

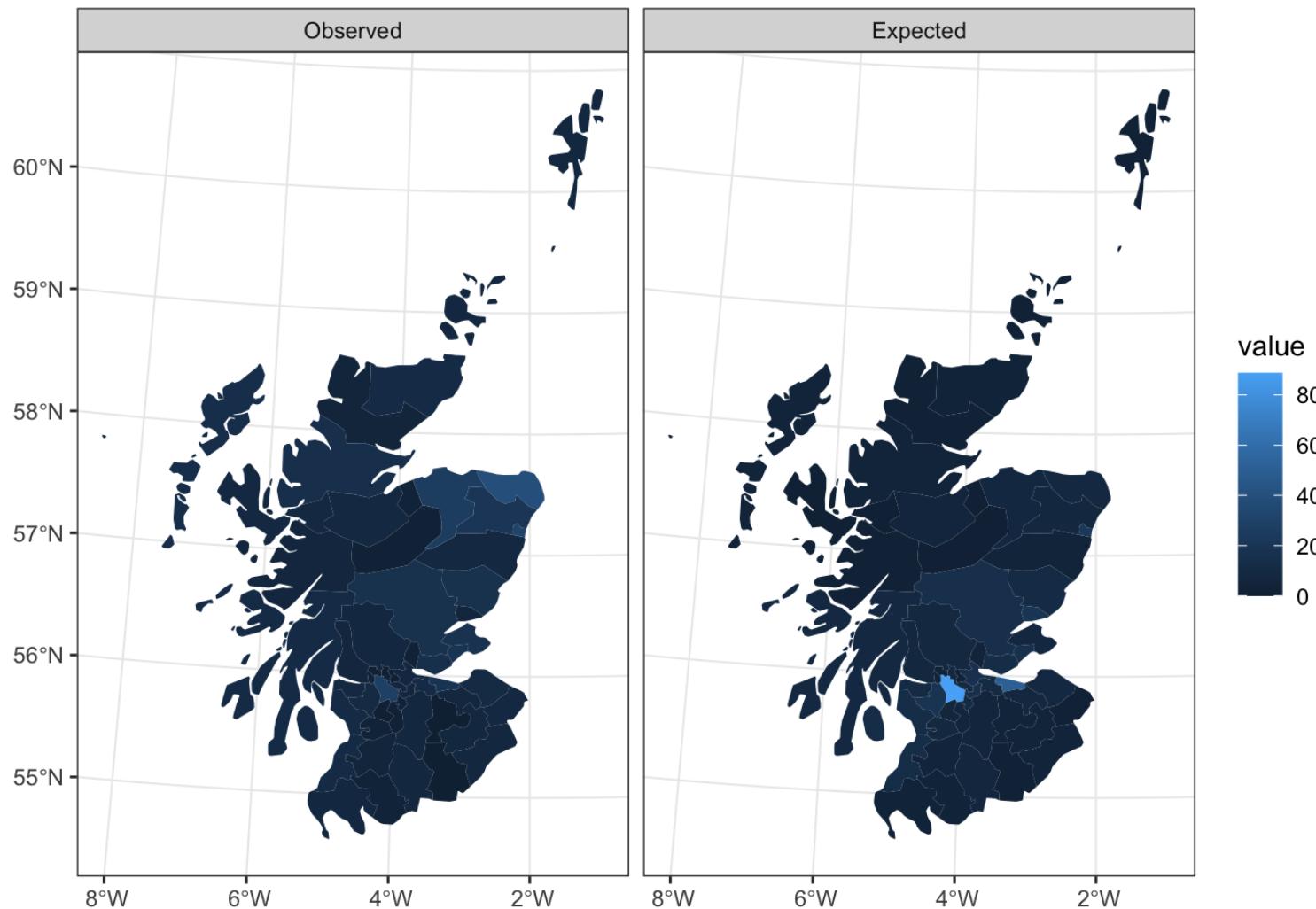
Spatial GLM + Point Reference Spatial Data

