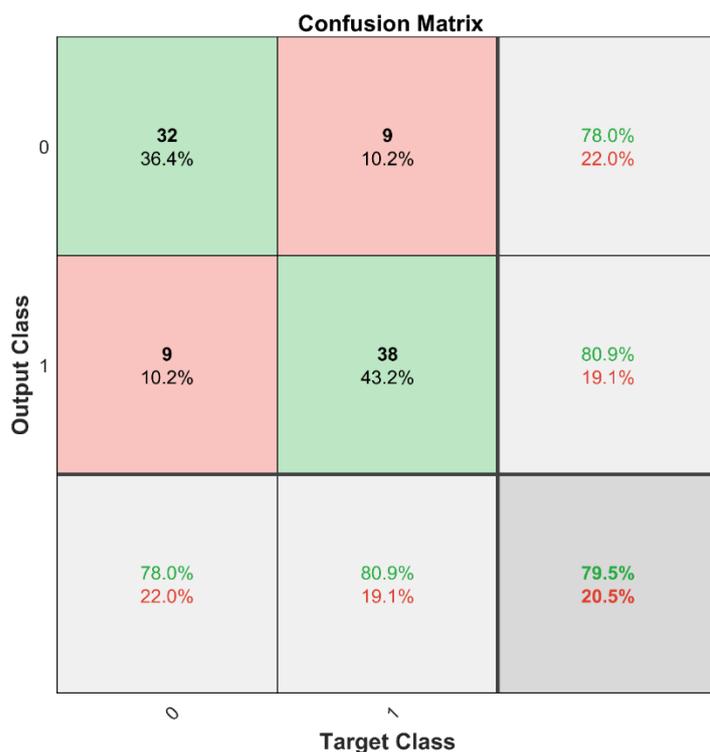


Fig. 1. Confusion matrix of spring SWD risk, where output class = predicted class, target class = true class, 0 = no captures in spring, and 1 = captures in spring. The grey column on the far right shows the positive predictive value in green and false discovery rate in red. The grey bottom row shows the true positive rate in green and the false negative rate in red. The cell in the bottom right shows the overall predictive accuracy.

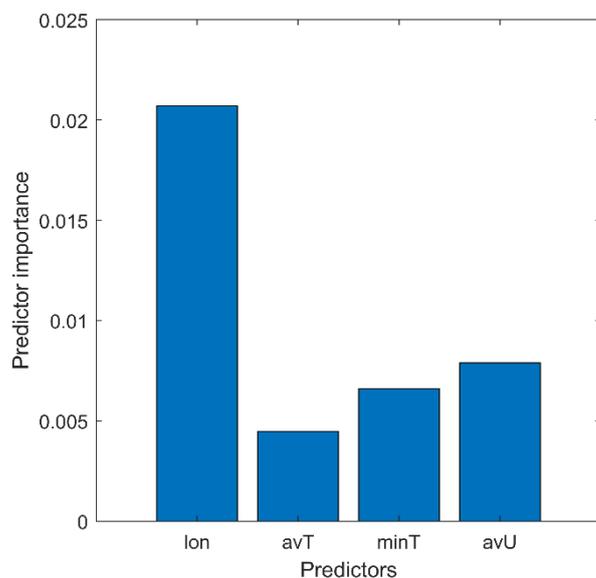


Fig. 2. Predictor importance scores for the Decision Tree model, where lon = longitude, avT = average temperature, mint = minimum temperature, and avU = average wind speed.

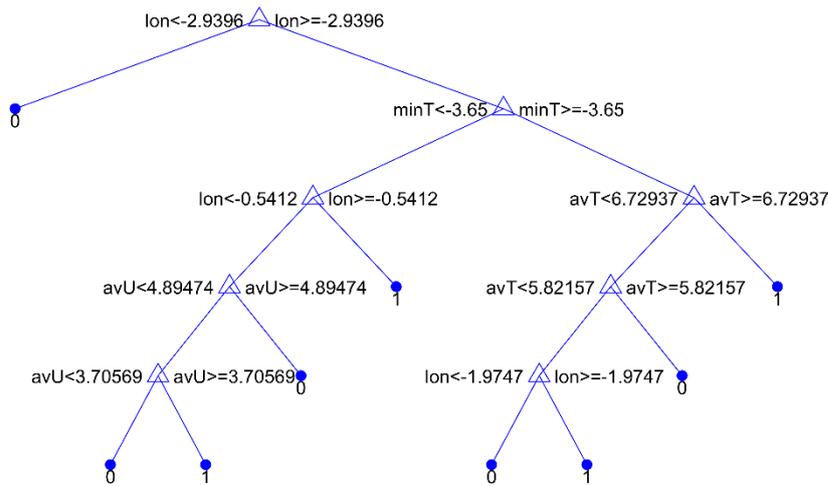


Fig. 3. Graphic description of the Decision Tree model, where 0 = no captures in spring, 1 = captures in spring, lon = longitude, avT = average temperature, mint = minimum temperature, and avU = average wind speed.

Predicting key percentiles of SWD activity

It was not possible to provide predictions for day-of-season of the 1-percentile of captures as the data for early captures were too variable among sites and between years. This was established previously using generalized linear mixed effects models in a modelling report for AHDB project SF/TF 145a. It is illustrated here using a plot of proportional SWD captures for the exponential phase of activity (0-5% of total captures) of each dataset (Fig. 4)

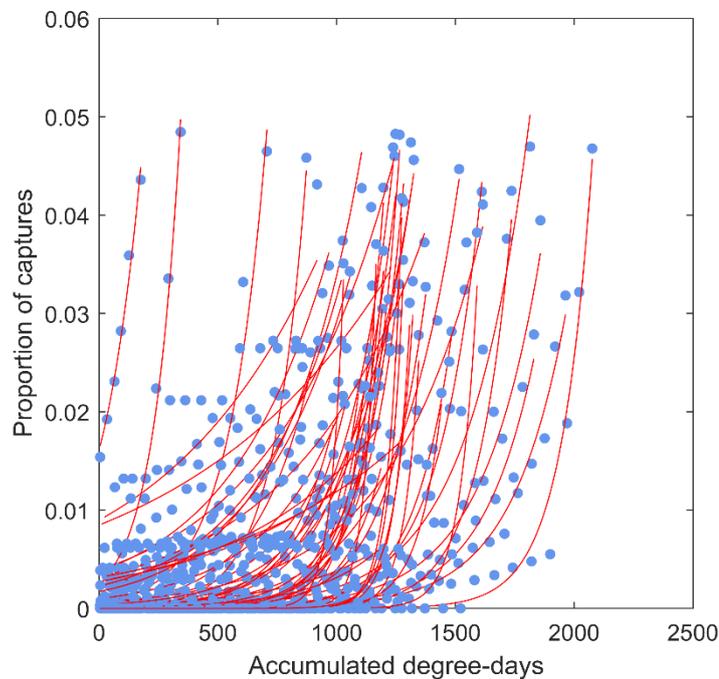


Fig. 4. Observed cumulative normalized SWD captures (blue markers) for the exponential phase of activity in individual site x year datasets. The red lines are fitted exponential curves.

