# West Nile Virus Forecast Report for 2018-08-15 South Dakota

DEMO for Arbovirus Modeling and Prediction (ArboMAP): Synthetic Data
Not for Epidemiological Use

Dawn M. Nekorchuk, Justin K. Davis, and Michael C. Wimberly
(mcwimberly@ou.edu)
Geography and Environmental Sustainability, University of Oklahoma

Report compiled on June 01, 2022

## Contents

# 1    Forecast results

The Arbovirus Monitoring and Prediction (ArboMAP) system produces a weekly, county-level forecast of human West Nile virus (WNV) cases using environmental data combined with entomological data.
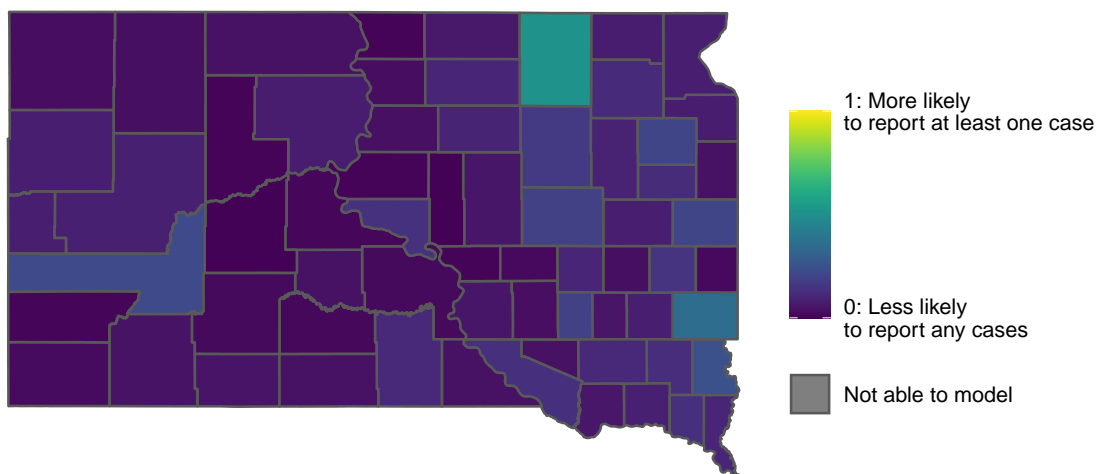
*Modeling overview*: The transmission of mosquito-borne diseases, such as WNV, is influenced by environmental conditions that affect many aspects of the disease transmission system. ArboMAP uses an ensemble of different mathematical models that each are predicting if a county will report at least one case in a given week ('positive county-week'). Results presented are an average of the models with ranges as appropriate. As part of the process, mosquito infection rate is also modeled based on the mosquito pool data, and is included in the default modeling. ArboMAP uses generalized additive models (GAMs) with smooths for seasonality, and also lagged weather data, which allows it to model the time-delayed effects of weather conditions. The appendix will expand all the results to show each individual model.

## 1.1    Forecast week WNV absolute risk

The following map displays the **absolute risk** of predicted positive counties during epidemiological week 33.

This map can be used in conjunction with the **relative risk** map. The absolute risk map shows the risk of a county reporting at least one WNV positive human case during this week, and the relative risk map shows if this risk is elevated (or not) as compared to previous years.
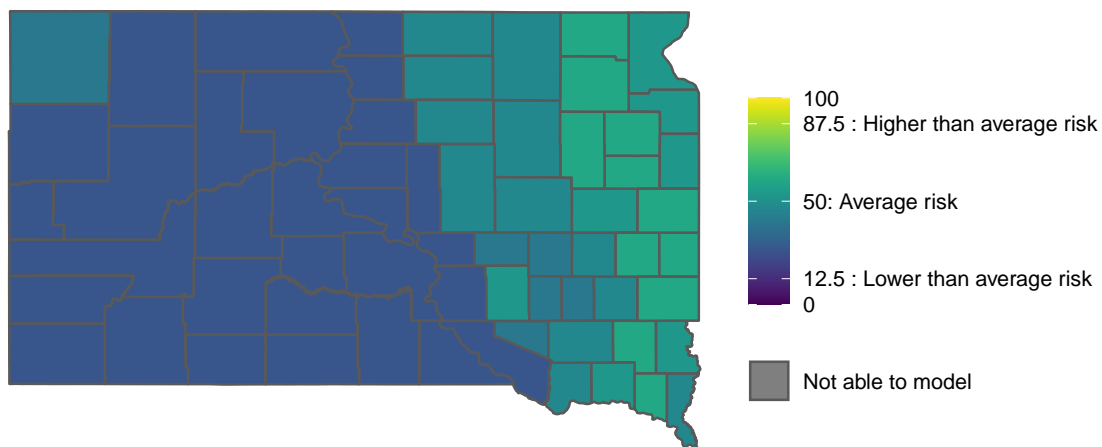
Absolute risk in forecast week

## 1.2 Forecast week WNV relative risk

In the forecast week there are 0 counties with higher than average risk as compared to the same epidemiological week in previous years (2004 through 2018).

This **relative risk** map may be used in conjunction with the **absolute risk** map. The absolute risk map shows the risk of a county reporting at least one WNV positive human case during this week, and the relative risk map shows if this risk is elevated (or not) as compared to previous years.
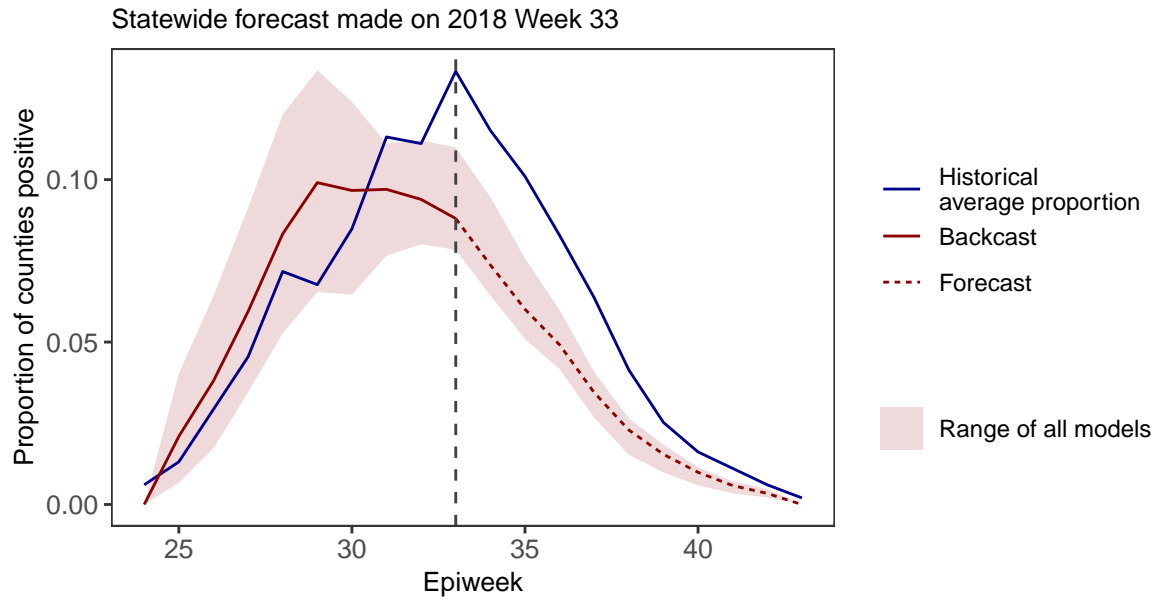
Risk in forecast week relative to the same epiweek in previous years



## 1.3 Forecast year

The following graph is the predicted epicurve of the forecast year: the average of all models is shown as a dark red line, with the range of all models in the shaded ribbon. Forecasts are shown as a dotted line and the predicted values from before the current forecast week ('backcast') are shown as a solid line. The appendix will have a version of this chart with a series for each model, rather than an average.

The historical **observed** proportion of counties positive, averaged over all known years, is also shown, here as a dark blue line. This is excluding human cases that occurred very early or very late in the season (temporal outliers), based on the percentage cut-off in the parameters, 0.02. This plotted curve allows a comparison between the timing and height of the predicted peak of cases as compared to averaged historical years. In the averaged year, 49% of the yearly cases would have been observed by this forecast week.

Statewide forecast made on 2018 Week 33



## 1.4 Case estimation

ArboMAP models are based on 'positive county-weeks', the probability that a county would have at least one human WNV case in a given week. These values can be used to predict a total number of **cases**, shown in the table below.

Table 1: Estimated number of WNV cases

| Year | Predicted positive county-weeks | Average estimated cases (standard dev) | Range of estimated cases |
|------|--------------------------------|----------------------------------------|--------------------------|
| 2018 | 63 | 75 (+/-15) | 56 - 89 |

## 1.5 Model fit statistics

The following table gives a summary of how well the model is fitting the historical years. The Area Under the ROC curve (AUC) is a statistic that ranges from 0 (model is right 0% of the time) to 1 (model is right 100% of the time). Scores above 0.5 are better than a random model, with >0.7 generally considered acceptable and >0.8 as good.
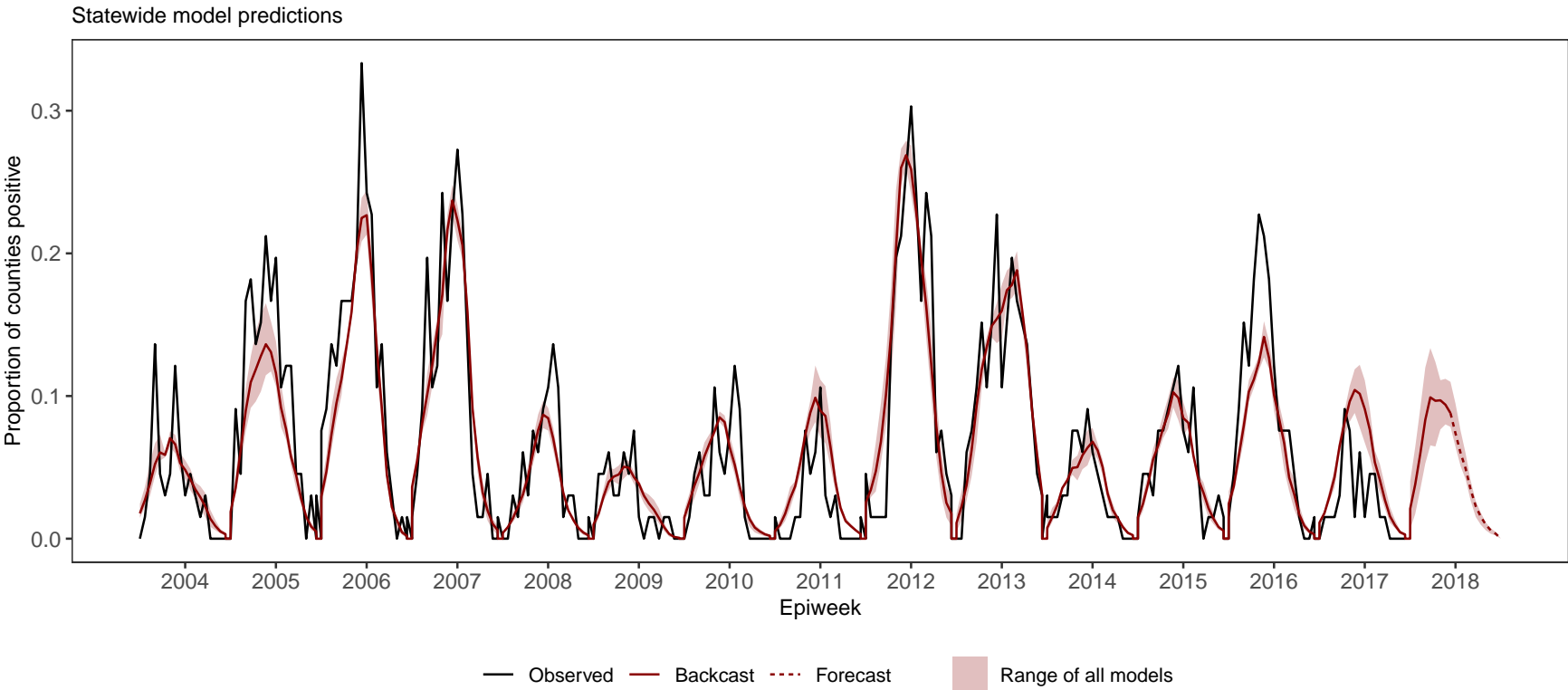
Table 2: Area Under Curve (AUC) statistics of all model fits

| Model | Average AUC | Min AUC | Max AUC |
|-------|-------------|---------|---------|
| Average of all models | 0.84 | 0.84 | 0.85 |

## 1.6   Multi-year forecast

The following chart shows the model results for the entire modeled period from 2004 through 2018. Years prior to the forecast year that had human case data were used for fitting the model.

Similar to the previous forecast year chart, the average of all models is shown as a dark red line, with the range of all models as the shaded ribbon. Forecasts are shown as a dotted line and predicted values from before the current forecast week ('backcast') are shown as a solid line. The historical **observed** values are shown in black. The appendix will have a version of this chart with a series for each model, rather than an average.



