

Ecological forecasting in R

Lecture 4: evaluating dynamic models

Nicholas Clark

School of Veterinary Science, University of Queensland

0900–1200 CET Wednesday 6th September, 2023

Workflow

Press the "o" key on your keyboard to navigate among slides

Access the [tutorial html here](#)

Download the data objects and exercise  script from the html file

Complete exercises and use Slack to ask questions

Relevant open-source materials include:

[Evaluating distributional forecasts](#)

[Approximate leave-future-out cross-validation for Bayesian time series models](#)

[The Marginal Effects Zoo \(0.14.0\)](#)

This lecture's topics

Forecasting from dynamic models

Bayesian posterior predictive checks

Point-based forecast evaluation

Probabilistic forecast evaluation

Forecasting from dynamic models

Forecasting in mvgam

Two options

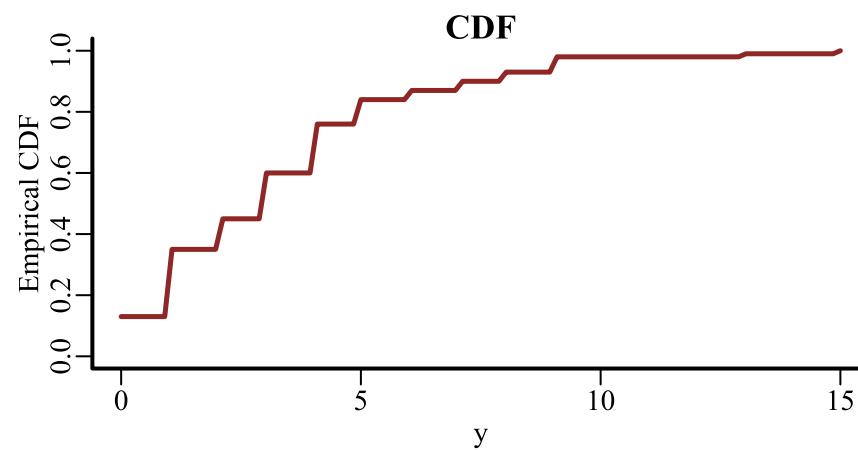
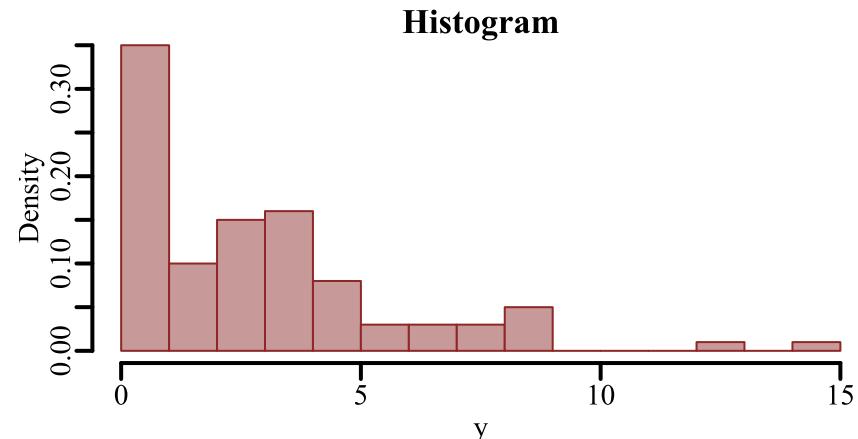
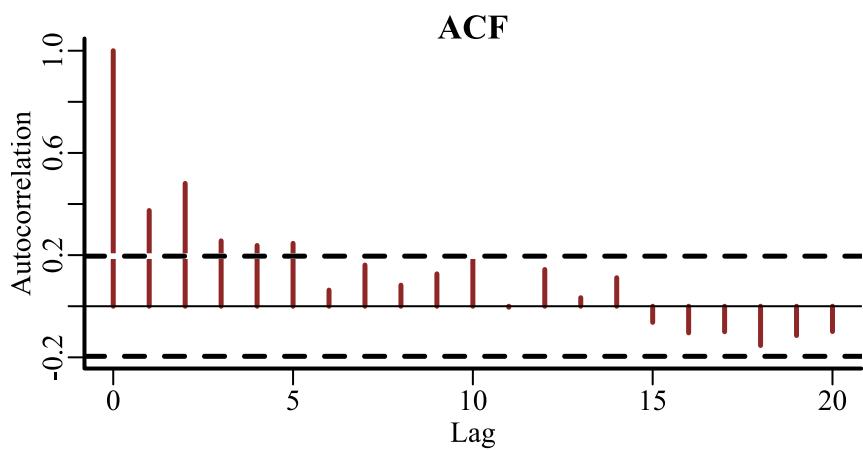
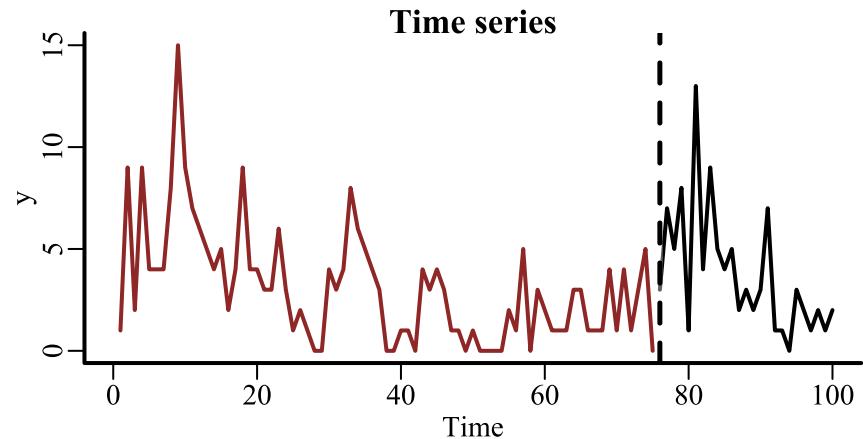
Feed `newdata` into the `mvgam` function for automatic probabilistic forecasts through `Stan`

Produce forecasts outside of `Stan` by feeding `newdata` and the fitted model into the `forecast.mvgam` function

Both require any out-of-sample covariates to be supplied

Both should give equivalent results

Simulated data



The model

```
library(mvgam)
model <- mvgam(y ~
  s(season, bs = 'cc', k = 8),
  data = data_train,
  newdata = data_test,
  trend_model = 'GP',
  family = poisson())
```

A cyclic smooth of `season` to capture repeated periodic variation

The model

```
library(mvgam)
model ← mvgam(y ~
  s(season, bs = 'cc', k = 8),
  data = data_train,
  newdata = data_test,
  trend_model = 'GP',
  family = poisson())
```

A Gaussian Process trend (approximated with Hilbert basis functions)

The model

```
library(mvgam)
model <- mvgam(y ~
  s(season, bs = 'cc', k = 8),
  data = data_train,
  newdata = data_test,
  trend_model = 'GP',
  family = poisson())
```

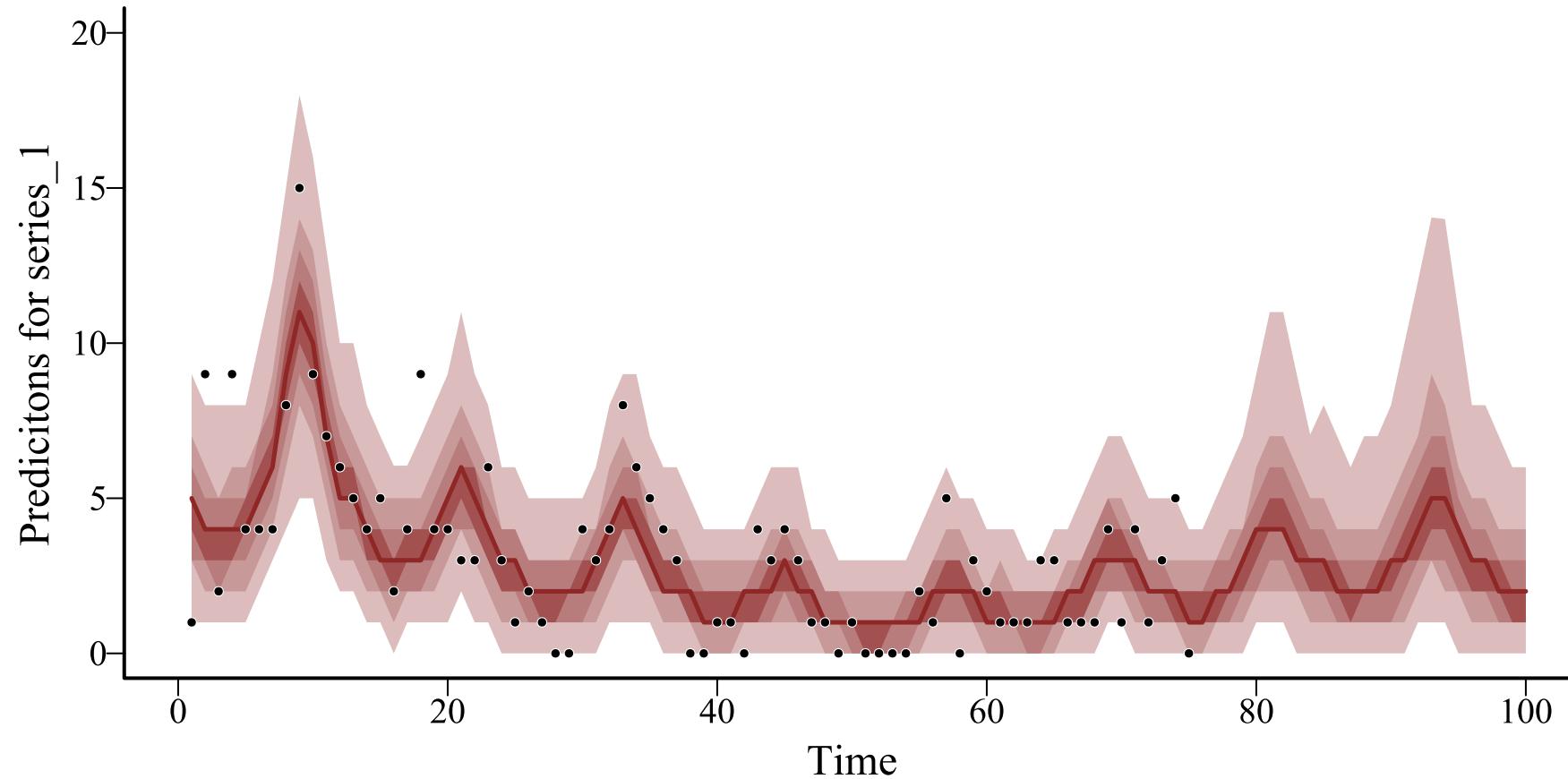
Forecasts will be computed automatically using the generated quantities block in Stan

Dropping newdata

```
model2 ← mvgam(y ~  
                 s(season, bs = 'cc', k = 8),  
                 data = data_train,  
                 trend_model = 'GP',  
                 family = poisson())
```

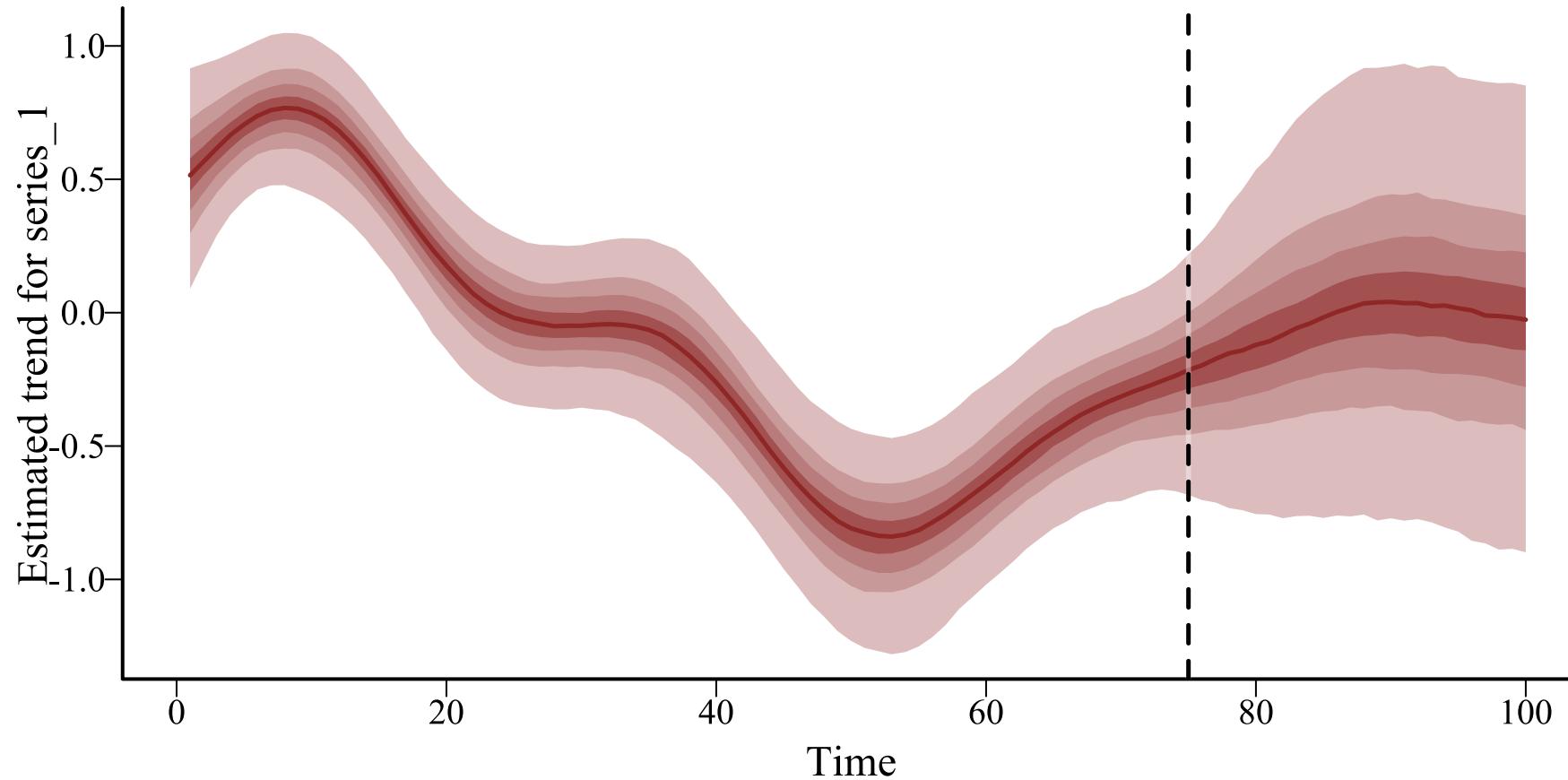
Predictions will only be calculated for the training data if no testing data (i.e. **newdata**) are supplied

```
plot(model, type = 'forecast')
```



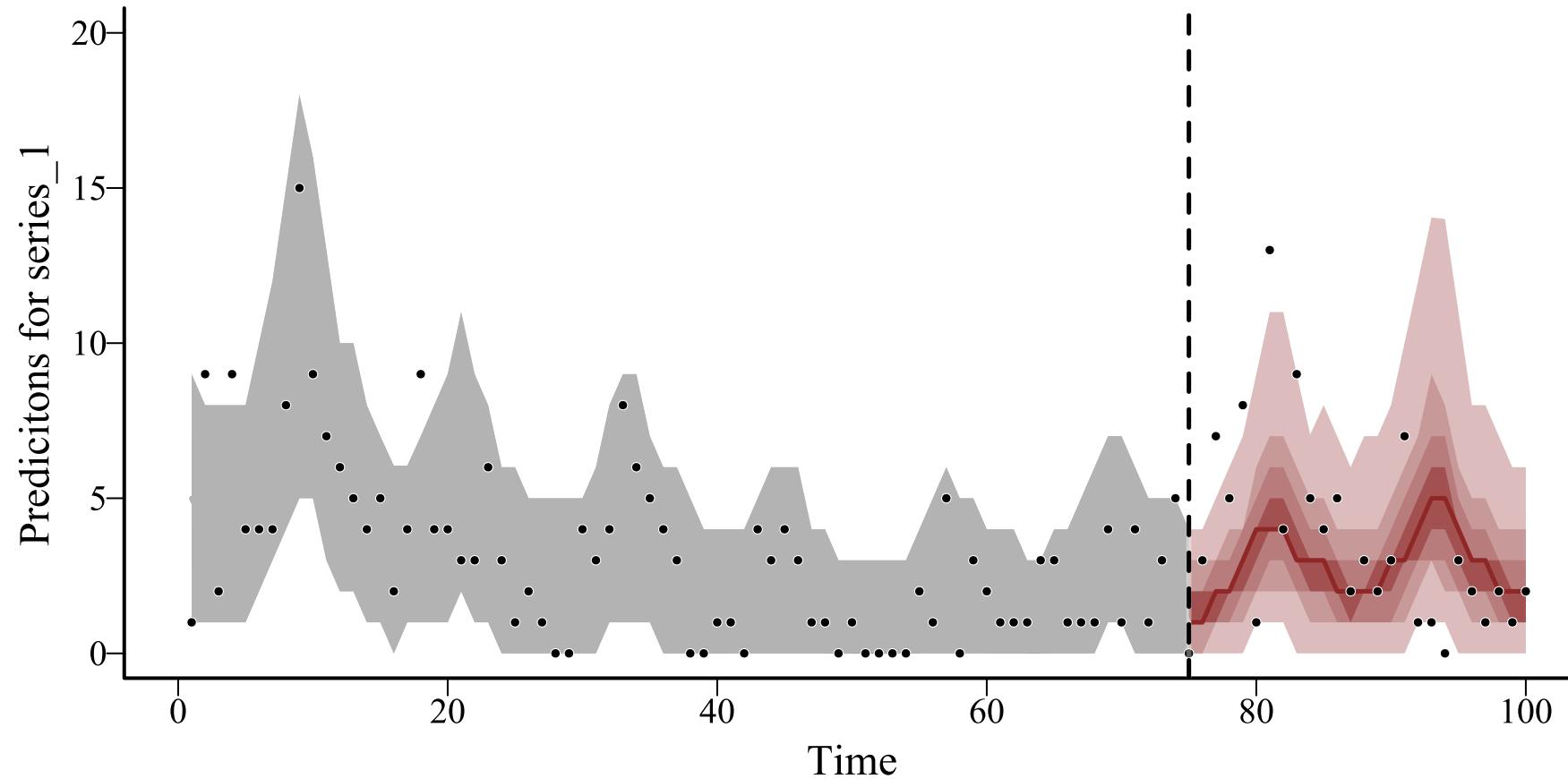
Automatic forecasts because `newdata` were supplied

```
plot(model, type = 'trend', newdata = data_test)
```