The MAE of the model on the test set (2018) was 1.8 days for day-of-season of the 5-percentile of captures, and 2.6 days for day-of-season of the 50-percentile of captures when data was grouped at a regional-scale. The optimal biofix-baseline settings were 18/07 and 0 degrees. Similar predictive accuracy was achieved when the model was refit to all years at the region-scale using optimal biofix-baseline settings (Fig. 5; Table S1).

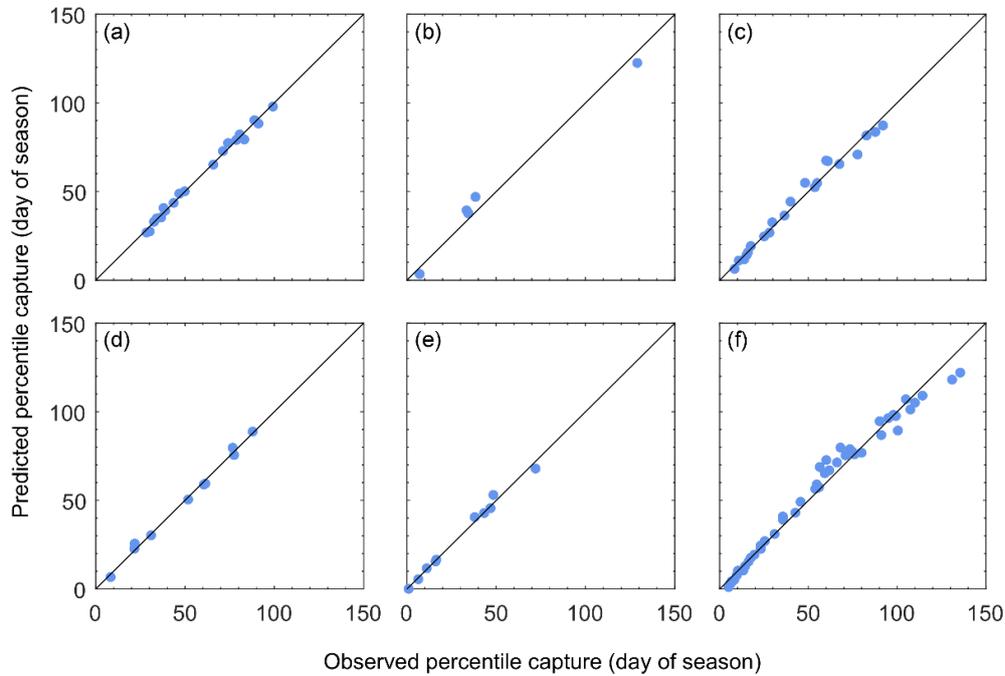


Fig. 5. Observed and predicted day of season of 5- and 50-percentile captures: (a) East Scotland, (b) Yorkshire & Humber, (c) West Midlands, (d) East Midlands, (e) East of England, (f) South East England. Diagonal lines are 1:1 lines.

The MAE of the model on the test set (2018) at the national-scale was 3.15 days for day-of-season of the 5-percentile of captures, and 5.96 days for day-of-season of the 50-percentile of captures. Similar predictive accuracy was achieved when the model was refit to all years at the national-scale using optimal biofix-baseline settings (Fig. 6, Table S2).

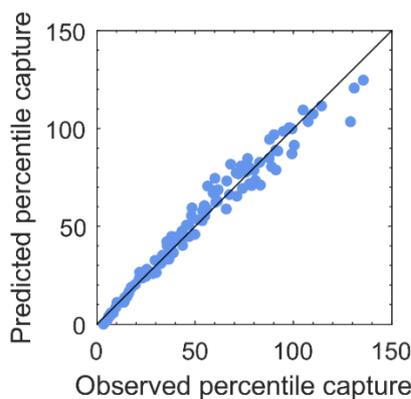


Fig. 6. Observed and predicted day of season of 5- and 50-percentile capture for GB. Diagonal line is a 1:1 line.

Machine learning

Support Vector Machine achieved the highest accuracy (rmse = 4.8) on the hold-out test data for prediction of day-of-season of percentile captures using degree-days as a predictor (Fig. 7).

1.1 Linear Regression Last change: Linear	RMSE (Validation): 5.4028 1/1 features	1.13 SVM Last change: Coarse Gaussian SVM	RMSE (Validation): 5.985 1/1 features
1.2 Linear Regression Last change: Interactions Linear	RMSE (Validation): 5.4028 1/1 features	1.14 Ensemble Last change: Boosted Trees	RMSE (Validation): 8.104 1/1 features
1.3 Linear Regression Last change: Robust Linear	RMSE (Validation): 5.2589 1/1 features	1.15 Ensemble Last change: Bagged Trees	RMSE (Validation): 8.4521 1/1 features
1.4 Stepwise Linear Regression Last change: Stepwise Linear	RMSE (Validation): 5.4028 1/1 features	1.16 Gaussian Process Regression Last change: Squared Exponential GPR	RMSE (Validation): 5.2078 1/1 features
1.5 Tree Last change: Fine Tree	RMSE (Validation): 6.719 1/1 features	1.17 Gaussian Process Regression Last change: Matern 5/2 GPR	RMSE (Validation): 5.195 1/1 features
1.6 Tree Last change: Medium Tree	RMSE (Validation): 12.416 1/1 features	1.18 Gaussian Process Regression Last change: Exponential GPR	RMSE (Validation): 5.498 1/1 features
1.7 Tree Last change: Coarse Tree	RMSE (Validation): 18.566 1/1 features	1.19 Gaussian Process Regression Last change: Rational Quadratic GPR	RMSE (Validation): 5.2078 1/1 features
1.8 SVM Last change: Linear SVM	RMSE (Validation): 5.3468 1/1 features	1.20 Neural Network Last change: Narrow Neural Network	RMSE (Validation): 5.2345 1/1 features
1.9 SVM Last change: Quadratic SVM	RMSE (Validation): 4.7885 1/1 features	1.21 Neural Network Last change: Medium Neural Network	RMSE (Validation): 5.239 1/1 features
1.10 SVM Last change: Cubic SVM	RMSE (Validation): 4.854 1/1 features	1.22 Neural Network Last change: Wide Neural Network	RMSE (Validation): 5.3446 1/1 features
1.11 SVM Last change: Fine Gaussian SVM	RMSE (Validation): 9.3962 1/1 features	1.23 Neural Network Last change: Bilayered Neural Network	RMSE (Validation): 5.2778 1/1 features
1.12 SVM Last change: Medium Gaussian SVM	RMSE (Validation): 6.9832 1/1 features	1.24 Neural Network Last change: Trilayered Neural Network	RMSE (Validation): 5.2914 1/1 features

Fig. 7. Accuracy of the suite of 24 machine learning algorithms for predicting day-of-season of percentile captures using degree-days as a predictor.