Lecture 21

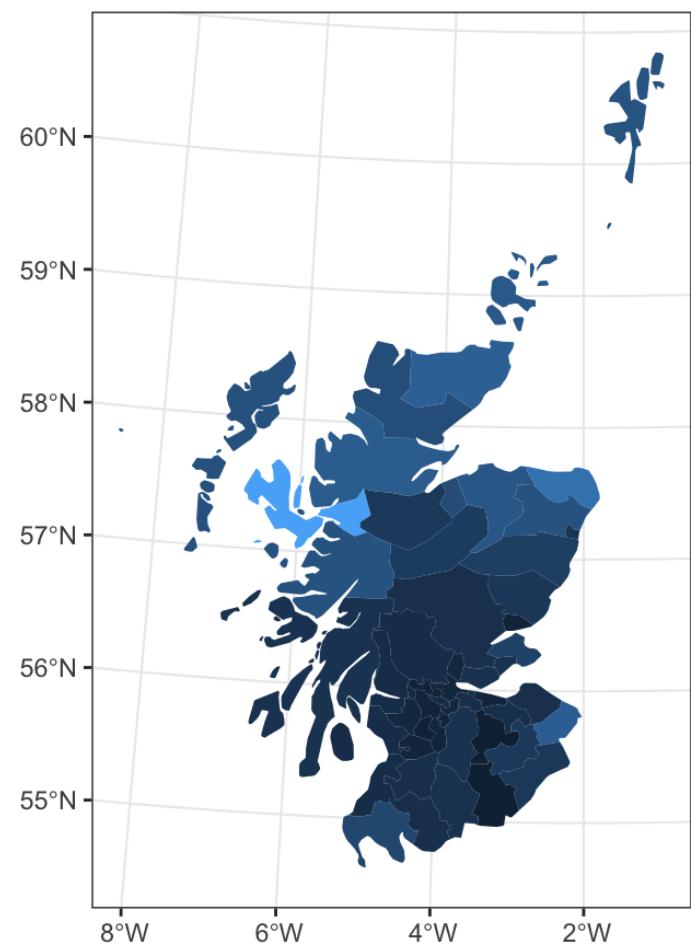
Dr. Colin Rundel

Spatial GLM Models

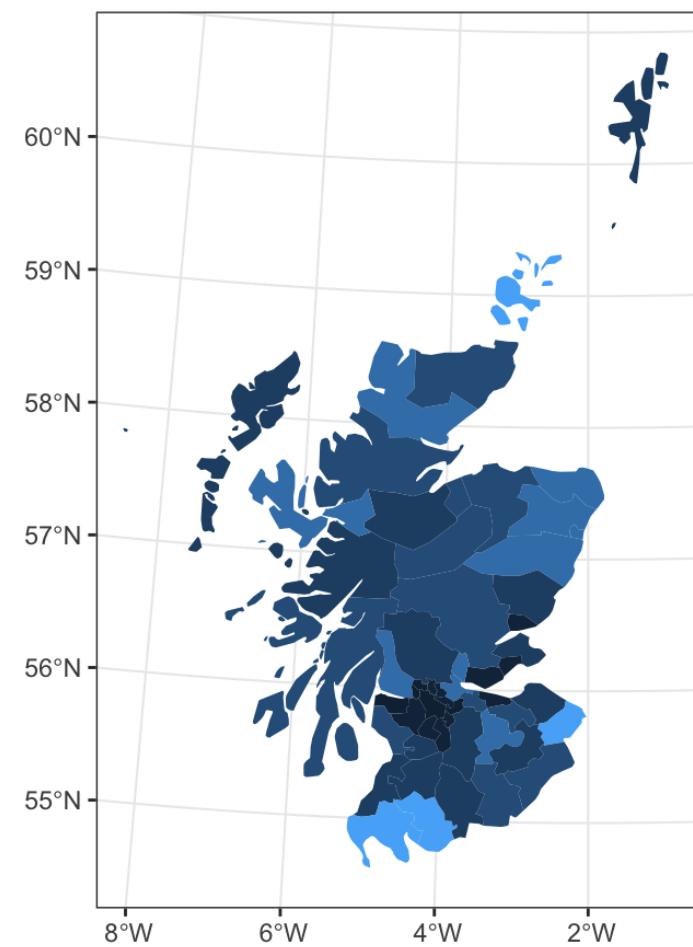
Scottish Lip Cancer Data



Obs/Exp

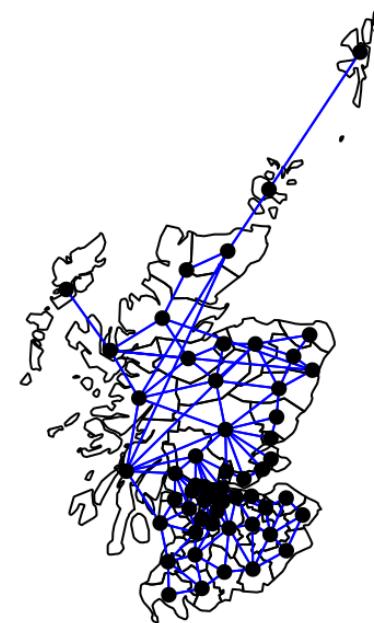


% Agg Fish Forest



Neighborhood / weight matrix

```
1 A = (st_distance(lip_cancer) |> unclass()) < 1e-6  
2 listw = spdep::mat2listw(A)
```



