

Gaussian Process Models

Part 2

Lecture 14

Dr. Colin Rundel

EDA and GPs

Variogram

When fitting a Gaussian process model, it is often difficult to fit the covariance parameters (they are often correlated and hard to identify).

Today we will discuss some EDA approaches for getting a sense of the values for the scale, lengthscale / effective range and nugget parameters.

From the spatial modeling literature the typical approach is to examine an *empirical variogram*, first we will define the *theoretical variogram* and its connection to the covariance.

Variogram & semivariogram

Variogram

$$2\gamma(t_i, t_j) = \text{Var} (y(t_i) - y(t_j))$$

SemiVariogram

$$\gamma(t_i, t_j) = \frac{1}{2} \text{Var} (y(t_i) - y(t_j))$$

Properties of the Variogram / Semivariogram

- both are non-negative

$$\gamma(t_i, t_j) \geq 0$$

- both are equal to 0 at distance 0

$$\gamma(t_i, t_i) = 0$$

- both are symmetric -

$$\gamma(t_i, t_j) = \gamma(t_j, t_i)$$

- if observations are independent

$$2\gamma(t_i, t_j) = \text{Var}(y(t_i)) + \text{Var}(y(t_j)) \quad \text{for all } i \neq j$$

- if the process *is not* stationary

$$2\gamma(t_i, t_j) = \text{Var}(y(t_i)) + \text{Var}(y(t_j)) - 2 \text{Cov}(y(t_i), y(t_j))$$

- if the process *is* stationary

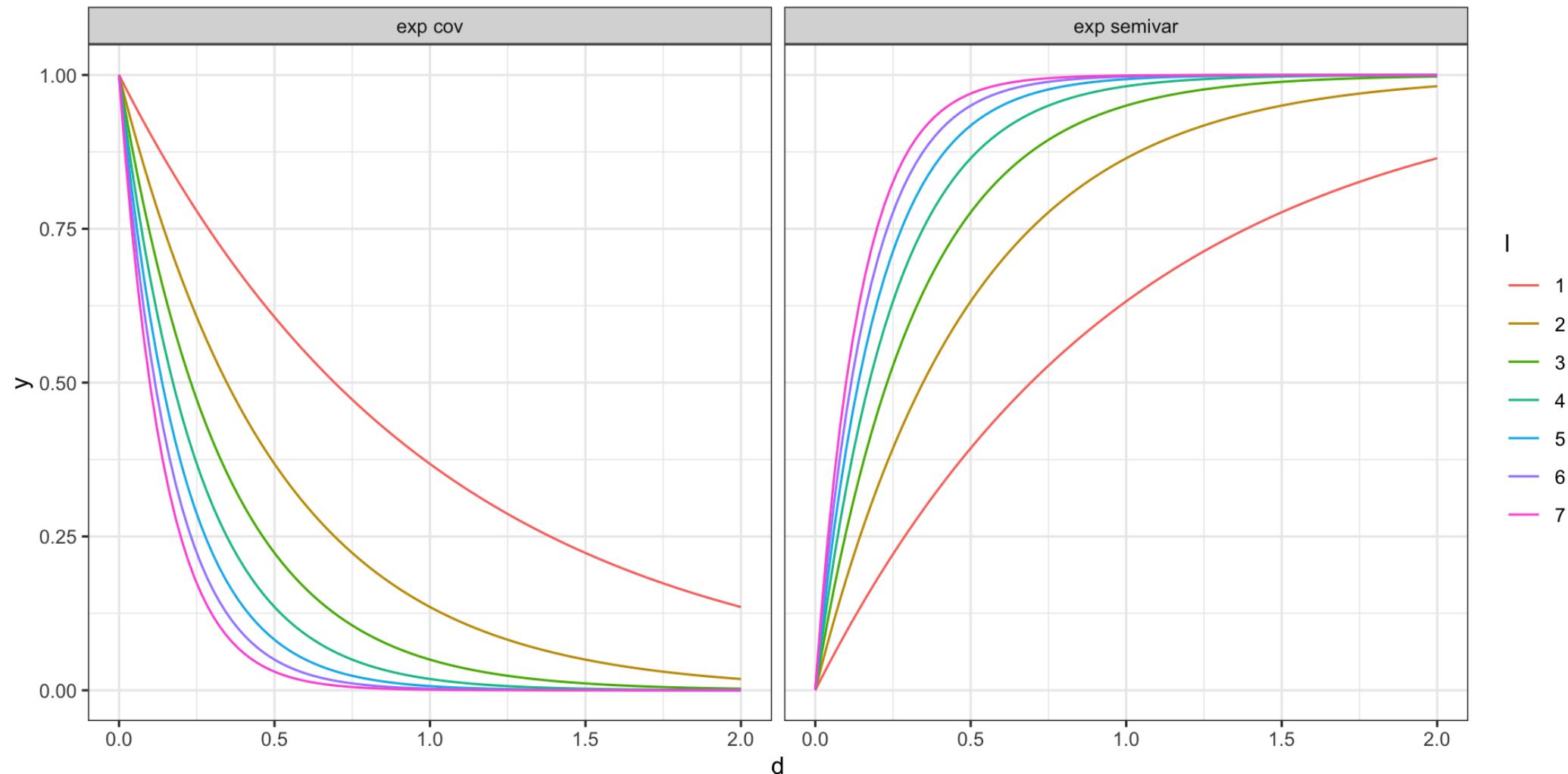
$$2\gamma(t_i, t_j) = 2 \text{Var}(y(t_i)) - 2 \text{Cov}(y(t_i), y(t_j))$$

Connection to Covariance

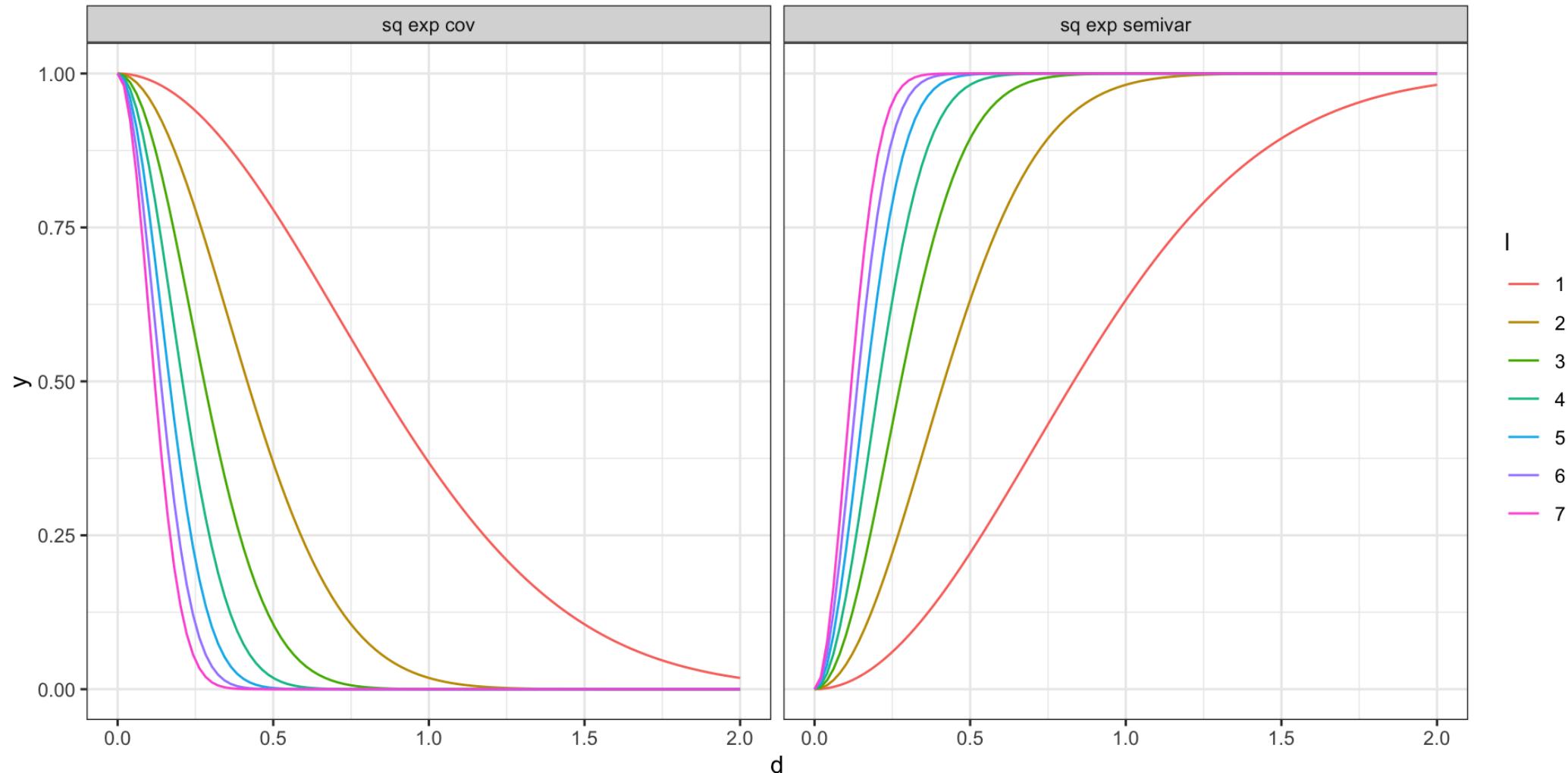
Assuming a squared exponential covariance structure and stationarity,

$$\begin{aligned}2\gamma(t_i, t_j) &= 2\text{Var}(y(t_i)) - 2\text{Cov}(y(t_i), y(t_j)) \\ \gamma(t_i, t_j) &= \text{Var}(y(t_i)) - \text{Cov}(y(t_i), y(t_j)) \\ &= \sigma^2 - \sigma^2 \exp(-(|t_i - t_j|/\ell)^2)\end{aligned}$$

Covariance vs Semivariogram - Exponential



Covariance vs Semivariogram - Sq. Exp.



Nugget variance

Very often in the real world we will observe that $\gamma(t_i, t_i) = 0$ is not true - there will be an initial discontinuity in the semivariogram at $|t_i - t_j| = 0$.

Why is this?