When region ID (a code from 1 to 6) was added as a categorical predictor, Gaussian Process Regression was the superior algorithm and predictive performance increased to rmse = 3.25 (Fig. 8). Note that the linear regression described previously to test the validity of percentile capture charts achieved an overall rmse of 2.78 for prediction of day-of-season of percentile captures using data grouped by region.

1.1 Linear Regression Last change: Linear	RMSE (Validation): 5.395 2/2 features	1.13 SVM Last change: Coarse Gaussian SVM	RMSE (Validation): 5.1618 2/2 features
1.2 Linear Regression Last change: Interactions Linear	RMSE (Validation): 5.0844 2/2 features	1.14 Ensemble Last change: Boosted Trees	RMSE (Validation): 7.0869 2/2 features
1.3 Linear Regression Last change: Robust Linear	RMSE (Validation): 5.202 2/2 features	1.15 Ensemble Last change: Bagged Trees	RMSE (Validation): 10.174 2/2 features
1.4 Stepwise Linear Regression Last change: Stepwise Linear	RMSE (Validation): 5.0844 2/2 features	1.16 Gaussian Process Regression Last change: Squared Exponential GPR	RMSE (Validation): 3.7093 2/2 features
1.5 Tree Last change: Fine Tree	RMSE (Validation): 6.6717 2/2 features	1.17 Gaussian Process Regression Last change: Matern 5/2 GPR	RMSE (Validation): 3.524 2/2 features
1.6 Tree Last change: Medium Tree	RMSE (Validation): 10.541 2/2 features	1.18 Gaussian Process Regression Last change: Exponential GPR	RMSE (Validation): 3.2512 2/2 features
1.7 Tree Last change: Coarse Tree	RMSE (Validation): 17.68 2/2 features	1.19 Gaussian Process Regression Last change: Rational Quadratic GPR	RMSE (Validation): 3.7093 2/2 features
1.8 SVM Last change: Linear SVM	RMSE (Validation): 5.2081 2/2 features	1.20 Neural Network Last change: Narrow Neural Network	RMSE (Validation): 3.5487 2/2 features
1.9 SVM Last change: Quadratic SVM	RMSE (Validation): 4.2265 2/2 features	1.21 Neural Network Last change: Medium Neural Network	RMSE (Validation): 4.2355 2/2 features
1.10 SVM Last change: Cubic SVM	RMSE (Validation): 6.4082 2/2 features	1.22 Neural Network Last change: Wide Neural Network	RMSE (Validation): 21.112 2/2 features
1.11 SVM Last change: Fine Gaussian SVM	RMSE (Validation): 13.289 2/2 features	1.23 Neural Network Last change: Bilayered Neural Network	RMSE (Validation): 5.3609 2/2 features
1.12 SVM Last change: Medium Gaussian SVM	RMSE (Validation): 6.1995 2/2 features	1.24 Neural Network Last change: Trilayered Neural Network	RMSE (Validation): 4.0738 2/2 features

Fig. 8. Accuracy of the suite of 24 machine learning algorithms for predicting day-of-season of percentile captures using degree-days and region ID as predictors.

The difference in rmse between the linear regression model here (Fig. 8) and that described in the previous section is likely due to the different hold-out cross-validation procedures; here the test set was a random selection of 25% of the data, whereas previously all data for 2018 was used as a test set.

Percentile capture charts

There was a marked difference in the 5- and 50-percentile thresholds among the six regions, providing evidence of the need for a separate chart for each geographic region (Fig. 9). Note that the thresholds on the charts are generated from *observations* of accumulated degree-days at percentile capture points and not from model predictions; the purpose of the modelling was to establish if accumulated degree-days are a good predictor of the day-of-season of key capture percentiles. Also note that the thresholds are mean values for all sites and years within each region, and there was considerable variation in those values within regions.

An alternative approach is to also display a measure of uncertainty around the mean (within-region) percentile capture points, or to use a more robust estimator of central tendency that reduces the effects of outlier bias. Yorkshire & Humber had insufficient data to calculate a meaningful confidence interval or measure of variation, therefore one-sided trimmed means were used. A trimmed mean is a measure of central tendency that is less sensitive to outliers than the mean. It involves calculation of the mean after discarding values at the higher and lower end of a sample. When only values at one extreme are removed it is called a one-sided trimmed mean. The 'upper mean' is the average of a sample when low values are discarded (upper mean > mean) and the

'lower mean' is the average when high values are discarded (lower mean < mean). Lower means were calculated by discarding values > 75-percentile of the sample and upper means by discarding values < 25-percentile of the sample. The lower mean can be used to mark the action point for a risk averse strategy and the upper mean a risk tolerant strategy (Fig. 10).