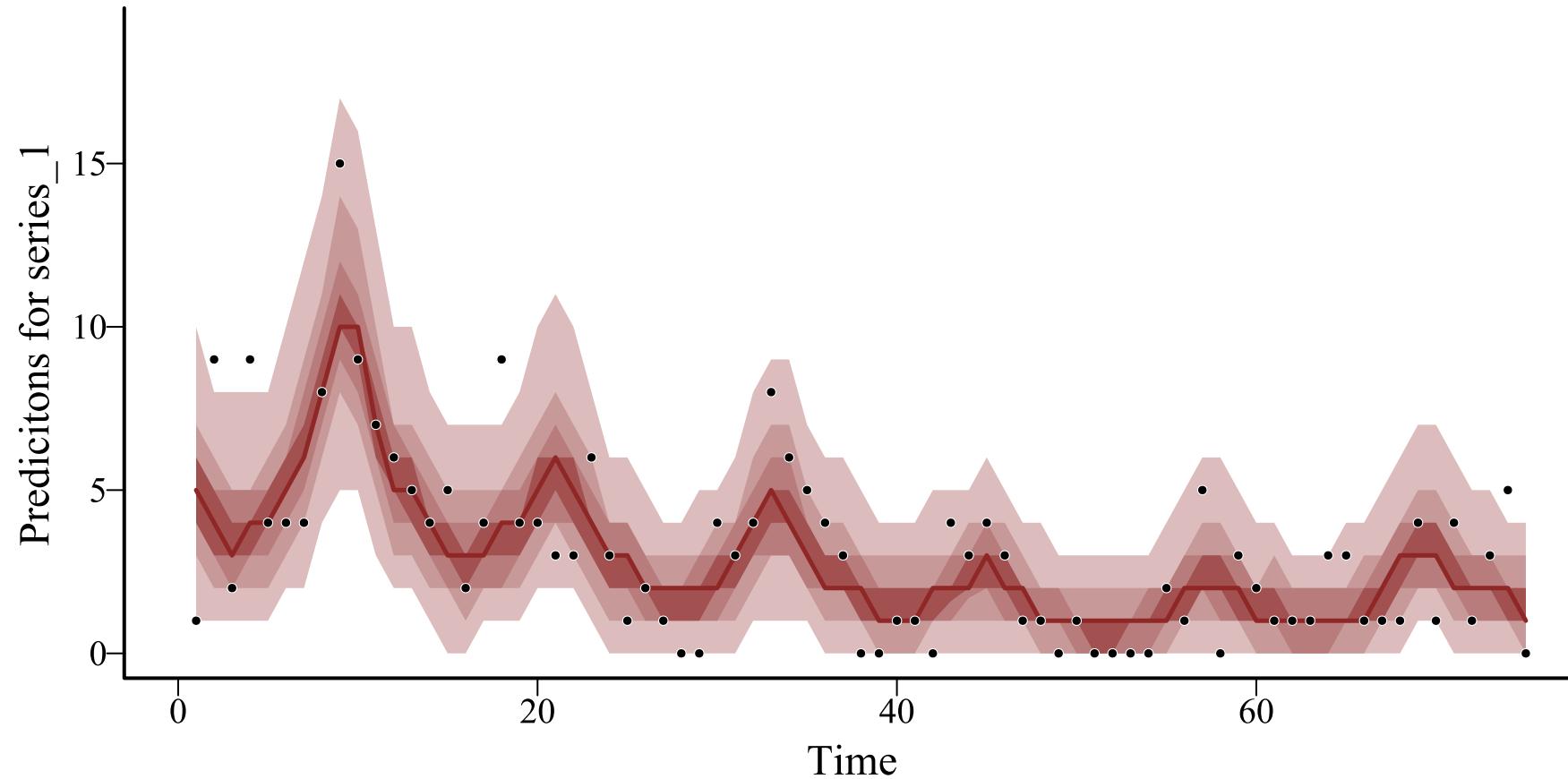
Trend extends into the future

```
plot(model, type = 'forecast', newdata = data_test)
```



Forecasts can be compared to truths quickly

```
plot(model2, type = 'forecast')
```



No forecasts in this case. Now what?

Posterior draws

dynamic `mvgam` models contain draws for many quantities

- β coefficients for linear predictor terms (called `b`)

- Any family-specific shape / scale parameters (i.e. ϕ for Negative Binomial; σ_{obs} for Normal / LogNormal etc...)

- Any trend-specific parameters (i.e. α and ρ for GP trends; σ and $ar1$ for AR trends etc...)

- In-sample posterior predictions (called `ypred`)

- In-sample posterior trend estimates (called `trend`)

All stored as MCMC draws in an object of class `stanfit` in the `model_output` slot

The stanfit object

```
summary(model$model_output)
```

```
##  Length Class Mode
## 1 stanfit S4
```

```
model$model_output@model_pars
```

```
## [1] "rho"          "b"           "ypred"        "mus"          "trend"        "alpha_gp"
"rho_gp"
## [8] "b_gp"
```

```
model$model_output@sim$chains
```

```
## [1] 4
```

```
model$model_output@sim$iter
```

```
## [1] 1000
```

Visualising chains

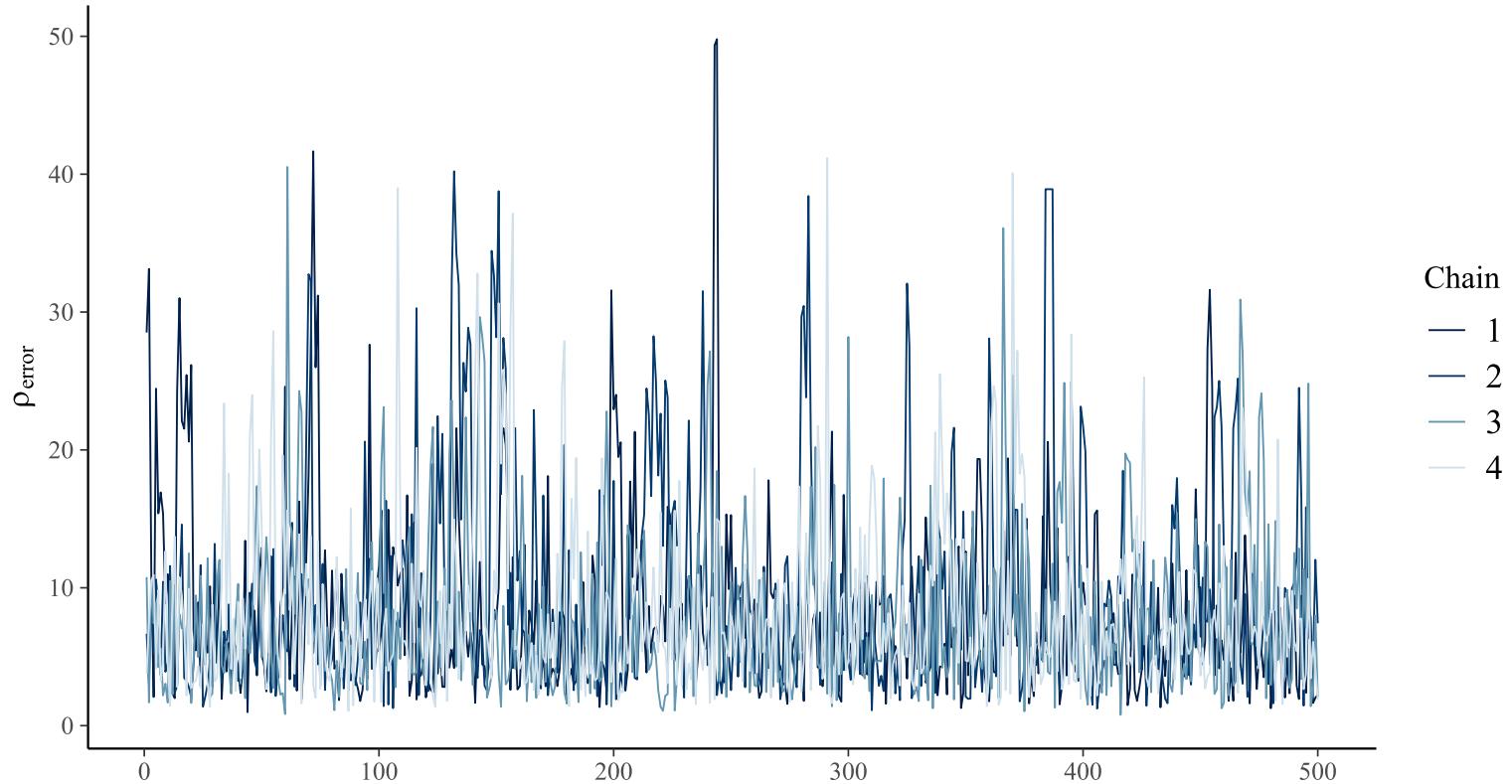
Code Plot

```
# use bayesplot to visualise posterior draws
library(ggplot2)
library(bayesplot)
mcmc_trace(model$model_output,
            # View draws of the GP length scale
            pars = 'rho_gp[1]' +
            # use ggplot2 to add informative labels
            ylab(expression(rho[error])))
```

Visualising chains

Code

Plot



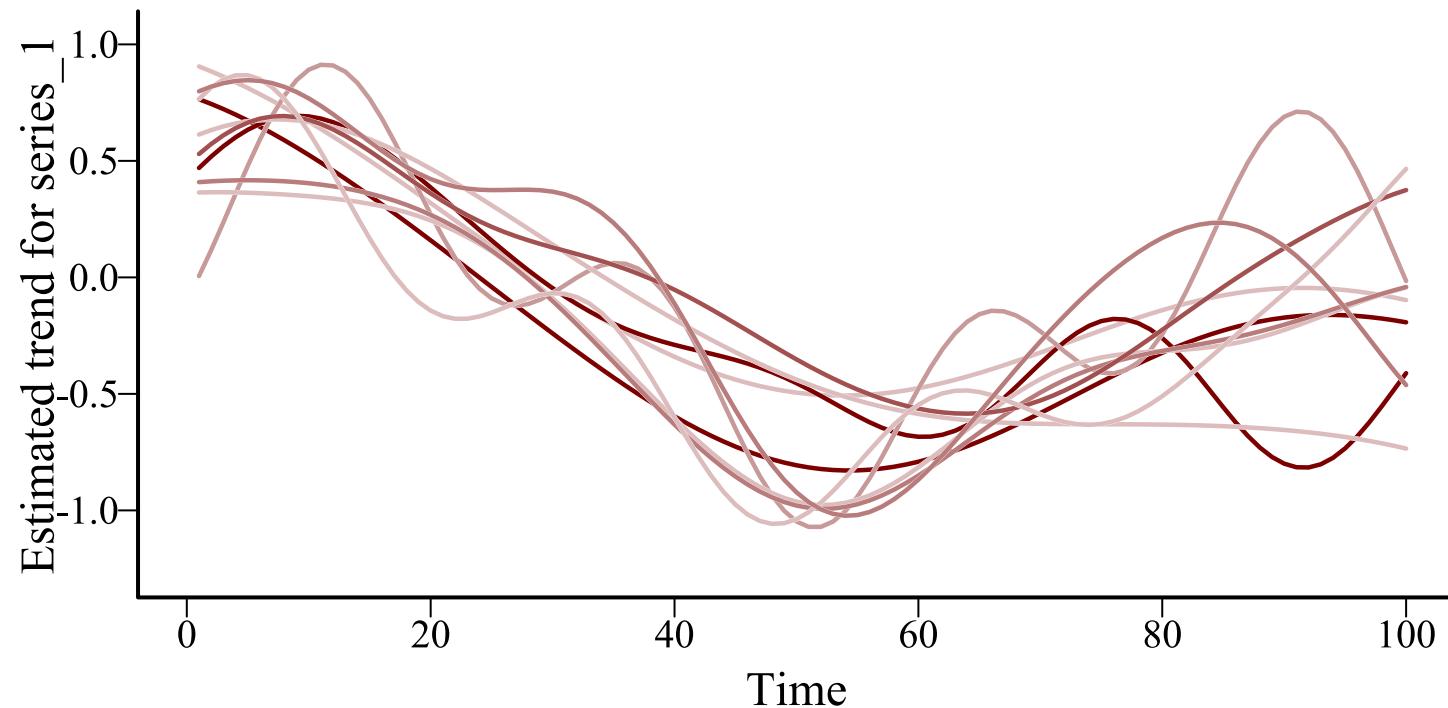
Draws of trend

Code Plot

```
# view posterior draws of the trend
plot(model, type = 'trend', realisations = TRUE,
     n_realisations = 10)
```

Draws of trend

Code Plot



But how can we extrapolate these to the future?

Ready for some multivariate statistical wizardry?

Ready





```
sim_gp = function(trend_draw, h, rho, alpha){
  # extract training and testing times
  t ← 1:length(trend_draw); t_new ← 1:(length(trend_draw) + h)
  # calculate training covariance
  Sigma ← alpha^2 * exp(-0.5 * ((outer(t, t, "-") / rho) ^ 2)) +
    diag(1e-9, length(t))
  # calculate training vs testing cross-covariance
  Sigma_new ← alpha^2 * exp(-0.5 * ((outer(t, t_new, "-") / rho) ^ 2))
  # calculate testing covariance
  Sigma_star ← alpha^2 * exp(-0.5 * ((outer(t_new, t_new, "-") / rho) ^ 2))
  +
    diag(1e-9, length(t_new))
  # draw one function realization of the stochastic Gaussian Process
  t(Sigma_new) %*% solve(Sigma, trend_draw) +
    MASS::mvrnorm(1, mu = rep(0, length(t_new)),
                  Sigma = Sigma_star - t(Sigma_new) %*% solve(Sigma,
Sigma_new))
}
```

Wizardize one trend draw

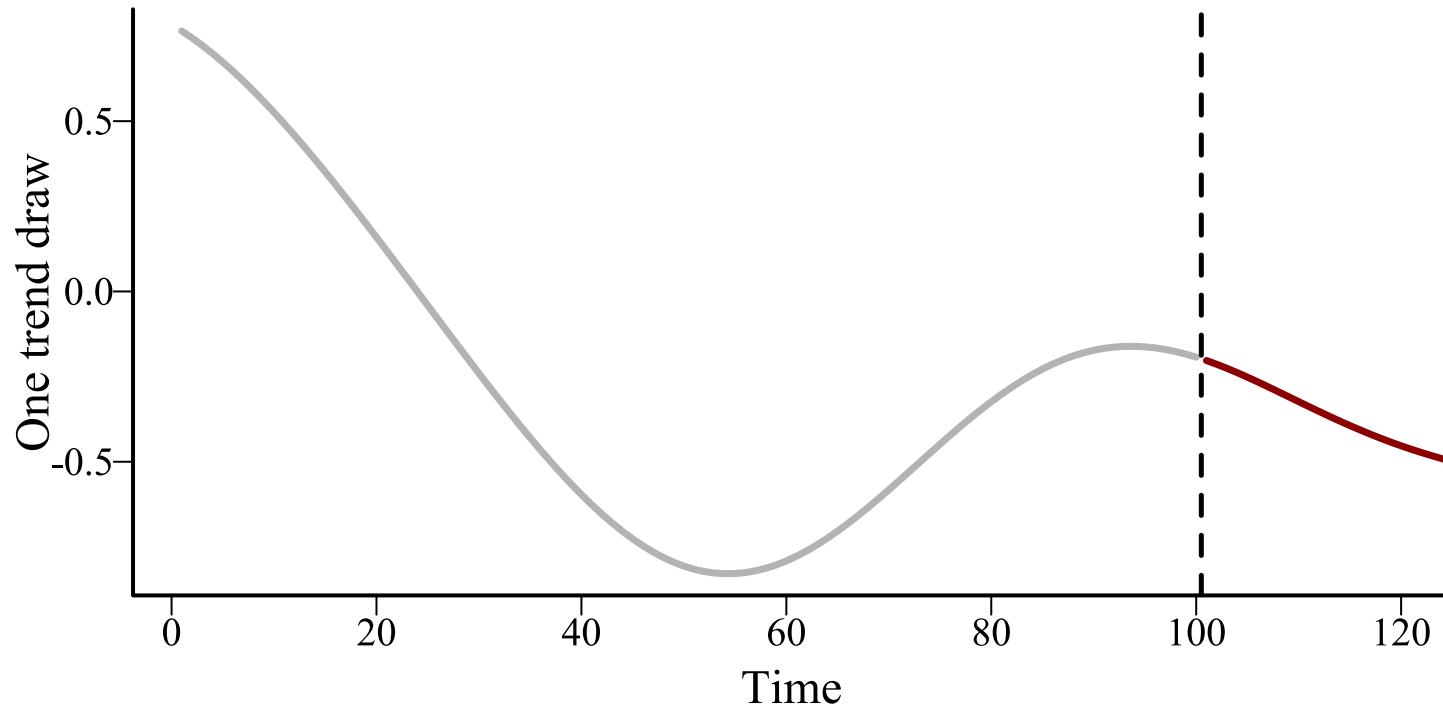
Wizardry Plot

```
# extract trend parameter draws and plot one draw
trend_draws ← as.matrix(model, variable = 'trend', regex = TRUE)
alpha_draws ← as.matrix(model, variable = 'alpha_gp', regex = TRUE)
rho_draws ← as.matrix(model, variable = 'rho_gp', regex = TRUE)
plot(1, type = 'n', bty = 'l',
      xlim = c(1, 130), ylim = range(trend_draws[1,]),
      ylab = 'One trend draw', xlab = 'Time')
lines(trend_draws[1,], col = 'gray70', lwd = 3.5)
# wizardize to extend draw forward 30 timesteps and plot
forecast_draw = sim_gp(trend_draw = trend_draws[1,], h = 30
                        alpha = alpha_draws[1,], rho = rho_draws[1,])
lines(x = 101:130, y = forecast_draw[101:130], lwd = 3.5, col = 'darkred')
abline(v = 100.5, lty = 'dashed', lwd = 2.5)
```