Statewide model predictions

# 2 Input data summaries

The report was requested for 2018-08-15, which is CDC/MMWR epiweek 33 in epiyear 2018.

## 2.1 Human cases

After data processing, the human case data contained a total of 1300 rows containing data from years: 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, and 2017. Parameters were set to include human data from 2004 through 2017. Data from all years were found in the data file.
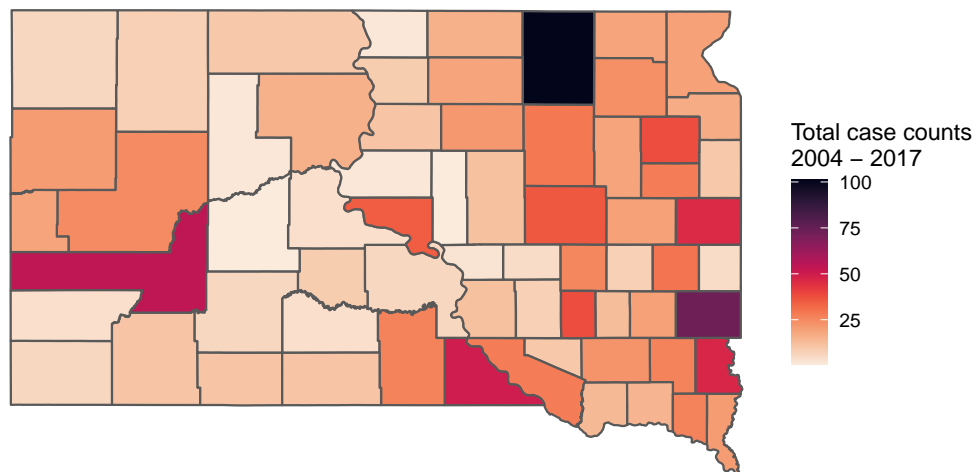
The human case data entries that were unmatched to spatial data during processing are in the table below. Please check for mispellings in the original file.
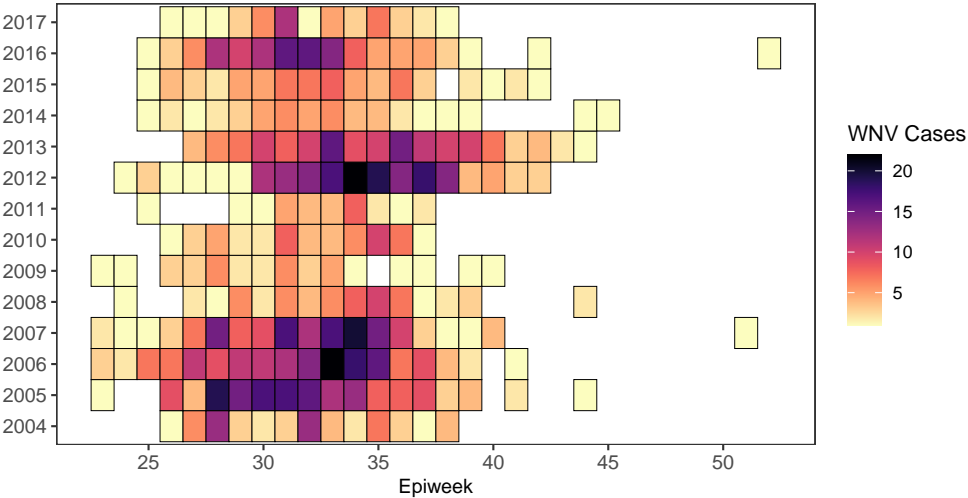
Table 3: Unmatched human case entries

| arbo_ID | date |
|---------|-----------|
| greg0ry | 7/4/2004 |

Over all years, the state saw a cumulative total of 1300 human cases representing 1141 positive county-weeks from a total of 66 counties.

Historical cumulative human cases



To compare the epicurve of human cases in each year, the heatmap below shows when in each year the cases occurred.

## 2.2   Mosquito pools

After data processing, the mosquito pool data contained a total of 16274 rows, containing data from years: 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, and 2018. Overall during this time frame, a total of 29 counties reported mosquito data.

Parameters were set to include mosquito data from 2004 through 2018. Data from all years were found in the data file.

Note that even if there were no positive pools in a given year, if there were any pools tested then the data will be useful; zero infection rates do predict low-risk years and should be used.

Parameters were set to include mosquito data from day of year 140 through 366. The mosquito infection rate modeling is very sensitive to early mosquito pool results, which is why a cut-off is used. Sensitivity analyses indicate that a start day of year of 140 is a reasonable cut-off for a high modeling accuracy.

Modeling was done from 2004 through 2018. Years without mosquito data during these years are assumed to have average mosquito infection rates. This allows us to estimate relationships with environmental data even when mosquito data are not available. There were 6 years where mosquito infection rates were imputed. These years are: 2004, 2007, 2009, 2010, 2012, and 2013. In modeling years where sufficient mosquito data were present, the mosquito infection statistic was created using the model specified in the input parameter: stratifiedMIGR.

In the forecast year to date, there were 1697 pools reported from 12 counties. Of these pools, 66 (4%) were reported WNV positive.

Pool statistics for the past two weeks are also included. If pool data exists for the forecast epiweek, then the two weeks will be the forecast week and the week prior. If data does not exist yet for the requested forecast epiweek, then the weeks shown will be the two epiweeks prior to the forecast week. In this report, the two weeks are 07/29/2018 through 08/05/2018 (epiweeks 31 & 32) with mosquito data existing between 07/29/2018 through 07/31/2018.

Count of WNV positive mosquito pools in past two weeks
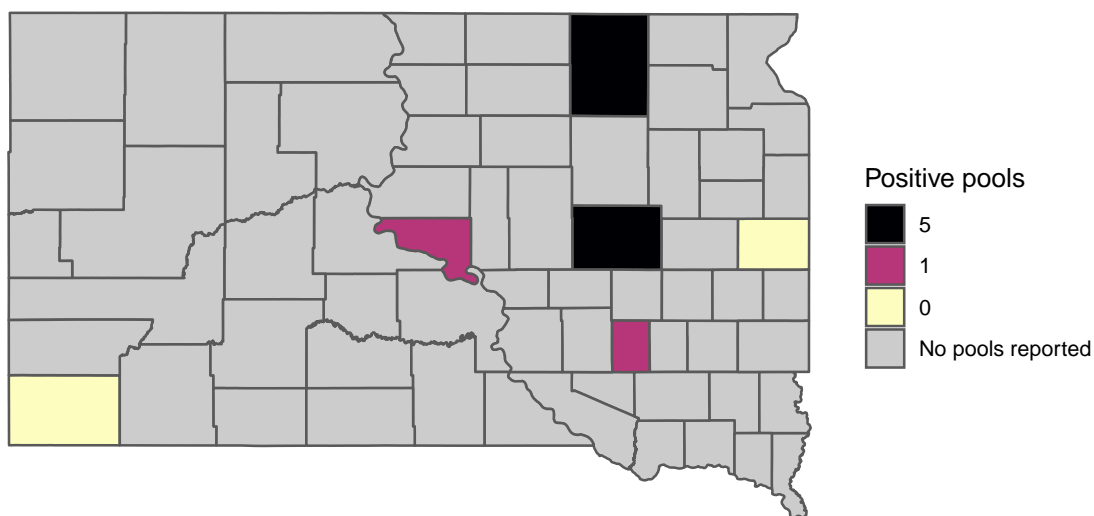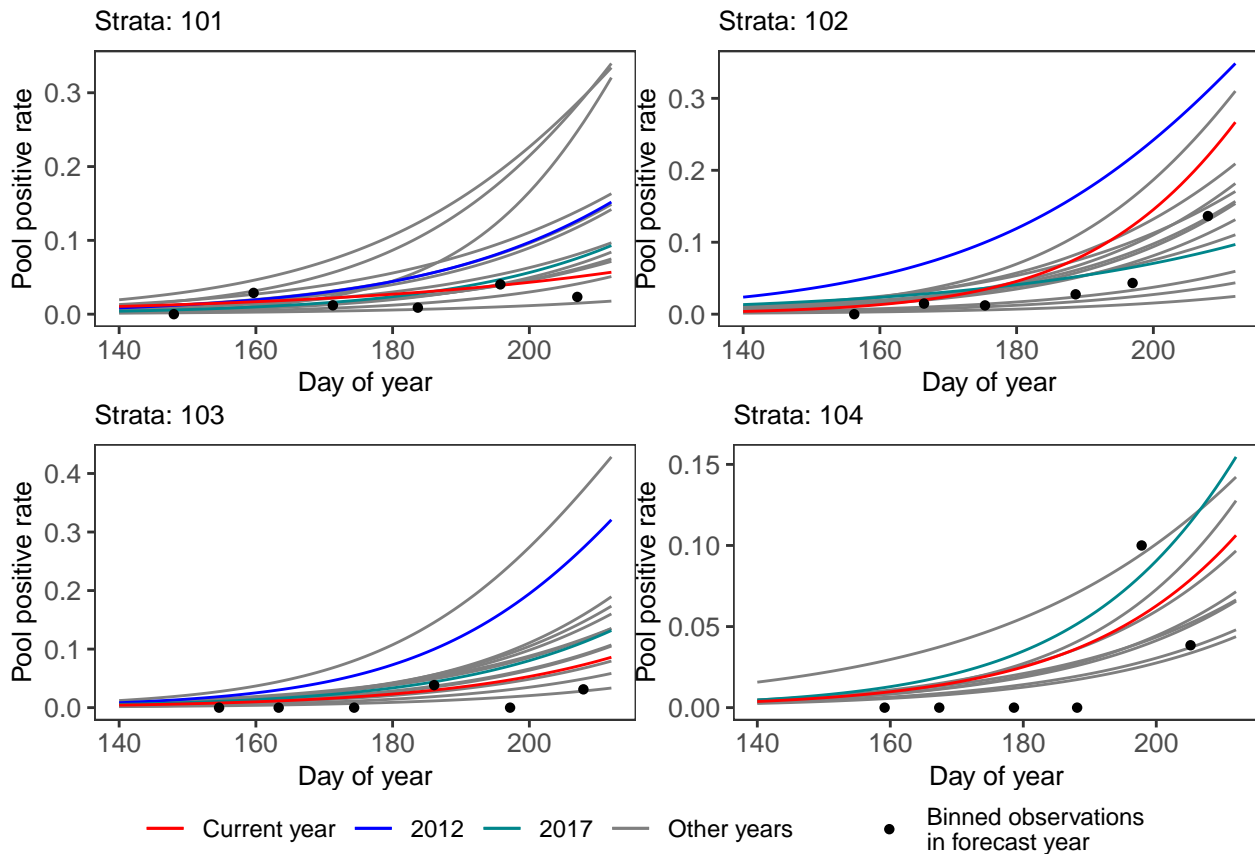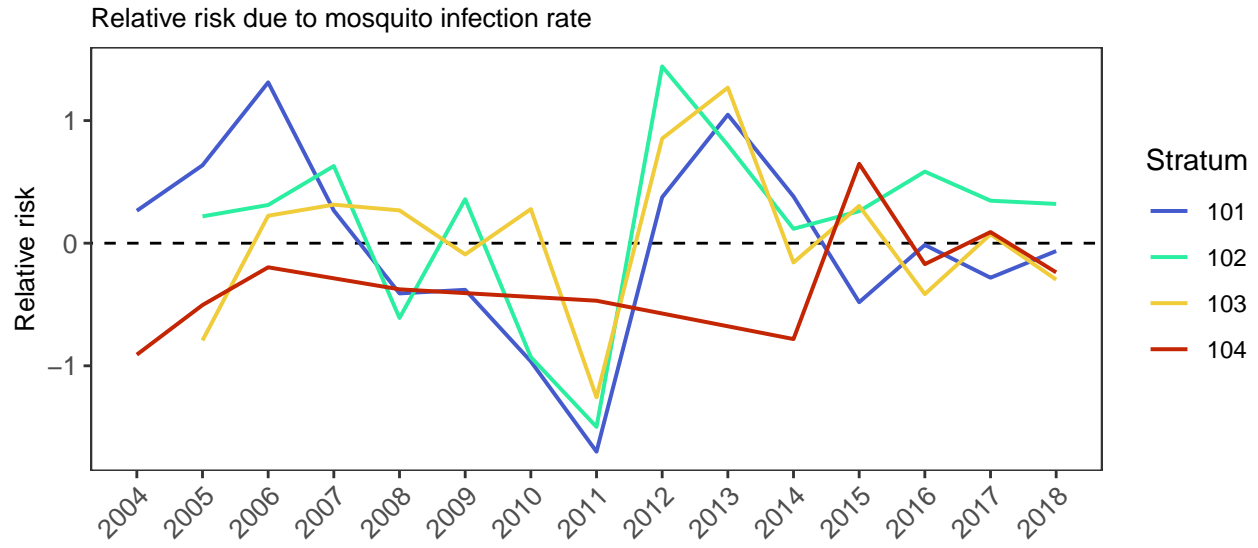Date range: 07/29/2018 through 08/05/2018

Table 4: Total reported and WNV positive mosquito pools: Past
two weeks, and year to forecast week (YTD)

| County | Pools reported last 2 weeks | Pools positive last 2 weeks | Pools reported YTD | Pools positive YTD |
|---|---|---|---|---|
| Beadle | 30 | 5 | 161 | 11 |
| Brookings | 13 | 0 | 290 | 7 |
| Brown | 38 | 5 | 556 | 30 |
| Custer | - | - | 5 | 0 |
| Davison | 5 | 1 | 38 | 2 |
| Edmunds | - | - | 21 | 1 |
| Fall River | 2 | 0 | 57 | 2 |
| Hughes | 16 | 1 | 76 | 1 |
| Lincoln | - | - | 47 | 1 |
| Minnehaha | - | - | 372 | 10 |
| Stanley | - | - | 29 | 1 |
| Todd | - | - | 45 | 0 |

The next graph shows the percentage of predicted positive pools by year comparing the forecast year (in
red) to the requested comparison years (shades of blue) and all other years (gray). If there is sufficient data
in the forecast year, the observed pools rates are shown as black dots, binned into six different time points.