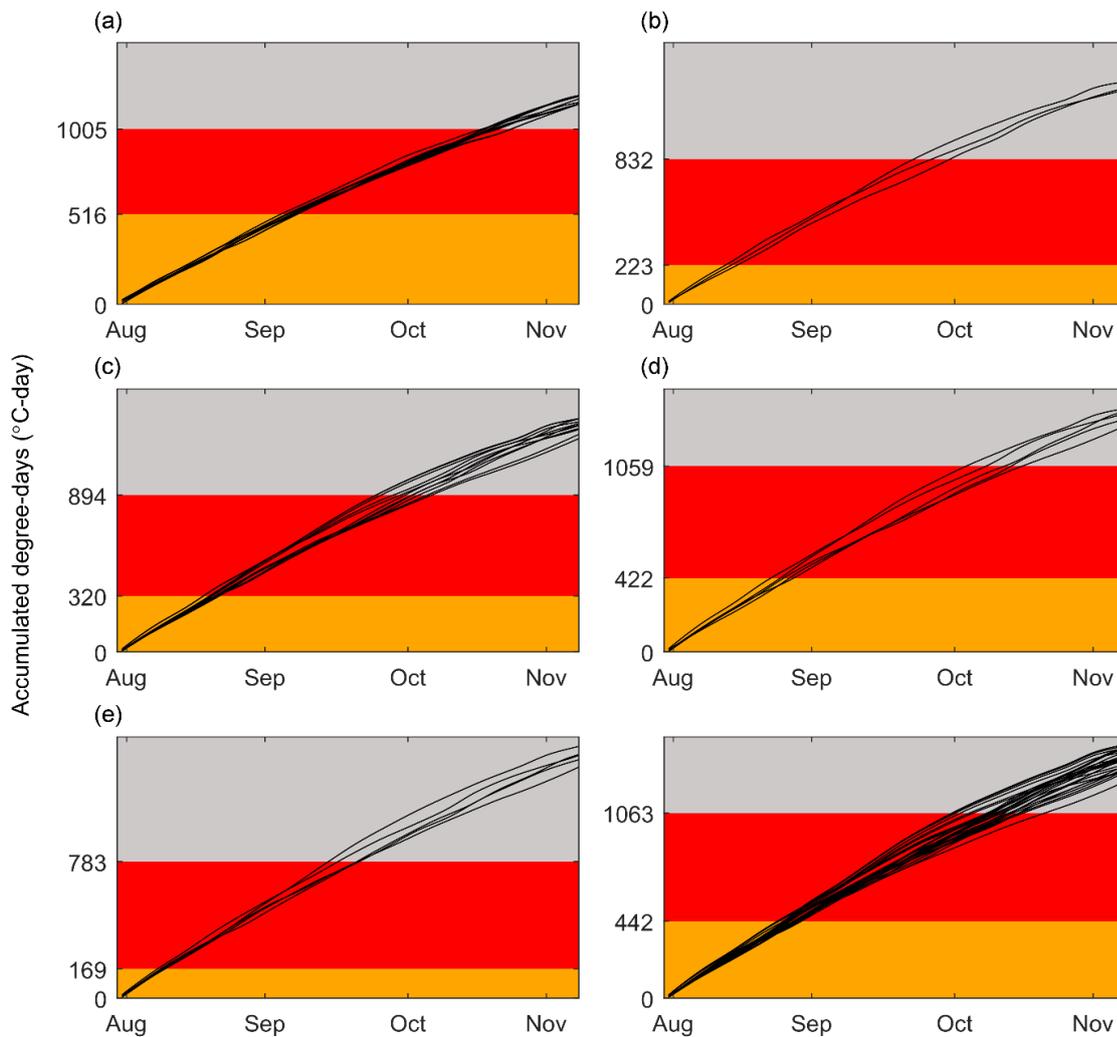


Fig. 9. Spotted wing drosophila percentile capture charts: (a) East Scotland, (b) Yorkshire & Humber, (c) West Midlands, (d) East Midlands, (e) East of England, (f) South East England. The orange zone marks 0 to 5-percentile of captures and the red zone the 5- to 50-percentile. The curves show accumulated degree days for the datasets used to generate the thresholds in each region.

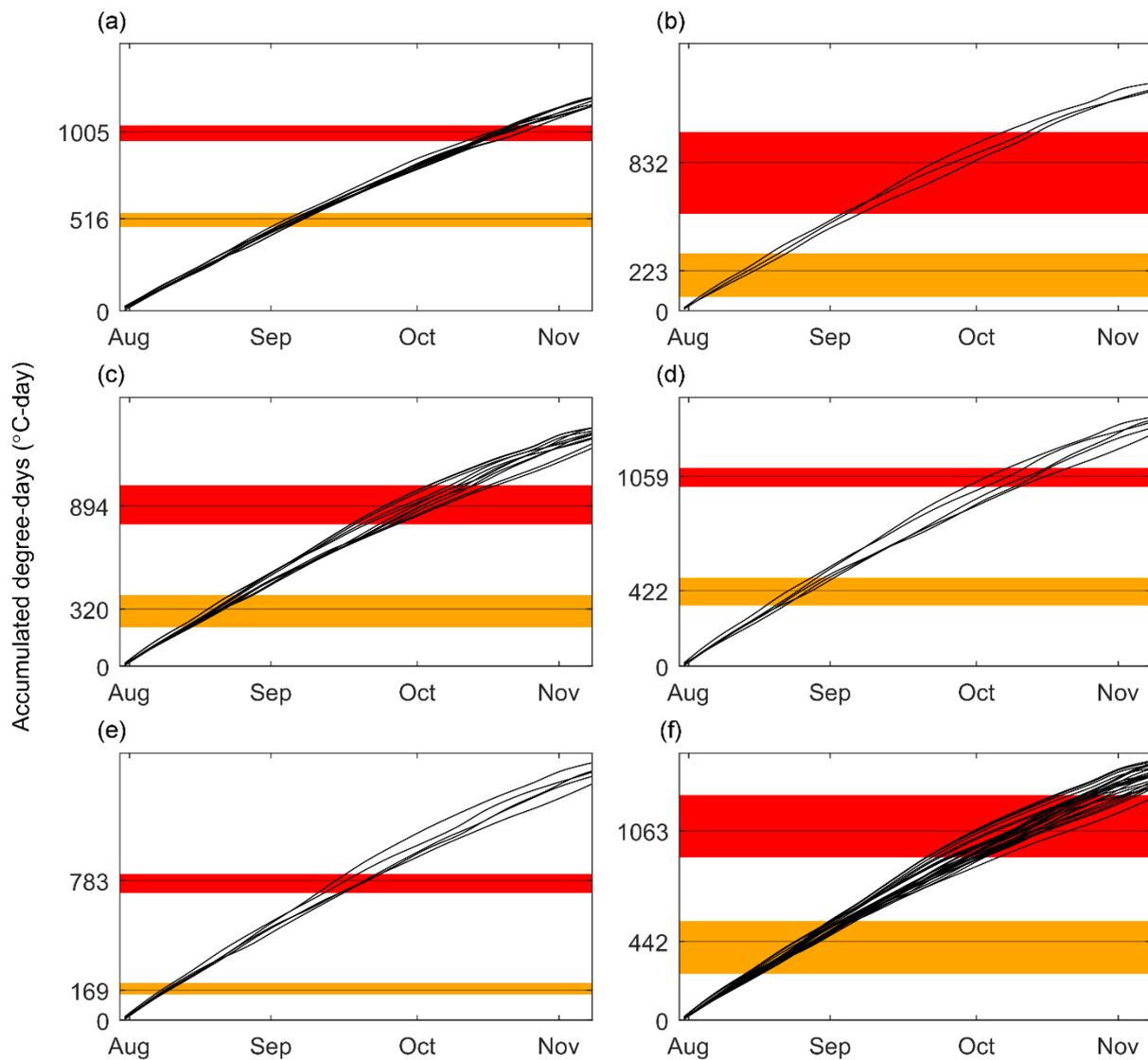


Fig. 10. Spotted wing drosophila percentile capture charts with risk tolerance bands: (a) East Scotland, (b) Yorkshire & Humber, (c) West Midlands, (d) East Midlands, (e) East of England, (f) South East England. Horizontal lines show the mean accumulated degree days for the 5- and 50-percentile capture points in each region. The coloured bands delineate lower (risk averse) and upper (risk tolerant) trimmed mean values. Curves show accumulated degree days for the datasets used to generate the thresholds in each region.

If confidence intervals are preferred, Yorkshire & Humber must be removed (Fig. 11). Alternatively, it could be merged with another region.