We can think about Gaussian process regression in the following way,

$$y(t) = \mu(t) + w(t) + \epsilon(t)$$

where

$$\mu(t) = X\beta$$

$$w(t) \sim N(\mathbf{0}, \Sigma)$$

$$\epsilon(t) \stackrel{\text{iid}}{\sim} N(0, \sigma_w^2)$$

Implications

With the inclusion of the $\epsilon(t)$ terms in the model we now have,

$$\text{Var}(y(t_i)) = \sigma_w^2 + \Sigma_{ii}$$

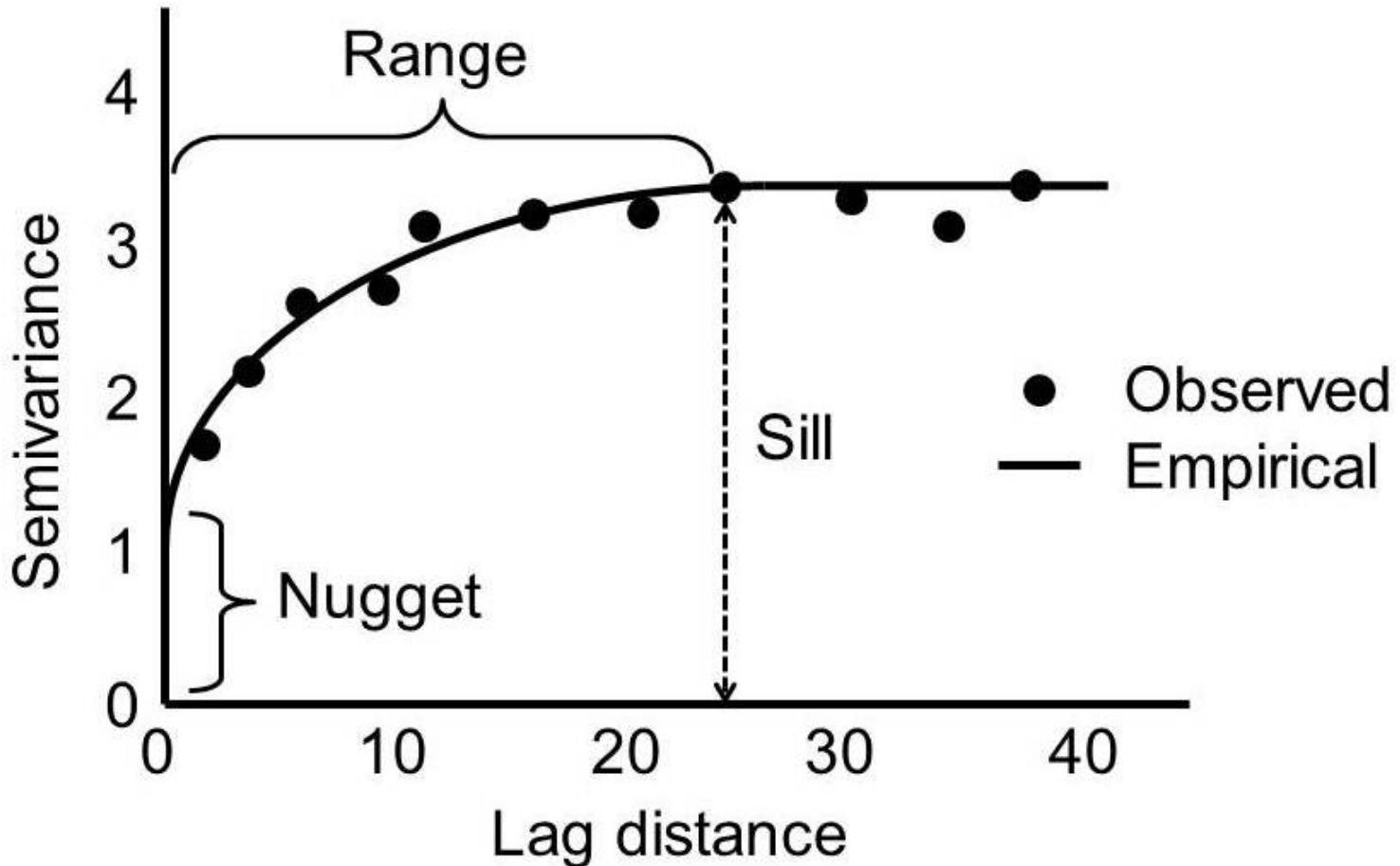
$$\text{Cov}(y(t_i), y(t_j)) = \Sigma_{ij}$$

Therefore, for a squared exponential covariance model with a nugget component the semivariogram is given by,

$$\gamma(t_i, t_j) = (\sigma^2 + \sigma_w^2) - \sigma^2 \exp(-(|t_i - t_j| l)^2)$$

$$\begin{aligned}\gamma(t_i, t_i) &= (\sigma^2 + \sigma_w^2) - \sigma^2 \exp(-(|t_i - t_i| l)^2) \\ &= (\sigma^2 + \sigma_w^2) - \sigma^2 \exp(0) = \sigma_w^2\end{aligned}$$

Semivariogram features



Empirical Semivariogram

We will assume that our process of interest is stationary, in which case we will parameterize the semivariogram in terms of $d = |t_i - t_j|$.

Empirical Semivariogram:

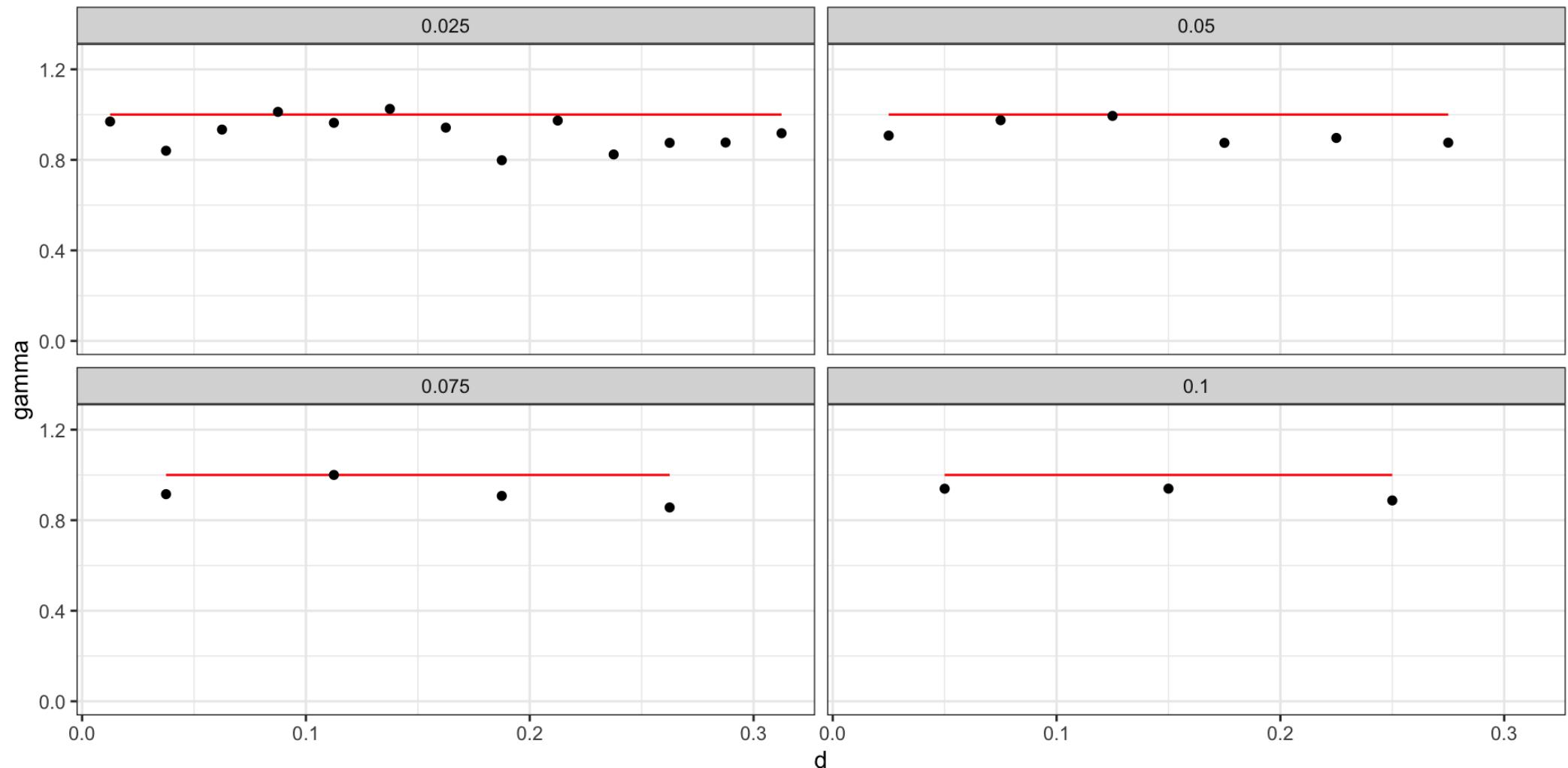
$$\hat{\gamma}(d) = \frac{1}{2N(d)} \sum_{|t_i - t_j| \in (d-\epsilon, d+\epsilon)} (y(t_i) - y(t_j))^2$$

...

Practically, for any data set with n observations there are $\binom{n}{2} + n$ possible data pairs to examine. Each individually is not very informative, so we aggregate into bins and calculate the empirical semivariogram for each bin.