The last mosquito graph shows the relative risk due to the mosquito infection rate as a time-series of all
known years. Mosquito strata are shown in different colors.

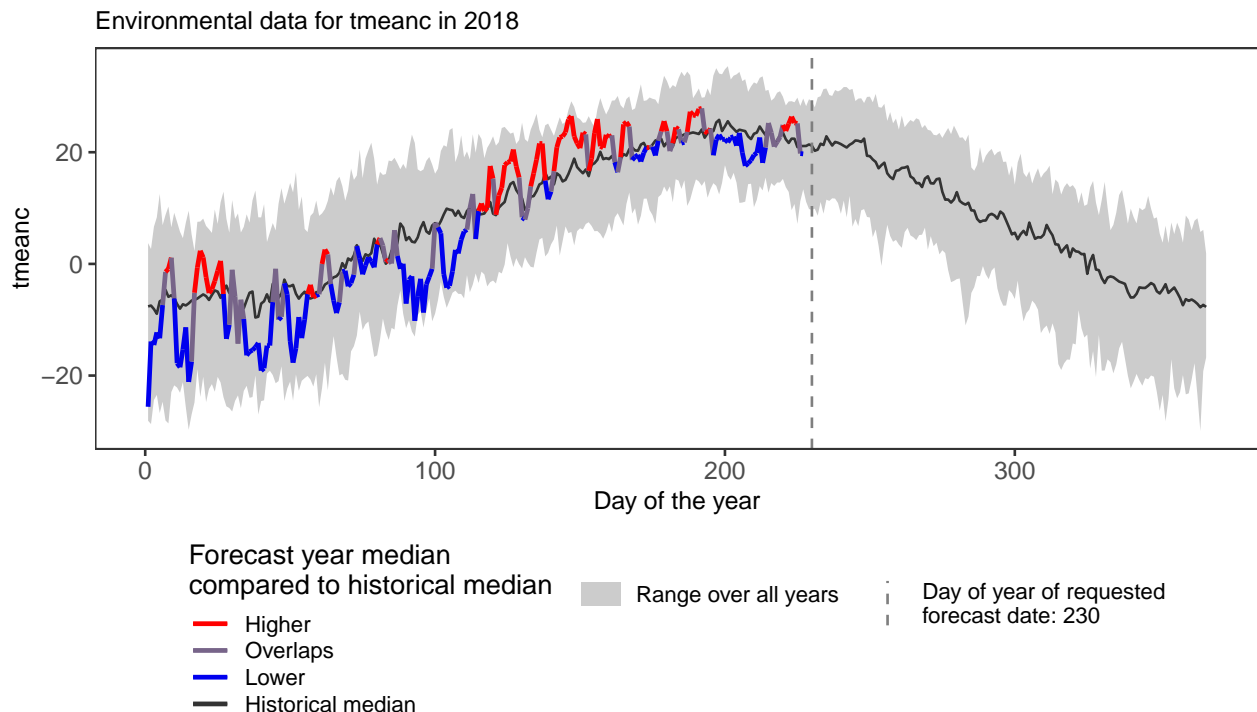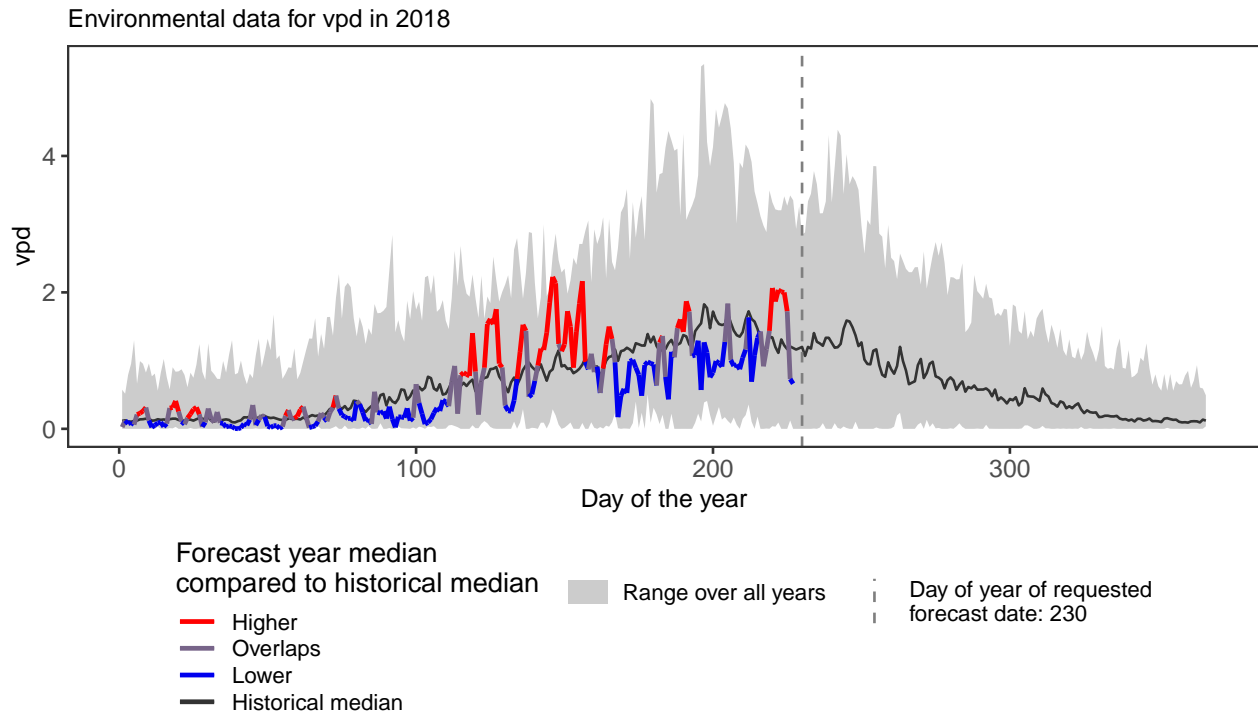Relative risk due to mosquito infection rate

## 2.3    Weather

After processing, weather data existed from 2001-01-01 through 2018-08-15 for 66 counties. Environmental data are read and collated from all files in the base `data_weather` folder and if there are duplicate data entries for any particular day, the value from the latest file is used (i.e. the latest updated value).

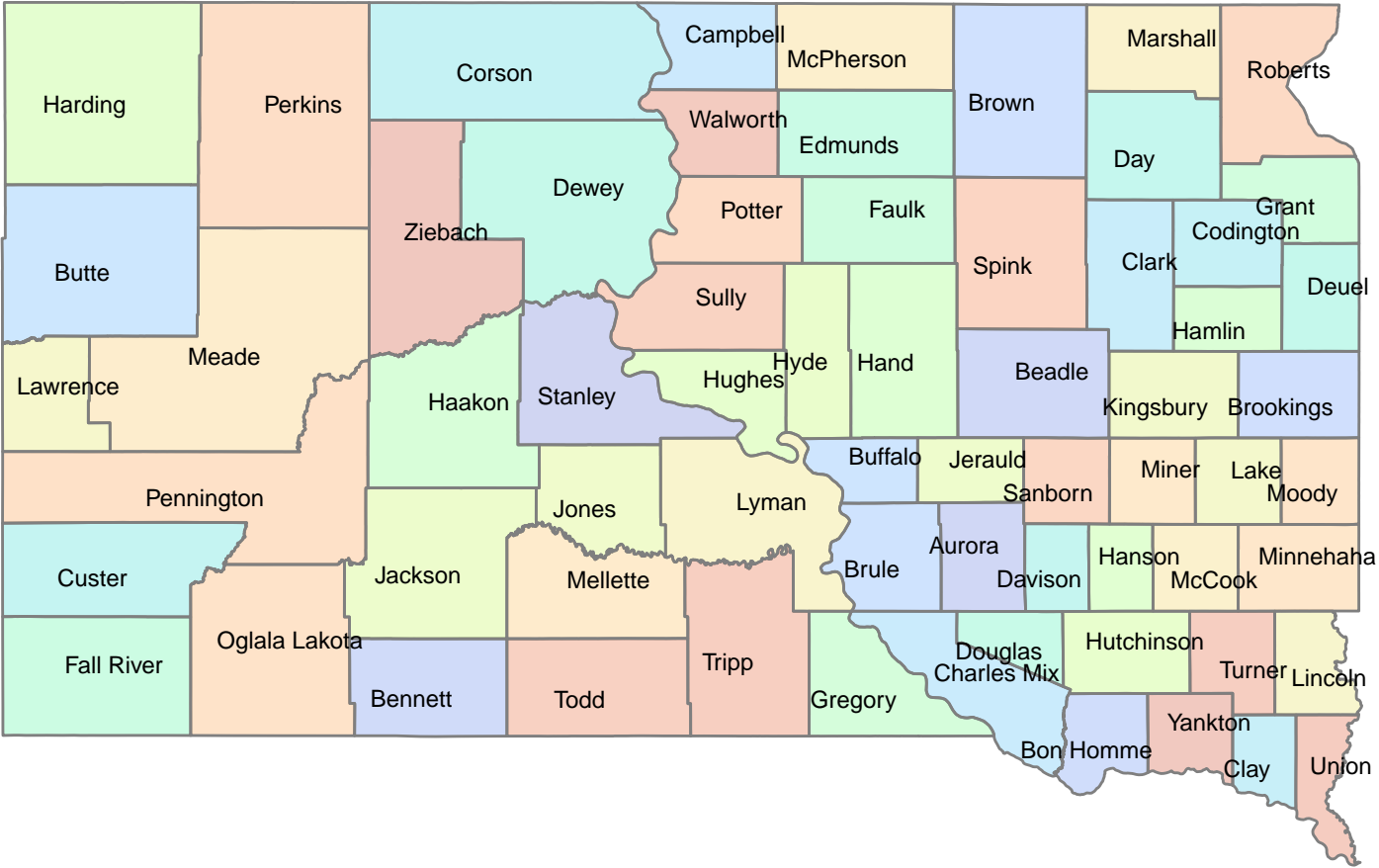Parameters were set to include weather data from 2000 through 2018. All necessary weather data were found in the data file.

The report parameters set the two environmental predictor variables as tmeanc and vpd. The following two graphs show the median state-wide observed weather variables for the forecast year, compared to the historical median. Two or more consecutive days that are **greater than** the historical median are drawn in red and consecutive days that are **less than** the historical median are drawn in blue. Consecutive days that overlap the historical median (i.e. one day above and the next below, or the opposite) are in purple. The gray shaded region is a ribbon showing the historical range (min to max). The appendix will show the anomaly graphs (same timeseries, but the weather variable has been anomalized).

Environmental data for vpd in 2018



Forecast year median
compared to historical median

— Higher
— Overlaps
— Lower
— Historical median

▮ Range over all years

┊ Day of year of requested
forecast date: 230

## 2.4 Reference map

For South Dakota, the spatial data (shapefile) contained 66 counties: Aurora, Beadle, Bennett, Bon Homme, Brookings, Brown, Brule, Buffalo, Butte, Campbell, Charles Mix, Clark, Clay, Codington, Corson, Custer, Davison, Day, Deuel, Dewey, Douglas, Edmunds, Fall River, Faulk, Grant, Gregory, Haakon, Hamlin, Hand, Hanson, Harding, Hughes, Hutchinson, Hyde, Jackson, Jerauld, Jones, Kingsbury, Lake, Lawrence, Lincoln, Lyman, Marshall, McCook, McPherson, Meade, Mellette, Miner, Minnehaha, Moody, Oglala Lakota, Pennington, Perkins, Potter, Roberts, Sanborn, Spink, Stanley, Sully, Todd, Tripp, Turner, Union, Walworth, Yankton, and Ziebach

## 2.5  Parameters used

The report was run with the following parameters set.