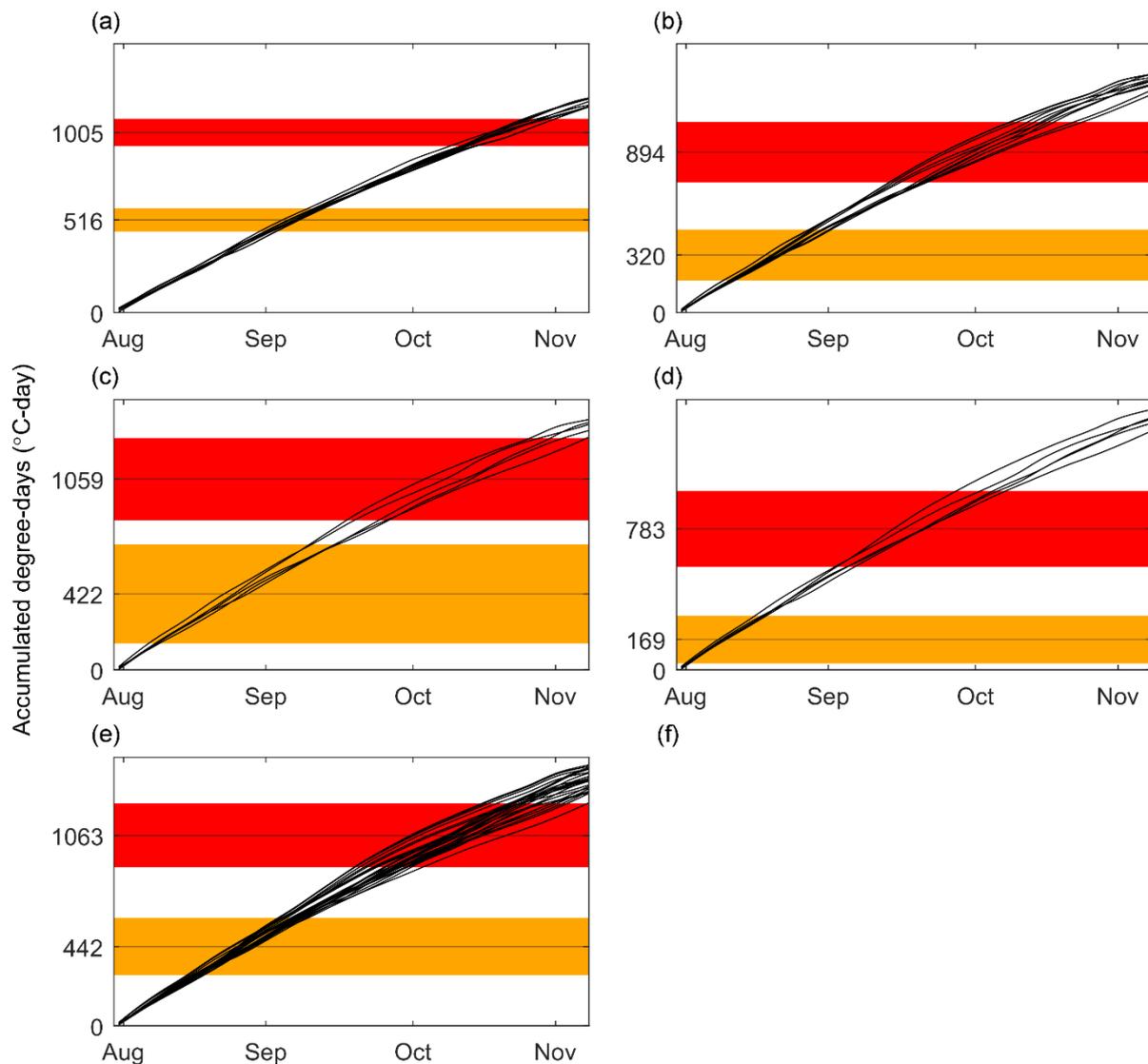


Fig. 11. Spotted wing drosophila percentile capture charts with risk tolerance bands: (a) East Scotland, (b) Yorkshire & Humber, (c) West Midlands, (d) East Midlands, (e) East of England, (f) South East England. Horizontal lines show the mean accumulated degree days for the 5- and 50-percentile capture points in each region. The coloured bands represent the 95% confidence interval of the mean. Curves show accumulated degree days for the datasets used to generate the thresholds in each region.

To plot a capture chart for GB the median accumulated degree days for the 5- and 50-percentile capture points is used, as the data for the 5-percentile points were skewed. The interquartile range is therefore used as a measure of uncertainty around median values, as opposed to a confidence interval or a trimmed mean, to give 'risk averse vs risk tolerant' action points (Fig. 12).

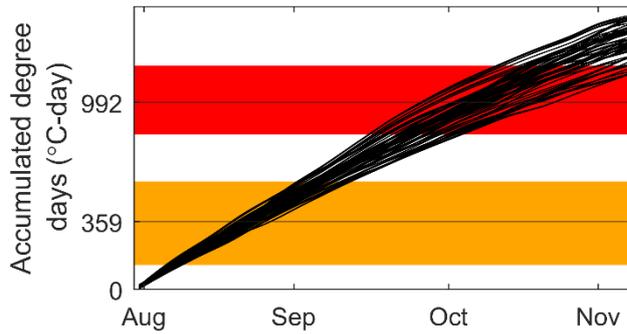


Fig. 12. Spotted wing drosophila percentile capture chart with risk tolerance bands for GB. Horizontal lines show the median accumulated degree days for the 5- and 50-percentile capture points across all trap locations. The coloured bands represent the interquartile range. Curves show accumulated degree days for all datasets used to generate the thresholds.

It can be seen in Figs. 9-12 that large differences in SWD trap captures between sites and years resulted in variation in the timing of percentile capture points, leading to uncertainty in the dates of action points.

Modelling population dynamics over the whole season

The fit results for the site-, region-, country- and national-scales of analysis are given in Figs. 13-16 and Tables S3-6 (Appendix), respectively. The logistic model was the best overall for grouping by site, region, and country, whereas the Richards model was superior at the national-scale. The optimal biofix-baseline settings were 18/06 and 0 degrees for all geographic groupings.

Machine learning

Gaussian Process Regression achieved the highest accuracy (rmse = 0.14) on the hold-out test data for prediction of proportional captures using degree-days as a predictor (Fig. 17). Note that the Logistic model fit to all the (ungrouped) proportional capture data also achieved an overall rmse of 0.14.

When site ID (a code from 1 to 16) was added as a categorical predictor, predictive performance on the test set increased to rmse = 0.1 (Fig. 18). The logistic population model fit to all the proportional captured data grouped by site achieved an overall rmse of 0.09.