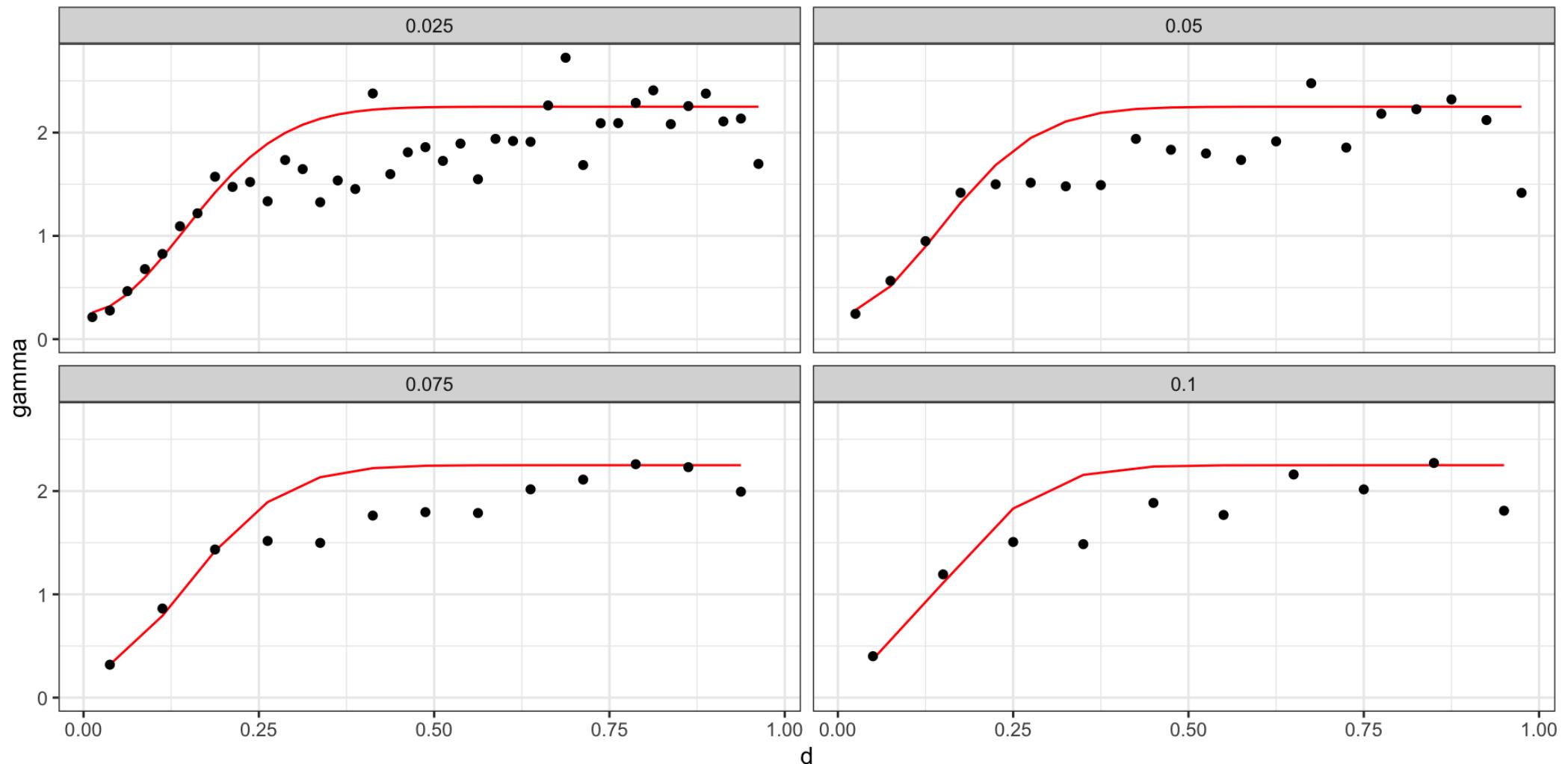
Empirical semivariogram of white noise

Where $\sigma_w^2 = 1$,



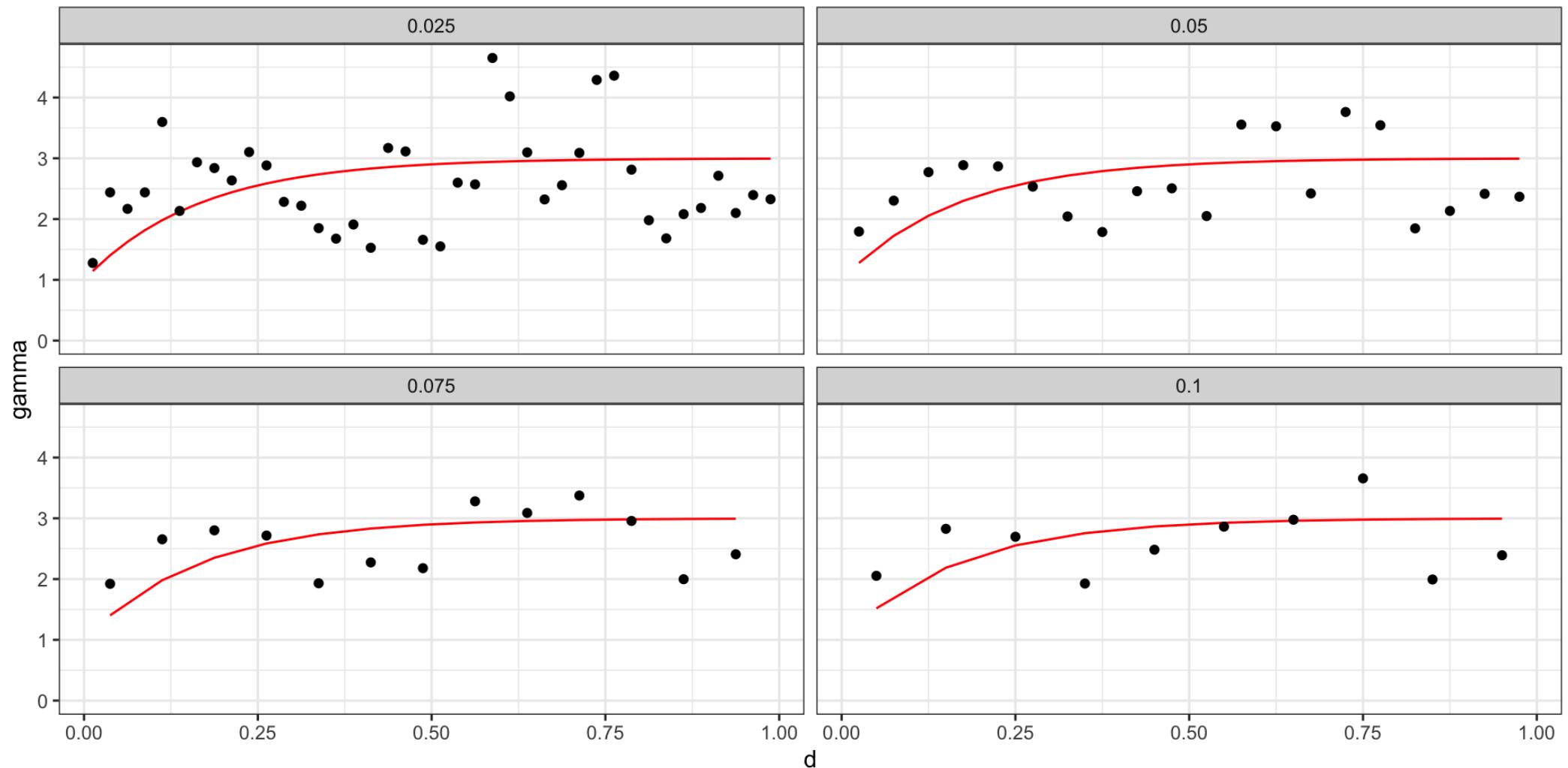
Empirical Variogram of GP w/ Sq Exp cov

Where $\sigma^2 = 2$, $l = 5$, and $\sigma_w^2 = 0.25$,

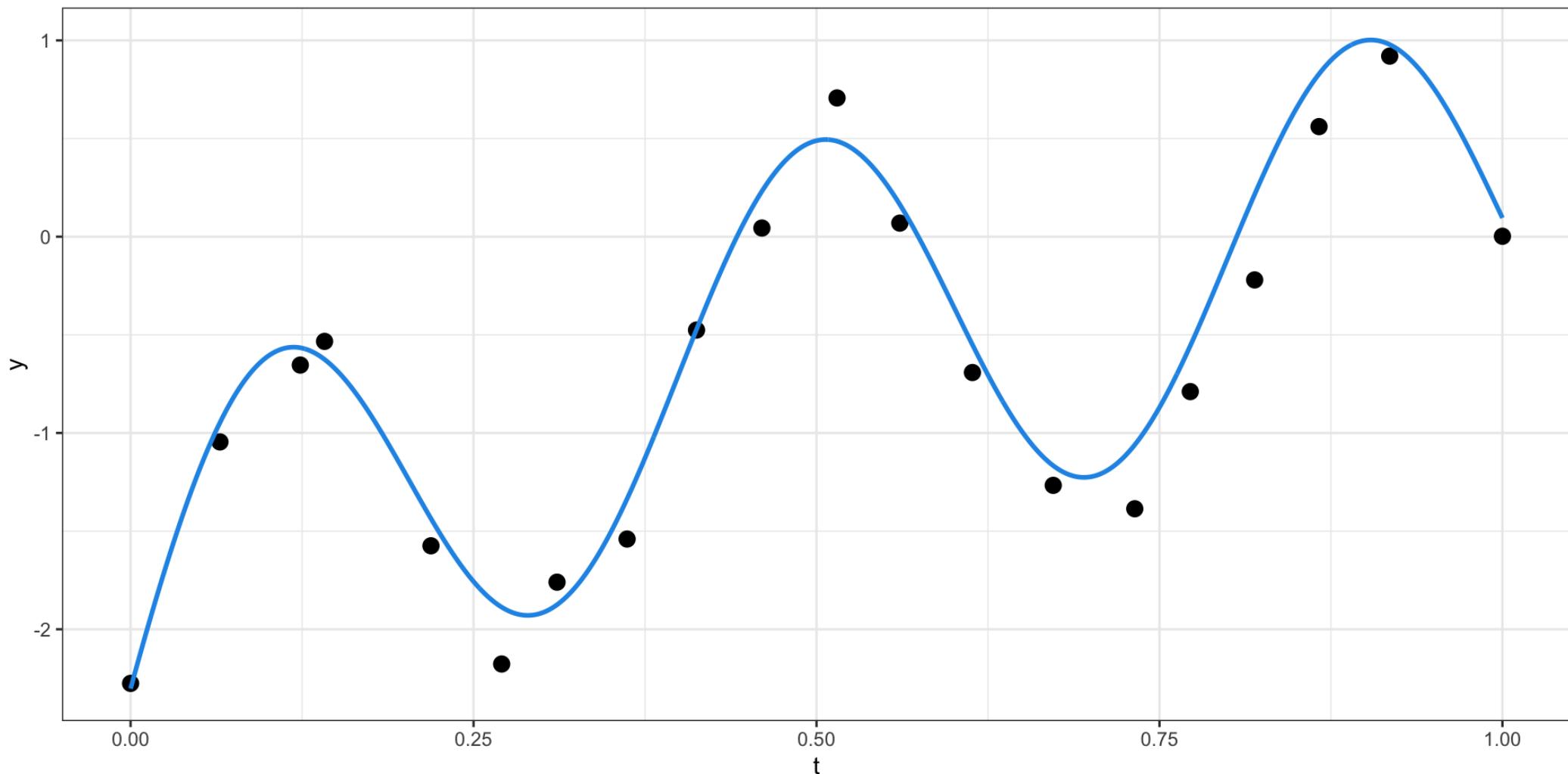