Table 5: Parameters used

| Parameter | Value |
| --- | --- |
| forecast_date | 2018-08-15 |
| state_name | South Dakota |
| state_code | SD |
| predictor_var1 | tmeanc |
| predictor_var2 | vpd |
| mosquito_model | stratifiedMIGR |
| mosquito_doy_start | 140 |
| mosquito_doy_end | 366 |
| file_human | data_human/simulated_human_data.csv |
| file_mosquito | data_mosquito/simulated_mosquito_data.csv |
| file_strata | data_strata/example_strata_SD.csv |
| file_county_sf | data_spatial/sd_counties.RDS |
| file_models | data_models/models.txt |
| folder_weather | data_weather |
| year_human_start | 2004 |
| year_human_end | 2017 |
| year_mosquito_start | 2004 |
| year_mosquito_end | 2018 |
| year_weather_start | 2000 |
| year_weather_end | 2018 |
| year_compare_vis1 | 2012 |
| year_compare_vis2 | 2017 |
| create_appendix | TRUE |
| lag_length | 121 |
| case_trim_alpha | 0.02 |

# 3    Appendix

This appendix will provide more details into some of the underlying forecast modeling and break out the results per model, rather than an average of all models run (as in the main report).

## 3.1    Forecast results

### 3.1.1    Current-week WNV absolute risk

Following are the absolute risk maps generated by **each** model:

Absolute risk in forecast week
Non−anomalized weather with fixed cubic splines: "cub−fx−nonanom"

## Absolute risk in forecast week

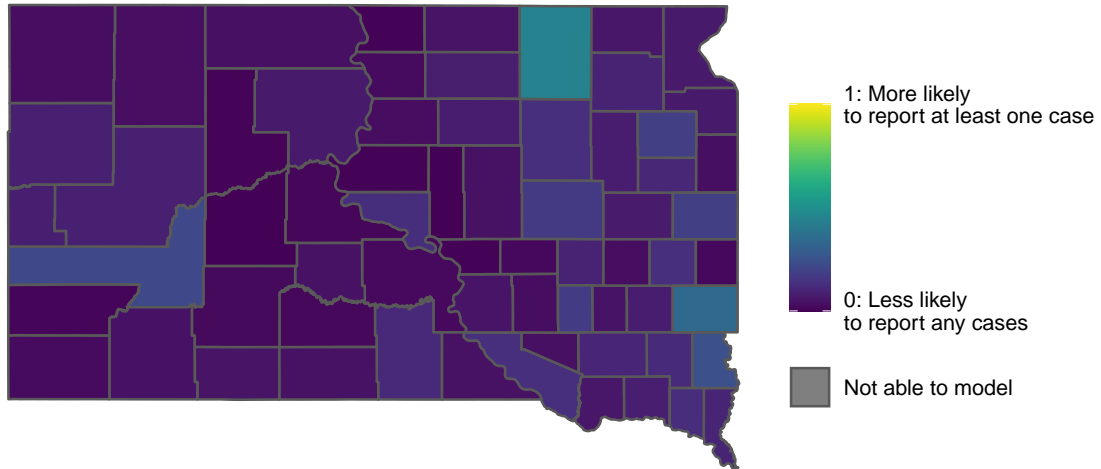Anomalized weather with fixed cubic splines: "cub–fx–anom"



## Absolute risk in forecast week

Non–anomalized weather with seasonally–varying cubic splines: "cub–sv–nonanom"

## Absolute risk in forecast week

Non–anomalized weather with seasonally–varying cubic splines: "cub–sv–anom"



## Absolute risk in forecast week

Non–anomalized weather with fixed thin plate splines: "tp–fx–nonanom"

## Absolute risk in forecast week
Anomalized weather with fixed thin plate splines: "tp–fx–anom"



## Absolute risk in forecast week
Non–anomalized weather with seasonally–varying thin plate splines: "tp–sv–nonanom"

## Absolute risk in forecast week

Anomalized weather with seasonally–varying thin plate splines: "tp–sv–anom"

### 3.1.2 Current-week WNV relative risk

Following are the relative risk maps generated by **each** model:

## Risk in forecast week relative to the same epiweek in previous years
Non−anomalized weather with fixed cubic splines: "cub−fx−nonanom"



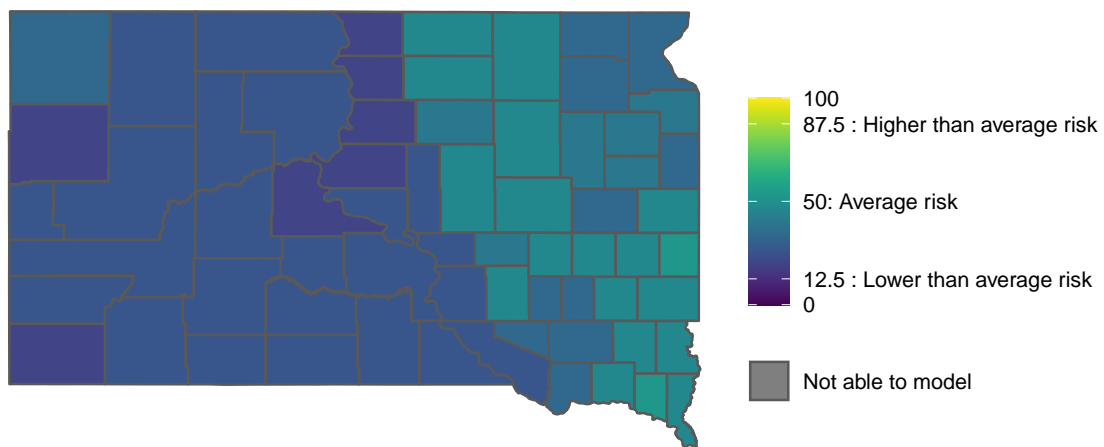## Risk in forecast week relative to the same epiweek in previous years
Anomalized weather with fixed cubic splines: "cub−fx−anom"

## Risk in forecast week relative to the same epiweek in previous years
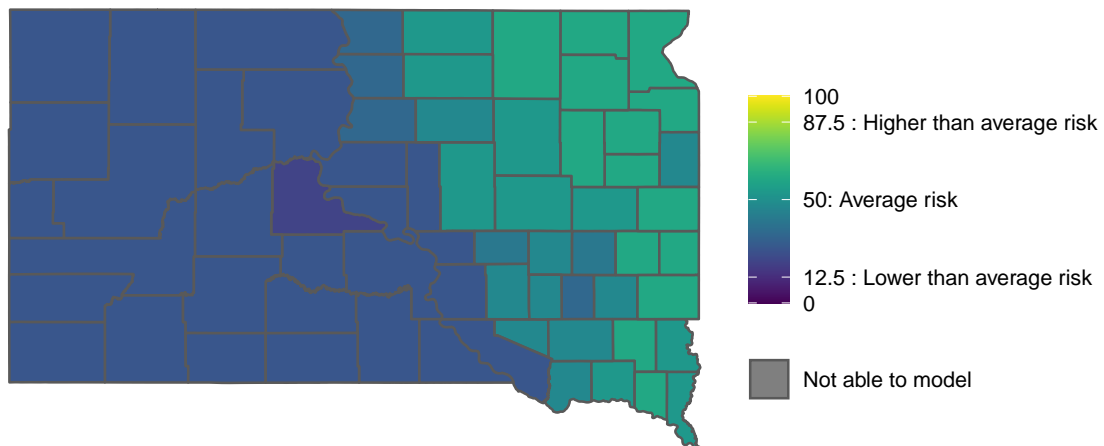Non–anomalized weather with seasonally–varying cubic splines: "cub–sv–nonanom"



## Risk in forecast week relative to the same epiweek in previous years
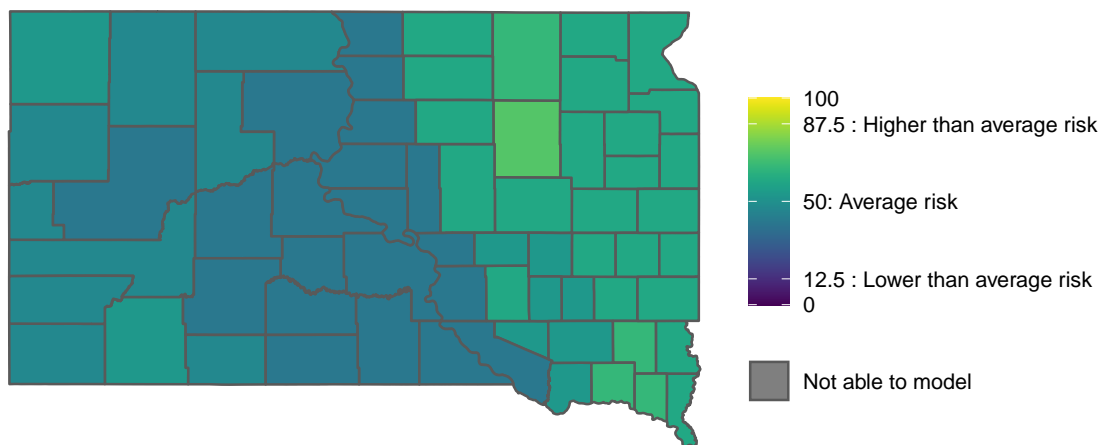Non–anomalized weather with seasonally–varying cubic splines: "cub–sv–anom"

## Risk in forecast week relative to the same epiweek in previous years
Non−anomalized weather with fixed thin plate splines: "tp−fx−nonanom"
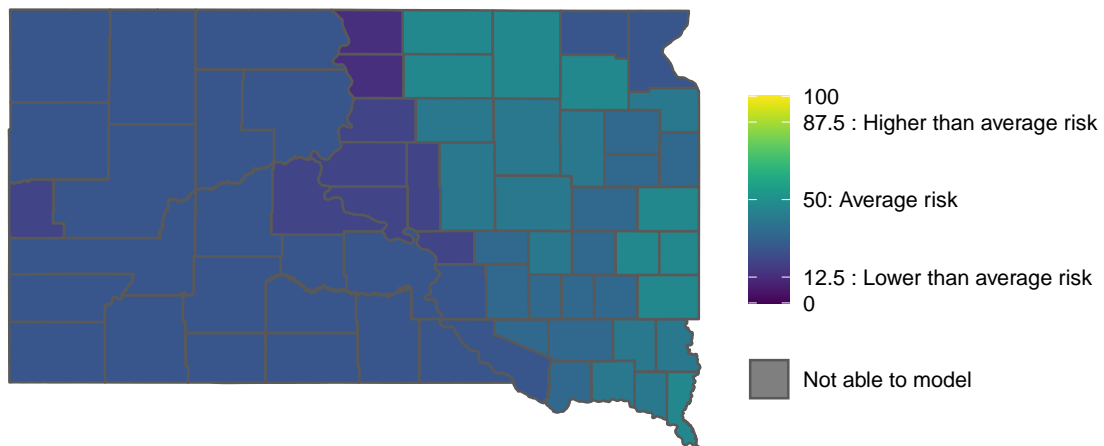


## Risk in forecast week relative to the same epiweek in previous years
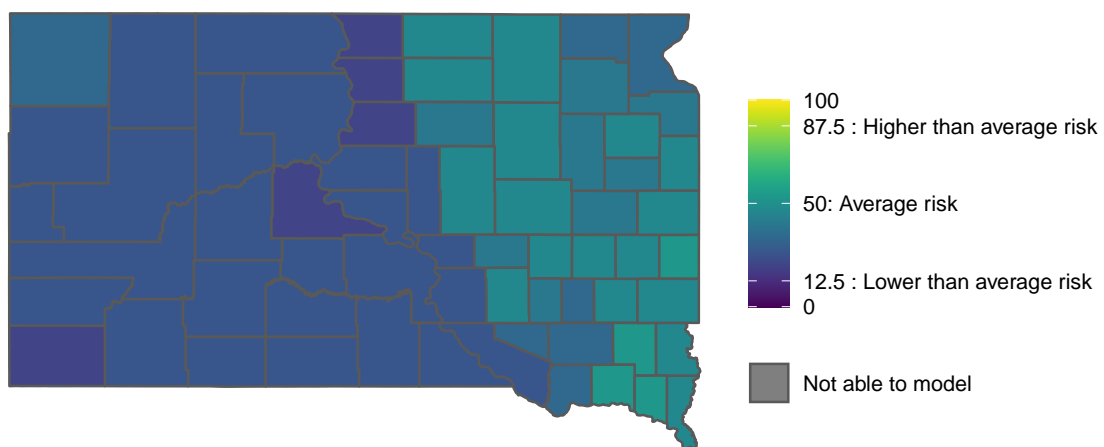Anomalized weather with fixed thin plate splines: "tp−fx−anom"