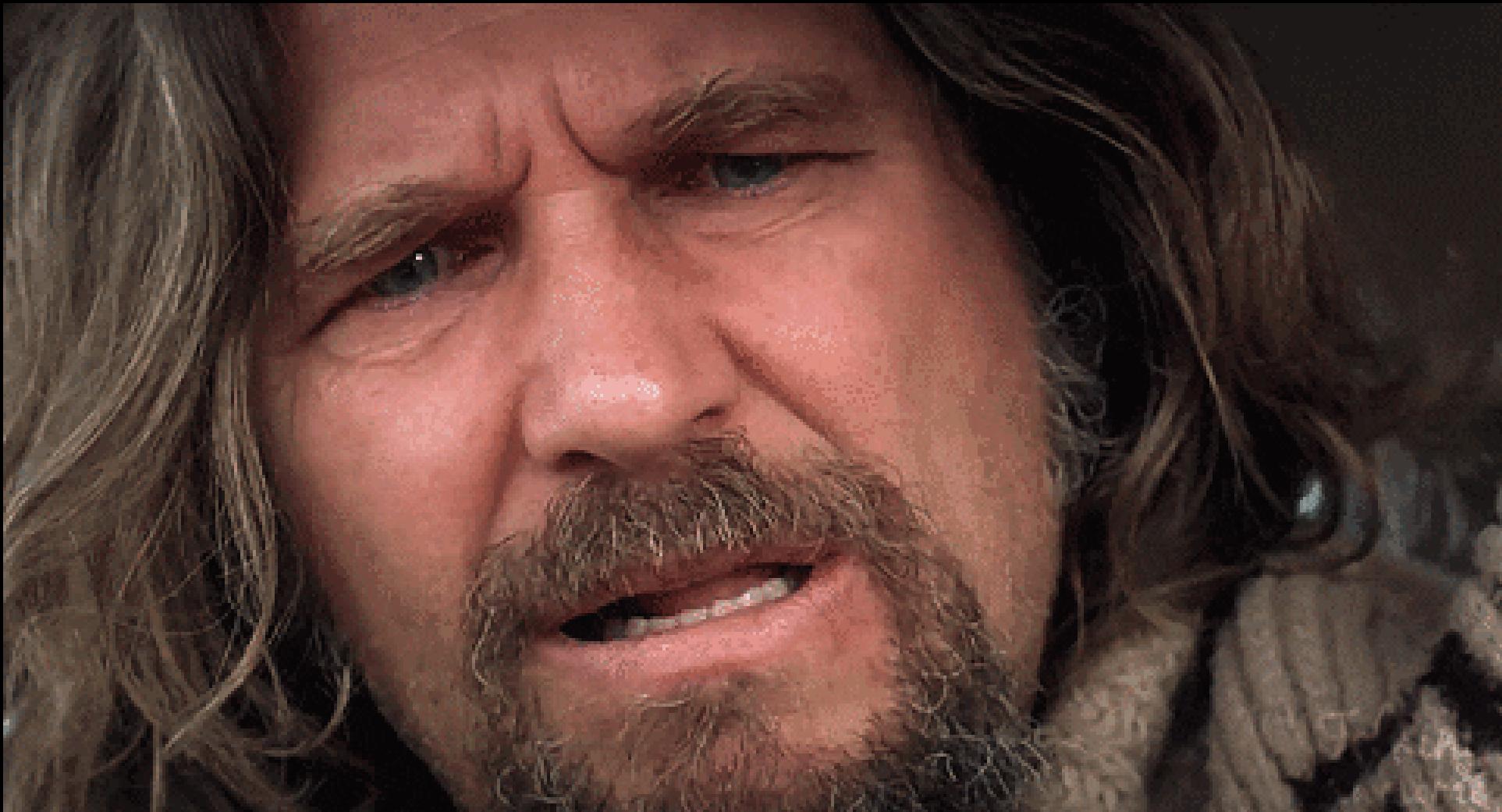
Wizardize one trend draw

Wizardry

Plot



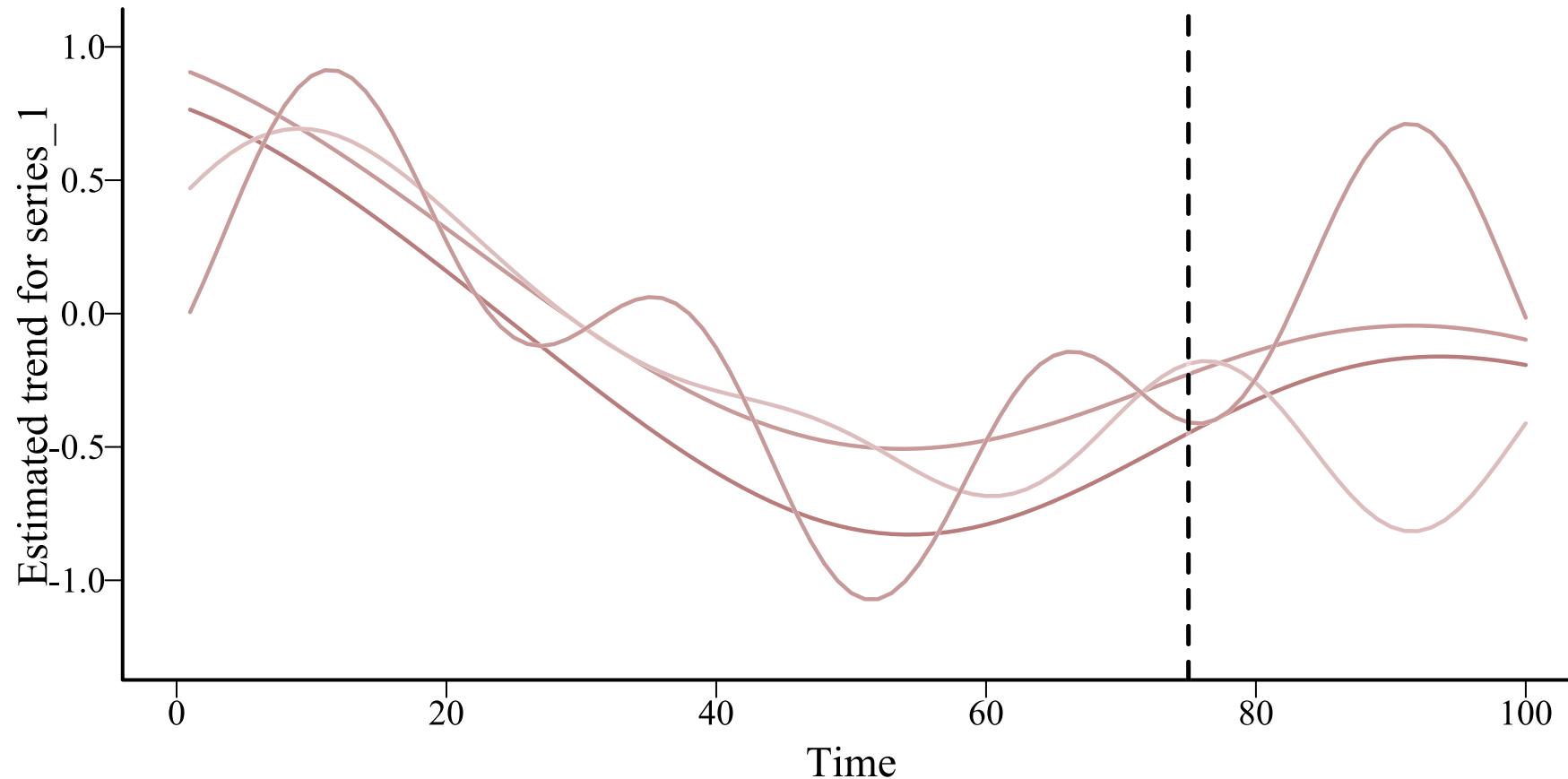
Piece of cake?



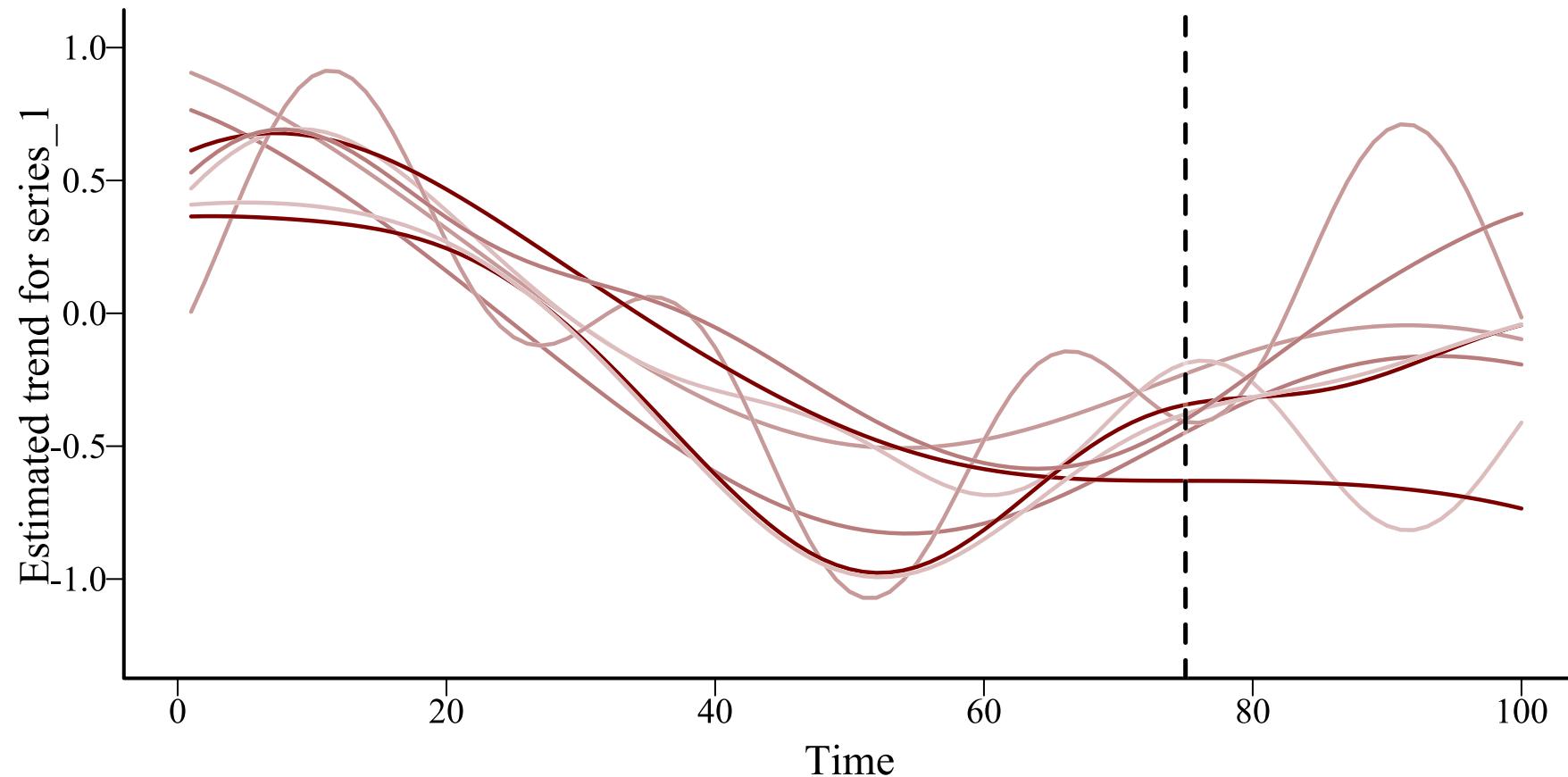
There is no wizardry 😊. Rather, each kind of trend (AR, GP etc...) has an underlying stochastic equation that can be used to extrapolate draws to the future

But doing this manually is slow and error-prone. `mvgam` does this *automatically* using `newdata`

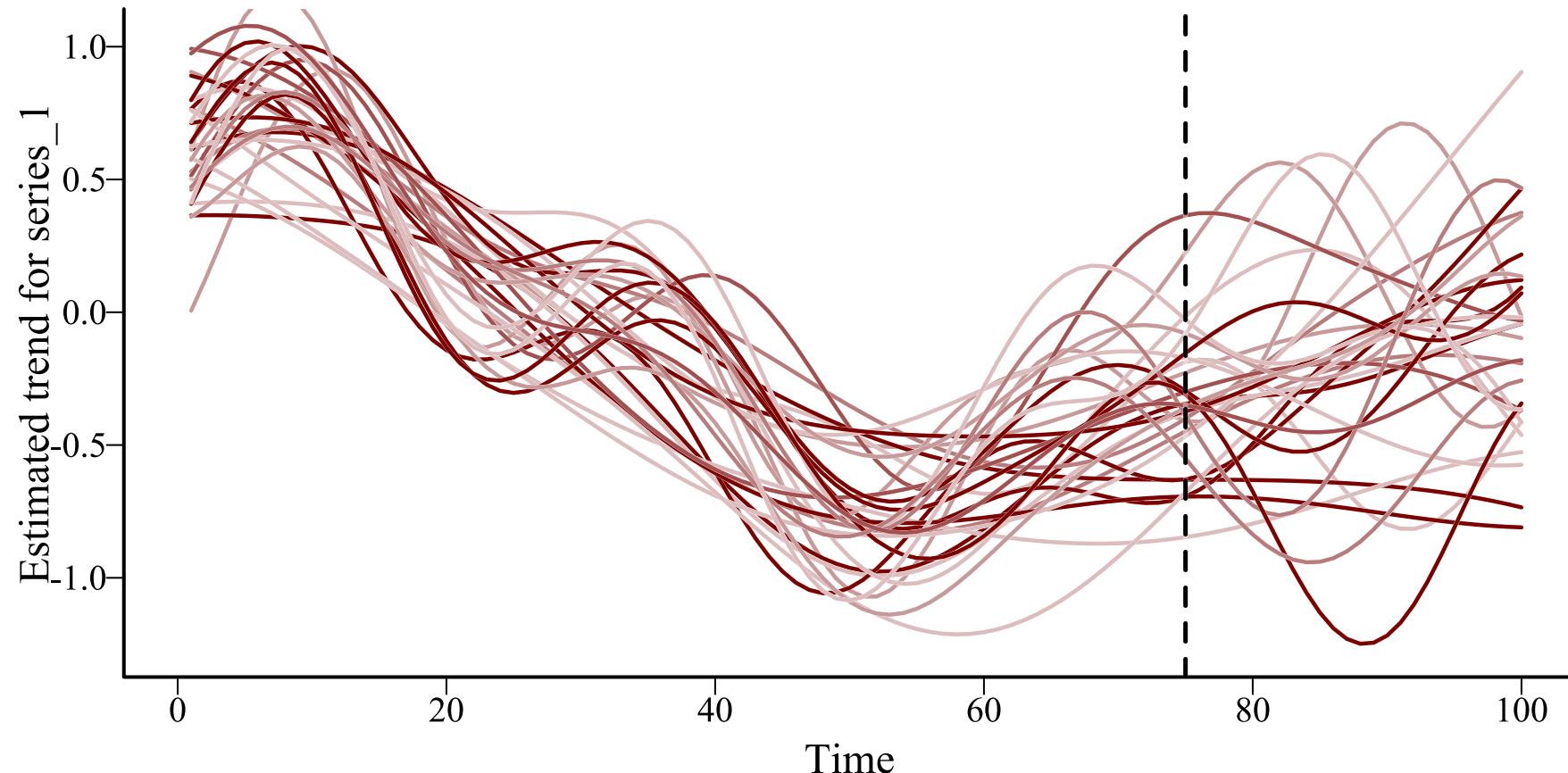
```
plot(model2, type = 'trend', newdata = data_test,  
realisations = TRUE, n_realisations = 4)
```



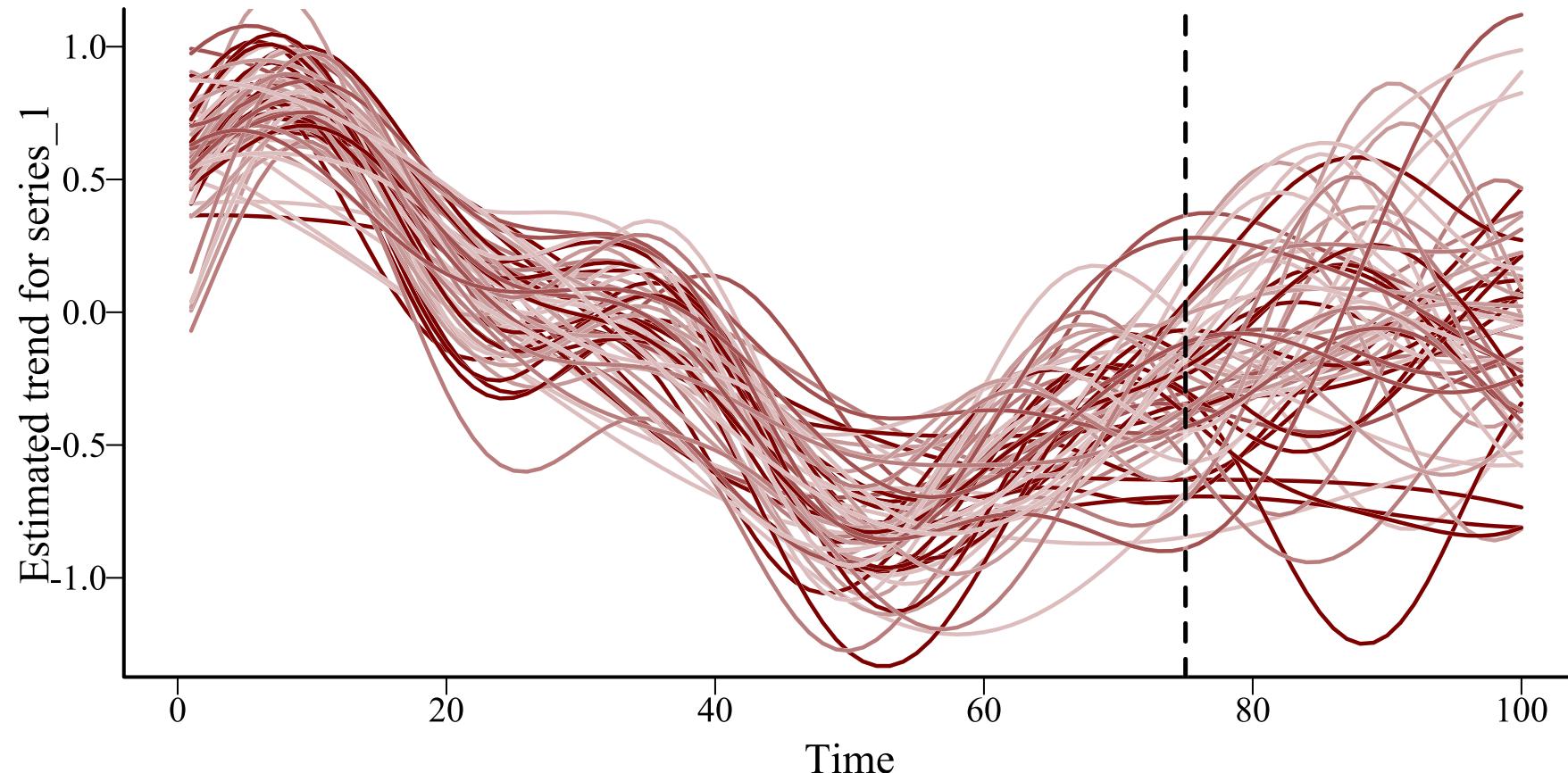
```
plot(model2, type = 'trend', newdata = data_test,  
realisations = TRUE, n_realisations = 8)
```