Empirical Variogram of GP w/ Exp cov

Where $\sigma^2 = 2$, $l = 6$, and $\sigma_w^2 = 1$,

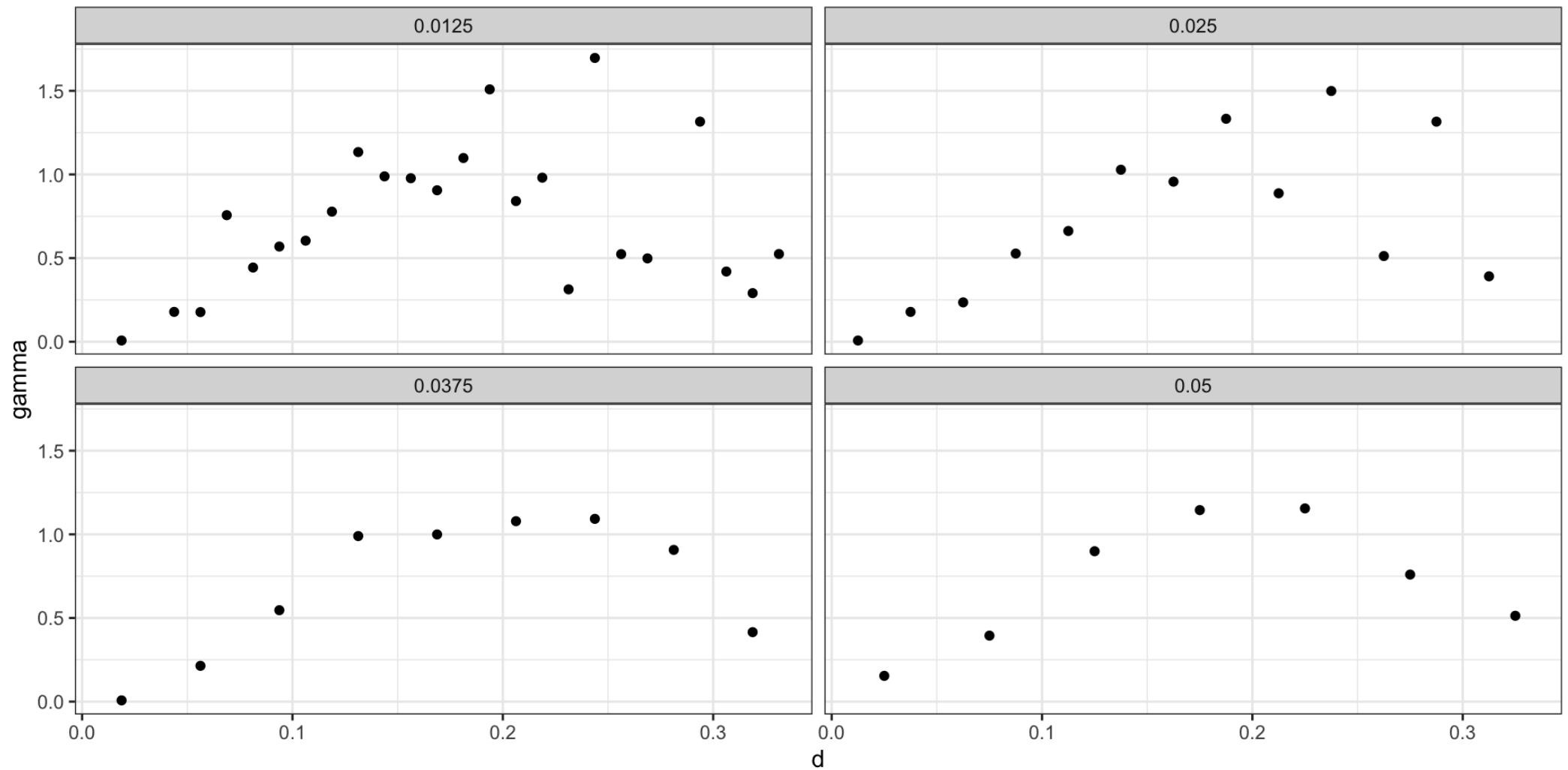


From last time

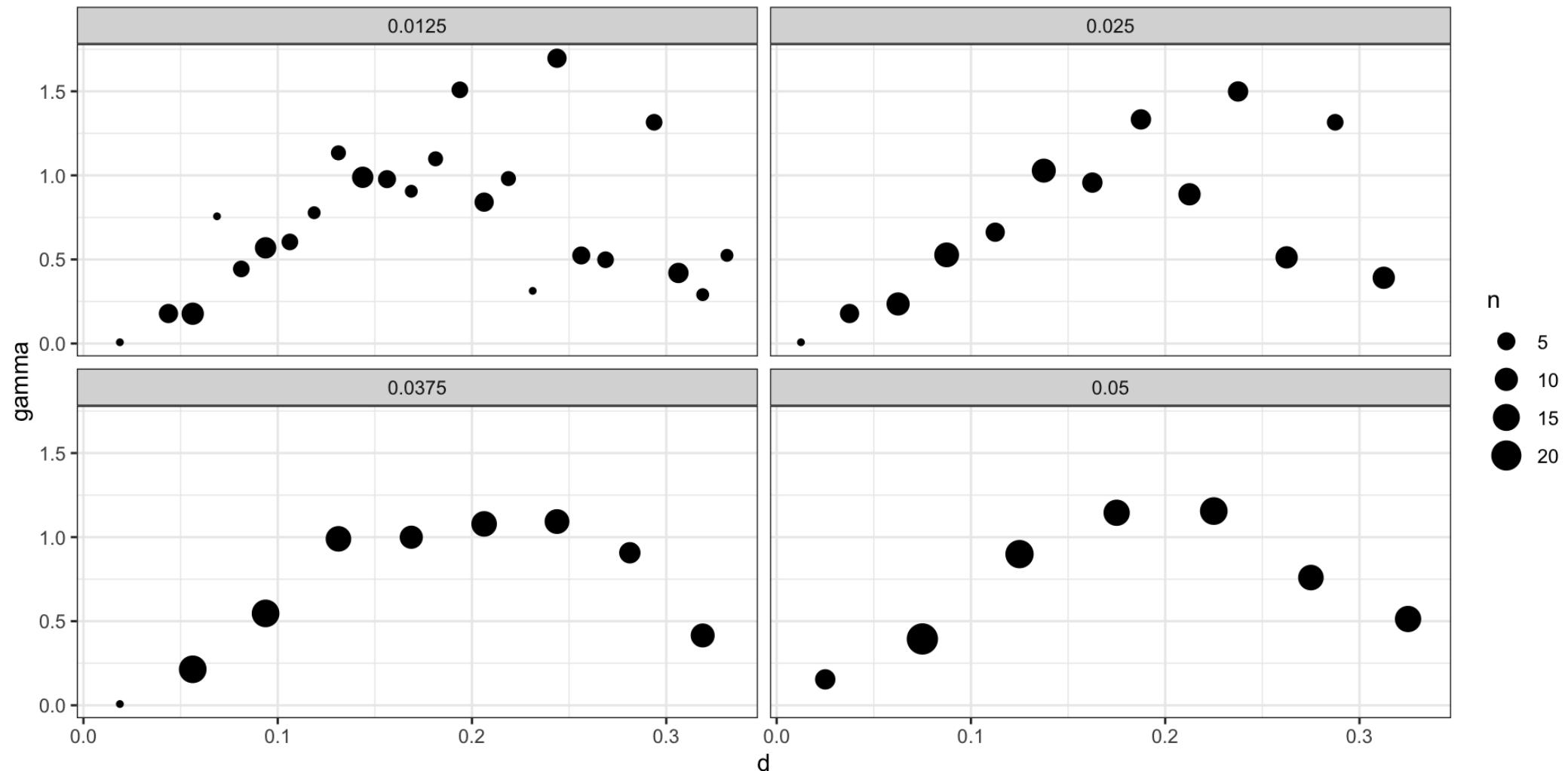


Empirical semivariogram - no bins / cloud

Empirical semivariogram (binned)



Empirical semivariogram (binned w/ size)