## Risk in forecast week relative to the same epiweek in previous years
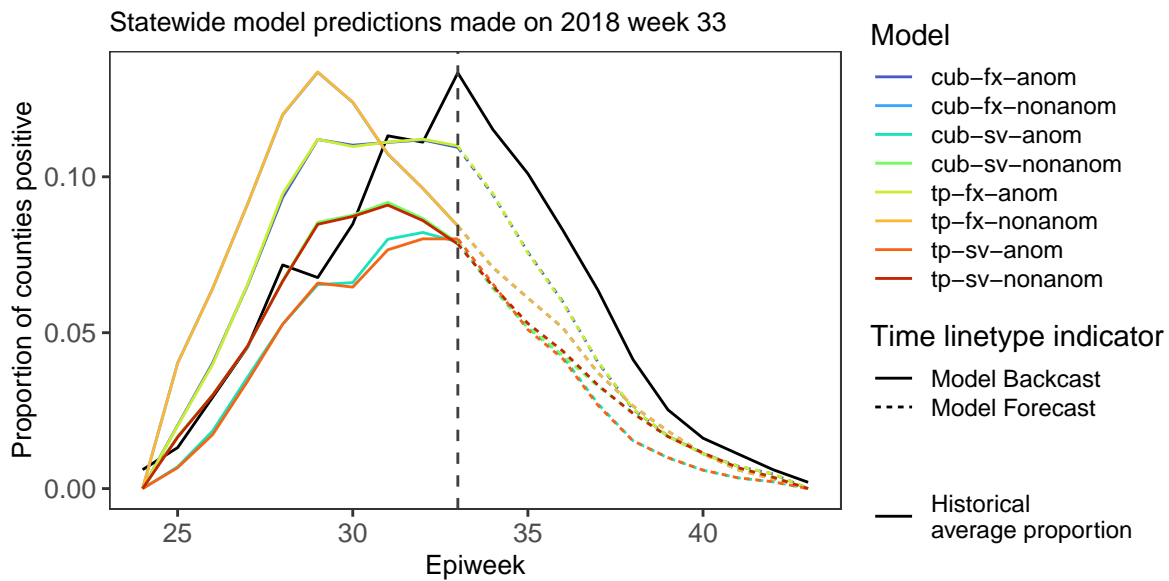Non−anomalized weather with seasonally−varying thin plate splines: "tp−sv−nonanom"



## Risk in forecast week relative to the same epiweek in previous years
Anomalized weather with seasonally−varying thin plate splines: "tp−sv−anom"

### 3.1.3 Current-year forecasts

The graph below shows the current year forecast, with lines for **each** model:



### 3.1.4 Case estimations

The table below lists the estimated case counts per model.

Table 6: Estimated number of WNV cases

| Year | Model | Predicted positive county-weeks | Estimated cases |
|------|-------|--------------------------------|-----------------|
| 2018 | cub-fx-anom | 73.2 | 86 |
| 2018 | cub-fx-nonanom | 75.7 | 89 |
| 2018 | cub-sv-anom | 46.8 | 56 |
| 2018 | cub-sv-nonanom | 55.7 | 66 |
| 2018 | tp-fx-anom | 73.3 | 86 |
| 2018 | tp-fx-nonanom | 75.7 | 89 |
| 2018 | tp-sv-anom | 46.2 | 55 |
| 2018 | tp-sv-nonanom | 55.7 | 66 |

### 3.1.5 Additional model fit statistics

The table below gives multiple model fit statistic per forecast model:

- AUC : Area Under ROC Curve, values range 0 - 1
- AIC : Akaike information criterion, relative fit statistic to other models
- Temporal MAE : Mean Average Error, mean of weeks (collapsed to state)
- Spatial MAE : Mean Average Error, mean of counties (collapsed all time)

Table 7: Fit statistcs by model

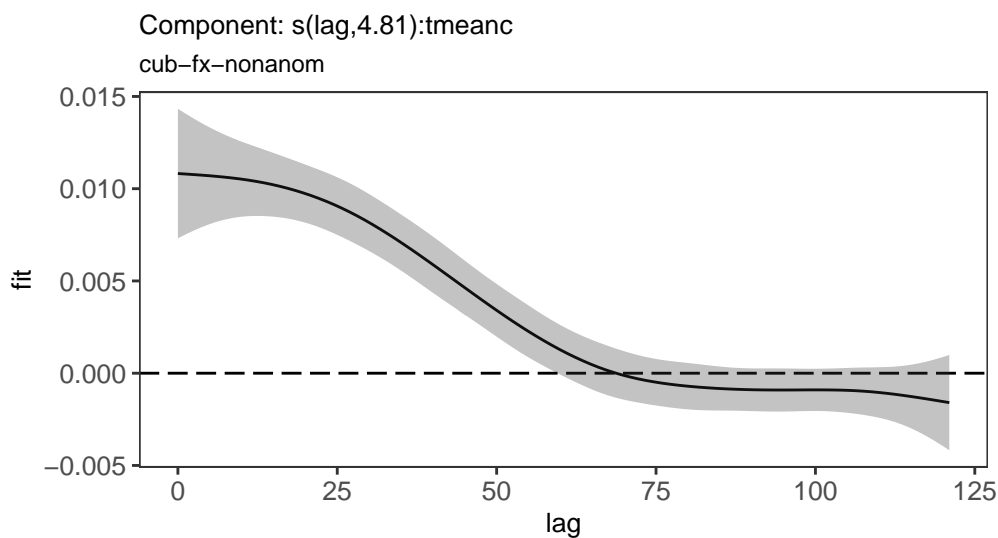| Model | AUC | AIC | MAE Temporal | MAE Spatial |
|-------|-----|-----|--------------|-------------|
| cub-fx-nonanom | 0.839 | 6728 | 1.837 | 0.772 |
| cub-fx-anom | 0.838 | 6749 | 1.791 | 0.766 |
| cub-sv-nonanom | 0.846 | 6634 | 1.622 | 0.743 |
| cub-sv-anom | 0.848 | 6645 | 1.539 | 0.729 |
| tp-fx-nonanom | 0.839 | 6728 | 1.837 | 0.772 |
| tp-fx-anom | 0.838 | 6749 | 1.790 | 0.766 |
| tp-sv-nonanom | 0.846 | 6634 | 1.620 | 0.743 |
| tp-sv-anom | 0.848 | 6644 | 1.538 | 0.729 |

### 3.1.6 Partial effects

ArboMAP allows the user to write custom model formulas, and as such the plots below are the partial effects of all the smooth terms for each model. These show the component effect of the term. All components (not just smooths) added together would be the overall prediction. For a table of all formulas, see the section on "Models and formulas".
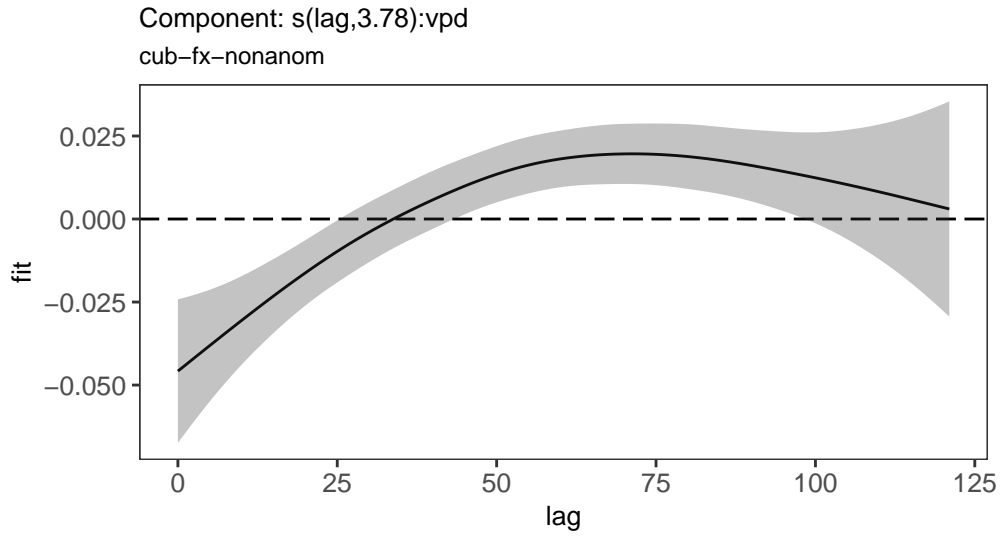
The number after the comma in the `s({item}, {number})` labels is the effective degrees of freedom (EDF). The EDF is a measurement of the complexity of the smooth term - a value of 1 is a straight line, higher values are more complex curves.

An easy way to check on the significance of the smooth term is if you cannot draw a horizontal line through the 95% confidence interval (value +/- se, shown in the gray shaded ribbon in the relevant graphs).
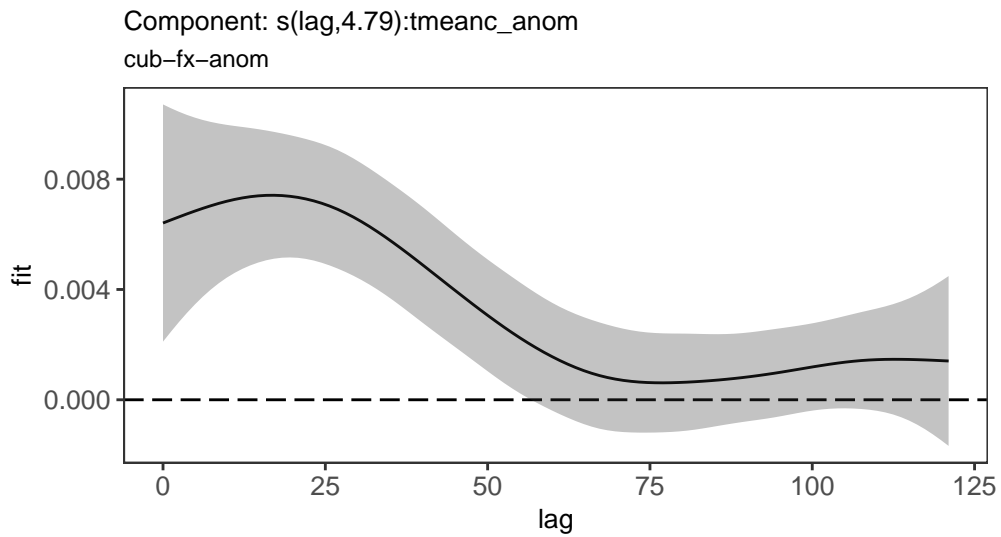
All models with smooths will have 1-D graphs. Seasonally-varying models will also have 2-D graphs (components with `doymat` in standard models), however a subset of y-values have been pulled out to plot as lines.
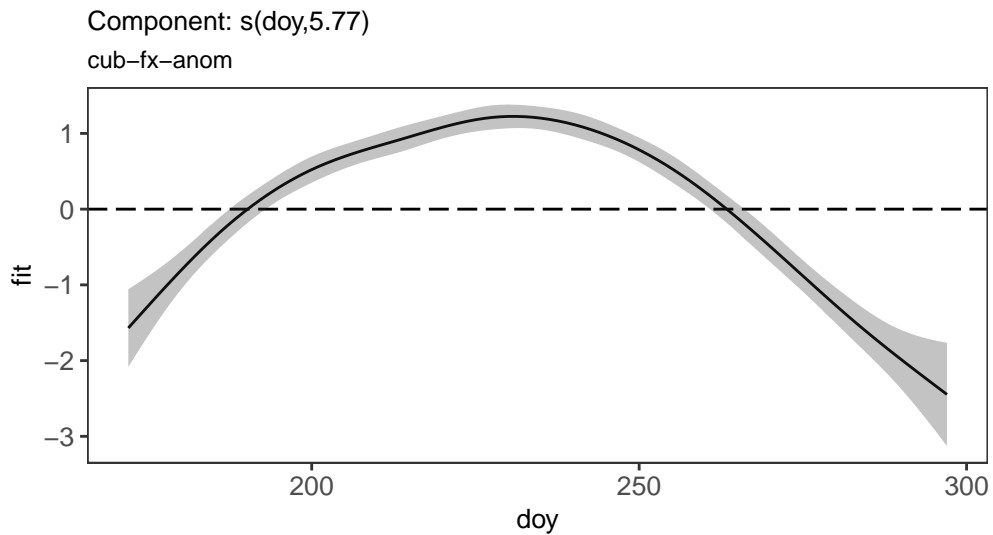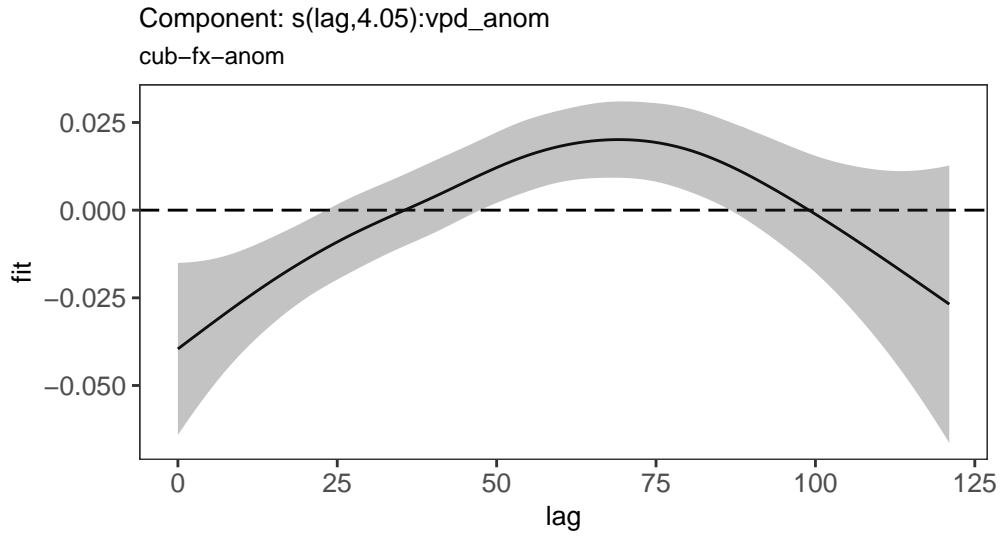
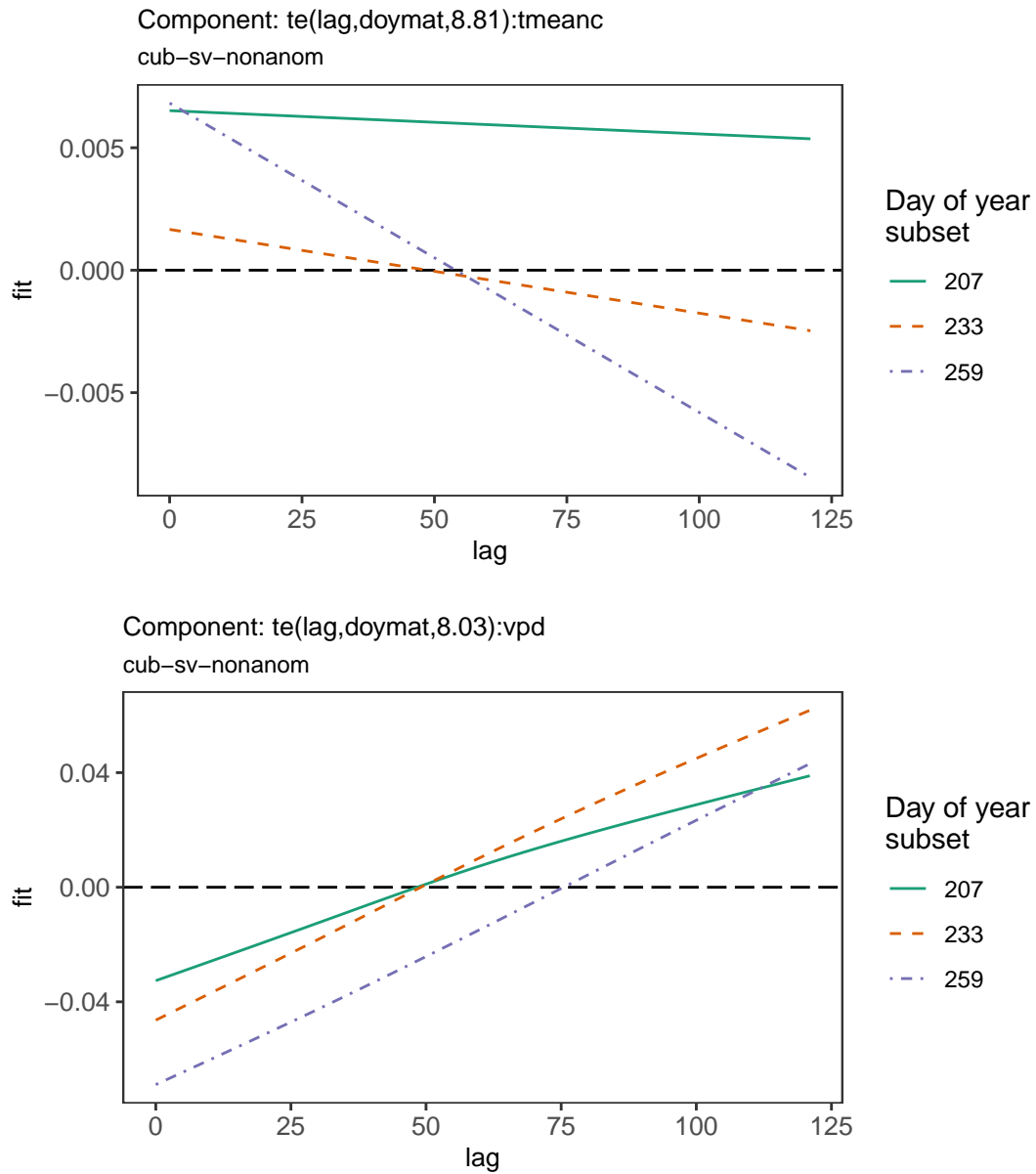#### 3.1.6.1 Non-anomalized weather with fixed cubic splines: "cub-fx-nonanom"