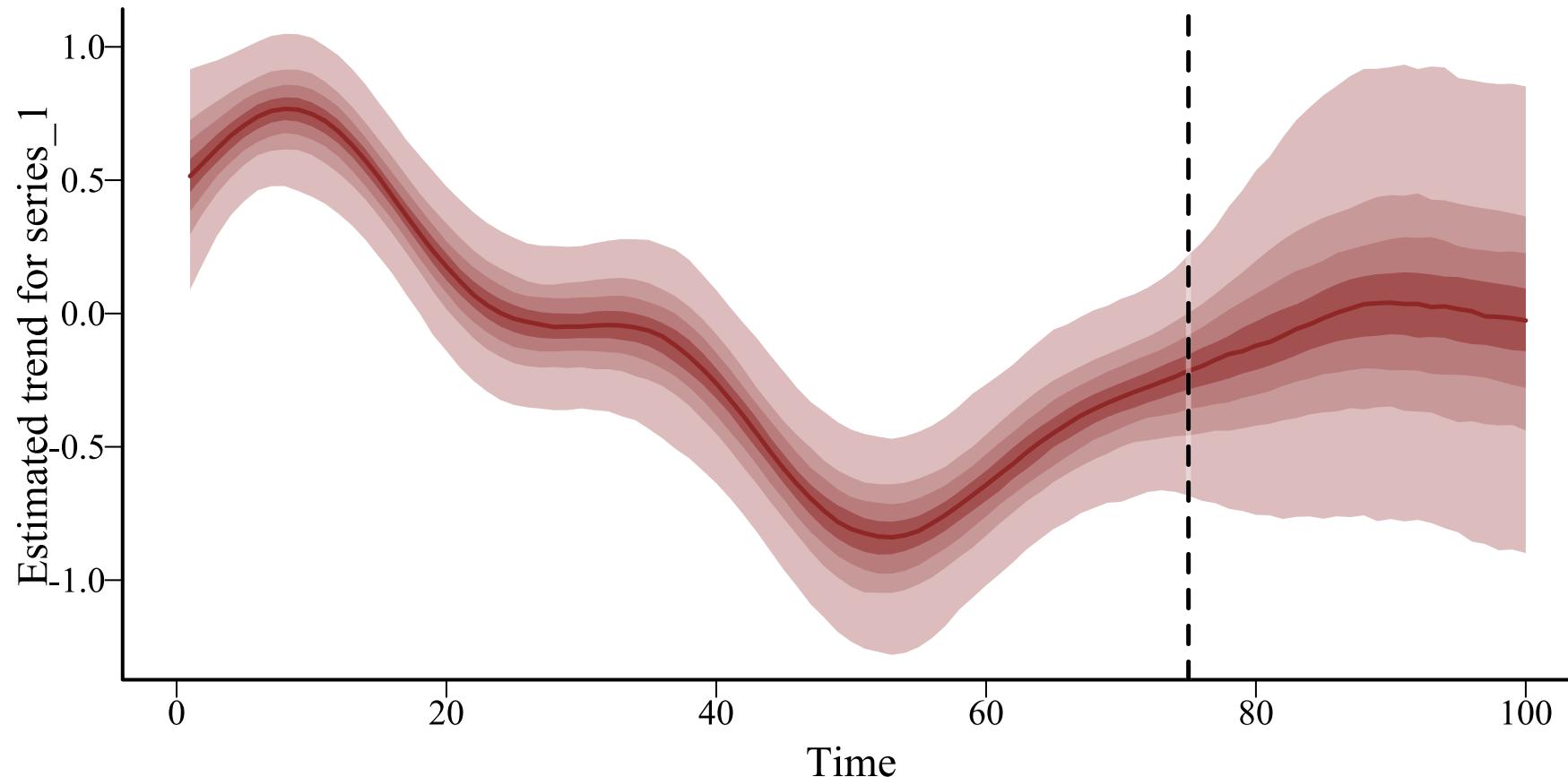
```
plot(model2, type = 'trend', newdata = data_test,  
realisations = TRUE, n_realisations = 30)
```



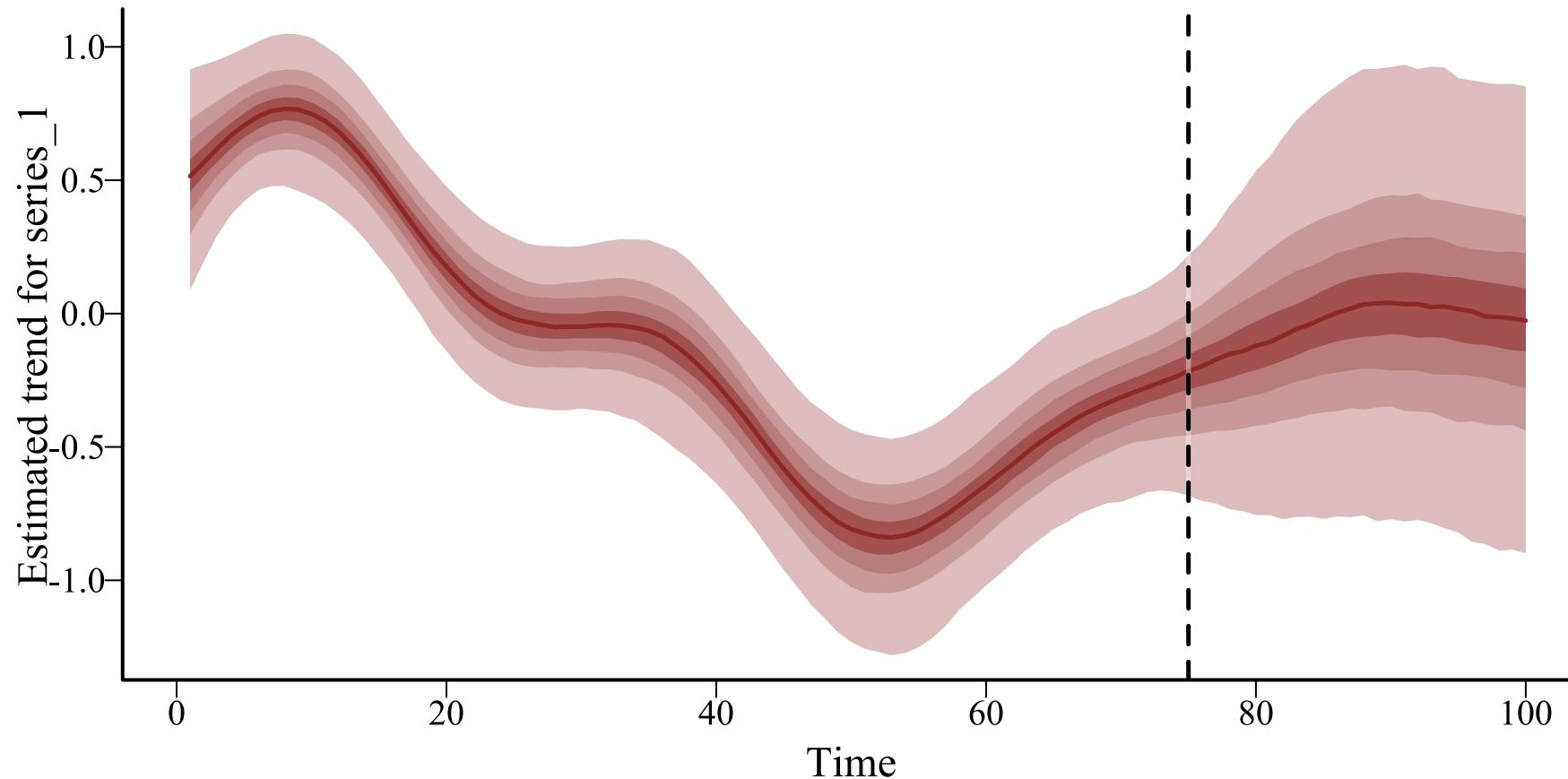
```
plot(model2, type = 'trend', newdata = data_test,  
realisations = TRUE, n_realisations = 60)
```



```
plot(model2, type = 'trend', newdata = data_test,  
realisations = FALSE)
```



```
Or: plot(forecast(model2, type = 'trend', newdata =  
data_test), realisations = FALSE)
```



Once dynamic trend is extrapolated, computing forecasts is easy

We only need any remaining "future" predictor values from covariates

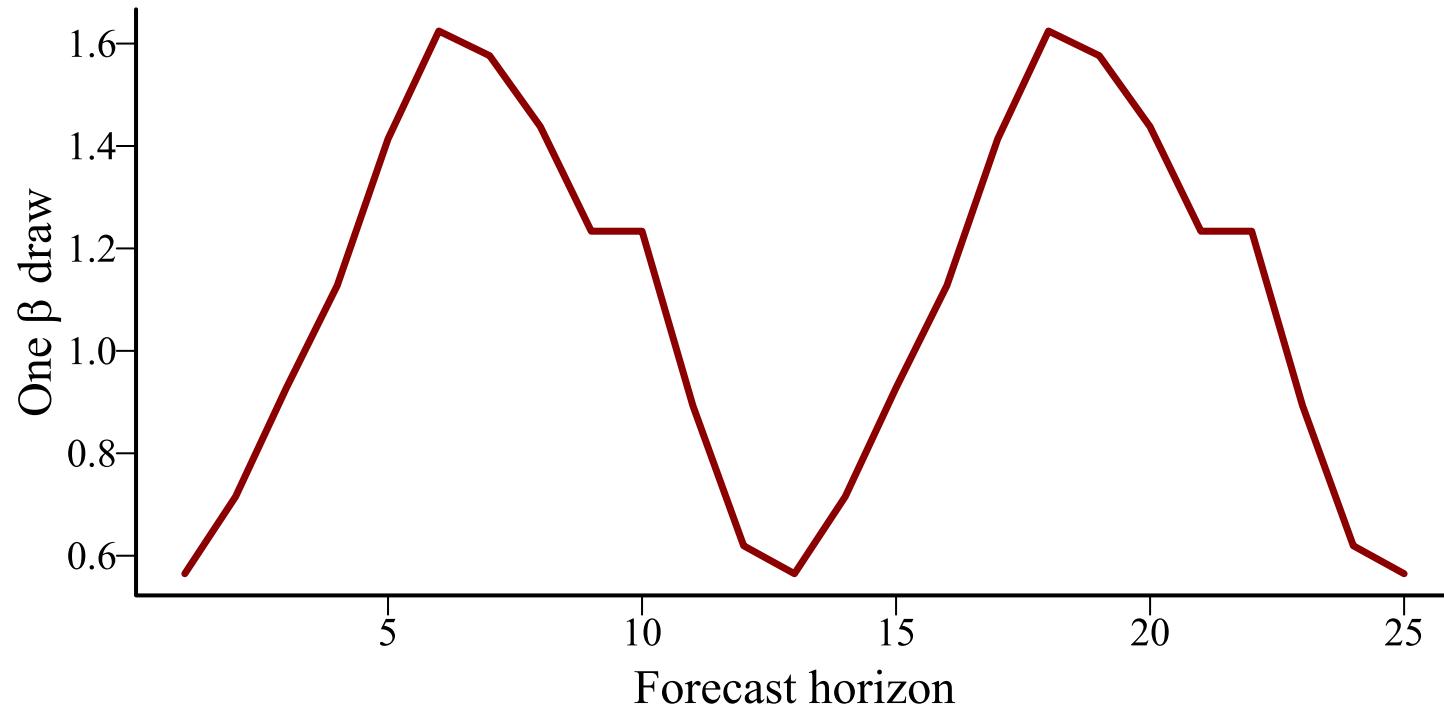
GAM covariate predictions

Code Plot

```
# extract beta regression coefficient draws
beta_draws <- as.matrix(model2, variable = 'betas')
# calculate the linear predictor matrix for the GAM component
lpmatrix <- mvgam:::obs_Xp_matrix(newdata = poisdat$data_test,
                                    mgcv_model = model$mgcv_model)
# calculate linear predictor (link-scale) predictions for one draw
linkpreds <- lpmatrix %*% beta_draws[1,] + attr(lpmatrix, 'model.offset')
# plot the linear predictor values
plot(1, type = 'n', bty = 'l',
      xlim = c(1, length(linkpreds)),
      ylim = range(linkpreds),
      ylab = expression(One~beta~draw), xlab = 'Forecast horizon')
lines(linkpreds, col = 'darkred', lwd = 3.5)
```

GAM covariate predictions

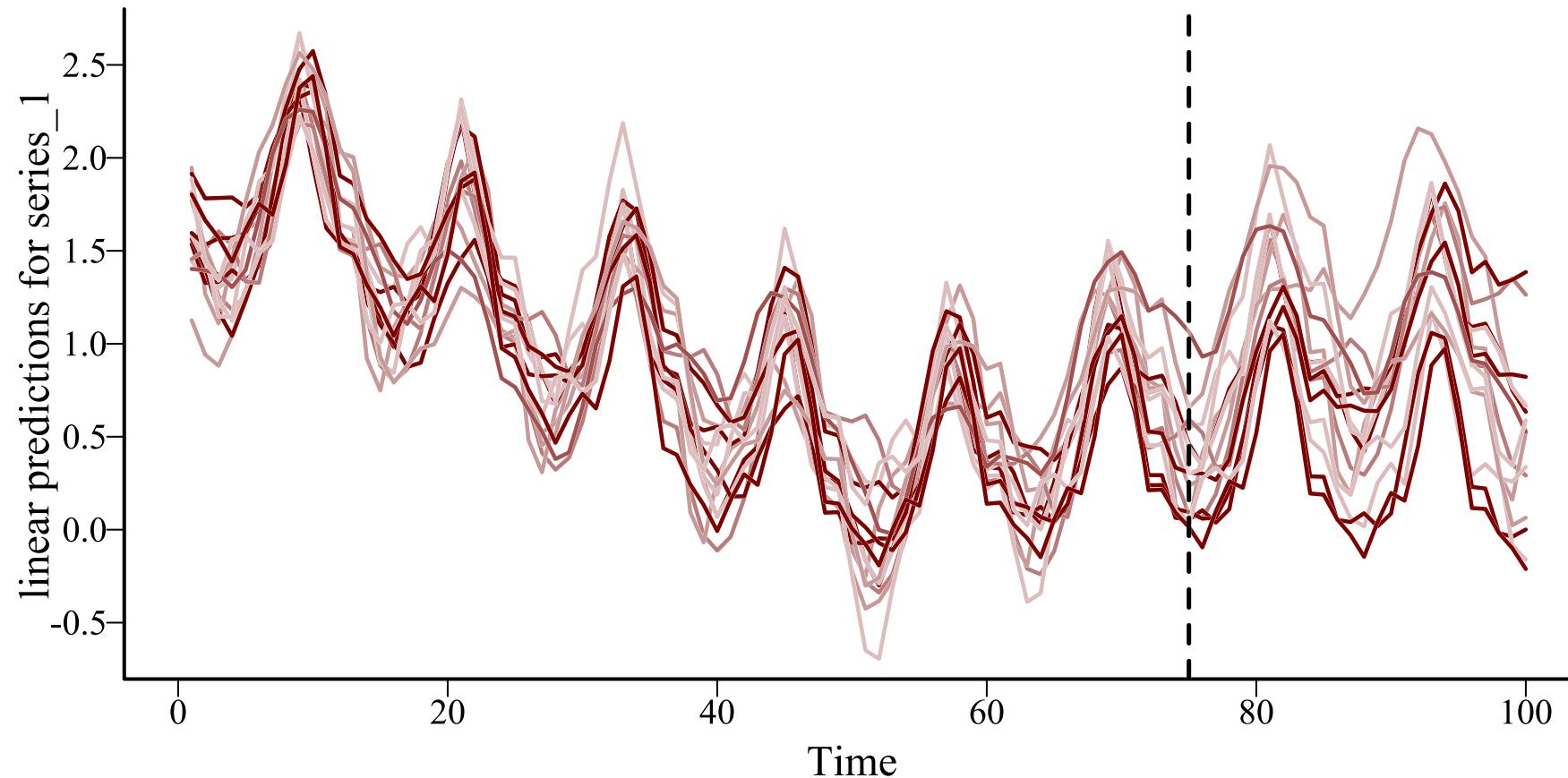
Code Plot



These covariate predictions are added to the trend predictions to give the full predictions *on the link scale*

Once again, `mvgam` does this *automatically* using the `forecast` function

```
plot(forecast(model2, type = 'link', newdata = data_test),  
realisations = TRUE)
```



Forecasting is easier if `newdata` are fed to `mvgam()`, but this results in a larger model object and requires test data be available now

When testing data not available, you can generate forecasts for new data later using `forecast.mvgam` (note, `time` values in `newdata` must follow immediately from `time` values in original training data)

But there are multiple *types* of predictions available. What are they?

Types of mvgam predictions

type = 'expected' gives the draws from the expected value of the posterior predictive distribution, or the average of each draw from type='response'

In logistic regression, this is π on the probability scale (or inverse logit)

$$\logit(\pi_i) = \alpha + \beta x_i$$
$$y_i \sim \text{Binomial}(1, \pi_i)$$
$$E(y_i)$$

type = 'response' gives the draws from a random binomial distribution with draws from the posterior distribution of π

type = 'link' gives the posterior draws of π on the logit or log odds scale

```
predict(object, type = 'link')
```

Gives the real-valued, unconstrained linear predictor

Takes into account uncertainty in GAM regression coefficients

Can include uncertainty in any dynamic trend components

Can be extracted from the fitted model as parameter `mus`

```
range(predict(model, type = 'link', process_error = FALSE))
```

```
## [1] -0.02145857 2.22647337
```

```
range(predict(model, type = 'link', process_error = TRUE))
```

```
## [1] -3.801796 5.165899
```

```
range(as.matrix(model, variable = 'mus', regex = TRUE))
```

```
## [1] -1.62582 4.99195
```

Hang on. Why do these differ?

```
range(predict(model, type = 'link', process_error = TRUE))
```

```
## [1] -3.801796 5.165899
```

```
range(as.matrix(model, variable = 'mus', regex = TRUE))
```

```
## [1] -1.62582 4.99195
```

`predict` assumes the dynamic process has reached stationarity to tell us what we might expect if we see these same covariate values *sometime in the future*

`mus` includes estimates for where the trend was *at each point in the training data* (hindcasts), so it is less uncertain

link predictions

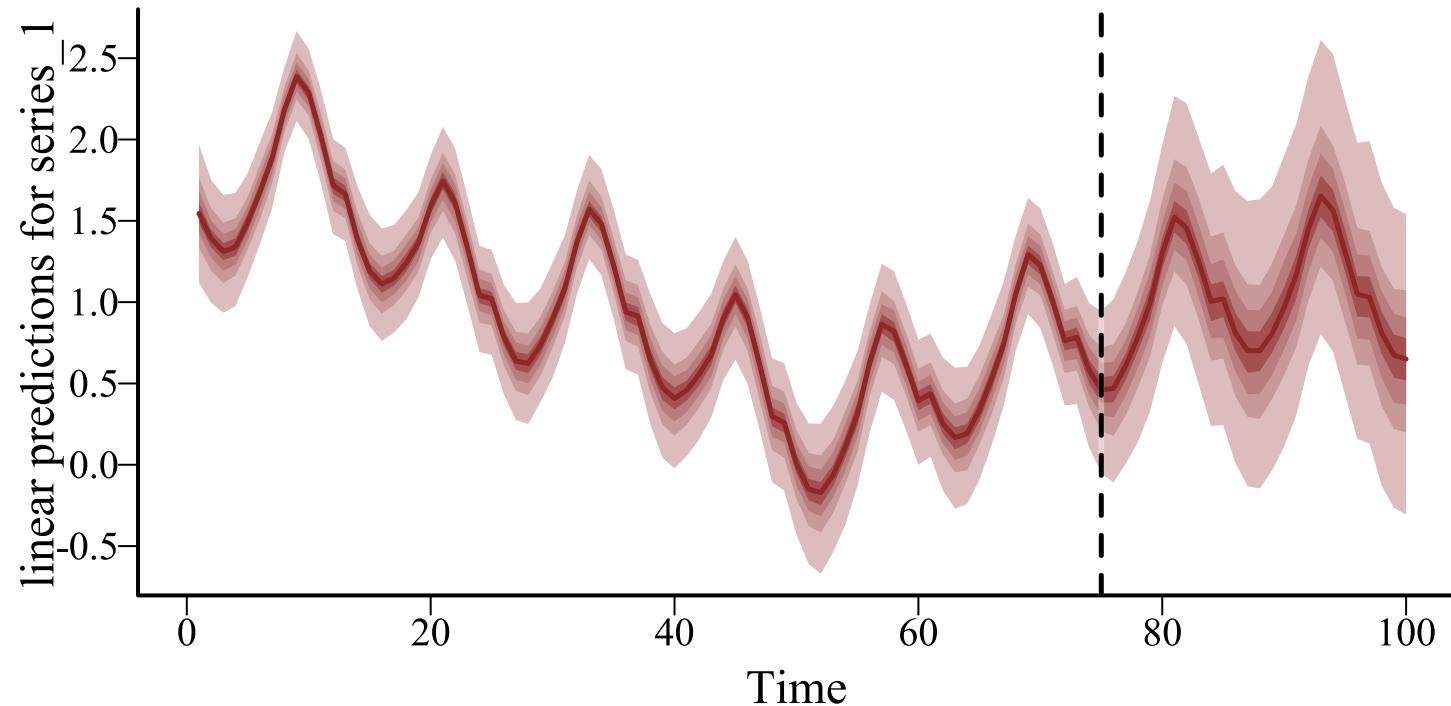
[Code](#) [Plot](#)

```
# extract link-scale forecasts from the model
fc <- forecast(model, type = 'link')

# plot using the available S3 plotting function
plot(fc)
```

link predictions

[Code](#) [Plot](#)



```
predict(object, type = 'expected')
```

Gives the **average** prediction on the observation (response) scale

Useful as we often want to get a sense of long-term averages for
guiding scenario analyses

Usually it is just the inverse link function applied to a prediction
from `type = link`

But not always!