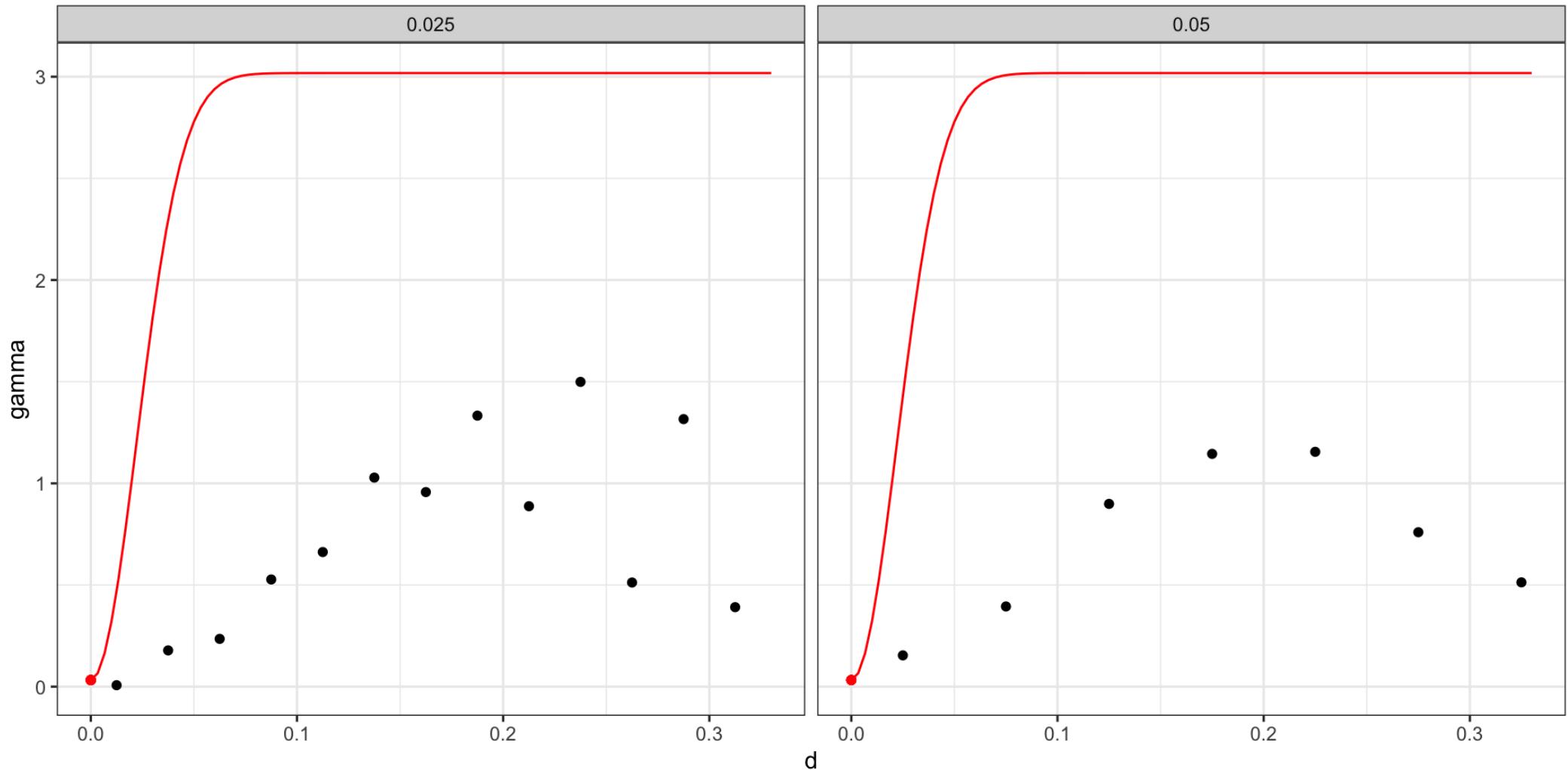
Theoretical vs empirical semivariogram

After fitting the model last time we came up with a posterior mean of $\sigma^2 = 2.99$, $l = 31.78$, and $\sigma_w^2 = 0.03$ for a square exponential covariance.

$$\text{Cov}(d) = \sigma^2 \exp(-d/l)^2 + \sigma_w^2 \mathbf{1}_{h=0}$$

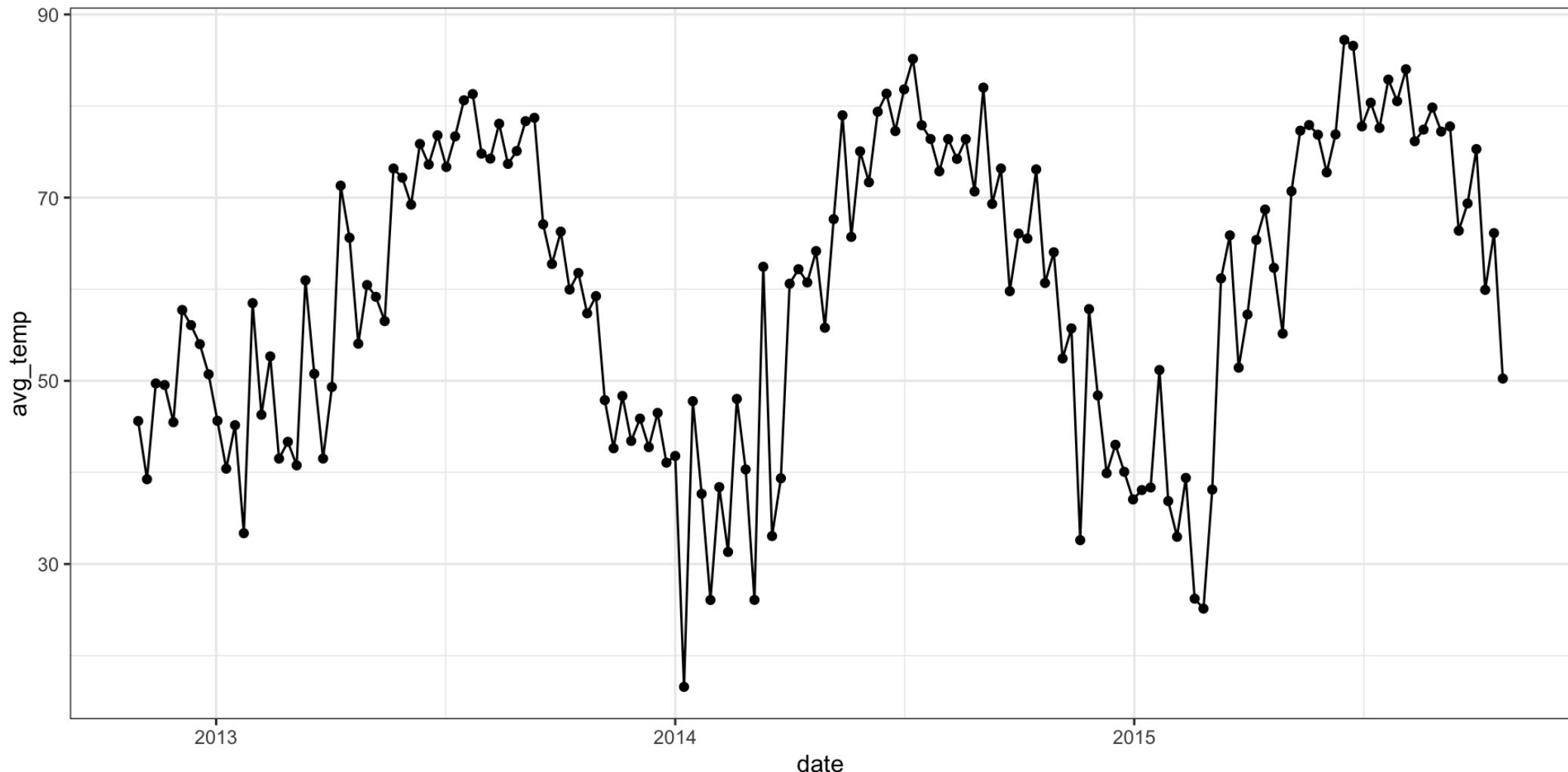
$$\gamma(h) = (\sigma^2 + \sigma_w^2) - \sigma^2 \exp(-h/l)^2$$

$$= (2.9850975 + 0.032637) - 2.9850975 \exp(-(31.7775906 h)^2)$$

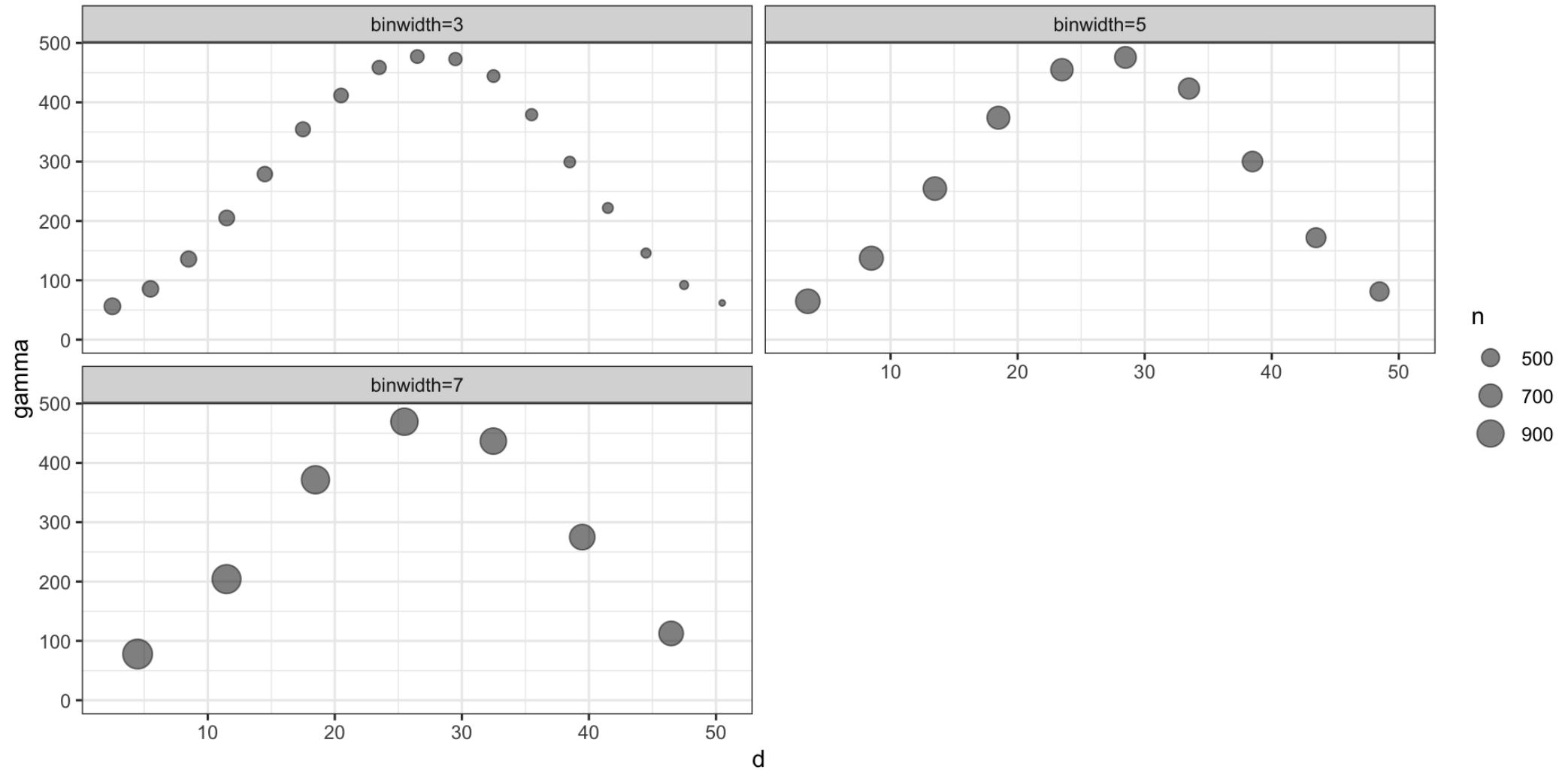


Durham Average Daily Temperature

Temp Data



Empirical semivariogram



Model

What does the model we are trying to fit actually look like?