Moran's I

```
1 spdep::moran.test(lip_cancer$Observed, listw)
```

Moran I test under randomisation

```
data: lip_cancer$Observed  
weights: listw
```

```
Moran I statistic standard deviate = 4.5416,  
p-value = 2.792e-06  
alternative hypothesis: greater  
sample estimates:  
Moran I statistic      Expectation  
0.311975396          -0.018181818  
Variance  
0.005284831
```

```
1 spdep::moran.test(lip_cancer$Observed / lip_c
```

Moran I test under randomisation

```
data: lip_cancer$Observed/lip_cancer$Expected  
weights: listw
```

```
Moran I statistic standard deviate = 8.2916,  
p-value < 2.2e-16  
alternative hypothesis: greater  
sample estimates:  
Moran I statistic      Expectation  
0.589795225          -0.018181818  
Variance  
0.005376506
```

GLM

```
1 l = glm(Observed ~ offset(log(Expected)) + pcaff,  
2         family="poisson", data=lip_cancer)  
3 summary(l)
```

Call:

```
glm(formula = Observed ~ offset(log(Expected)) + pcaff, family = "poisson",  
    data = lip_cancer)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-4.7632	-1.2156	0.0967	1.3362	4.7130

Coefficients:

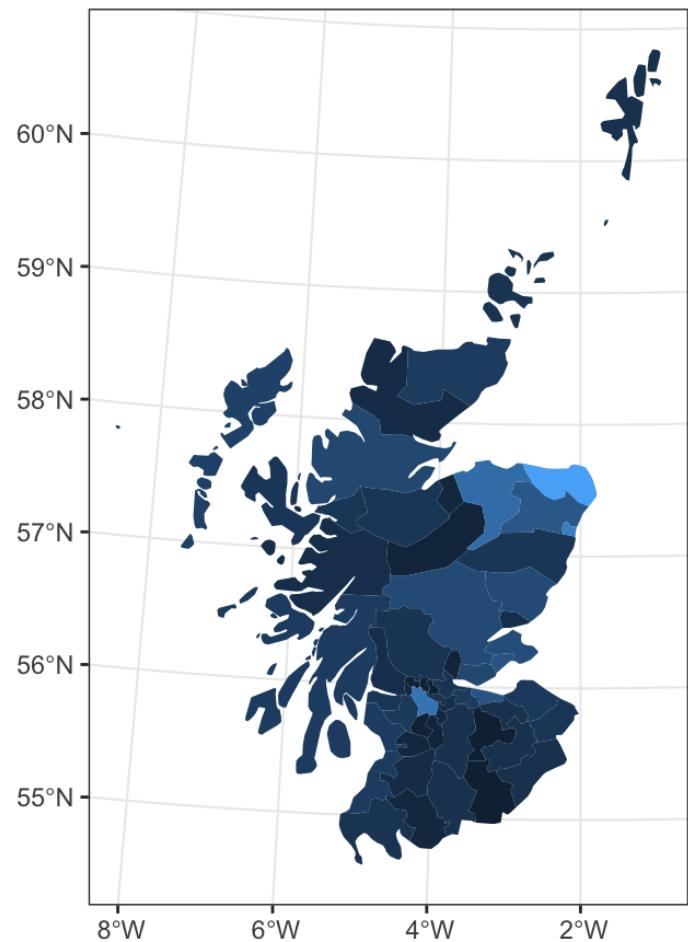
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.542268	0.069525	-7.80	6.21e-15
pcaff	0.073732	0.005956	12.38	< 2e-16

(Intercept) ***

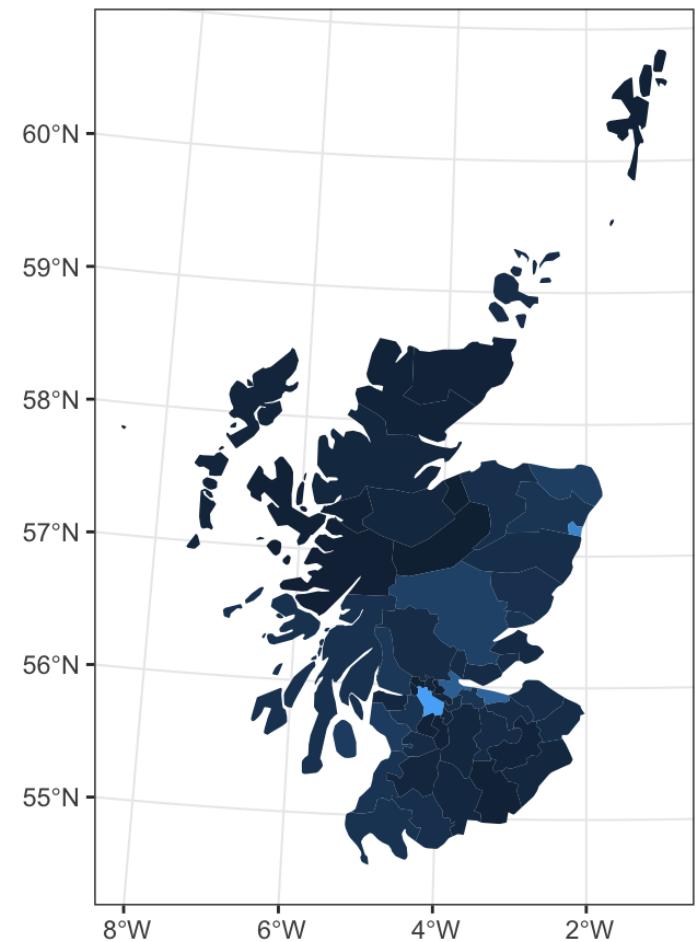
pcaff ***

GLM Fit

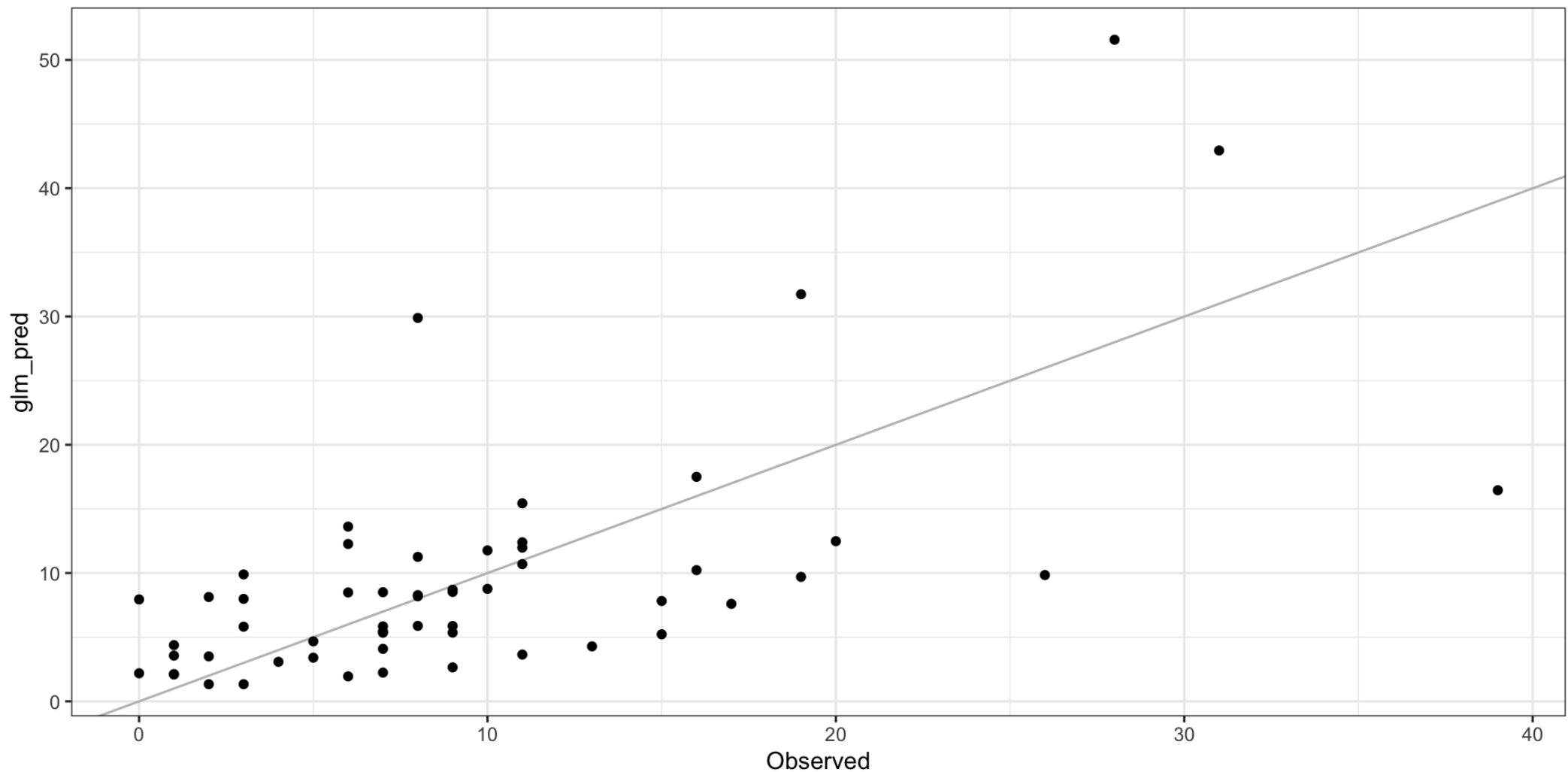
Observed Cases



GLM Predicted Cases

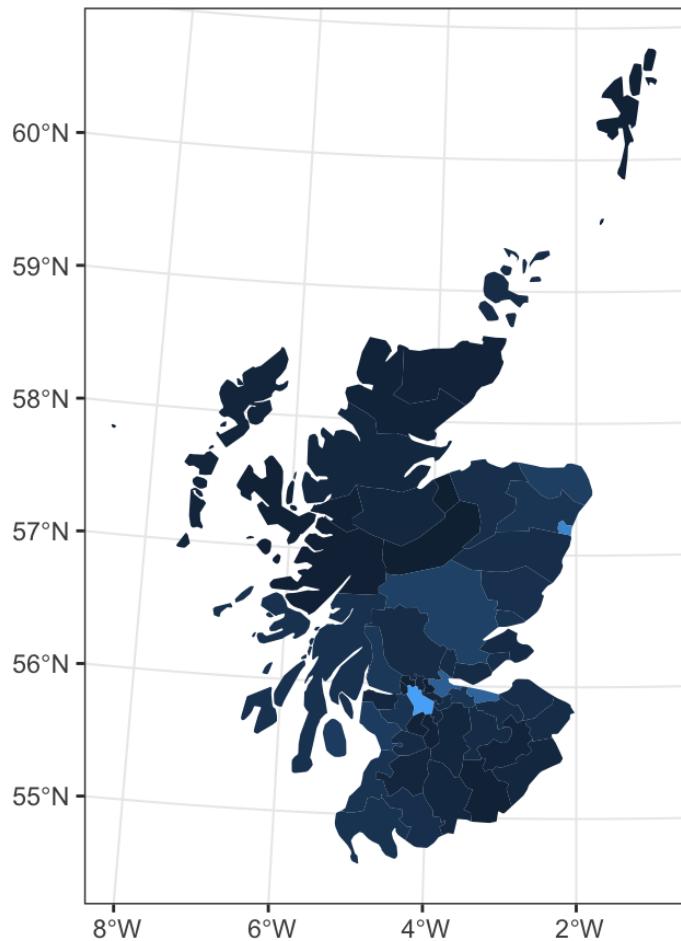


GLM Fit

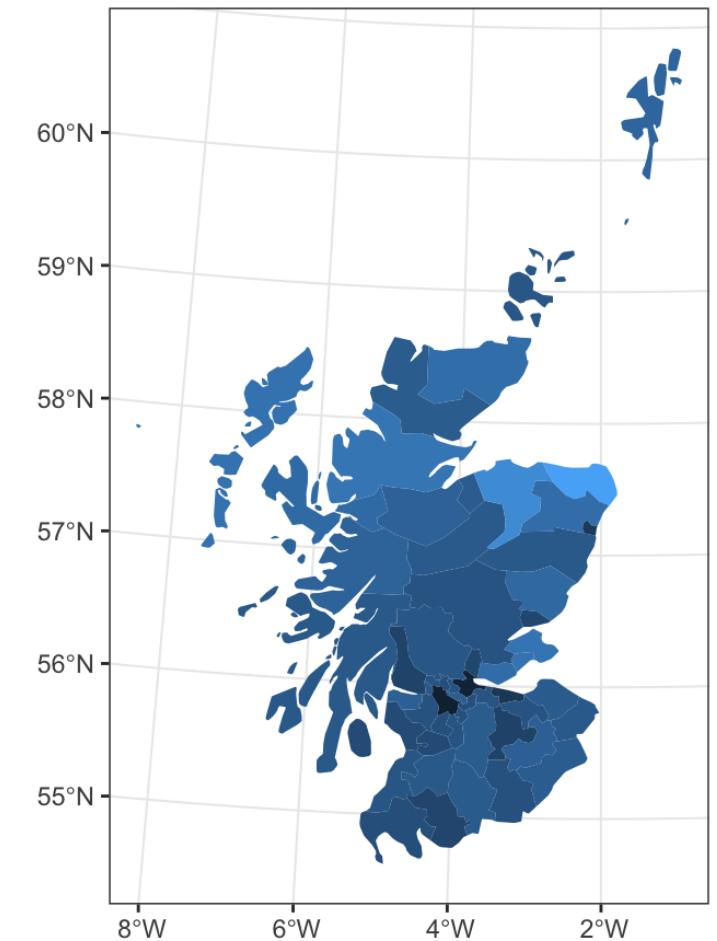


GLM Residuals

GLM Predicted Cases



GLM Residuals



Model Results