This is probably the most confusing type of prediction

Normal distribution

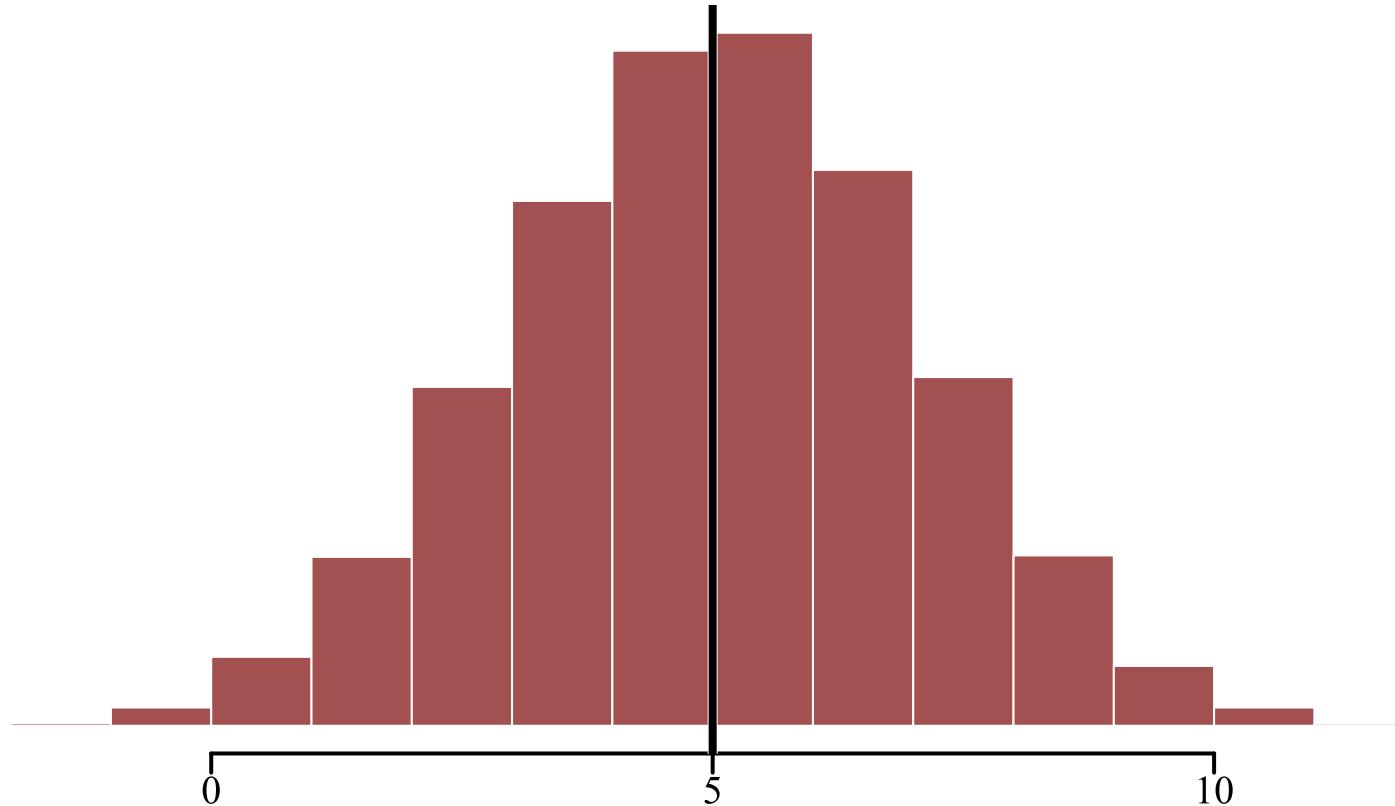
$$Y \sim \text{Normal}(\mu, \sigma)$$

$$\textit{link} = \mu = \alpha + \mathbf{X}\beta + z$$

$$\textit{expected} = \mathbb{E}(Y|\mu, \sigma) = \mu$$

Normal distribution

$Y \sim \text{Normal}(5, 2); \text{Mean}(Y) = 5$



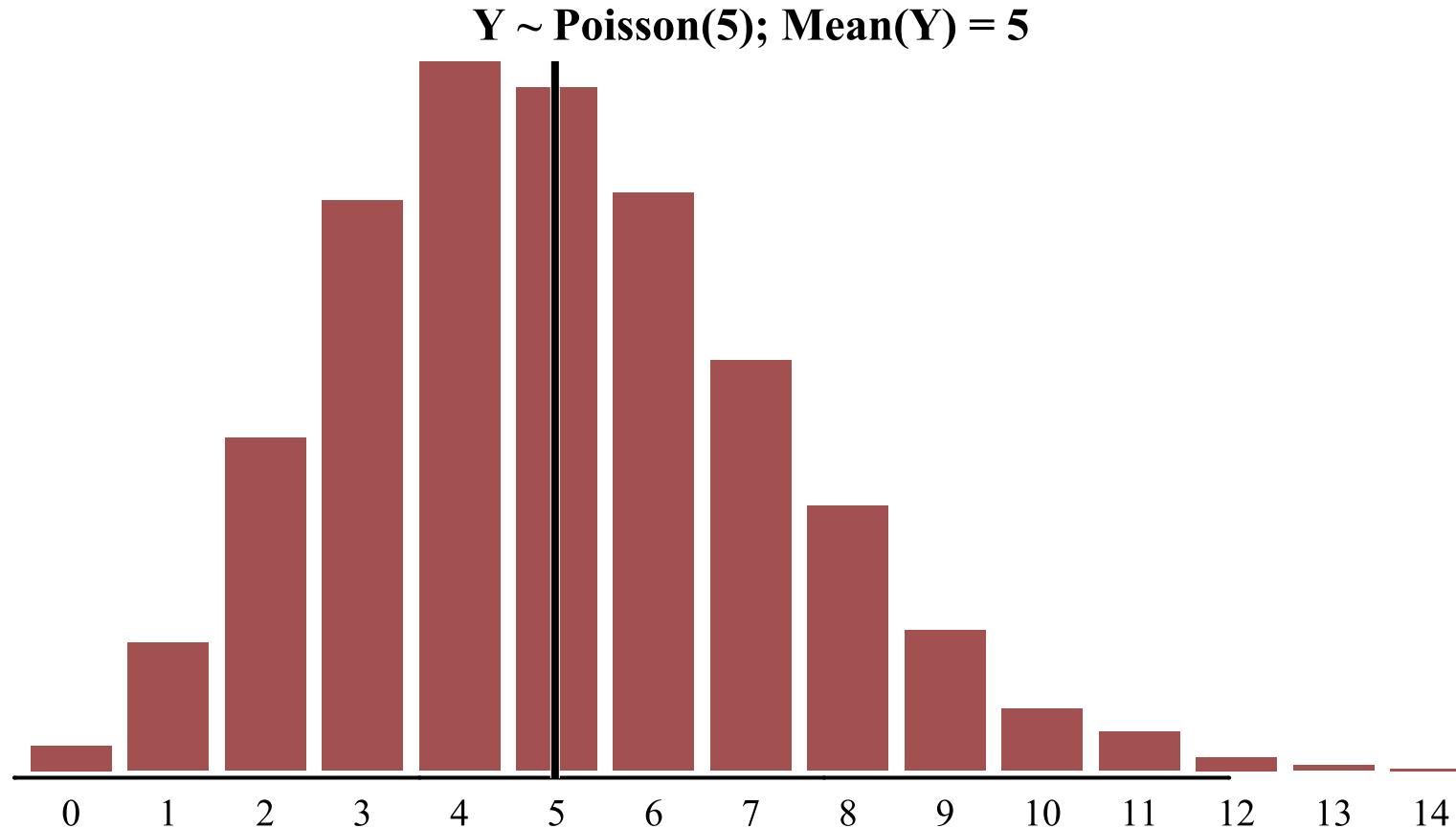
Poisson distribution

$$Y \sim \text{Poisson}(\lambda)$$

$$\textit{link} = \log(\lambda) = \alpha + \mathbf{X}\beta + z$$

$$\textit{expected} = \mathbb{E}(Y|\lambda) = \exp(\lambda)$$

Poisson distribution



LogNormal distribution

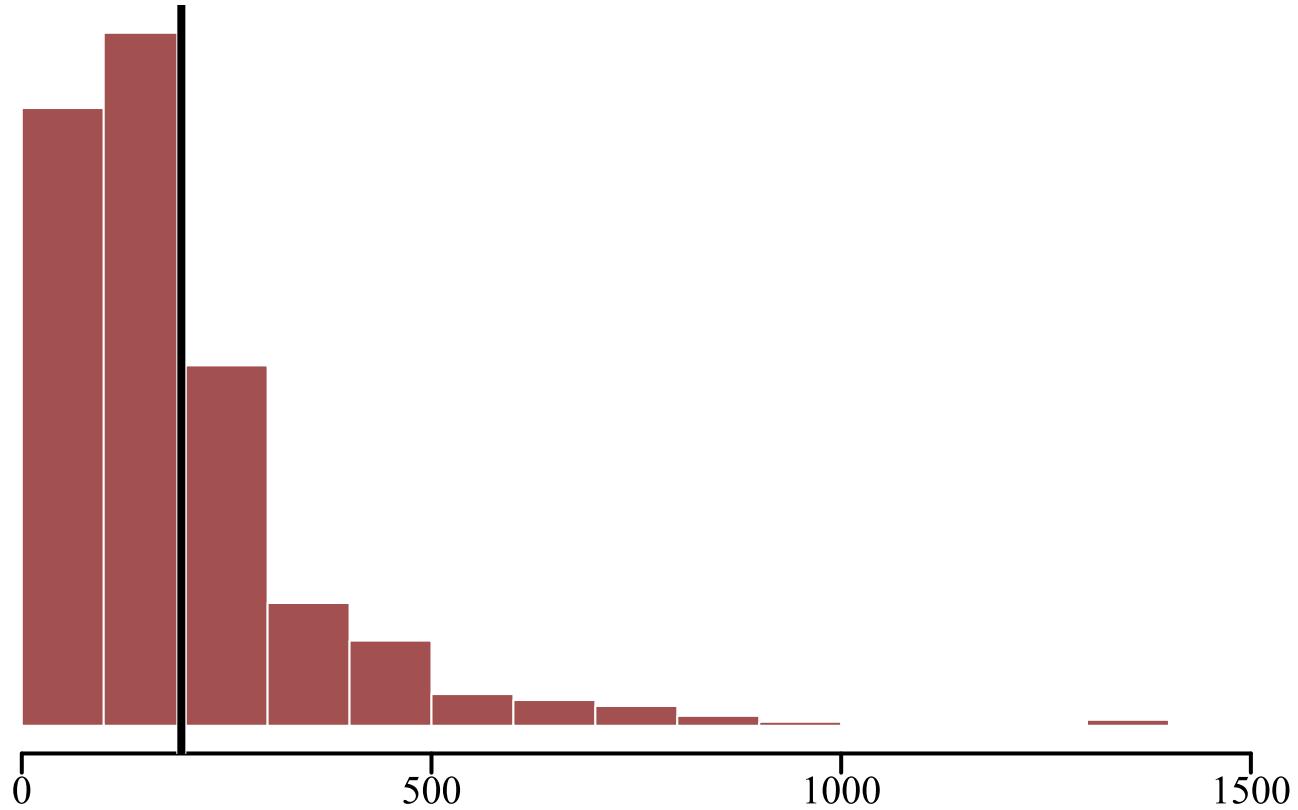
$$Y \sim \text{LogNormal}(\mu, \sigma)$$

$$\text{link} = \mu = \alpha + \mathbf{X}\beta + z$$

$$\text{expected} = \mathbb{E}(Y|\mu, \sigma) = \exp\left(\mu + \frac{\sigma^2}{2}\right)$$

LogNormal distribution

$Y \sim \text{LogNormal}(5, 0.75)$; $\text{Mean}(Y) = 195$



expected predictions

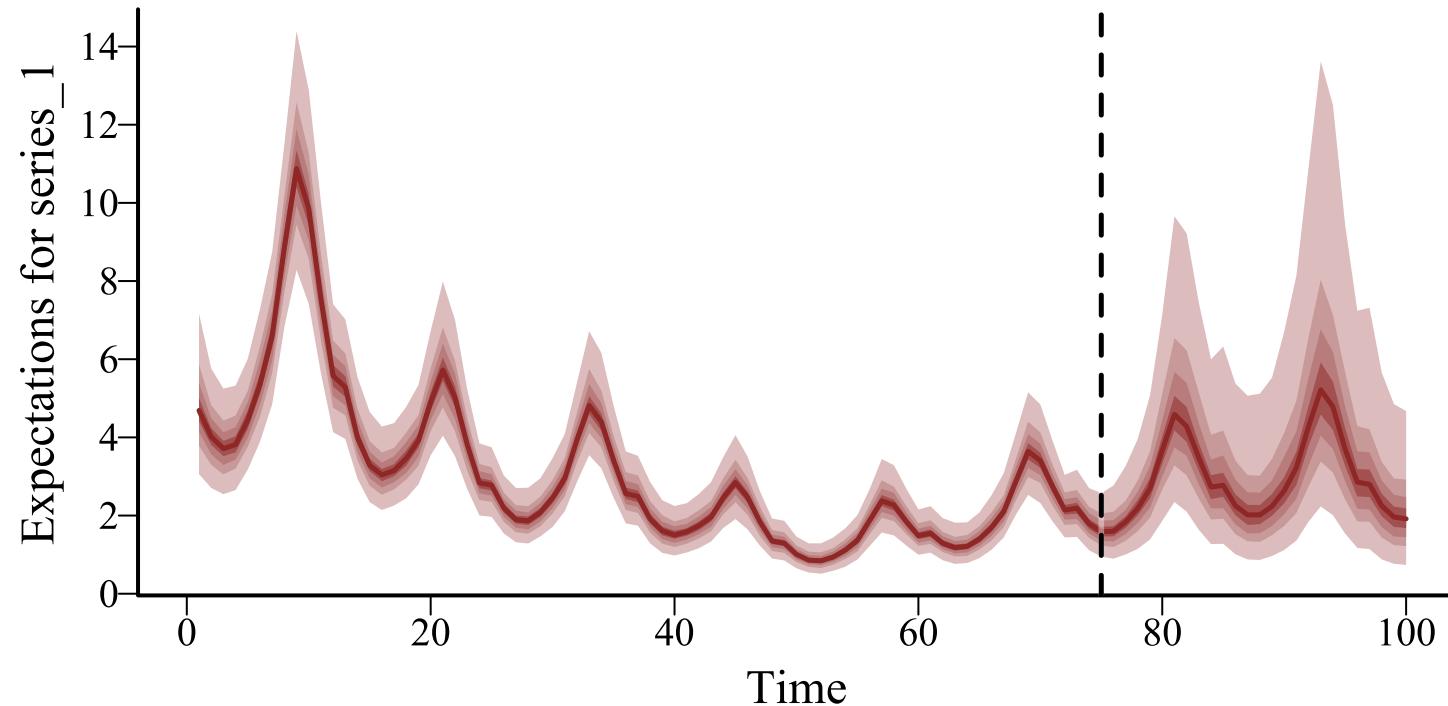
[Code](#) [Plot](#)

```
# extract expectation-scale forecasts from the model
fc <- forecast(model, type = 'expected')

# plot using the available S3 plotting function
plot(fc)
```

expected predictions

Code Plot



```
predict(object, type = 'response')
```

Gives the predictions on the observation (response) scale

Includes uncertainty in the linear predictor **and** any uncertainty arising from the observation process

Some distributions only depend on the inverse link of the linear predictor (i.e. $Poisson(\lambda)$ or $Bernoulli(\pi)$)

Others depend on additional shape / scale parameters (i.e. $Normal(\mu, \sigma)$ or $StudentT(\nu, \mu, \sigma)$)