$$y(t) = \mu(t) + w(t) + \epsilon(t)$$

where

$$\mu(t) = \beta_0$$

$$w(t) \sim \mathcal{N}(0, \Sigma)$$

$$\epsilon(t) \sim \mathcal{N}(0, \sigma_w^2)$$

$$\{\Sigma\}_{ij} = \text{Cov}(t_i, t_j) = \sigma^2 \exp(-(|t_i - t_j|/l)^2)$$

BRMS Model

```
1 library(brms)
2 ( m = brm(
3   avg_temp ~ 1 + gp(week), data=temp,
4   cores = 4, refresh=0
5 ) )
```

Family: gaussian

Links: mu = identity; sigma = identity

Formula: avg_temp ~ 1 + gp(week)

Data: temp (Number of observations: 156)

Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
total post-warmup draws = 4000

Gaussian Process Terms:

	Estimate	Est.Error	1-95% CI	u-95% CI	Rhat	Bulk_ESS	Tail_ESS
sdgp(gpweek)	10.64	5.41	4.20	16.01	3.82	4	15
lscale(gpweek)	1.26	1.20	0.06	2.69	4.34	4	NA

	Estimate	Est.Error	1-95% CI	u-95% CI	Rhat	Bulk_ESS	Tail_ESS
Intercept	6.32	7.27	-1.28	16.12	3.13	5	30

BRMS Alternatives

The BRMS model (and hence Stan) took between 5-10 minutes (per chain) to attempt to fit the model and failed spectacularly.

We could potentially improve things by tweaking the priors and increasing iterations but this won't solve the slowness issue.

The stop gap work around - using `spBayes`

- Interface is old and clunky (inputs and outputs)
- Designed for spatial GPs
- Super fast (~10 seconds for 20k iterations)
- `dukestm` has a wrapper, called `gplm()` to make the interface / usage not as terrible

Fitting the model

```
1 m = gplm(  
2   avg_temp~1,  
3   data = temp, coords = "week",  
4   cov_model = "gaussian",  
5   starting=list(  
6     "phi"=sqrt(3)/4, "sigma.sq"=1, "tau.sq"=1  
7   ),  
8   priors=list(  
9     "phi.unif"=c(sqrt(3)/52, sqrt(3)/1),  
10    "sigma.sq.ig"=c(2, 1),  
11    "tau.sq.ig"=c(2, 1)  
12  ),  
13  thin=10  
14 )
```

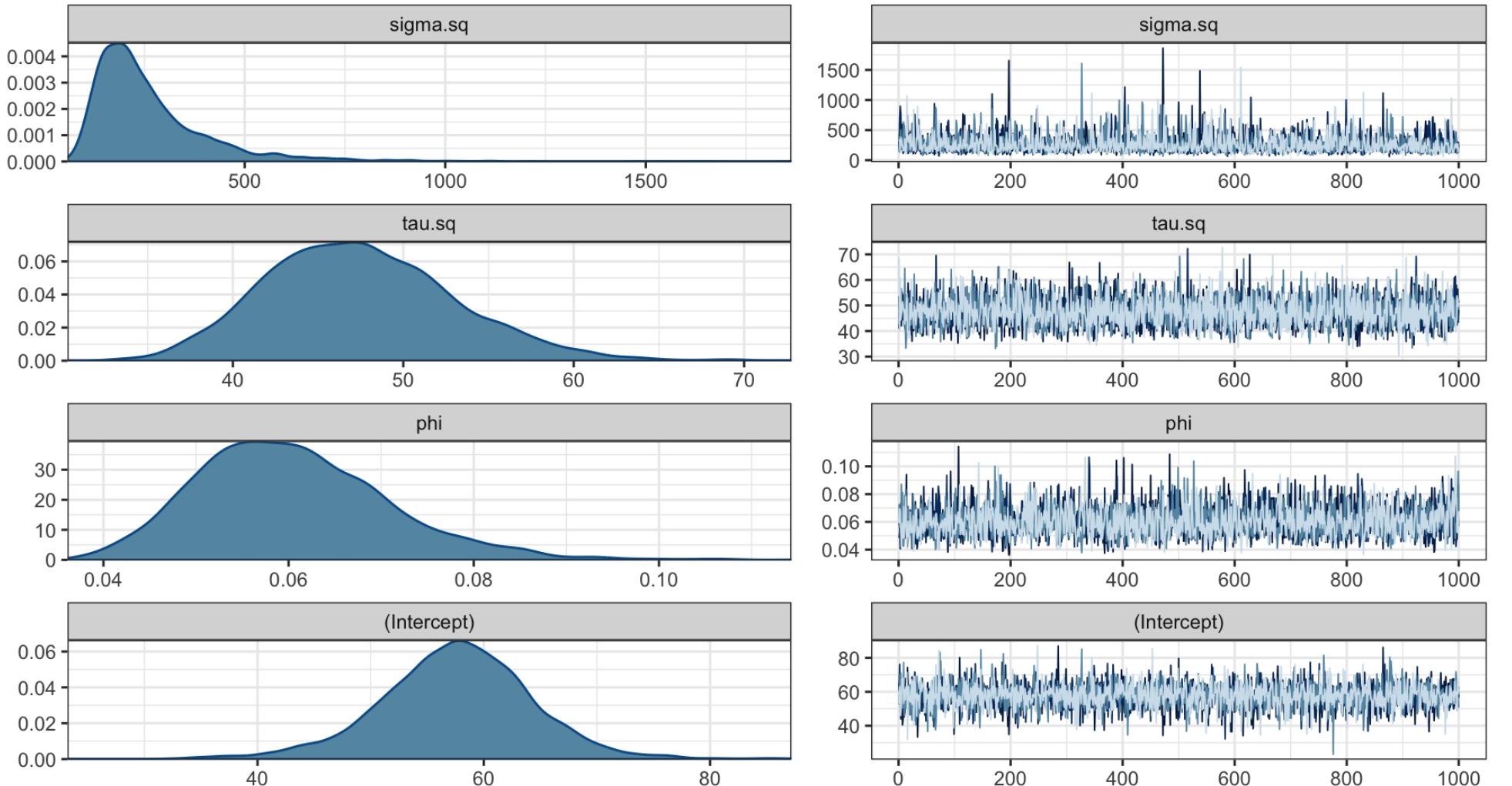
Model results

```
1 m
```

```
# A gplm model (spBayes spLM) with 4 chains, 4 variables, and 4000
iterations.
# A tibble: 4 × 10
  variable      mean    median      sd     mad     q5     q95   rhat
  <chr>        <num>    <num>    <num>    <num>    <num>    <num>  <num>
  ess_bulk
<num>
1 sigma.sq    263.     221.     150.     98.2    1.15e+2 5.54e+2 1.00
3075.
2 tau.sq      47.5     47.1     5.54     5.50    3.91e+1 5.72e+1 1.00
3704.
3 phi         0.0605   0.0595   0.0104   0.00997 4.55e-2 7.92e-2 1.00
3097.
4 (Intercept) 57.6     57.7     6.80     6.15    4.63e+1 6.81e+1 1.00
3955.
# i 1 more variable: ess_tail <num>
```