Component: s(lag,3.78):vpd

cub–fx–nonanom



### 3.1.6.2 Anomalized weather with fixed cubic splines: "cub-fx-anom"

Component: s(lag,4.79):tmeanc_anom

cub–fx–anom

Component: s(lag,4.05):vpd_anom
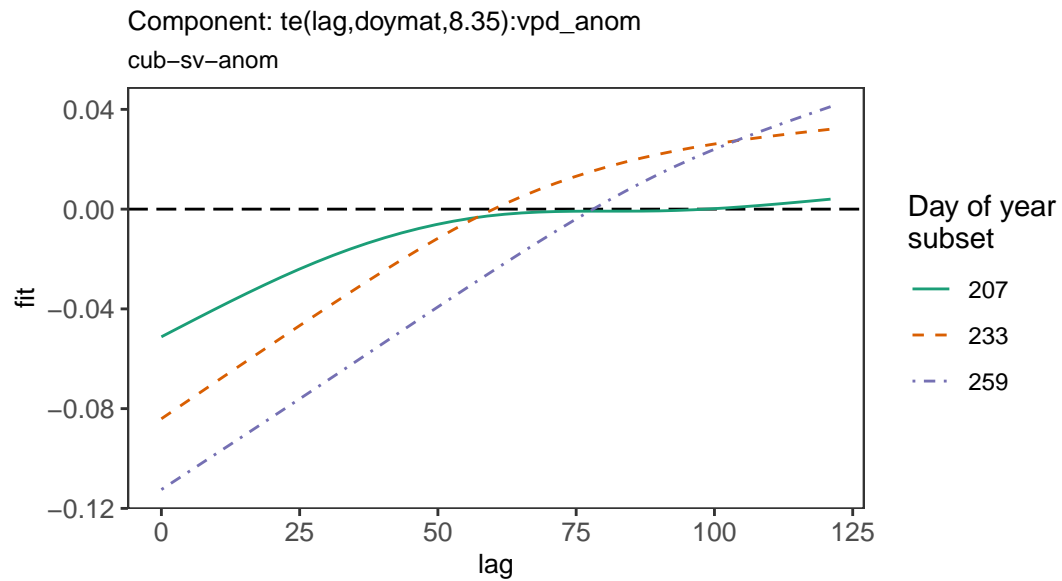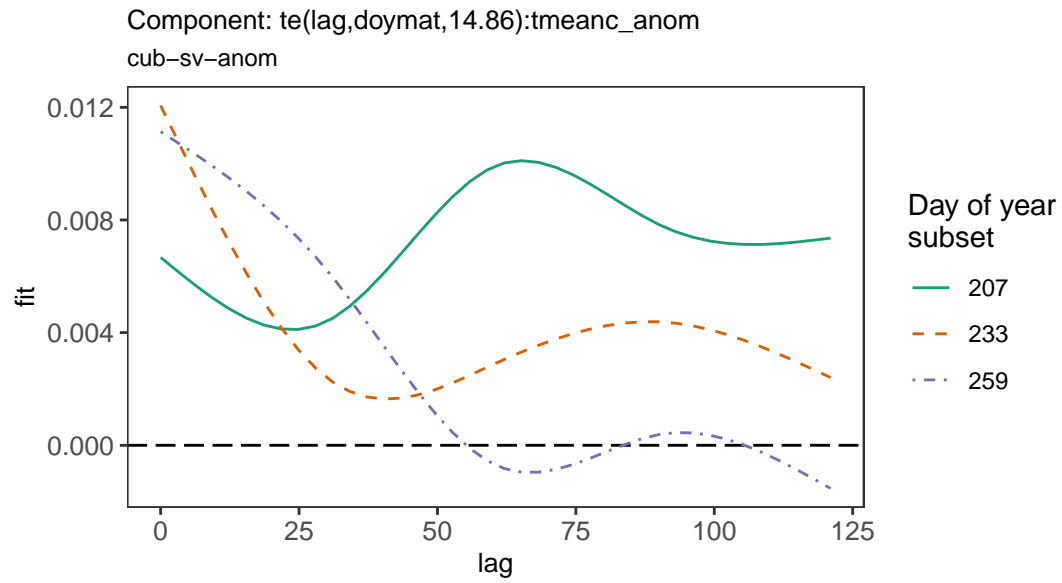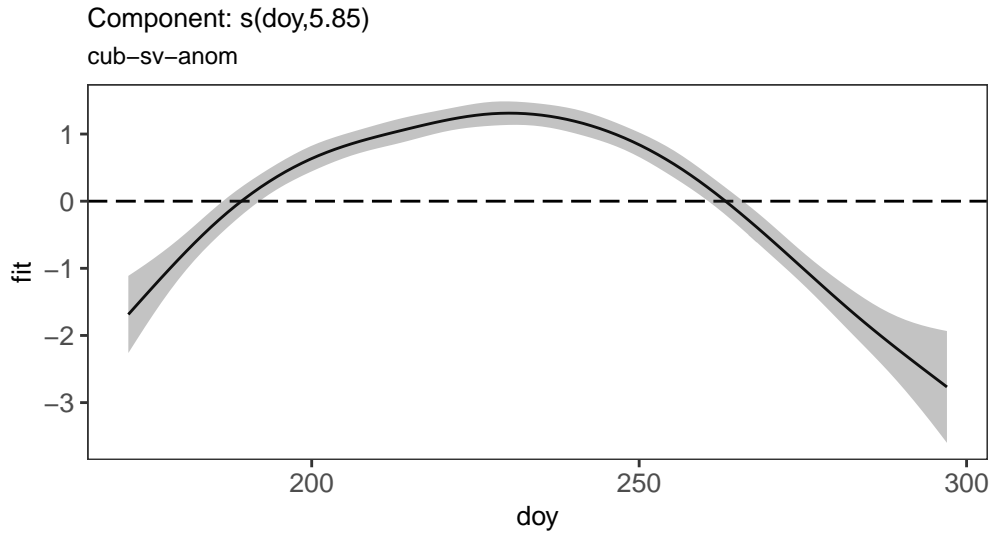
cub−fx−anom

Component: s(doy,5.77)

cub−fx−anom

### 3.1.6.3 Non-anomalized weather with seasonally-varying cubic splines: "cub-sv-nonanom"

Component: te(lag,doymat,8.81):tmeanc

cub−sv−nonanom



Component: te(lag,doymat,8.03):vpd

cub−sv−nonanom

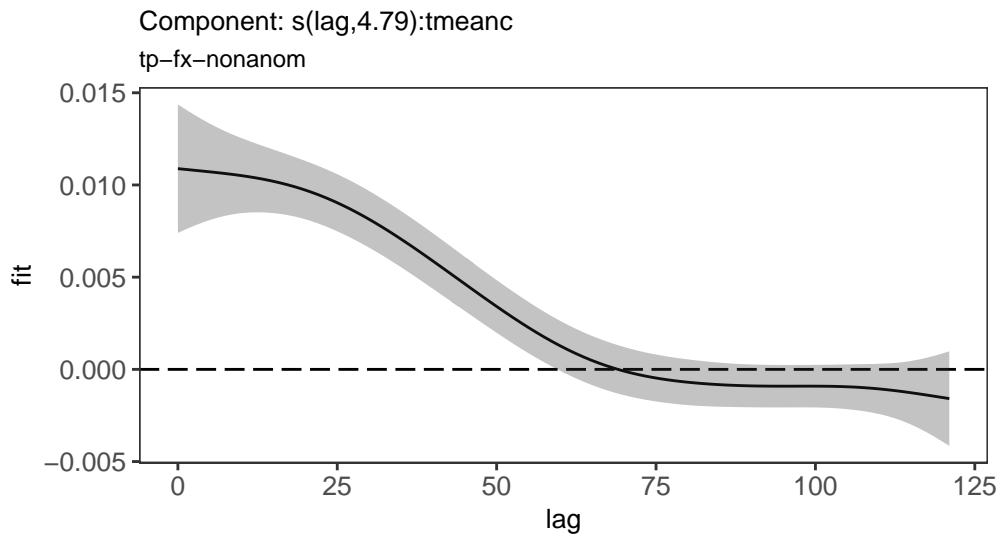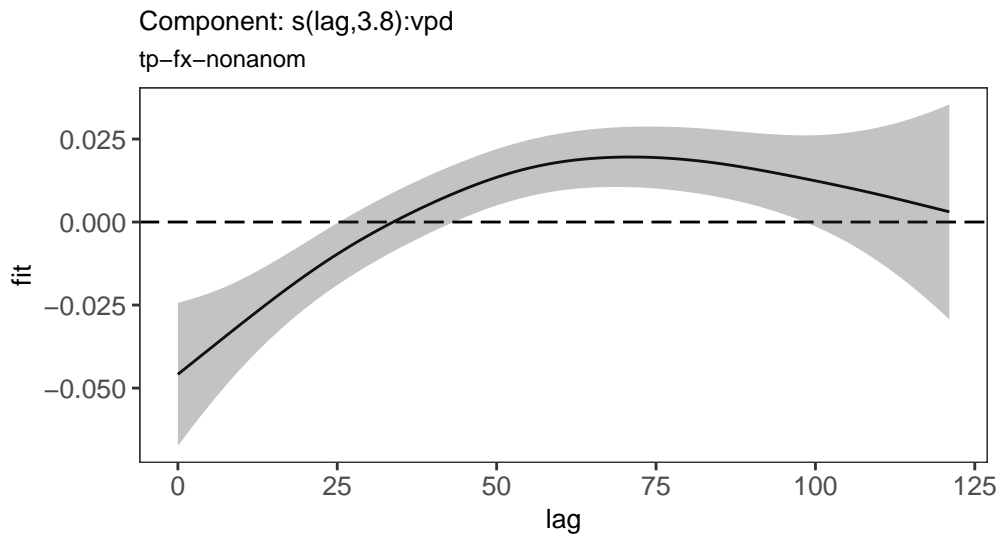**3.1.6.4   Non-anomalized weather with seasonally-varying cubic splines: "cub-sv-anom"**