These are the most often used type of predictions for evaluating forecasts

response predictions

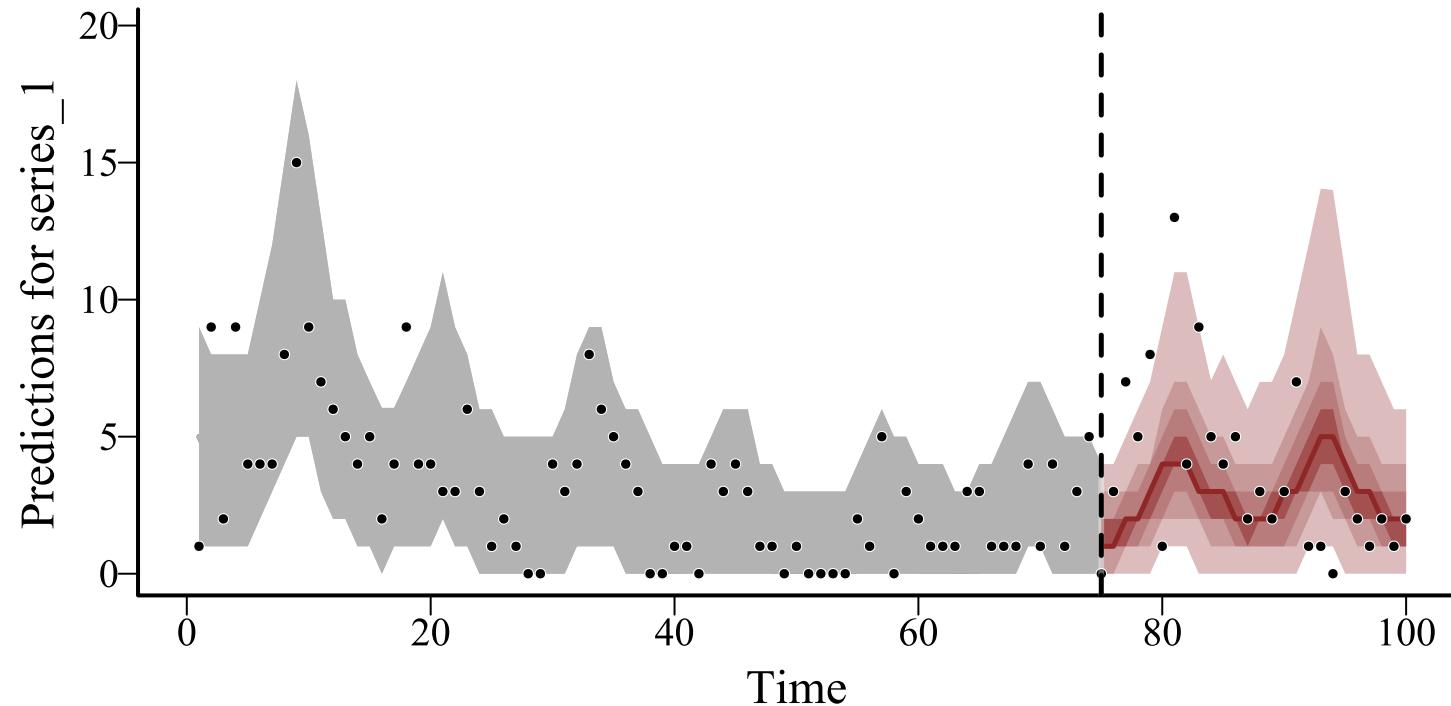
[Code](#) [Plot](#)

```
# extract response-scale forecasts from the model
fc <- forecast(model, type = 'response')

# plot using the available S3 plotting function
plot(fc)
```

response predictions

Code Plot



mvgam and brms



's

Type	mvgam	brms
link	<code>predict(type = 'link')</code>	<code>posterior_linpred()</code>
expected	<code>predict(type = 'expected')</code>	<code>posterior_epred()</code>
response	<code>predict(type = 'response')</code>	<code>posterior_predict()</code>

For all `mvgam` predictions, whether to include error in the dynamic process can be controlled using `process_error = TRUE` or `process_error = FALSE`

Posterior predictive checks

Fitted models yield coefficients

```
coef(model)
```

	2.5%	50%	97.5%	Rhat	n.eff
## (Intercept)	0.8489335	1.0693400	1.30013750	1	2610
## s(season).1	-0.6248299	-0.3361070	-0.03161323	1	2840
## s(season).2	-0.7103276	-0.3848765	-0.07808232	1	2976
## s(season).3	-0.4944772	-0.1698800	0.13811907	1	3297
## s(season).4	-0.1477989	0.1311880	0.41641965	1	2589
## s(season).5	0.2683289	0.5494570	0.83324085	1	2906
## s(season).6	0.1107273	0.3828205	0.65990857	1	2774

Interpret coefficients?

These coefficients are acting on the ***link scale***

Often result in nonlinear relationships on response scale

Very often, the coefficients are ***correlated somehow***

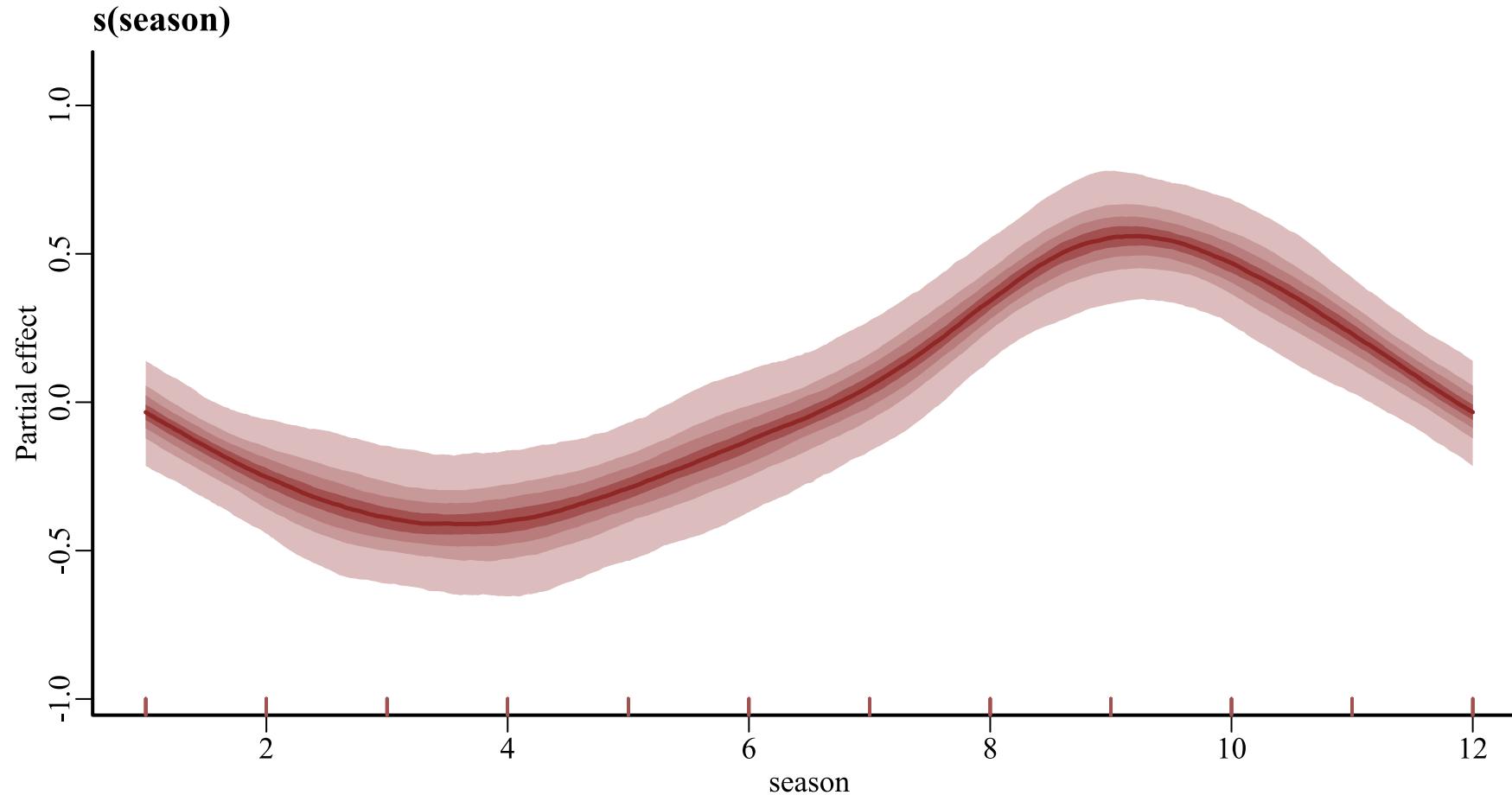
This is especially the case in GAMs!

Don't worry about *p*-values or intervals, use ***posterior predictions***
instead

Start with ***partial effects*** on link scale

These are conditional on all other effects being zero
negative values \Rightarrow covariate reduces the response
positive values \Rightarrow covariate increases the response

```
plot(model, type = 'smooths')
```



Look at partial residuals

Partial effect residuals can be thought of as **residuals that would be obtained** by dropping a specific term from the model

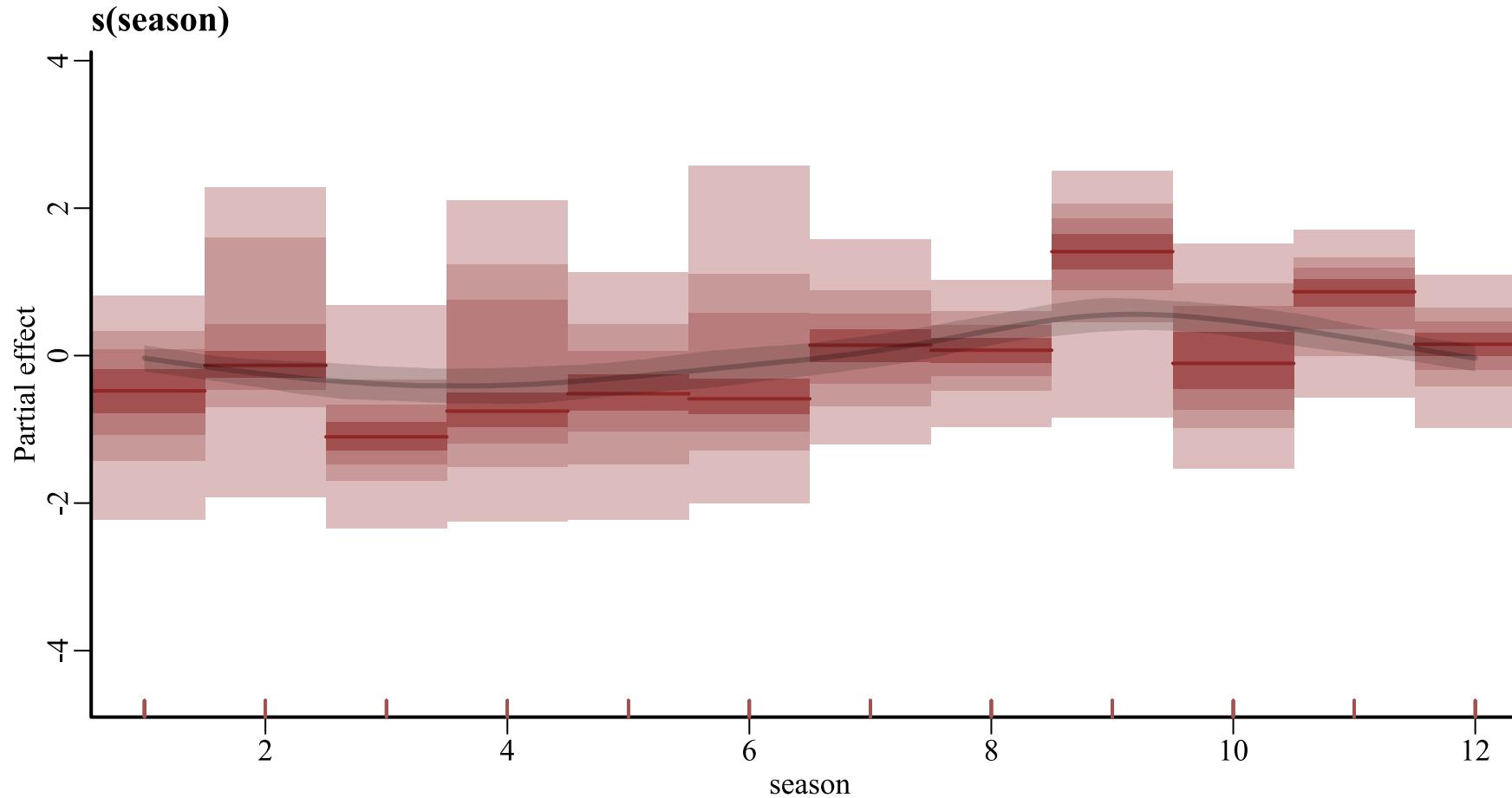
$$\hat{\epsilon}^{partial} = \hat{f}(x) + \hat{\epsilon}^{DS}$$

Where:

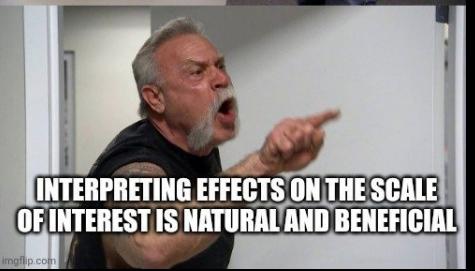
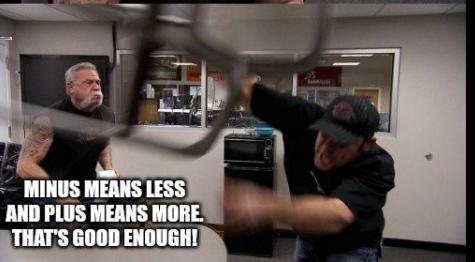
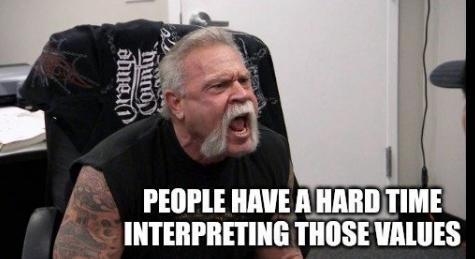
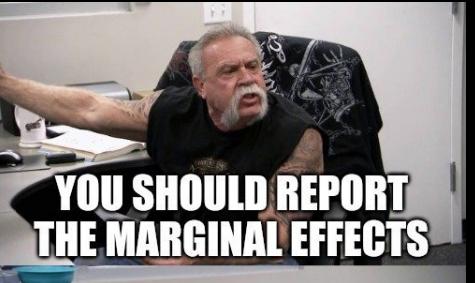
$\hat{f}(x)$ is estimated smooth function for the effect of covariate x
 $\hat{\epsilon}^{DS}$ is a draw of randomized quantile (Dunn-Smyth) residuals

We would expect these to be scattered evenly around the smooth for a well fitting model

```
plot(model, type = 'smooths', residuals = TRUE)
```



Ok. but what do these things actually, really *mean*?



Credit [@stephenjwild](#)

Interpreting on the *response scale*

Some key questions you should ask of a fitted model

Can the model simulate realistic data?

Does the model capture salient features of the data that you'd like to predict?

What criteria would you use to determine whether one model is more suitable than another?

Very often, these questions can only be answered by looking at what kinds of predictions a model makes ***on the response scale***

Posterior predictive checks

Statistical models can be used to generate (i.e. simulate) new outcome data

Can either use the same covariates used to train the model

Or can use `newdata` for scenario modelling (including forecasting)

To generate new outcome data we can simulate from the model's posterior predictive distribution

"The idea is simple: if a model is a good fit then we should be able to use it to generate data that looks a lot like the data we observed"

Gabry & Mahr

All mvgams use simulations

Uses the extremely efficient generated quantities block

```
## [1] "generated quantities {"
## [2] "vector[total_obs] eta;"
## [3] "matrix[n, n_series] mus;"
## [4] "vector[n_sp] rho;"
## [5] "array[n, n_series] int ypred;"
## [6] "rho = log(lambda);"
## [7] ""
## [8] "// posterior predictions"
## [9] "eta = X * b;"
## [10] "for(s in 1:n_series){ "
## [11] "mus[1:n, s] = eta[ytimes[1:n, s]] + trend[1:n, s];"
## [12] "ypred[1:n, s] = poisson_log_rng(mus[1:n, s]);"
## [13] "}"
## [14] "}"
```

A PPC histogram

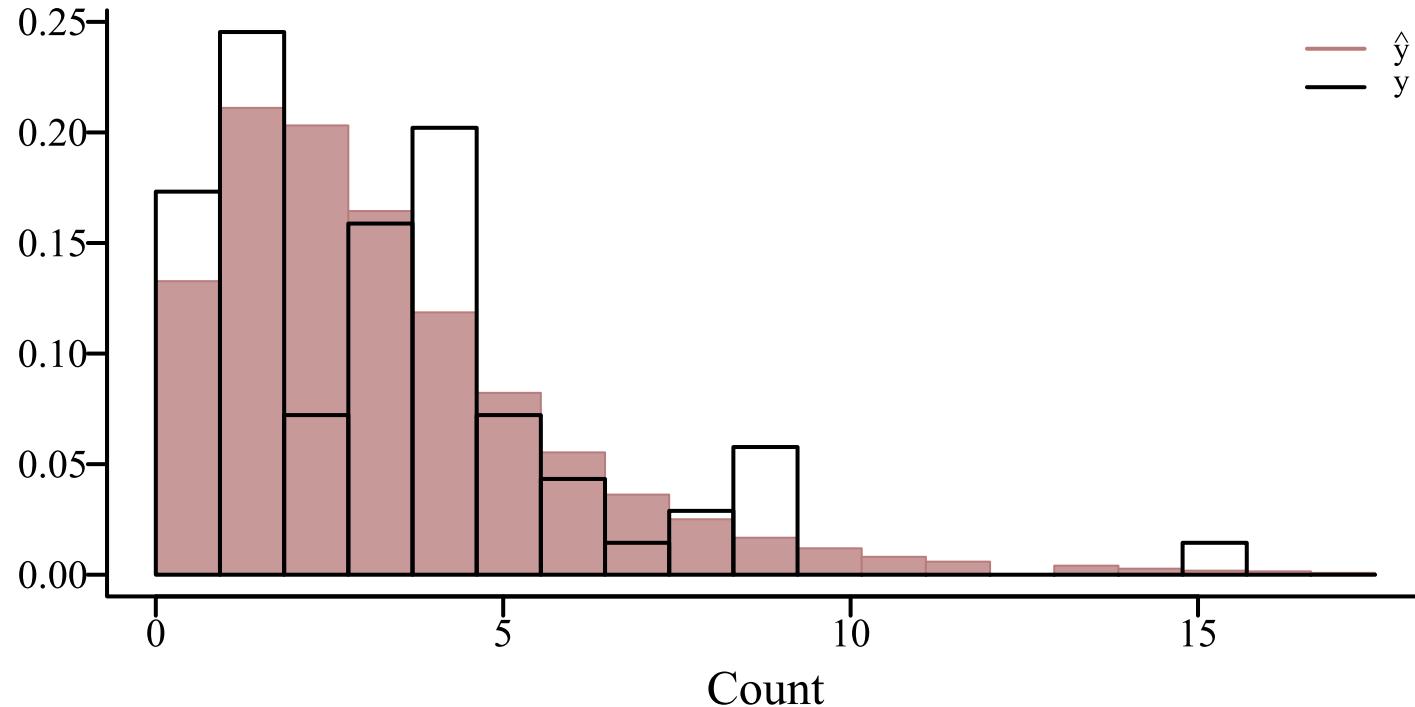
Code Plot

```
# view a histogram of true data vs simulated predictions
ppc(model, type = 'hist')
title(expression('Predictions ('*hat(y)*') vs true observations (y)'))
```