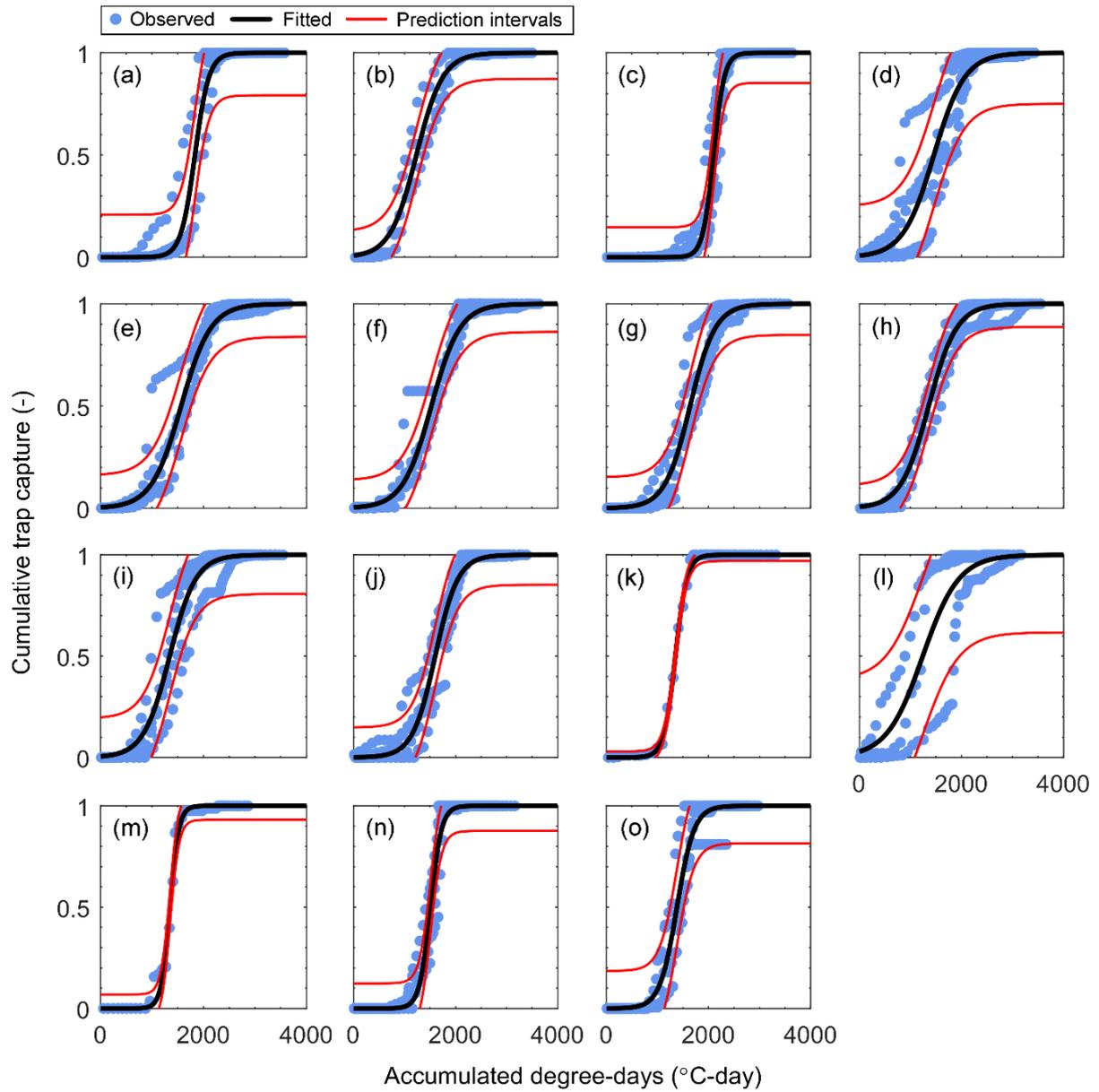


Fig. 13. Observed and predicted proportional Spotted wing drosophila captures grouped by site, where (a) to (j) are sites 1 to 10, and panels (k) to (o) are sites 10a, 10b, 1100, 1300, and 1400.

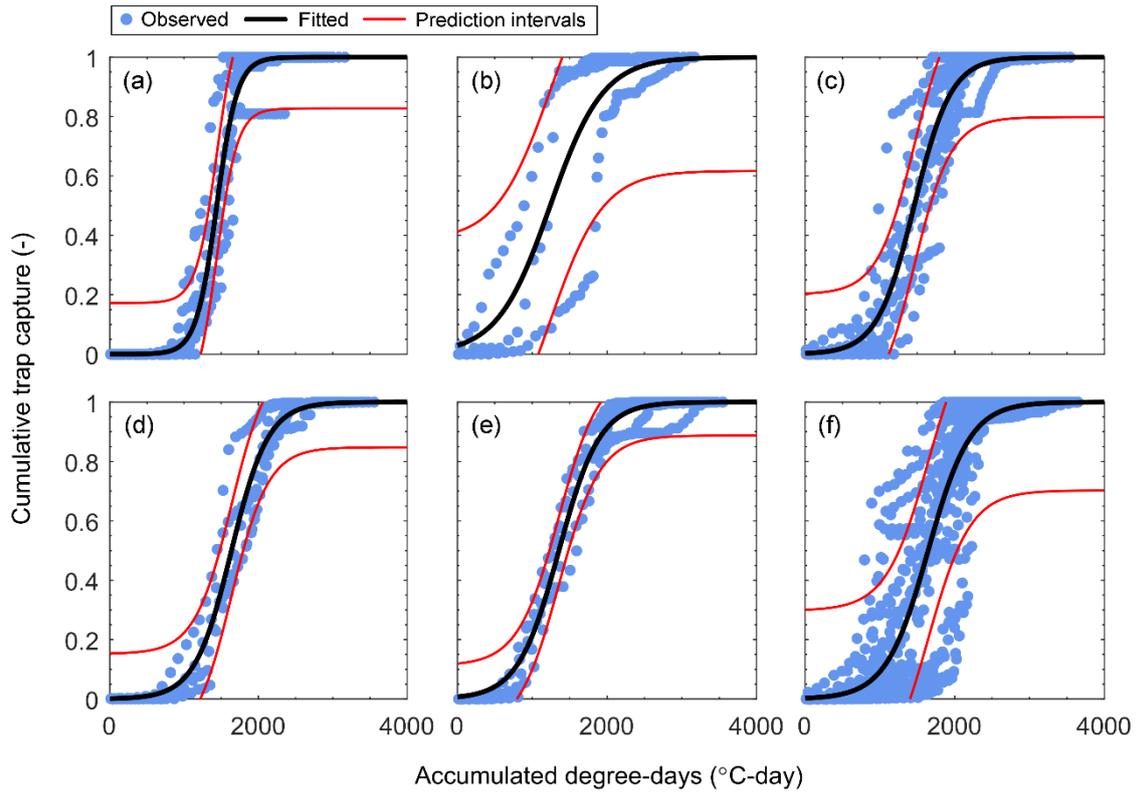


Fig. 14. Observed and predicted proportional *Spotted wing drosophila* captures grouped by region: (a) East Scotland, (b) Yorkshire & Humber, (c) West Midlands, (d) East Midlands, (e) East of England, (f) South East England.