Parameter posteriors

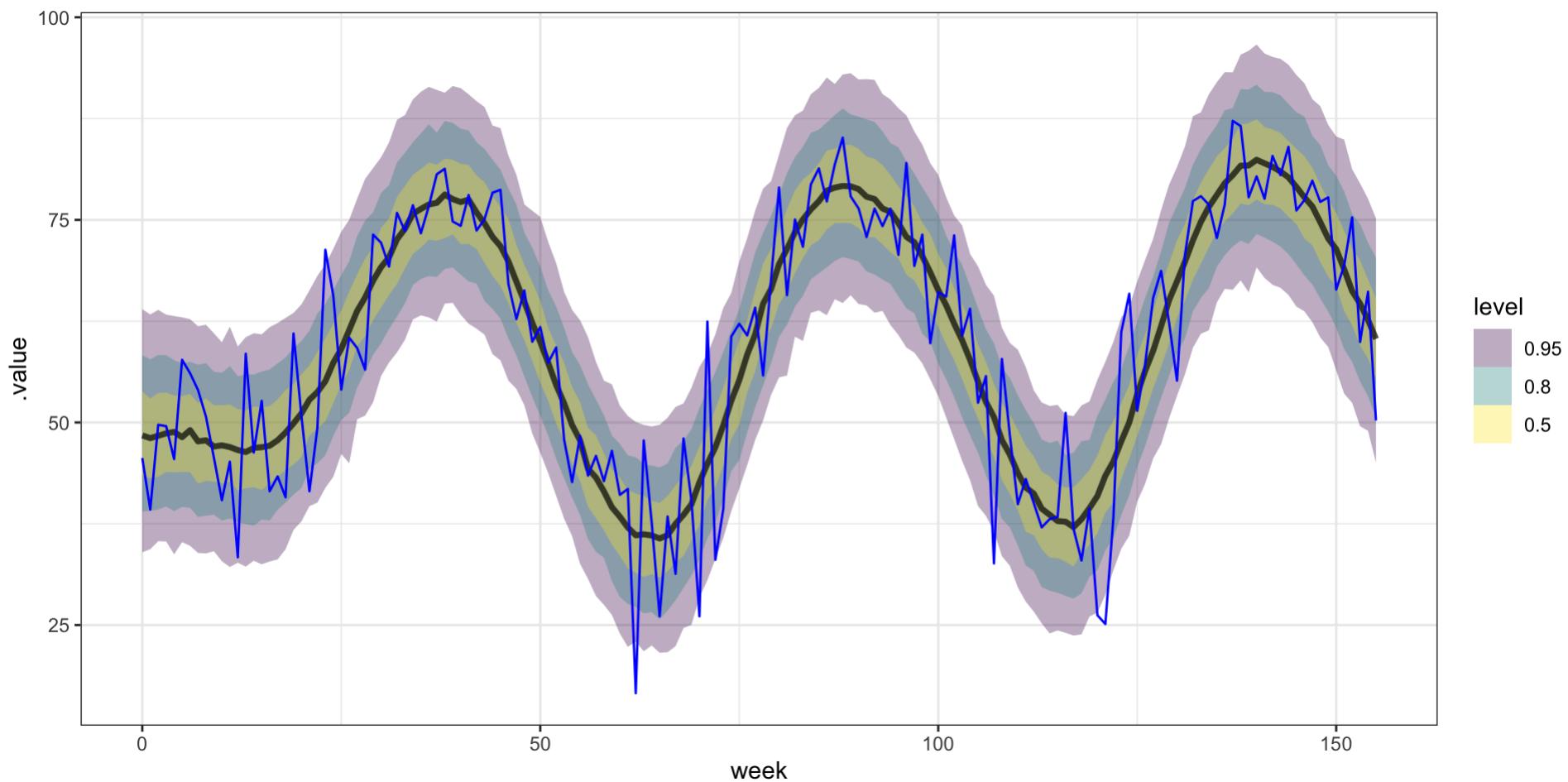
```
1 plot(m)
```



Fitted model

```
1 predict(m, newdata = tibble(week=1:(3*52)-1+1e-6), coords = "week") |>
2   tidybayes::gather_draws(y[i]) |>
3   mutate(week = i-1) |>
4   filter(.chain == 1, week <= 3*52) |>
5   ggplot(aes(x=week, y=.value)) +
6     tidybayes::stat_lineribbon(alpha=0.33) +
7     geom_line(data=temp, aes(y=avg_temp), color="blue")
```

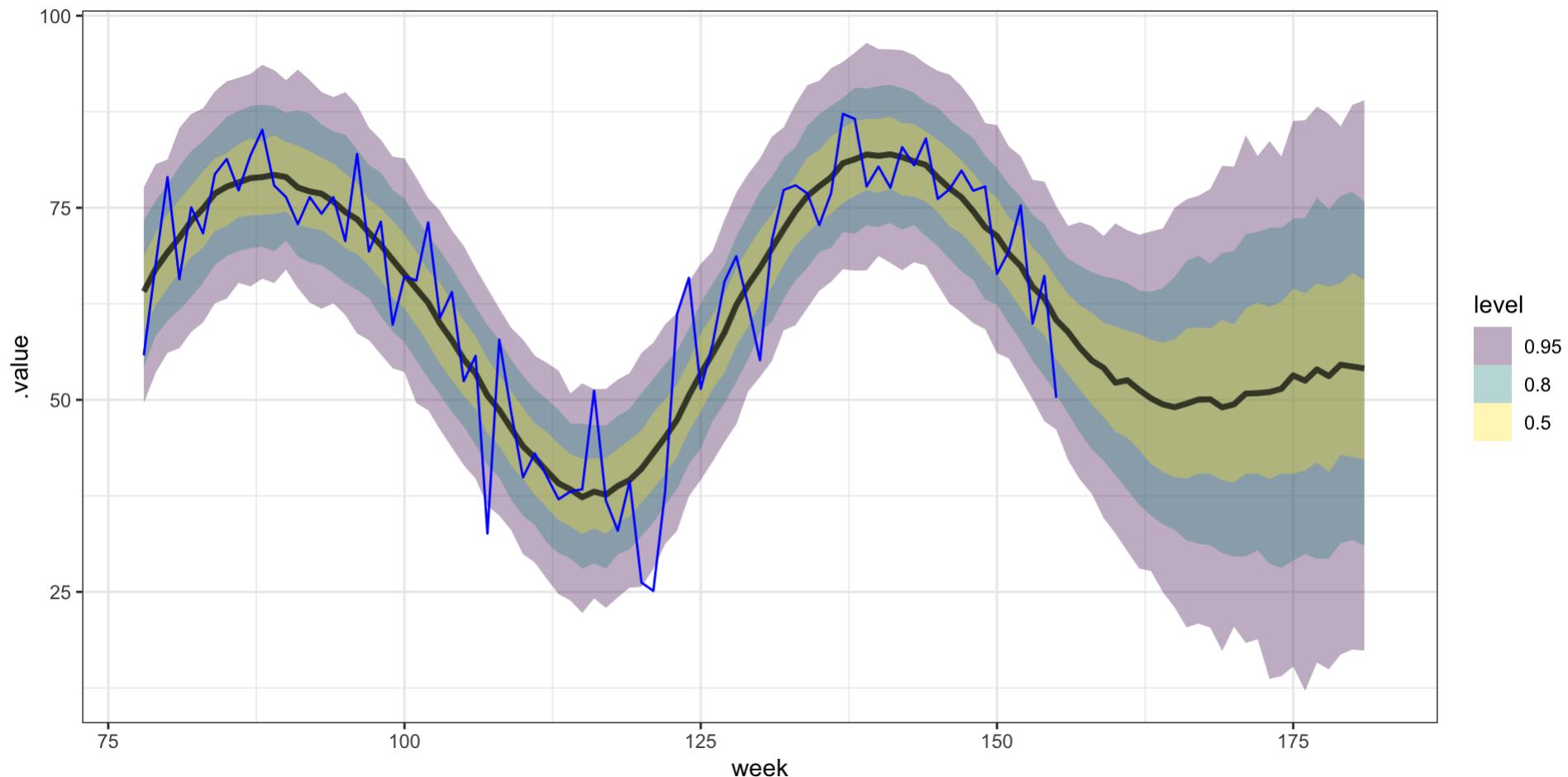
Fitted model



Forecasting

```
1 predict(m, newdata = tibble(week=1:(3.5*52)-1+1e-6), coords = "week") |>
2   tidybayes::gather_draws(y[i]) |>
3   mutate(week = i-1) |>
4   filter(.chain == 1) |>
5   ggplot(aes(x=week, y=.value)) +
6     tidybayes::stat_lineribbon(alpha=0.33) +
7     geom_line(data=temp, aes(y=avg_temp), color="blue") +
8     xlim(1.5*52, 3.5*52)
```

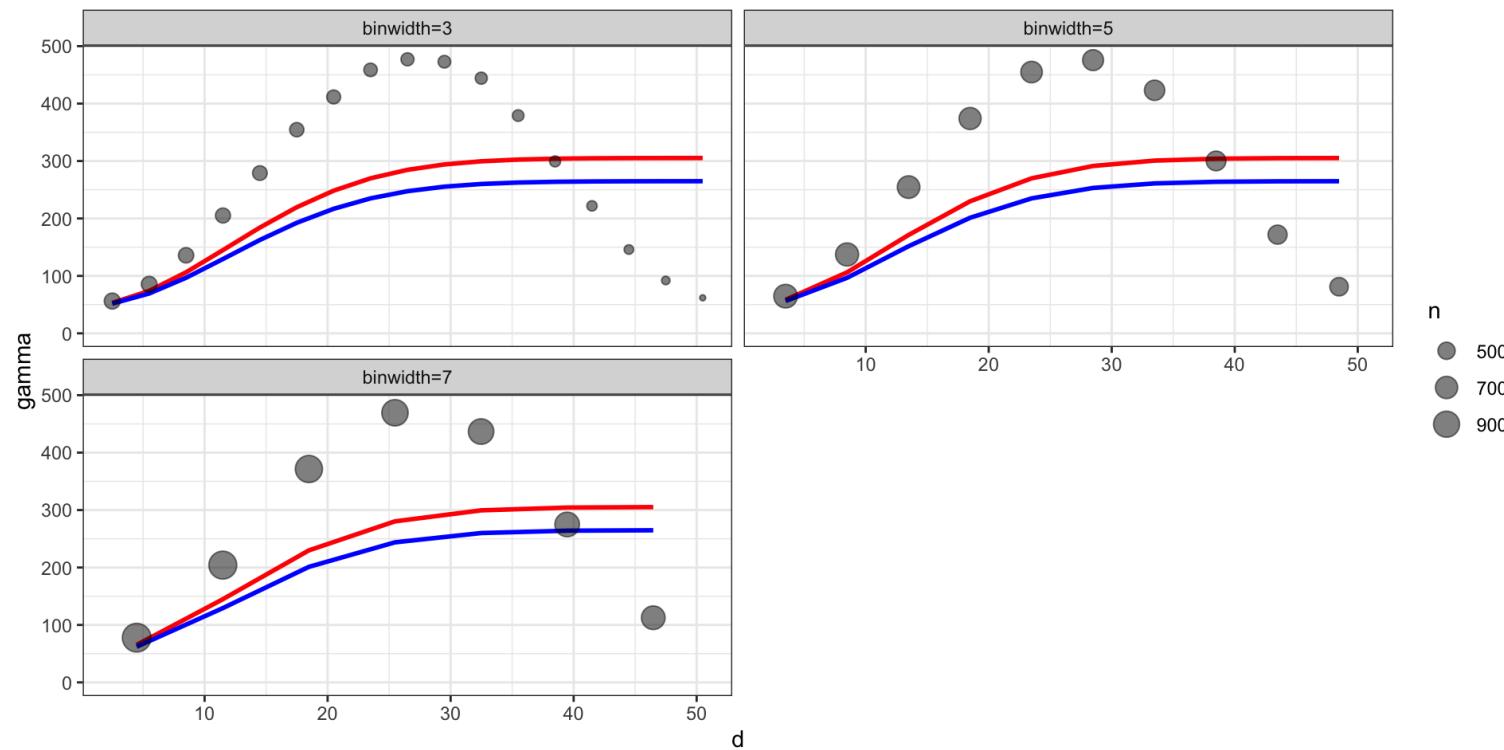
Forecasting



Empirical semivariogram vs. model

From the model summary we have the following,

- *posterior means*: $\sigma^2 = 258$, $\sigma_w^2 = 47.3$, $l = 0.06$
- *posterior medians*: $\sigma^2 = 218$, $\sigma_w^2 = 46.9$, $l = 0.06$