A PPC histogram

Code Plot

Predictions (\hat{y}) vs true observations (y)



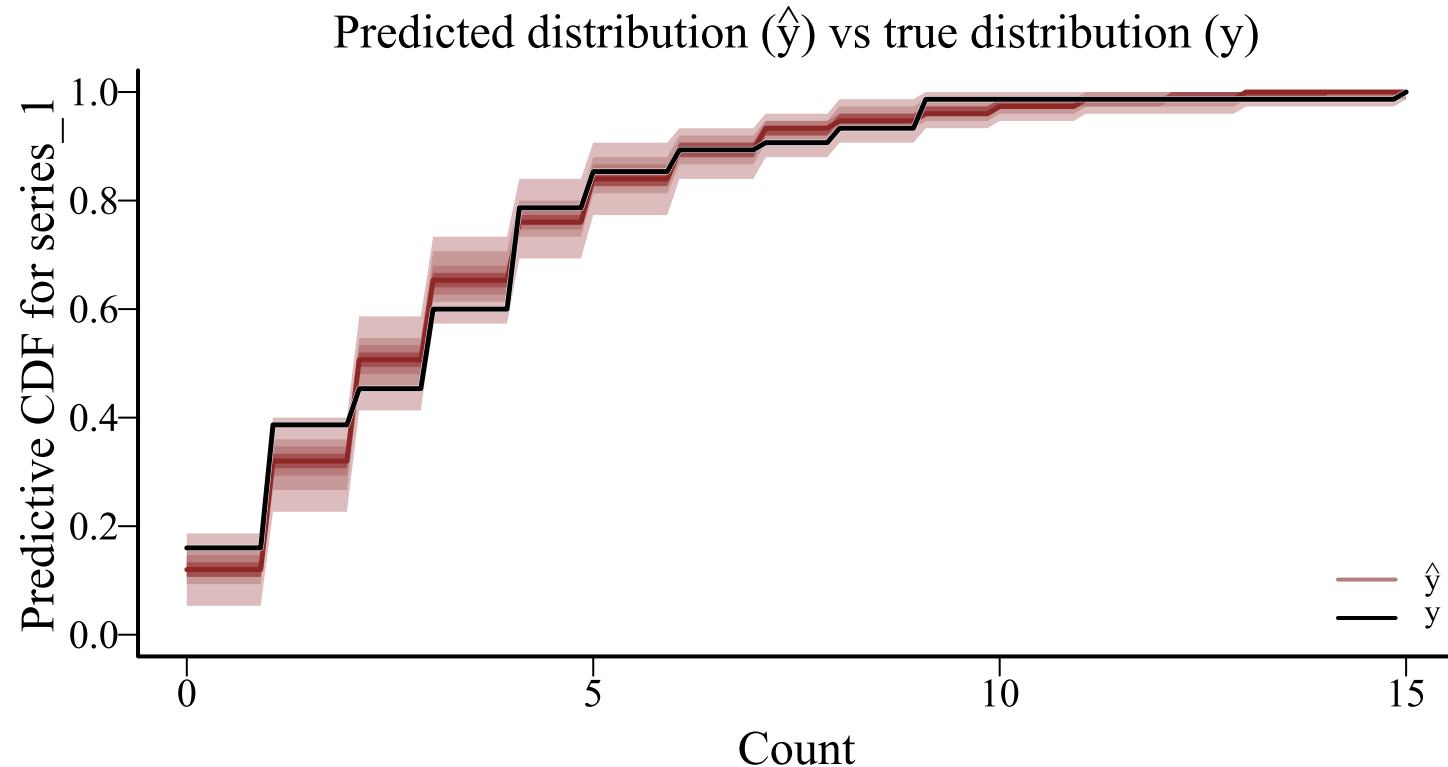
A PPC cumulative distribution

Code Plot

```
# view the simulated vs true cumulative distribution functions
ppc(model, type = 'cdf')
title(expression('Predicted distribution ('*hat(y)*') vs true distribution
(y)'))
title(xlab = 'Count')
```

A PPC cumulative distribution

Code Plot



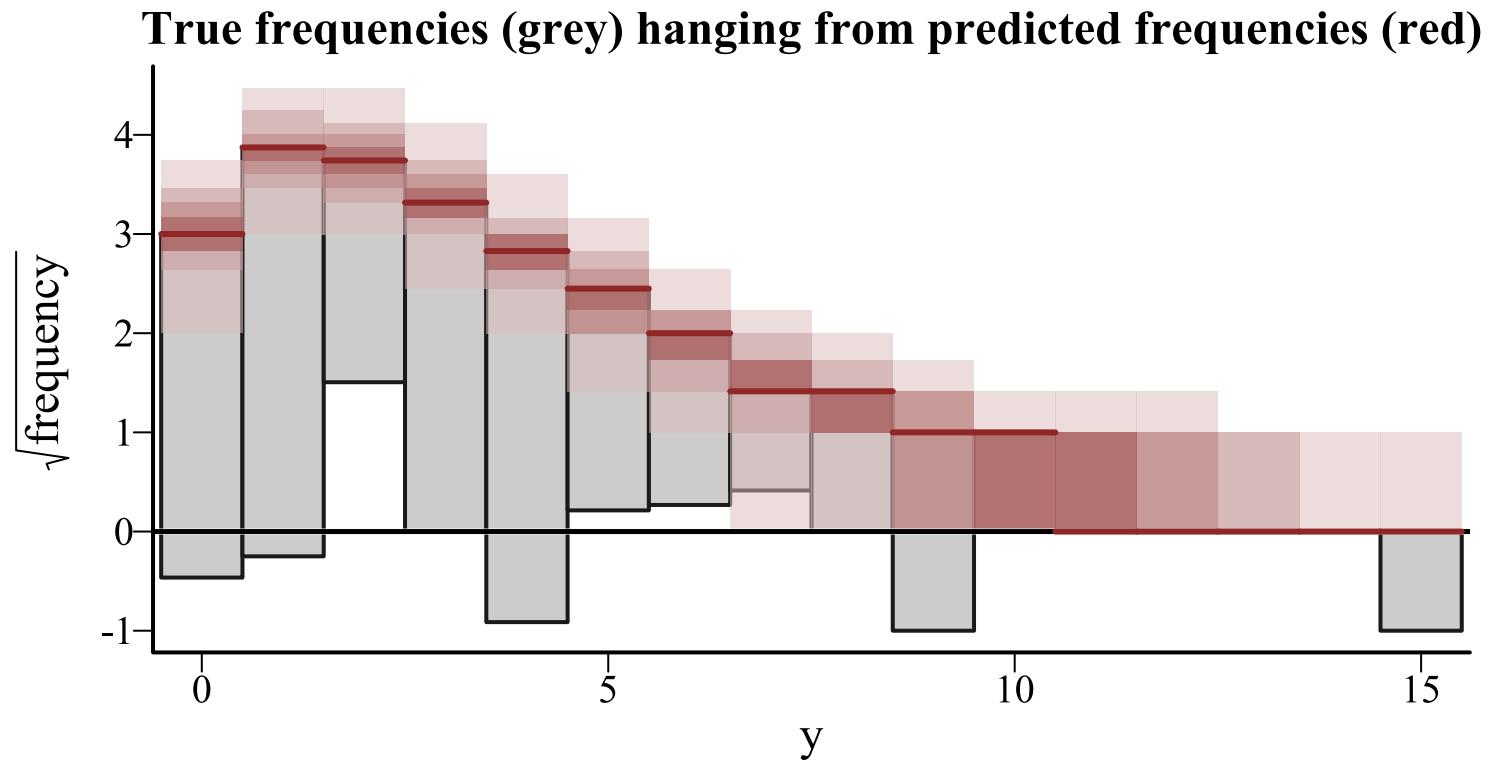
A PPC rootogram

Code Plot

```
# view the simulated vs true count frequencies
ppc(model, type = 'rootogram')
title('True frequencies (grey) hanging from predicted frequencies (red)')
```

A PPC rootogram

Code Plot



PPCs using the training covariates are a great first step to check model validity

But they can be misleading (very flexible dynamic processes mean predictions may be nearly perfect)

How else can we verify models? Using `newdata` for response predictions \Rightarrow counterfactual *scenarios*

Marginal & conditional predictions

"Applied researchers are keen to report simple quantities that carry clear scientific meaning" ([Arel-Bundock 2023](#))

This is often challenging because:

- Intuitive estimands and uncertainties are tedious to compute
- Nonlinear terms, nonlinear link functions, interaction effects and observation parameters all make these effects nearly impossible to gain from looking at coefficients alone
- Most software emphasizes coefficients and p -values over meaningful interpretations

`predict.mvgam`

Feed `newdata` consisting of particular covariate values that represent scenarios you'd like to explore

Can be simple: predict a smooth function along a fine-spaced grid to explore the smooth's shape and / or derivatives

Or can be complex: integrate over a high-dimensional grid of predictors to understand the average impact of a predictor on the response

Users can implement the wonderful `datagrid` function from `marginaleffects`  to effortlessly generate a `data.frame` of covariate values for scenario predictions

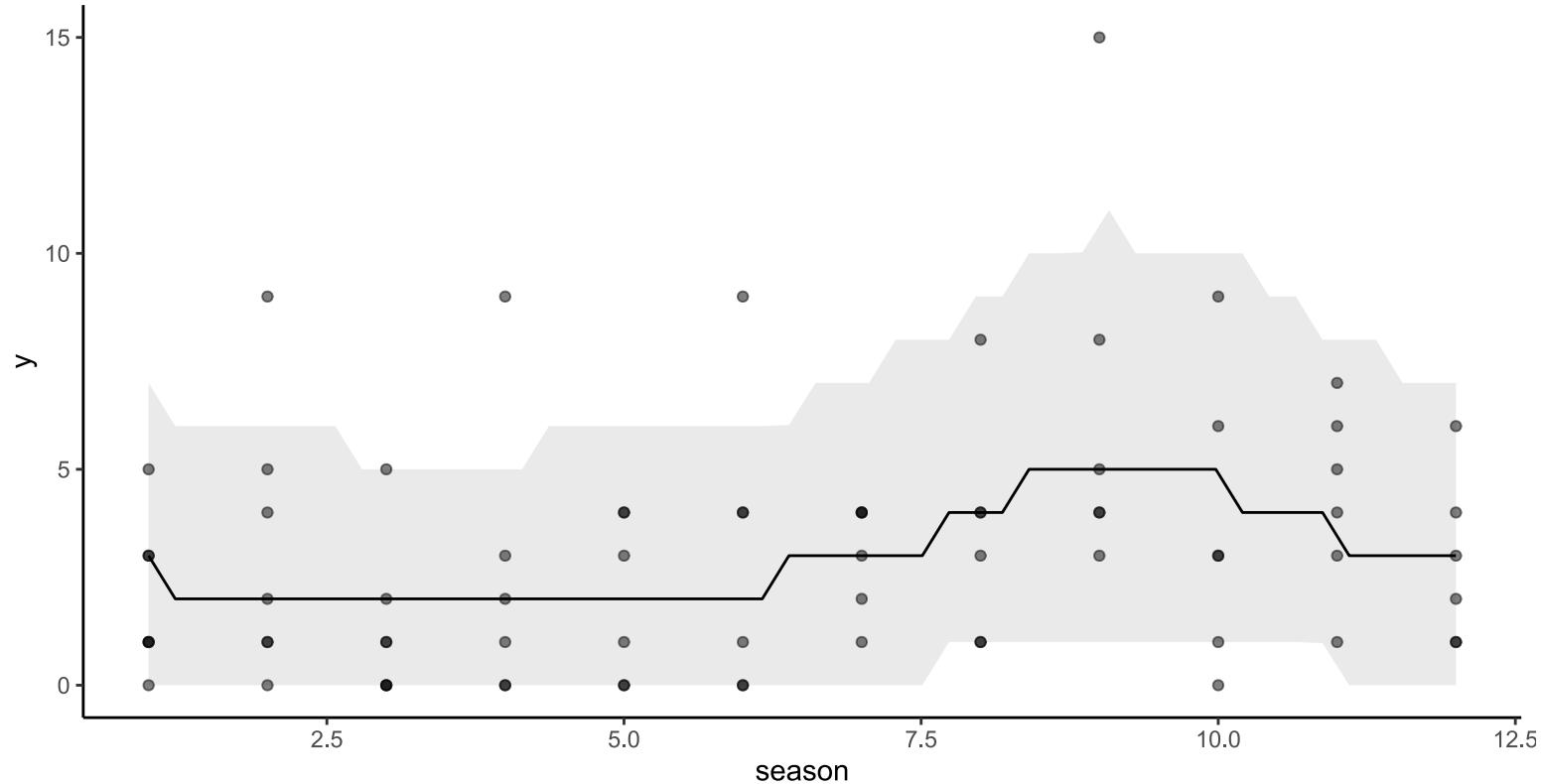
Conditional smooths

Code Plot

```
# use plot_predictions to visualise conditional effects
# on the scale of the response
library(ggplot2)
plot_predictions(model, condition = 'season',
                  points = 0.5, process_error = FALSE) +
  theme_classic()
```

Conditional smooths

Code Plot



Posterior contrasts

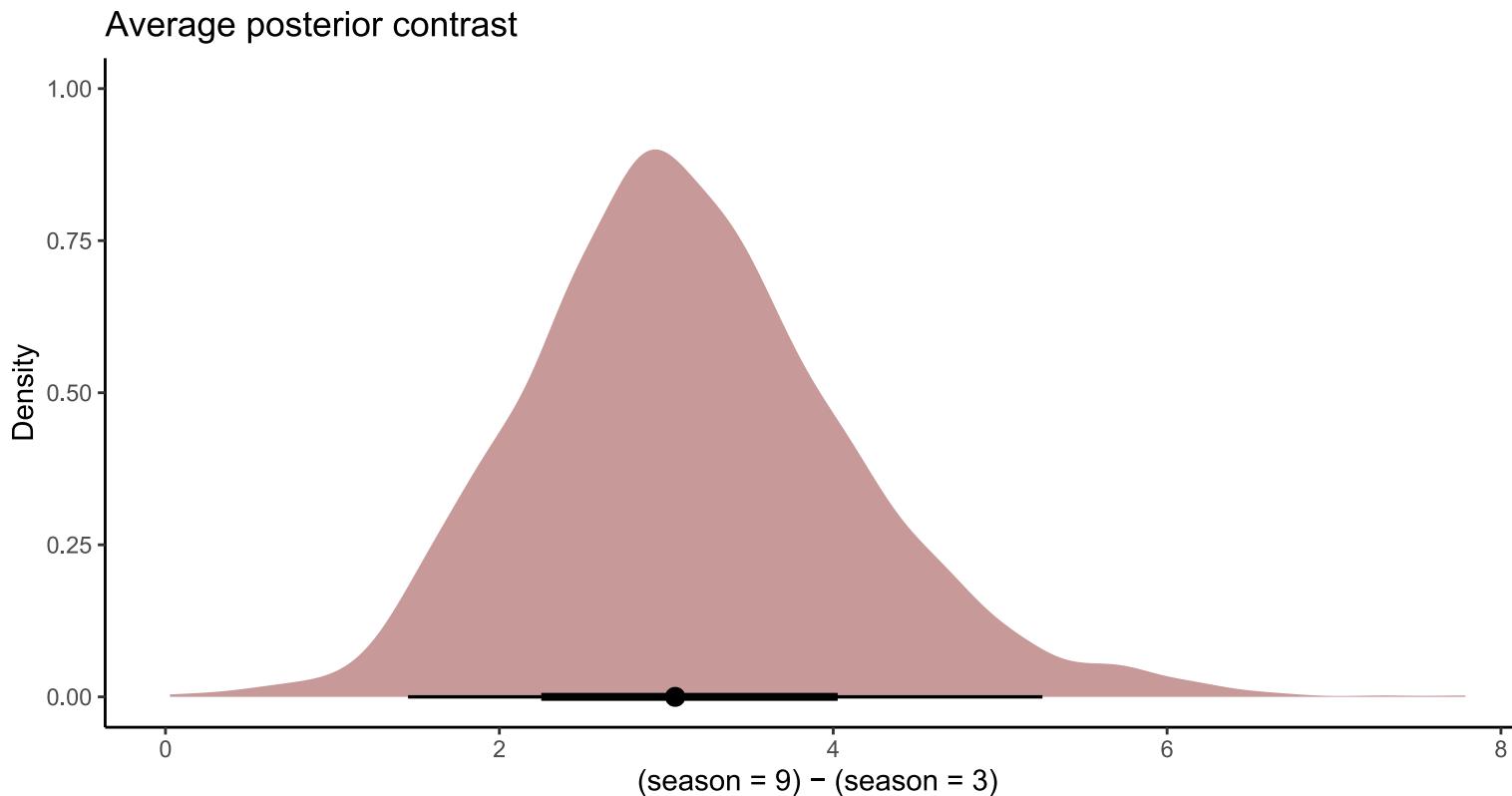
Code Plot

```
# take draws of average comparison between season = 9 vs season = 3
post_contrasts ← avg_comparisons(model,
                                    variables = list(season = c(9, 3)),
                                    proces_error = FALSE) %>%
  posteriordraws()

# use the resulting posterior draw object to plot a density of the
# posterior contrasts
library(tidybayes)
post_contrasts %>% ggplot(aes(x = draw)) +
  # use the stat_halfeye function from tidybayes for a nice visual
  stat_halfeye(fill = "#C79999") +
  labs(x = "(season = 9) - (season = 3)", y = "Density",
       title = "Average posterior contrast") + theme_classic()
```

Posterior contrasts

[Code](#) [Plot](#)



The ability to readily interpret models from `mgvam` and `brms` 's
is a *huge advantage* over traditional time series models



But this is a forecasting course. So how can we evaluate forecast distributions?

The forecasting workflow

"The accuracy of forecasts can only be determined by considering how well a model performs on new data that were not used when fitting the model." Hyndman and Athanasopoulos

We must evaluate on data that was not used to train the model (i.e. ***leave-future-out cross-validation***) because:

Models that fit training data well do not always provide good forecasts

We can easily engineer a model that perfectly fits the training data, leading to overfitting

Leave-future-out CV

Important to train the model on some portion of data and use a hold-out portion (test data) to evaluate forecasts:

$$p(y_{T+H} | y_{1:T})$$

Some points to consider:

The test set should ideally be at least as large as the maximum forecast horizon required for decision-making

Ideally, this process would be repeated many times to incorporate variation in forecast performance

Usually good to compare models against simpler **benchmark** models to ensure added complexity improves forecasts

Approximate leave-future-out CV

Re-fitting Bayesian models to obtain exact forecasts on different training sets can be computationally infeasible

Some models fit with `mvgam` and `brms` 's can take several hours to compute posterior distributions

Doing this even 5 or 6 times may be impractical, especially when comparing multiple models

There are ways we can **approximate** the leave-future-out forecast distribution using importance sampling of the posterior distribution

Won't discuss details here, but see [Bürkner et al 2020](#) for in-depth information

We must obtain leave-future-out forecasts (ideally for many different training / testing splits) to compare ecological forecasting models

But how do we *evaluate* forecasts?