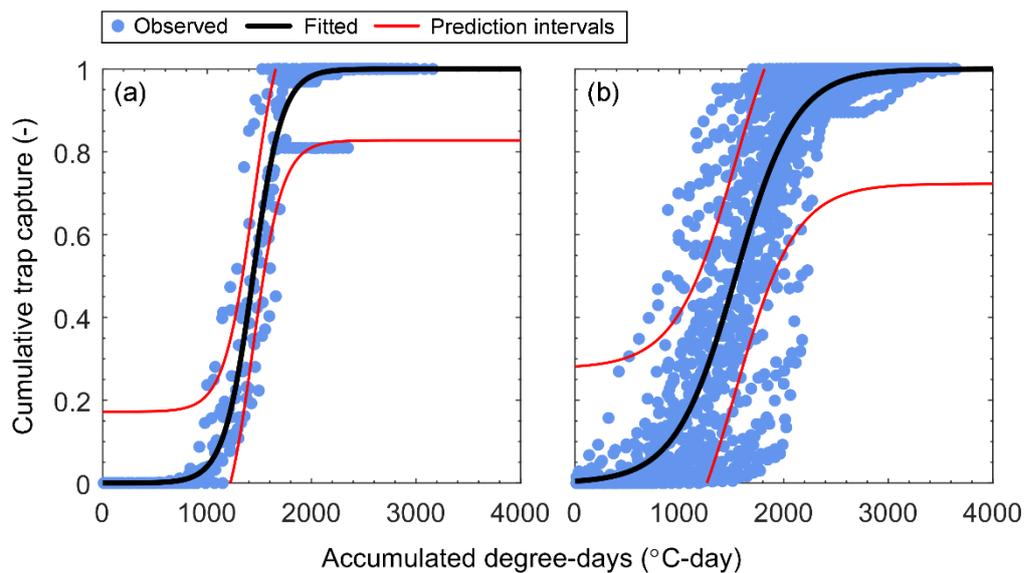


Fig. 15. Observed and predicted proportional *Spotted wing drosophila* captures grouped by country: (a) Scotland, (b) England.

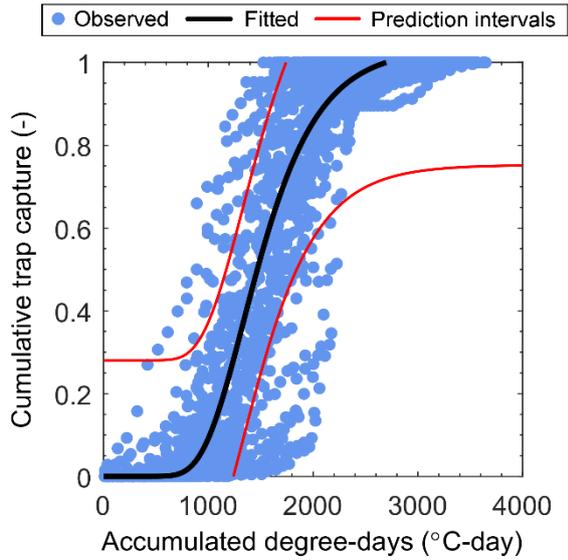


Fig. 16. Ungrouped observed and predicted proportional Spotted wing drosophila captures

1.1 Linear Regression Last change: Linear	RMSE (Validation): 0.18883 1/1 features	1.13 SVM Last change: Coarse Gaussian SVM	RMSE (Validation): 0.14921 1/1 features
1.2 Linear Regression Last change: Interactions Linear	RMSE (Validation): 0.18883 1/1 features	1.14 Ensemble Last change: Boosted Trees	RMSE (Validation): 0.14805 1/1 features
1.3 Linear Regression Last change: Robust Linear	RMSE (Validation): 0.18876 1/1 features	1.15 Ensemble Last change: Bagged Trees	RMSE (Validation): 0.15282 1/1 features
1.4 Stepwise Linear Regression Last change: Stepwise Linear	RMSE (Validation): 0.18883 1/1 features	1.16 Gaussian Process Regression Last change: Squared Exponential GPR	RMSE (Validation): 0.14057 1/1 features
1.5 Tree Last change: Fine Tree	RMSE (Validation): 0.16377 1/1 features	1.17 Gaussian Process Regression Last change: Matern 5/2 GPR	RMSE (Validation): 0.14064 1/1 features
1.6 Tree Last change: Medium Tree	RMSE (Validation): 0.15293 1/1 features	1.18 Gaussian Process Regression Last change: Exponential GPR	RMSE (Validation): 0.1421 1/1 features
1.7 Tree Last change: Coarse Tree	RMSE (Validation): 0.1464 1/1 features	1.19 Gaussian Process Regression Last change: Rational Quadratic GPR	RMSE (Validation): 0.14056 1/1 features
1.8 SVM Last change: Linear SVM	RMSE (Validation): 0.18896 1/1 features	1.20 Neural Network Last change: Narrow Neural Network	RMSE (Validation): 0.14141 1/1 features
1.9 SVM Last change: Quadratic SVM	RMSE (Validation): 0.17629 1/1 features	1.21 Neural Network Last change: Medium Neural Network	RMSE (Validation): 0.14262 1/1 features
1.10 SVM Last change: Cubic SVM	RMSE (Validation): 0.1713 1/1 features	1.22 Neural Network Last change: Wide Neural Network	RMSE (Validation): 0.14318 1/1 features
1.11 SVM Last change: Fine Gaussian SVM	RMSE (Validation): 0.14982 1/1 features	1.23 Neural Network Last change: Bilayered Neural Network	RMSE (Validation): 0.14193 1/1 features
1.12 SVM Last change: Medium Gaussian SVM	RMSE (Validation): 0.14612 1/1 features	1.24 Neural Network Last change: Trilayered Neural Network	RMSE (Validation): 0.14182 1/1 features

Fig. 17. Accuracy of the suite of 24 machine learning algorithms for predicting proportional captures using degree-days as a predictor.