```
1 #RMSE  
2 yardstick::rmse_vec(lip_cancer$Observed, lip_cancer$glm_pred)
```

```
[1] 7.480889
```

```
1 #Moran's I  
2 spdep::moran.test(lip_cancer$glm_resid, listw)
```

Moran I test under randomisation

```
data: lip_cancer$glm_resid  
weights: listw  
  
Moran I statistic standard deviate = 4.8186,  
p-value = 7.228e-07  
alternative hypothesis: greater  
sample estimates:  
Moran I statistic      Expectation  
0.333403223      -0.018181818  
Variance  
0.005323717
```

A hierarchical model for lip cancer

We have observed counts of lip cancer for 56 districts in Scotland. Let y_i represent the number of lip cancer for district i .

$$y_i \sim \text{Poisson}(\lambda_i)$$

$$\log(\lambda_i) = \log(E_i) + x_i\beta + \omega_i$$

$$\omega \sim (\mathbf{0}, \sigma^2(D - \phi A)^{-1})$$

where E_i is the expected counts for each region (and serves as an offset).

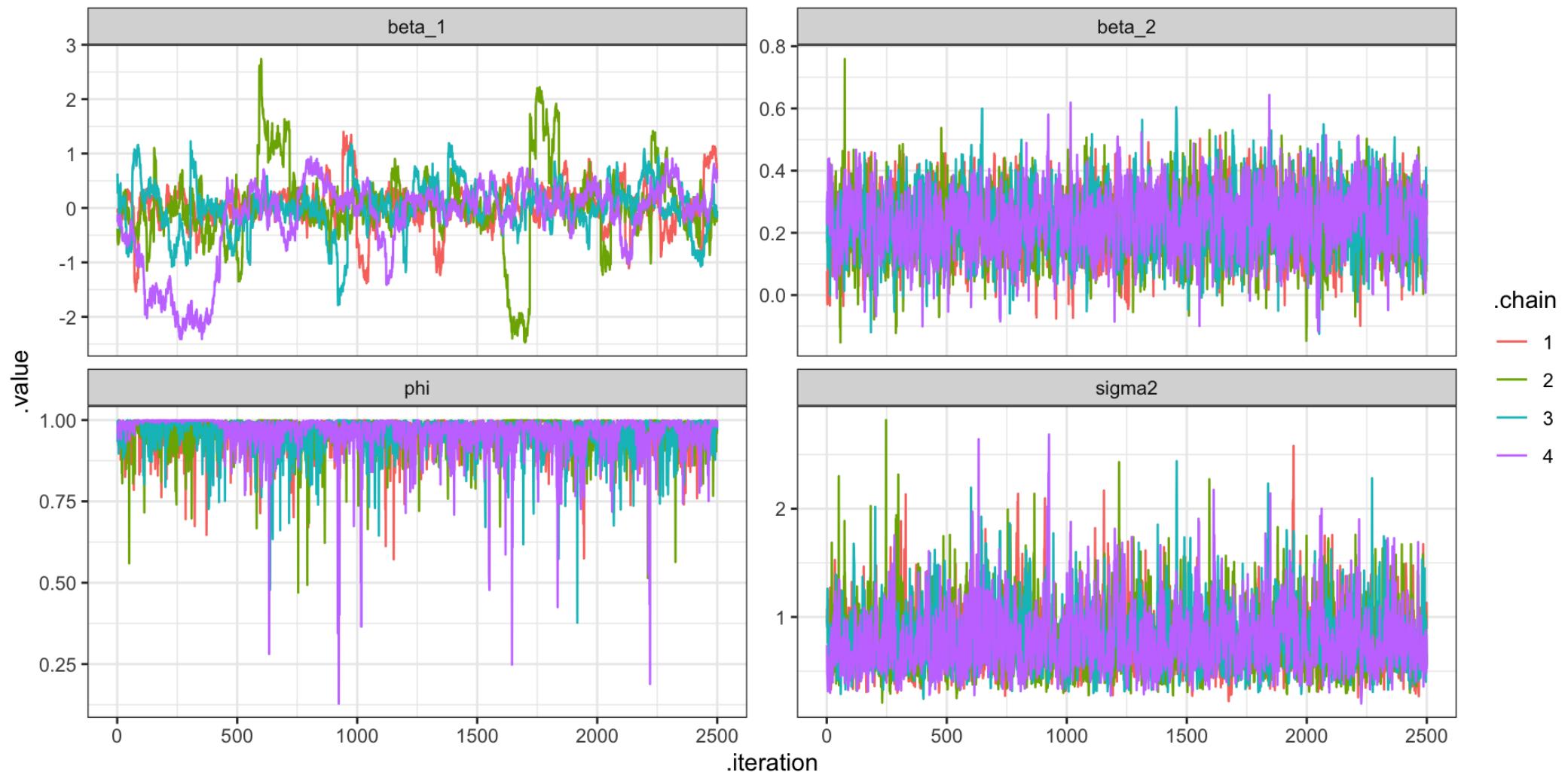
Data prep & CAR model

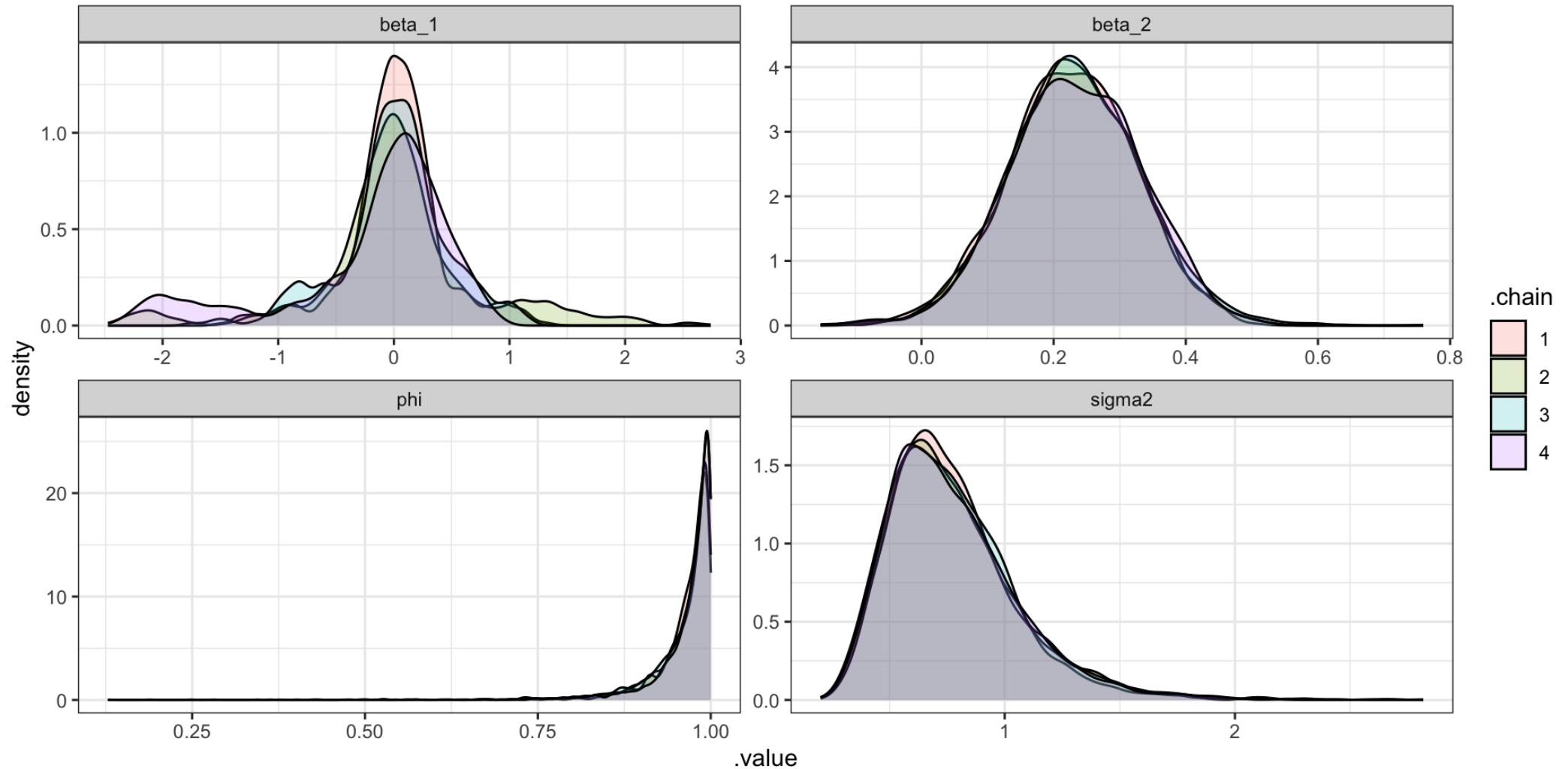
```
1 X = model.matrix(~scale(lip_cancer$pcaff))
2 offset = lip_cancer$Expected
3 y = lip_cancer$Observed
```