The most common evaluation practice in forecasting tasks is to evaluate point predictions

Point-based
forecast
evaluation

Forecast errors

A forecast error (or forecast residual) is the difference between the true value in an out-of-sample set and the predicted response value:

$$\epsilon_{T+H} = \mathbf{y}_{T+H} - \hat{\mathbf{y}}_{T+H}$$

Where:

T is the total length of the training set

H is the forecast horizon

$\hat{\mathbf{y}}_{T+H}$ is the prediction at time $T + H$

Point-based measures use these errors in different ways

Common point-based measures

Scale-dependent measures

Mean Absolute Error: $\text{mean}(|\epsilon_t|)$

Root Mean Squared Error: $\sqrt{\text{mean}(\epsilon_t^2)}$

Scale-independent measures

Mean Absolute Percentage Error: $\text{mean}(|p_t|)$, where $p_t = 100\epsilon_t/y_t$

Mean Absolute Scaled Error: $\text{mean}(|q_t|)$, where q_t is the error

scaled against errors from an appropriate **benchmark** forecast

Lower values are better for all these measures

We won't dwell much on point-based measures because ecological predictions and their associated management decisions are inherently (but see this video for more details) *uncertain*

Point-based measures ignore far too much information in the forecast distribution

It is better to evaluate the *entire forecast distribution*

Probabilistic
forecast
evaluation

Interval scores

A common step to evaluate a forecast distribution is to compute how well its prediction intervals perform:

$$MSIS = (U_t - L_t) + \frac{2}{\alpha} (L_t - y_t) \mathbf{1}(y_t < L_t) + \frac{2}{\alpha} (y_t - U_t) \mathbf{1}(y_t > U_t)$$

Where:

y_t is the true observed value at horizon H

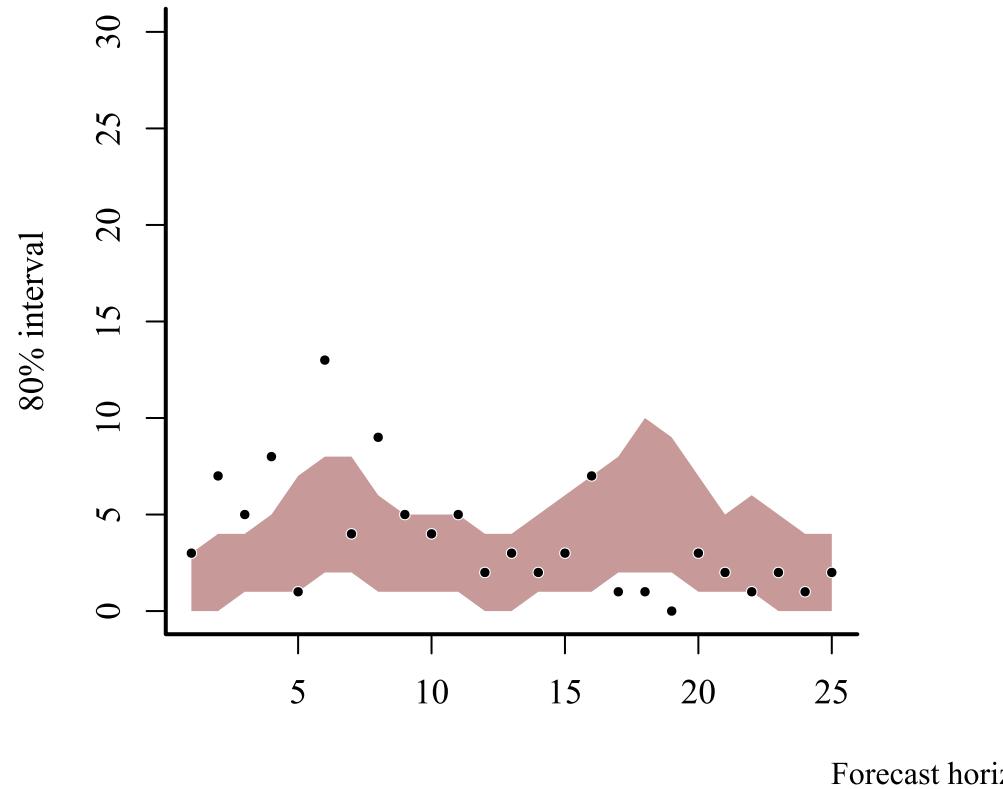
α is 1 – interval width

The $100(1 - \alpha)\%$ interval for horizon H is $[L_t, U_t]$

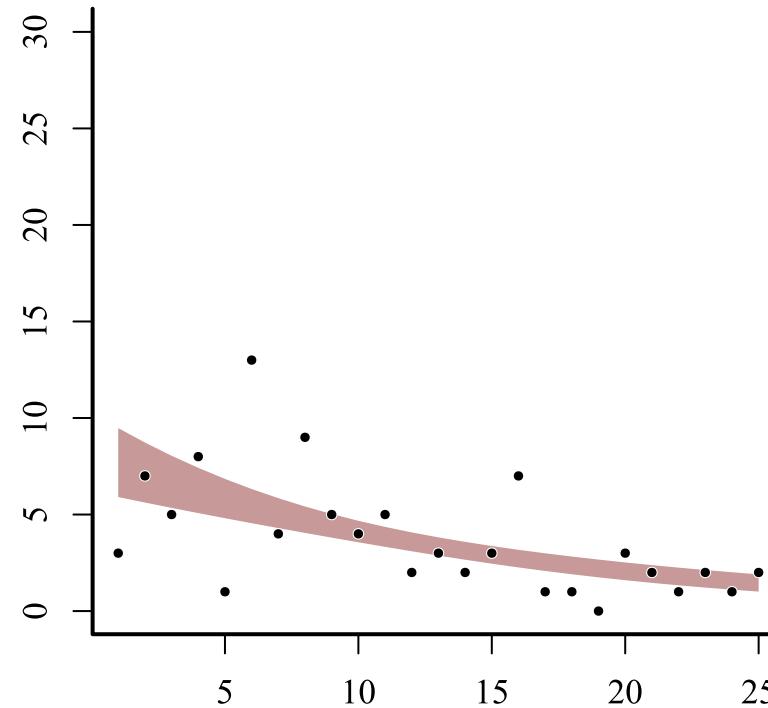
$\mathbf{1}$ is a binary indicator function

Penalize *overly precise* forecasts

$\text{MSIS}_{80} = 12.44$



$\text{MSIS}_{80} = 13.06$



Evaluating the full distribution

Interval scores are very useful when we want to target a particular interval or if we don't have the full distribution

Allows different teams to submit a few intervals rather than thousands of posterior samples

Can compare forecasts from many different algorithms / models

But if we do have a full distribution, we have other options

"Scoring rules provide summary measures for the evaluation of probabilistic forecasts, by assigning a numerical score based on the predictive distribution and on the event or value that materializes"
(Gneiting and Raftery 2007)

What is a good forecast?

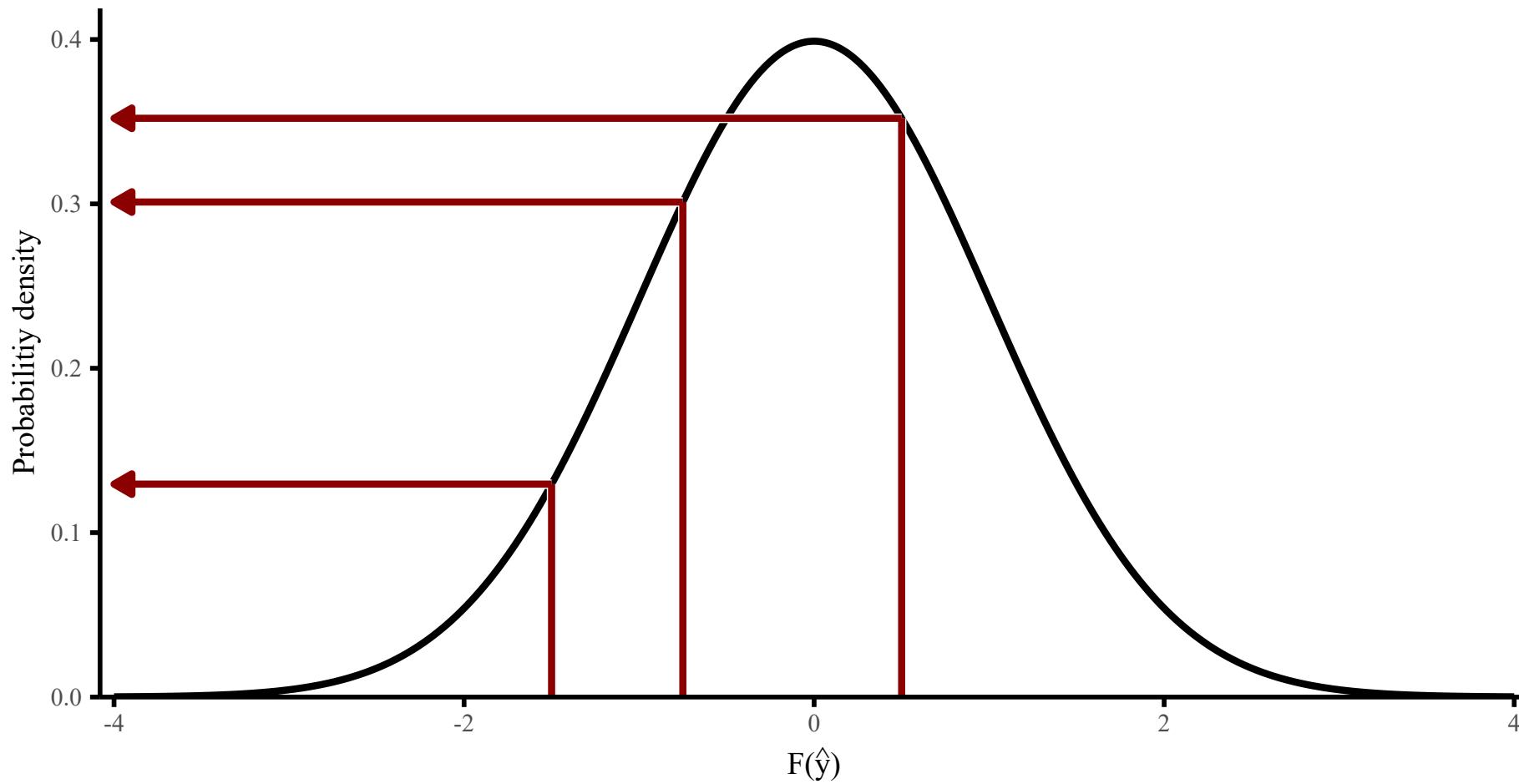
Reliable: good probabilistic calibration

Sharp: informative, with tight enough intervals to guide decisions

Skilled: performs better overall than simpler benchmark forecasts

Proper scoring rules attempt to address each of these goals using the full forecast distribution

Predictive density



Log predictive density

Compute $\log(\text{probability})$ of a given truth given distributional assumptions:

$$\log p(y_{T+H} | y_{t:T}, \theta)$$

Use density functions in , such as `dnorm` or `dnbino`; higher values are better

θ captures all unknown parameters:

Regression coefficients β

Dynamic parameters; α or ρ for GP; σ_{error} for RW

Observation parameters; ν for StudentT or σ_{obs} for Normal

logging brings numerical stability and makes joint calculations easier

But the log score can severely penalize over-confidence and is sensitive to outliers

Other proper scoring rules can provide more robust comparisons, without needing to rely on distributional assumptions

CRPS

Continuous Ranked Probability Score compares true Cumulative Distribution Function (CDF) to forecast CDF

$$CRPS(F, y) = \int_{-\infty}^{\infty} (F(\hat{y}) - \mathbf{1}(\hat{y} \geq y))^2 dy$$

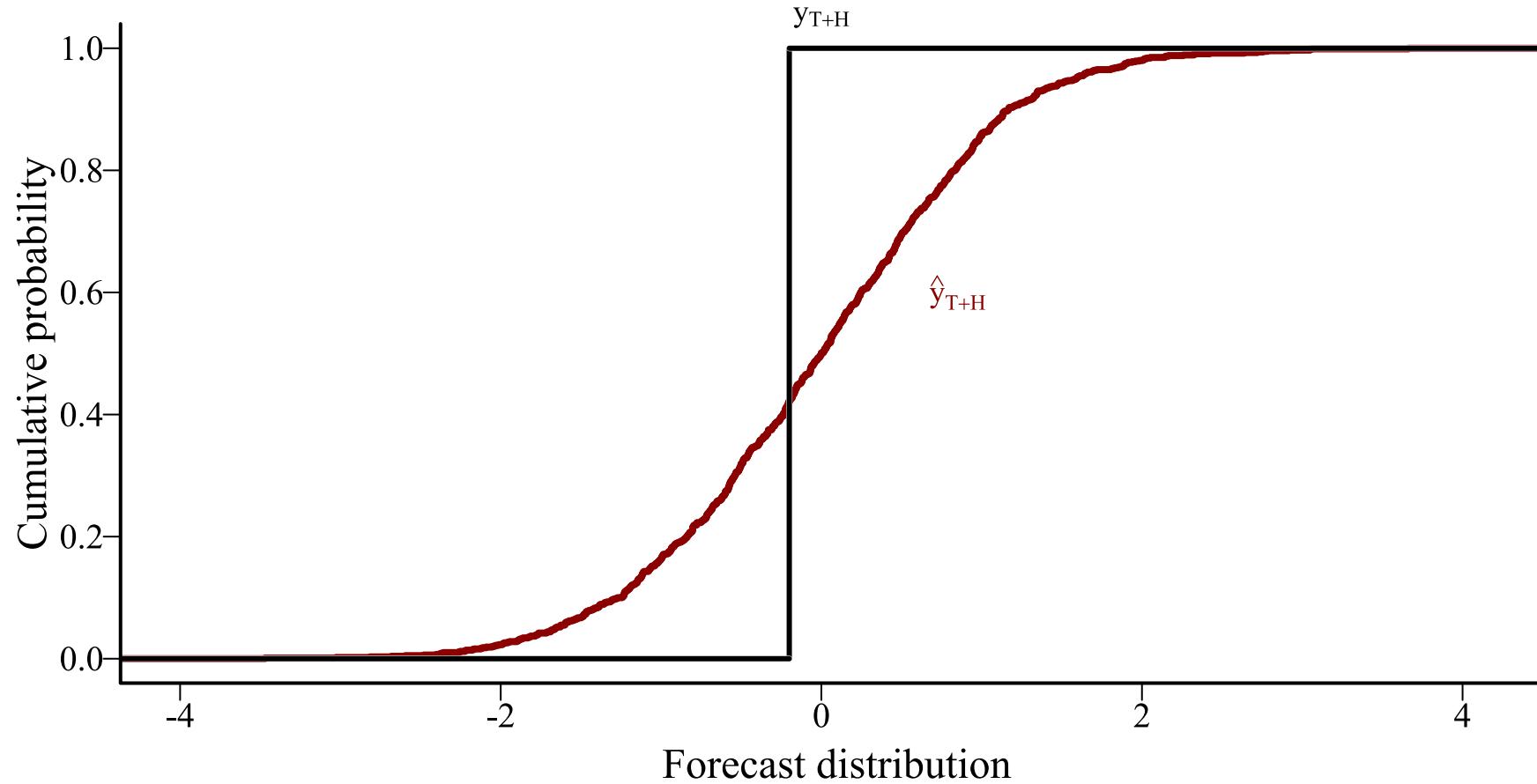
Where:

$F(\hat{y})$ is the forecast CDF evaluated at many points

$\mathbf{1}(\hat{y} \geq y)$ gives the true observed CDF

MSIS converges to CRPS when evaluating an increasing number of equally spaced intervals

CRPS



CRPS useful for both parametric and non-parametric predictions because we just need to calculate the CDF of the forecast distribution

Penalises over- and under-confidence similarly, and gives more stable handling of outliers

Score is on the scale of the outcome variable being forecasted, so is somewhat intuitive (a lower score is better)

DRPS

Similar to CRPS, the discrete version (DRPS) can be used to evaluate a forecast that is composed only of integers

Uses an approximation of the forecast and true CDFs at a range of possible count values

Interpretation is similar

score.mvgam_forecast

Once forecasts are computed and stored in an object of class `mvgam_forecast`, scores can be directly applied

User chooses among the Scaled Interval Score (`sis`), log score (`elpd`), CRPS (`crps`), DRPS (`drps`) and two multivariate scores (`energy` or `variogram`; more on this in the next lecture)

User also specifies an interval for calculating coverage and/or which interval to use for the Scaled Interval Score

`return` is a `list` with scores for each series in the data and an overall score (usually just the sum of series-level scores)

```
sc ← score(forecast(model),  
           score = 'crps',  
           interval = 0.90)  
sc$series_1[1:10,]
```

```
##          score in_interval interval_width eval_horizon score_type  
## 1  0.9019320            1             0.9            1      crps  
## 2  4.2066147            0             0.9            2      crps  
## 3  1.9790223            1             0.9            3      crps  
## 4  4.0763472            0             0.9            4      crps  
## 5  1.6685138            1             0.9            5      crps  
## 6  6.4072417            0             0.9            6      crps  
## 7  0.6793927            1             0.9            7      crps  
## 8  3.8851392            1             0.9            8      crps  
## 9  1.3764830            1             0.9            9      crps  
## 10 0.7970892            1             0.9           10      crps
```

Calculating the CRPS using the previously generated forecasts

```
sc ← score(forecast(model),  
           score = 'sis',  
           interval = 0.90)  
sc$series_1[1:10,]
```

```
##   score in_interval interval_width eval_horizon score_type  
## 1 4.00          1        0.9          1       sis  
## 2 45.00         0        0.9          2       sis  
## 3 6.00          1        0.9          3       sis  
## 4 27.00         0        0.9          4       sis  
## 5 8.00          1        0.9          5       sis  
## 6 50.00         0        0.9          6       sis  
## 7 10.00         1        0.9          7       sis  
## 8 9.00          1        0.9          8       sis  
## 9 7.05          1        0.9          9       sis  
## 10 8.00         1        0.9         10      sis
```

Calculating the SIS using the previously generated forecasts; values outside interval are more heavily penalized

We have seen how to produce out-of-sample forecasts from `mgvam` models and evaluate them against new observations

We have also investigated other ways that models can be critiqued, particularly making use of conditional predictions using `newdata`

But so far we have only considered univariate investigations. What happens if we want to forecast *multiple time series*?

In the next lecture, we will cover

Multivariate ecological time series

Vector autoregressive processes

Dynamic factor models

Multivariate forecast evaluation