Component: te(lag,doymat,14.86):tmeanc_anom

cub–sv–anom



Component: te(lag,doymat,8.35):vpd_anom
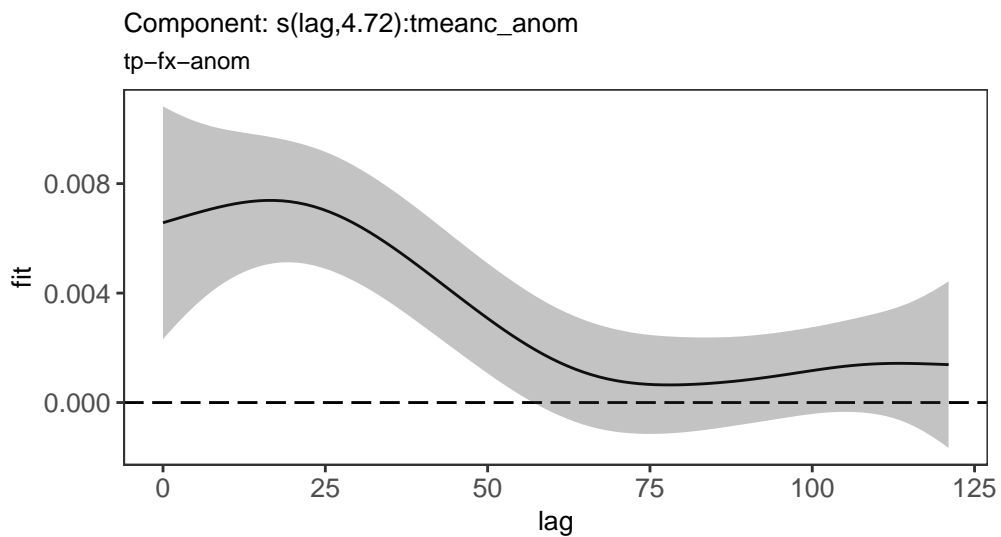
cub–sv–anom

Component: s(doy,5.85)

cub−sv−anom



### 3.1.6.5  Non-anomalized weather with fixed thin plate splines: "tp-fx-nonanom"

Component: s(lag,4.79):tmeanc

tp−fx−nonanom

Component: s(lag,3.8):vpd

tp–fx–nonanom

### 3.1.6.6 Anomalized weather with fixed thin plate splines: "tp-fx-anom"



Component: s(lag,4.72):tmeanc_anom

tp–fx–anom

Component: s(lag,4.05):vpd_anom

tp–fx–anom



Component: s(doy,5.86)

tp–fx–anom



**3.1.6.7  Non-anomalized weather with seasonally-varying thin plate splines: "tp-sv-nonanom"**
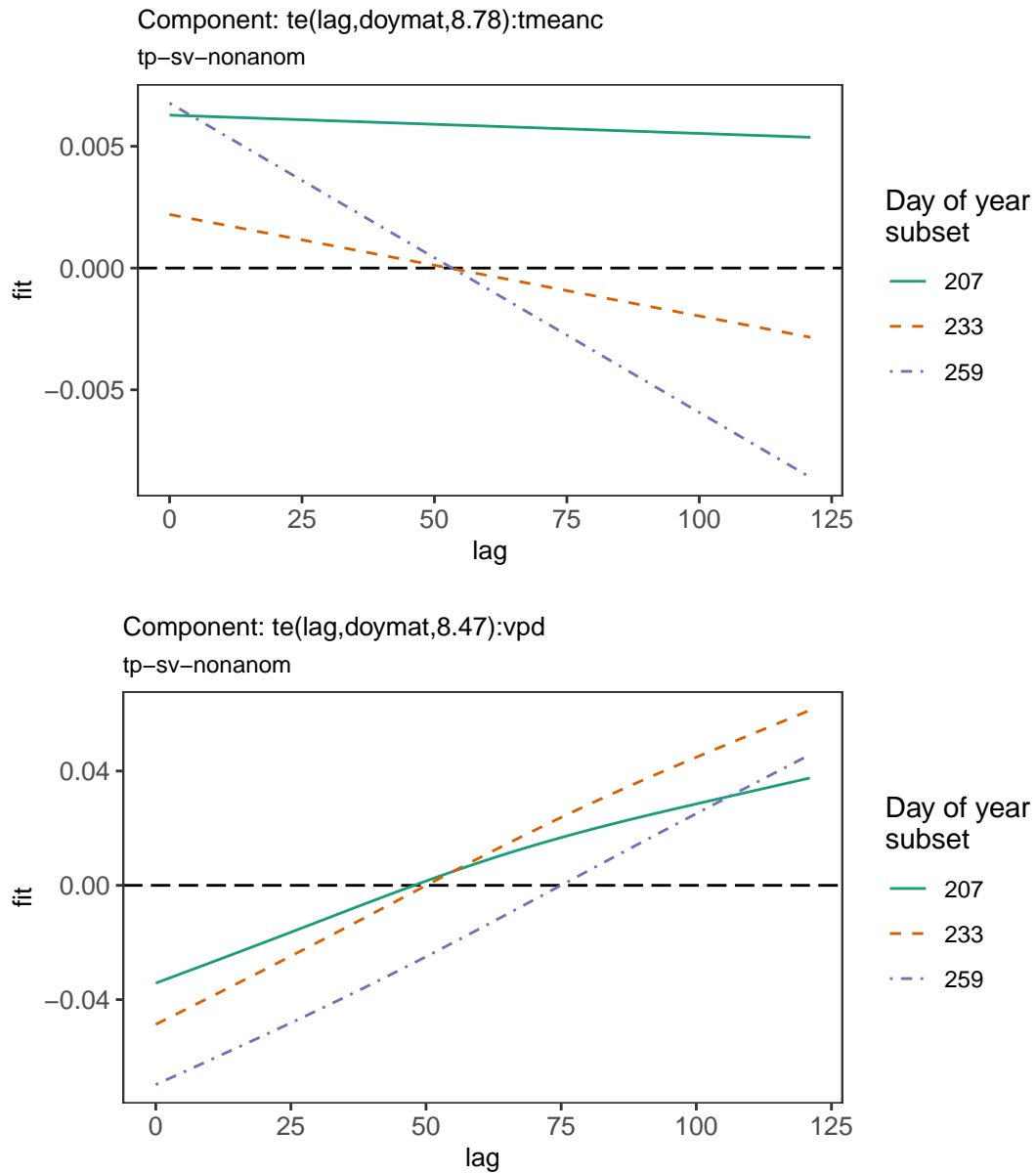
Component: te(lag,doymat,8.78):tmeanc
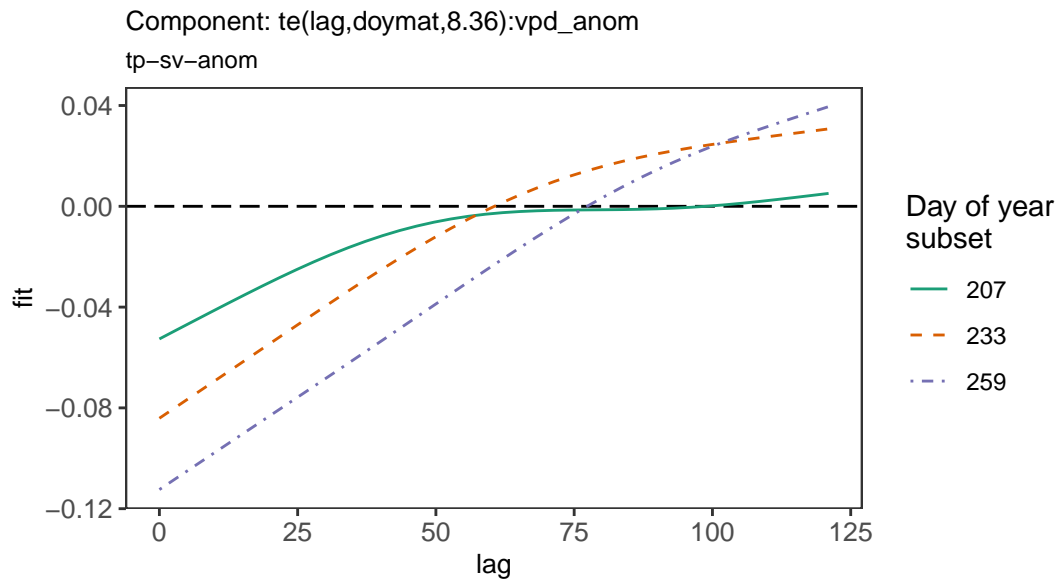tp−sv−nonanom



Component: te(lag,doymat,8.47):vpd
tp−sv−nonanom

**3.1.6.8   Anomalized weather with seasonally-varying thin plate splines: "tp-sv-anom"**

Component: te(lag,doymat,15.08):tmeanc_anom



Component: te(lag,doymat,8.36):vpd_anom

Component: s(doy,5.94)

tp−sv−anom

FOR DEMONSTRATION ONLY

### 3.1.7 Multi-year forecasts

The graph below shows the full forecast for all years, with lines for **each** model:



Statewide model predictions

### 3.1.8 Models and formulas

The table below lists the models that were found in the data_models/models.txt file. Standard models will have a text description, but all models run should appear in the table along with their formula.

The following fields may be present:

- `any_cases` : positive county-week
- `arbo_ID` : internal field for identifying counties
- `mir_stat` : the mosquito infection rate statistic
- `s(lag, by=var...)` : fixed smooth term for the environmental variable over the distributed lag period
- `te(lag, doymat, by=var...)` : seasonally-varying smooth term for the environmental variable over the distributed lag period
- `var1` : variable for tmeanc, observed value
- `var2` : variable for vpd, observed value
- `var1_anom` : variable for tmeanc, anomalized value
- `var2_anom` : variable for vpd, anomalized value
- `s(doy,...)` : smooth term for day of year, for seasonality