Exponential model

Fitting the model

```
1 m = gplm(  
2   avg_temp~1,  
3   data = temp, coords = "week",  
4   cov_model = "exponential",  
5   starting=list(  
6     "phi"=3/4, "sigma.sq"=1, "tau.sq"=1  
7   ),  
8   priors=list(  
9     "phi.unif"=c(3/52, 3/1),  
10    "sigma.sq.ig"=c(2, 1),  
11    "tau.sq.ig"=c(2, 1)  
12  ),  
13  thin=10  
14 )
```

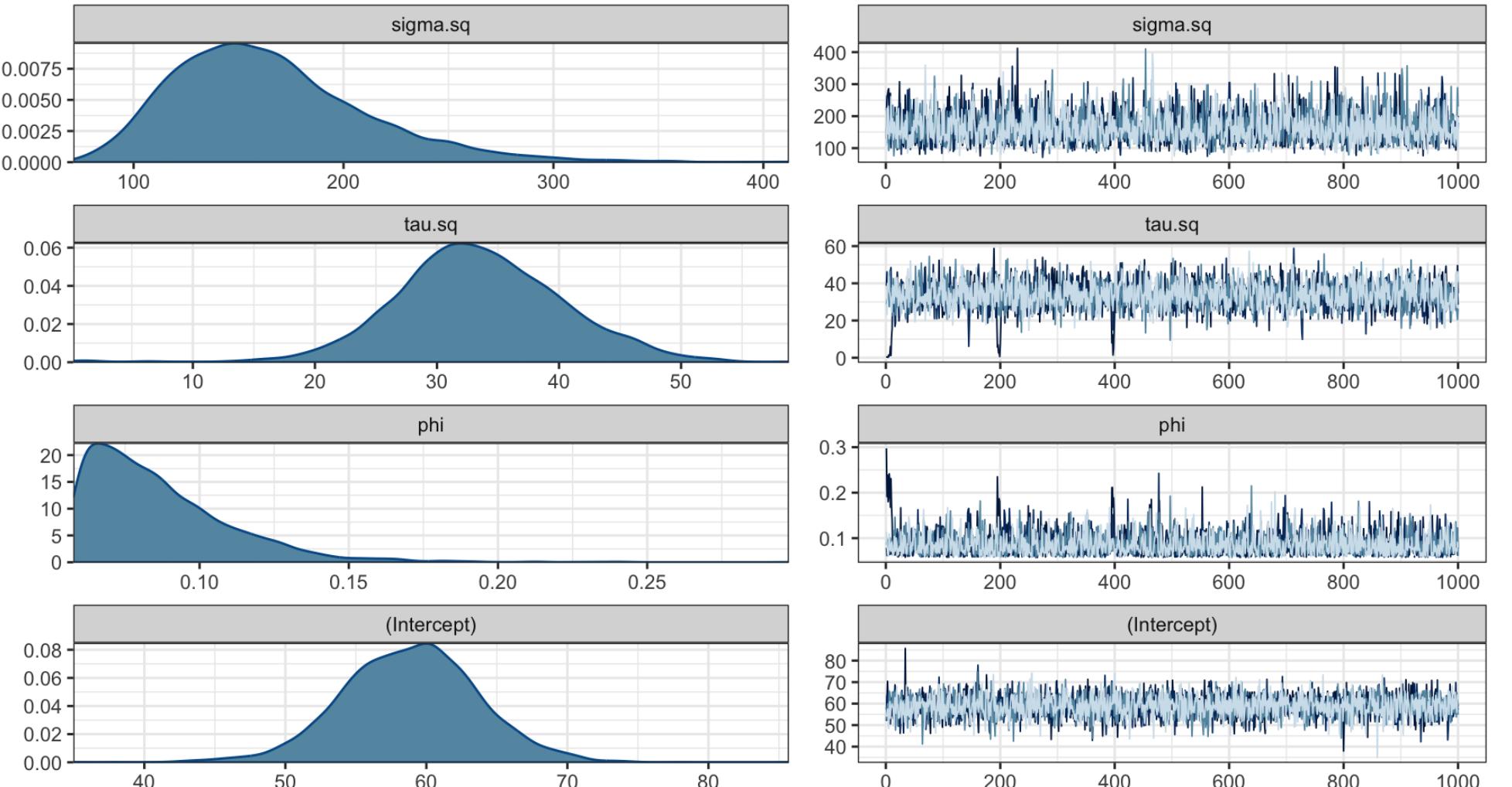
Model results

```
1 m
```

```
# A gplm model (spBayes spLM) with 4 chains, 4 variables, and 4000
iterations.
# A tibble: 4 × 10
  variable      mean    median      sd     mad      q5     q95   rhat
  <chr>        <num>    <num>    <num>    <num>    <num>    <num> <num>
  ess_bulk
  <num>
  1 sigma.sq    165.     159.     46.1     42.7     103.     251.    1.00
  2871.
  2 tau.sq      33.5     33.2     6.97     6.48     22.7     45.3    1.00
  2766.
  3 phi         0.0865   0.0809   0.0244   0.0210   0.0597   0.132    1.00
  2858.
  4 (Intercept) 58.8     59.0     4.80     4.68     51.2     66.7    1.00
  3910.
# i 1 more variable: ess_tail <num>
```

Parameter posteriors

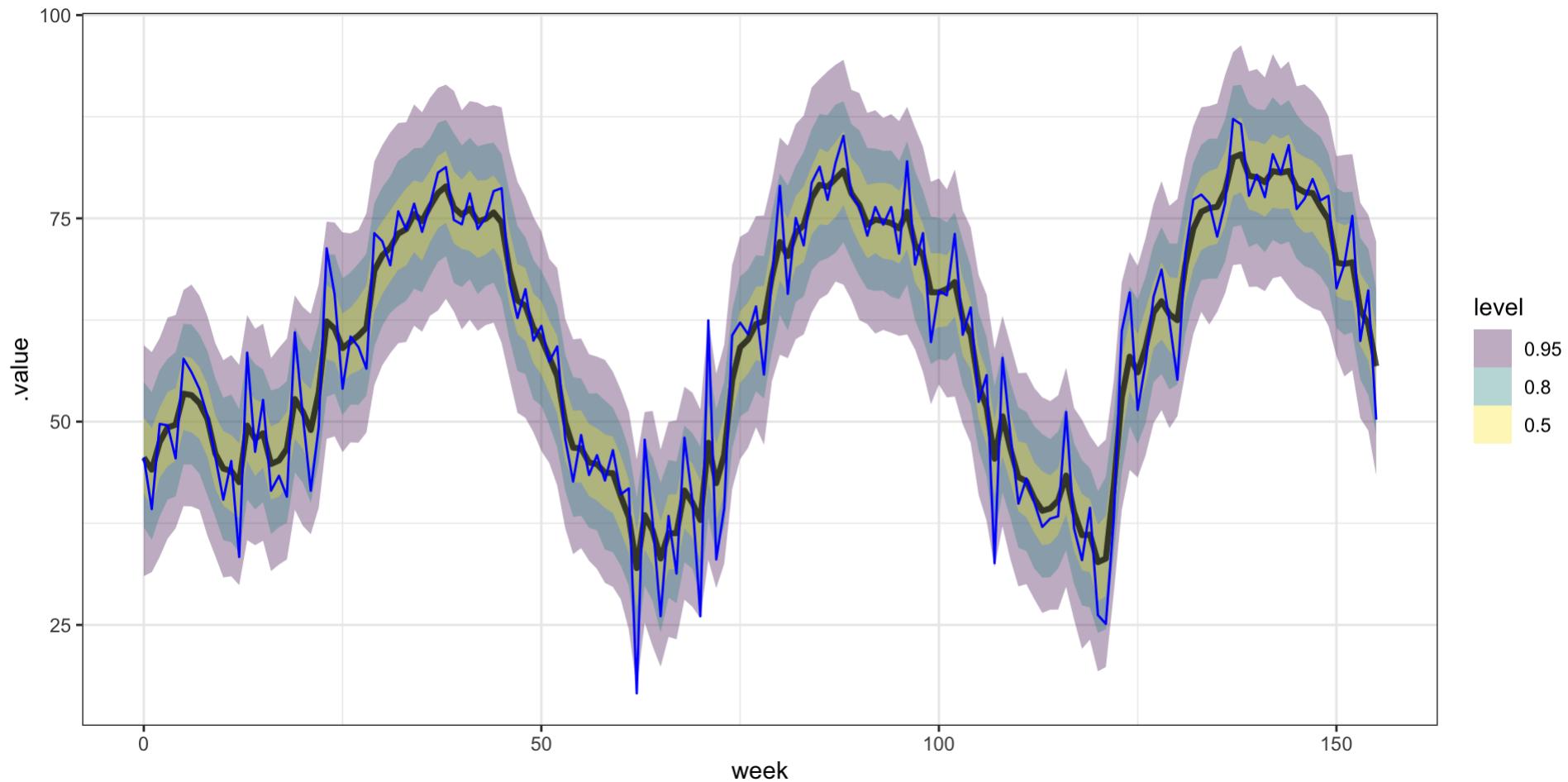
```
1 plot(m)
```



Fitted model

```
1 predict(m, newdata = tibble(week=1:(3*52)-1+1e-6), coords = "week") |>
2   tidybayes::gather_draws(y[i]) |>
3   mutate(week = i-1) |>
4   filter(.chain == 1, week <= 3*52) |>
5   ggplot(aes(x=week, y=.value)) +
6     tidybayes::stat_lineribbon(alpha=0.33) +
7     geom_line(data=temp, aes(y=avg_temp), color="blue")
```

Fitted model



Forecasting

```
1 predict(m, newdata = tibble(week=1:(3.5*52)-1+1e-6), coords = "week") |>
2   tidybayes::gather_draws(y[i]) |>
3   mutate(week = i-1) |>
4   filter(.chain == 1) |>
5   ggplot(aes(x=week, y=.value)) +
6     tidybayes::stat_lineribbon(alpha=0.33) +
7     geom_line(data=temp, aes(y=avg_temp), color="blue") +
8     xlim(1.5*52, 3.5*52)
```

Forecasting