1.1 Linear Regression Last change: Linear	RMSE (Validation): 0.17673 2/2 features	1.13 SVM Last change: Coarse Gaussian SVM	RMSE (Validation): 0.15074 2/2 features
1.2 Linear Regression Last change: Interactions Linear	RMSE (Validation): 0.1757 2/2 features	1.14 Ensemble Last change: Boosted Trees	RMSE (Validation): 0.10543 2/2 features
1.3 Linear Regression Last change: Robust Linear	RMSE (Validation): 0.17683 2/2 features	1.15 Ensemble Last change: Bagged Trees	RMSE (Validation): 0.11039 2/2 features
1.4 Stepwise Linear Regression Last change: Stepwise Linear	RMSE (Validation): 0.1757 2/2 features	1.16 Gaussian Process Regression Last change: Squared Exponential GPR	RMSE (Validation): 0.098798 2/2 features
1.5 Tree Last change: Fine Tree	RMSE (Validation): 0.11517 2/2 features	1.17 Gaussian Process Regression Last change: Matern 5/2 GPR	RMSE (Validation): 0.098416 2/2 features
1.6 Tree Last change: Medium Tree	RMSE (Validation): 0.109 2/2 features	1.18 Gaussian Process Regression Last change: Exponential GPR	RMSE (Validation): 0.10308 2/2 features
1.7 Tree Last change: Coarse Tree	RMSE (Validation): 0.11365 2/2 features	1.19 Gaussian Process Regression Last change: Rational Quadratic GPR	RMSE (Validation): 0.098482 2/2 features
1.8 SVM Last change: Linear SVM	RMSE (Validation): 0.17718 2/2 features	1.20 Neural Network Last change: Narrow Neural Network	RMSE (Validation): 0.10749 2/2 features
1.9 SVM Last change: Quadratic SVM	RMSE (Validation): 0.16741 2/2 features	1.21 Neural Network Last change: Medium Neural Network	RMSE (Validation): 0.10133 2/2 features
1.10 SVM Last change: Cubic SVM	RMSE (Validation): 0.25256 2/2 features	1.22 Neural Network Last change: Wide Neural Network	RMSE (Validation): 0.10513 2/2 features
1.11 SVM Last change: Fine Gaussian SVM	RMSE (Validation): 0.10684 2/2 features	1.23 Neural Network Last change: Bilayered Neural Network	RMSE (Validation): 0.10343 2/2 features
1.12 SVM Last change: Medium Gaussian SVM	RMSE (Validation): 0.10253 2/2 features	1.24 Neural Network Last change: Trilayered Neural Network	RMSE (Validation): 0.10416 2/2 features

Fig. 18. Accuracy of the suite of 24 machine learning algorithms for predicting proportional captures using degree-days and site ID as predictors.

Discussion

All modelling approaches were successful with good predictive performance. The classifier for spring risk can be used to provide a prediction on March 1 regarding the need to begin crop protection measures. The ‘percentile capture charts’ can be used in the same way as the Blackcurrant Gall Mite Emergence charts, where the degree-day curve of the current season is plotted on the chart. The percentile capture charts, however, have the added advantage of separate region-specific thresholds and, if desired, zones defining risk averse and risk tolerant strategies. The population models and machine learning algorithms, which performed equally well, can also be used to provide predictions of when certain thresholds of activity are reached. Note that a host of additional statistical models and machine learning algorithms for various aspects of the SWD life cycle are available from AHDB project SF/TF 145a, as detailed in the project reports. These include an Adaptively Boosted Decision Tree algorithm (adaboostm1) that can predict SWD flight activity on any given day with 91.8% accuracy using a range of weather variables, and a Fine *K*-Nearest Neighbor algorithm that can predict the first spring peaks of female activity with 93.3% accuracy.

Supplement

Table S1. Estimated coefficients (95% CI) and goodness of fit for day-of-season of percentile capture points regressed on accumulated degree-days, grouped by region.

Region	Parameter a	Parameter b	MAE	rmse	R ²
1	0.0886 ± 0.0034	-7.724 ± 2.757	1.46	1.91	0.99
2	0.0931 ± 0.0194	-8.377 ± 13.41	5.92	7.63	0.98
3	0.0747 ± 0.0046	-3.156 ± 3.267	2.72	3.76	0.98
4	0.0711 ± 0.0044	-2.832 ± 3.659	1.77	2.23	0.99
5	0.0653 ± 0.0054	-1.166 ± 3.153	1.59	2.47	0.99
6	0.0751 ± 0.0030	-4.475 ± 2.755	3.96	5.42	0.98

Region 1 = E Scotland, 2 = Yorkshire & Humber, 3 = W Midlands, 4 = E Midlands, 5 = E England, 6 = SE England

Table S2. Estimated coefficients (95% CI) and goodness of fit for day-of-season of percentile capture points regressed on accumulated degree-days, with data pooled for the UK.

Parameter a	Parameter b	MAE	rmse	R ²
0.0763 ± 0.0024	-3.880 ± 1.972	4.13	5.72	0.97

Region 1 = E Scotland, 2 = Yorkshire & Humber, 3 = W Midlands, 4 = E Midlands, 5 = E England, 6 = SE England

Table S3. Estimated coefficients (95% CI) and goodness of fit for the logistic model fit to Spotted wing drosophila capture data grouped by site.

Site	Parameter a	Parameter b	rmse	R ²
1	0.0073 ± 0.0013	1833 ± 29	0.11	0.94
2	0.0039 ± 0.0004	1233 ± 28	0.06	0.97
3	0.0094 ± 0.0009	2097 ± 12	0.07	0.97
4	0.0033 ± 0.0004	1459 ± 42	0.13	0.89
5	0.0035 ± 0.0003	1565 ± 27	0.08	0.95
6	0.0036 ± 0.0003	1525 ± 26	0.07	0.97
7	0.0040 ± 0.0003	1637 ± 23	0.08	0.96
8	0.0037 ± 0.0002	1357 ± 19	0.06	0.98
9	0.0039 ± 0.0004	1342 ± 31	0.10	0.93
10	0.0045 ± 0.0003	1601 ± 20	0.07	0.97
11	0.0088 ± 0.0005	1350 ± 8	0.01	1.00
12	0.0028 ± 0.0006	1241 ± 92	0.19	0.74
13	0.0116 ± 0.0015	1349 ± 13	0.03	0.99
15	0.0093 ± 0.0007	1508 ± 11	0.06	0.98
16	0.0059 ± 0.0007	1380 ± 24	0.09	0.95

Table S4. Estimated coefficients (95% CI) and goodness of fit for the logistic model fit to Spotted wing drosophila capture data grouped by region.

Region	Parameter a	Parameter b	rmse	R ²
1	0.0072 ± 0.0005	1439 ± 12	0.09	0.96
2	0.0028 ± 0.0006	1241 ± 92	0.19	0.74
3	0.0041 ± 0.0003	1459 ± 21	0.10	0.93
4	0.0040 ± 0.0003	1637 ± 23	0.08	0.96
5	0.0037 ± 0.0002	1357 ± 19	0.06	0.98
6	0.0035 ± 0.0002	1645 ± 21	0.15	0.86

Region 1 = E Scotland, 2 = Yorkshire & Humber, 3 = W Midlands, 4 = E Midlands, 5 = E England, 6 = SE England

Table S5. Estimated coefficients (95% CI) and goodness of fit for the logistic model fit to Spotted wing drosophila capture data grouped by country.

Country	Parameter a	Parameter b	rmse	R ²
Scotland	0.0072 ± 0.0005	1439 ± 13	0.09	0.96
England	0.0034 ± 0.0001	1538 ± 15	0.14	0.87

Table S6. Estimated coefficients (95% CI) and goodness of fit for the Richards model fit to Spotted wing drosophila capture data grouped by nation.

Nation	Parameter a	Parameter b	Parameter c	rmse	R ²
GB	1.033 ± 0.012	0.0025 ± 0.0001	28.79 ± 5.58	0.14	0.87