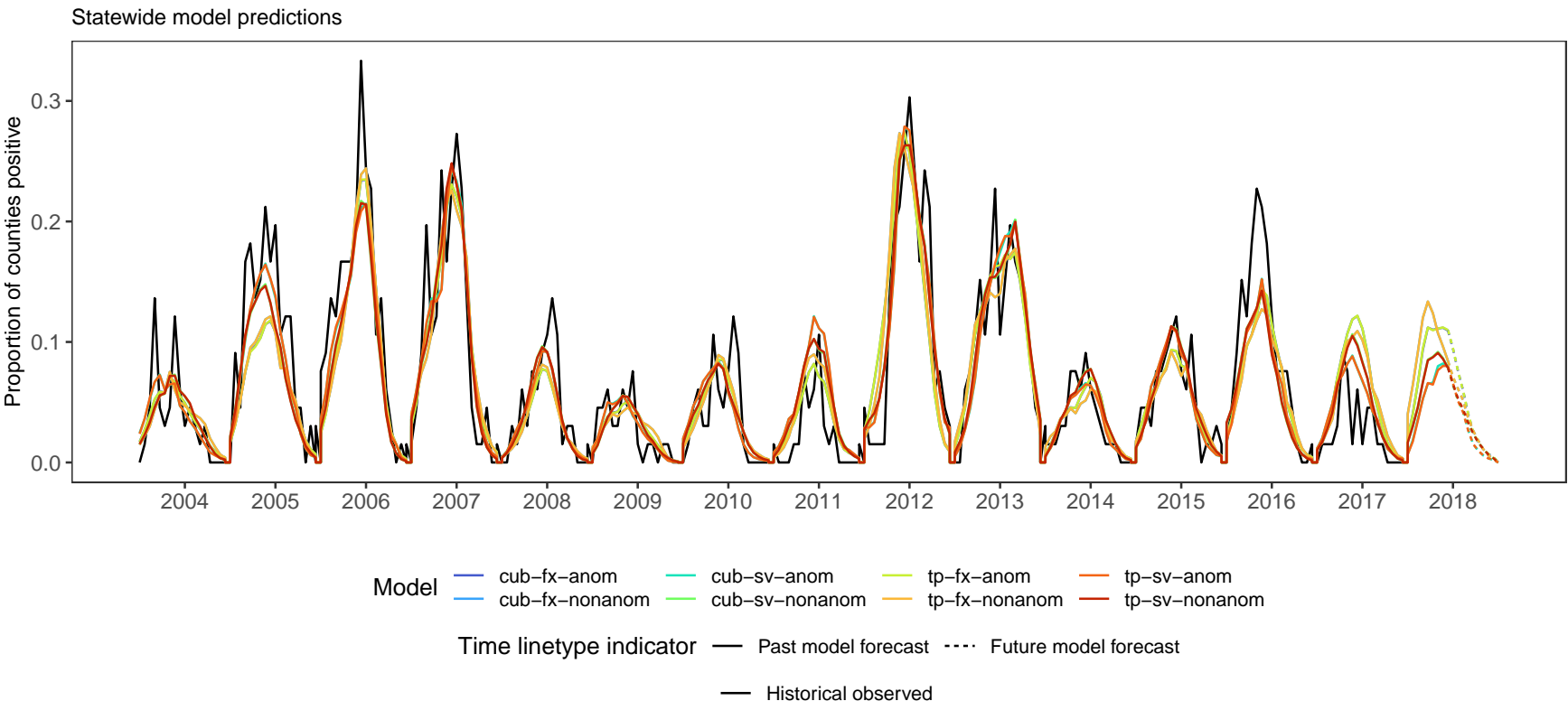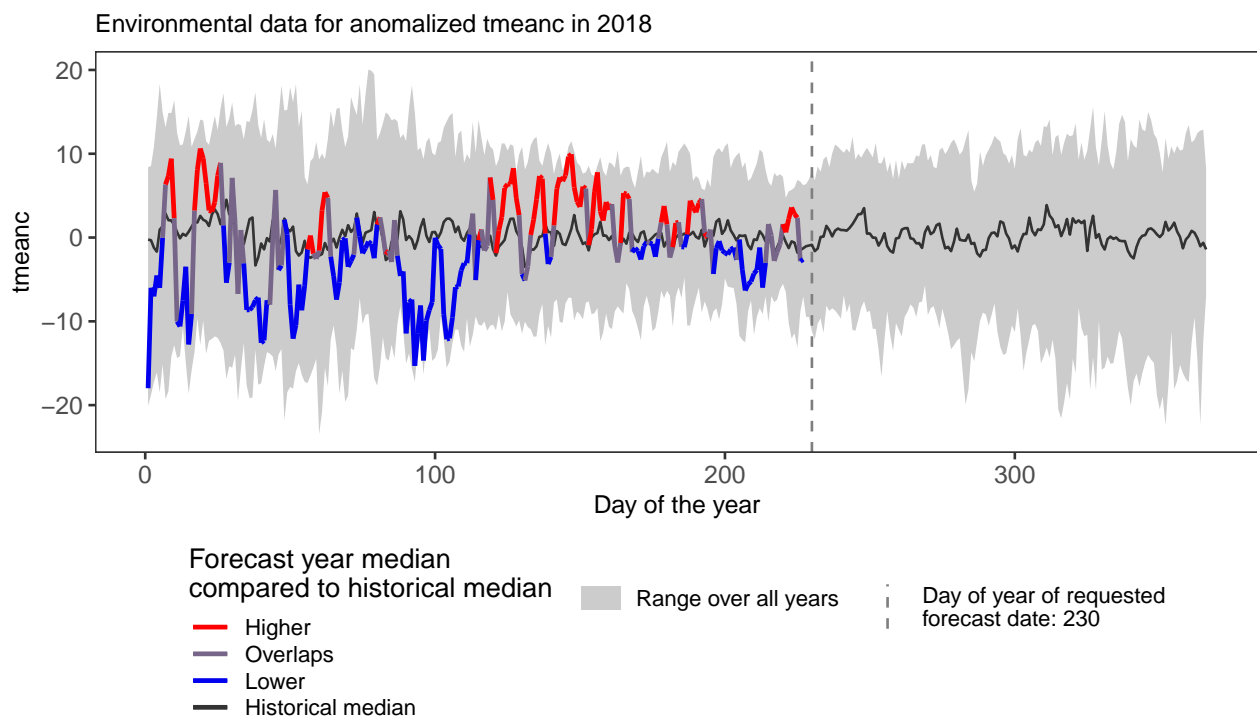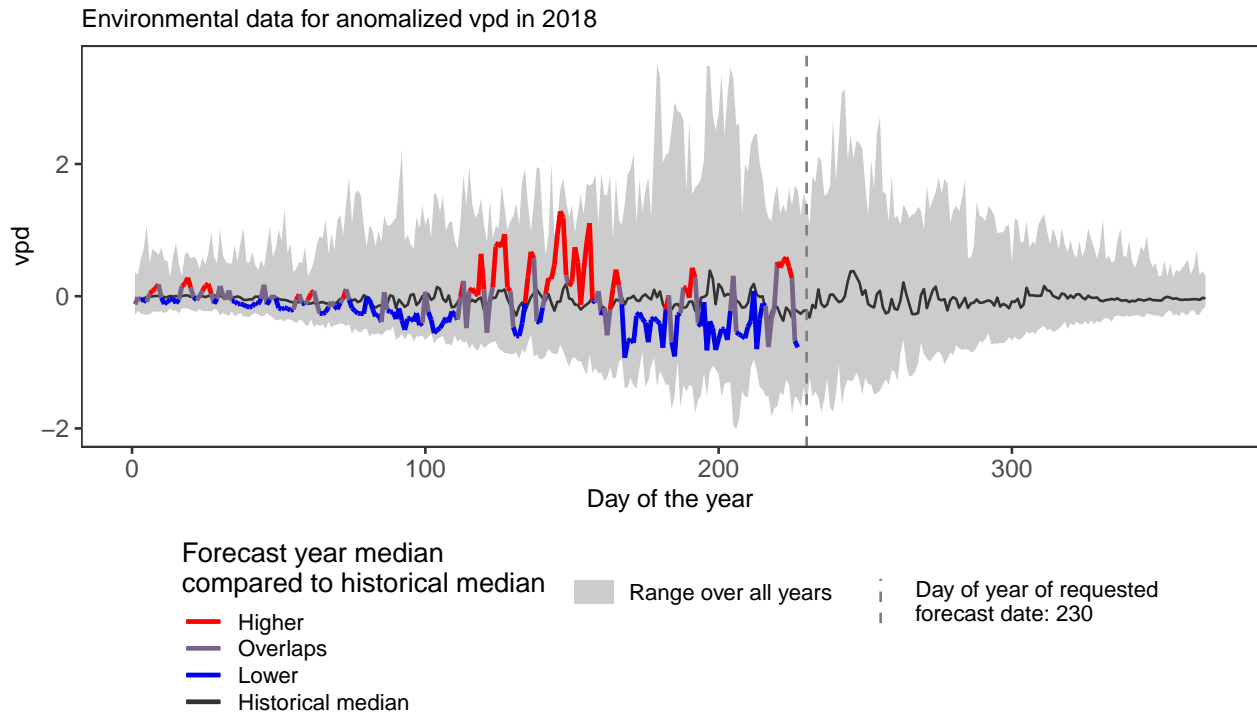| Model | Description | Formula |
|---|---|---|
| cub-fx-nonanom | Non-anomalized weather with fixed cubic splines | any_cases ~ 0 + arbo_ID + mir_stat + s(lag, by=var1, bs='cr') + s(lag, by=var2, bs='cr') |
| cub-fx-anom | Anomalized weather with fixed cubic splines | any_cases ~ 0 + arbo_ID + mir_stat + s(lag, by=var1_anom, bs='cr') + s(lag, by=var2_anom, bs='cr') + s(doy, bs='cr') |
| cub-sv-nonanom | Non-anomalized weather with seasonally-varying cubic splines | any_cases ~ 0 + arbo_ID + mir_stat + te(lag, doymat, by=var1, bs='cr') + te(lag, doymat, by=var2, bs='cr') |
| cub-sv-anom | Non-anomalized weather with seasonally-varying cubic splines | any_cases ~ 0 + arbo_ID + mir_stat + te(lag, doymat, by=var1_anom, bs='cr') + te(lag, doymat, by=var2_anom, bs='cr') + s(doy, bs='cr') |
| tp-fx-nonanom | Non-anomalized weather with fixed thin plate splines | any_cases ~ 0 + arbo_ID + mir_stat + s(lag, by=var1, bs='tp') + s(lag, by=var2, bs='tp') |
| tp-fx-anom | Anomalized weather with fixed thin plate splines | any_cases ~ 0 + arbo_ID + mir_stat + s(lag, by=var1_anom, bs='tp') + s(lag, by=var2_anom, bs='tp') + s(doy, bs='tp') |
| tp-sv-nonanom | Non-anomalized weather with seasonally-varying thin plate splines | any_cases ~ 0 + arbo_ID + mir_stat + te(lag, doymat, by=var1, bs='tp') + te(lag, doymat, by=var2, bs='tp') |
| tp-sv-anom | Anomalized weather with seasonally-varying thin plate splines | any_cases ~ 0 + arbo_ID + mir_stat + te(lag, doymat, by=var1_anom, bs='tp') + te(lag, doymat, by=var2_anom, bs='tp') + s(doy, bs='tp') |

## 3.2 Data summaries

### 3.2.1 Anomalized environmental variables

The report parameters set the two environmental predictor variables as tmeanc and vpd. The following two graphs show the median state-wide anomalized weather variables for the forecast year, compared to the historical median. Anomalies are calculated using deviance between the observed value and the predicted value from a GAM regression model using county and a smooth on day of year (seasonality) and county. An anomaly is the observed minus the predicted.

Two or more consecutive days that have anomalized values **greater** than the anomalized historical median are drawn in red and consecutive days that are **less** than the historical median are drawn in blue. Consecutive days that overlap the historical median (i.e. one day above and the next below, or the opposite) are in purple. The gray shaded region is a ribbon showing the historical range (min to max).



Environmental data for anomalized tmeanc in 2018

Environmental data for anomalized vpd in 2018



### 3.2.2 Modeled mosquito infection rate

In modeling years where sufficient mosquito data were present, the mosquito infection rate (MIR) statistic was created using the model specified in the input parameter: stratifiedMIGR. The following table presents the calculated centered MIR values that were used in the forecast modeling.

Table 9: Mosquito model summary statistic

| Year | Stratum | Centered MIR stat |
|------|---------|-------------------|
| 2004 | 101 | 0.265 |
| 2004 | 104 | -0.909 |
| 2005 | 101 | 0.635 |
| 2005 | 102 | 0.219 |
| 2005 | 103 | -0.793 |
| 2005 | 104 | -0.505 |
| 2006 | 101 | 1.312 |
| 2006 | 102 | 0.312 |
| 2006 | 103 | 0.224 |
| 2006 | 104 | -0.197 |
| 2007 | 101 | 0.266 |
| 2007 | 102 | 0.629 |
| 2007 | 103 | 0.314 |
| 2008 | 101 | -0.410 |
| 2008 | 102 | -0.610 |
| 2008 | 103 | 0.268 |
| 2008 | 104 | -0.376 |
| 2009 | 101 | -0.380 |
| 2009 | 102 | 0.359 |
| 2009 | 103 | -0.092 |

| Year | Stratum | Centered MIR stat |
|------|---------|-------------------|
| 2010 | 101 | -0.964 |
| 2010 | 102 | -0.929 |
| 2010 | 103 | 0.278 |
| 2011 | 101 | -1.700 |
| 2011 | 102 | -1.498 |
| 2011 | 103 | -1.257 |
| 2011 | 104 | -0.469 |
| 2012 | 101 | 0.375 |
| 2012 | 102 | 1.442 |
| 2012 | 103 | 0.853 |
| 2013 | 101 | 1.048 |
| 2013 | 102 | 0.799 |
| 2013 | 103 | 1.268 |
| 2014 | 101 | 0.381 |
| 2014 | 102 | 0.117 |
| 2014 | 103 | -0.158 |
| 2014 | 104 | -0.782 |
| 2015 | 101 | -0.481 |
| 2015 | 102 | 0.261 |
| 2015 | 103 | 0.304 |
| 2015 | 104 | 0.647 |
| 2016 | 101 | -0.014 |
| 2016 | 102 | 0.584 |
| 2016 | 103 | -0.415 |
| 2016 | 104 | -0.171 |
| 2017 | 101 | -0.282 |
| 2017 | 102 | 0.347 |
| 2017 | 103 | 0.073 |
| 2017 | 104 | 0.091 |
| 2018 | 101 | -0.064 |
| 2018 | 102 | 0.320 |
| 2018 | 103 | -0.297 |
| 2018 | 104 | -0.238 |