```
1 car_model = "
2 data {
3     int<lower=0> N;
4     int<lower=0> p;
5     int<lower=0> y[N];
6     matrix[N,N] A;
7     matrix[N,p] X;
8     vector[N] offset;
9 }
10 transformed data {
11     vector[N] nb = A * rep_vector(1, N);
12     matrix[N,N] D = diag_matrix(nb);
13 }
14 parameters {
15     vector[N] w_s;
16     vector[p] beta;
17     real<lower=0> sigma2;
18     real<lower=0,upper=1> phi;
19 }
20 transformed parameters {
21     vector[N] eta = log(offset) + X * beta + w_s;
```

CAR Fitting

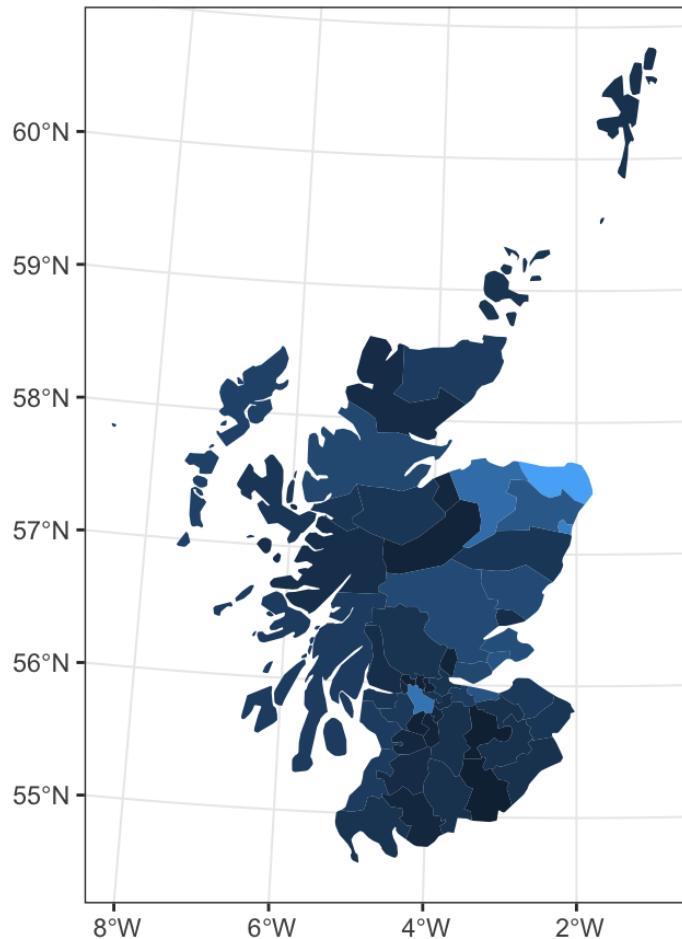
```
1 car = rstan::stan_model(model_code = car_model)
2
3 car_m = rstan::sampling(
4   car, iter=5000, cores=4,
5   data = list(N=nrow(X), A=A, X=X, p=ncol(X), offset=offset, y=y)
6 )
```



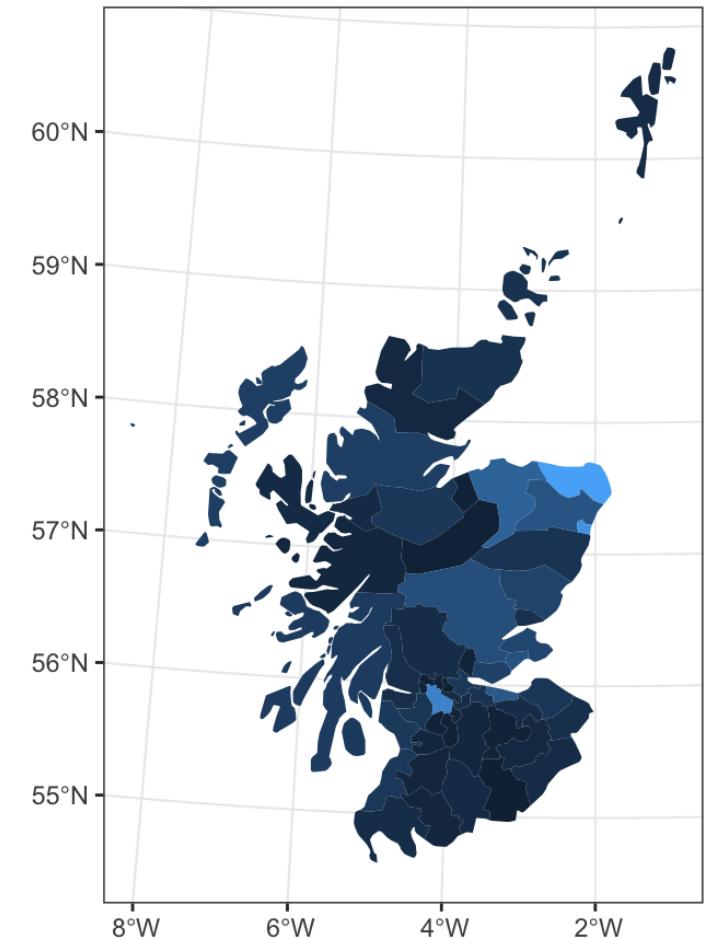


CAR Predictions ($\hat{\lambda}$)

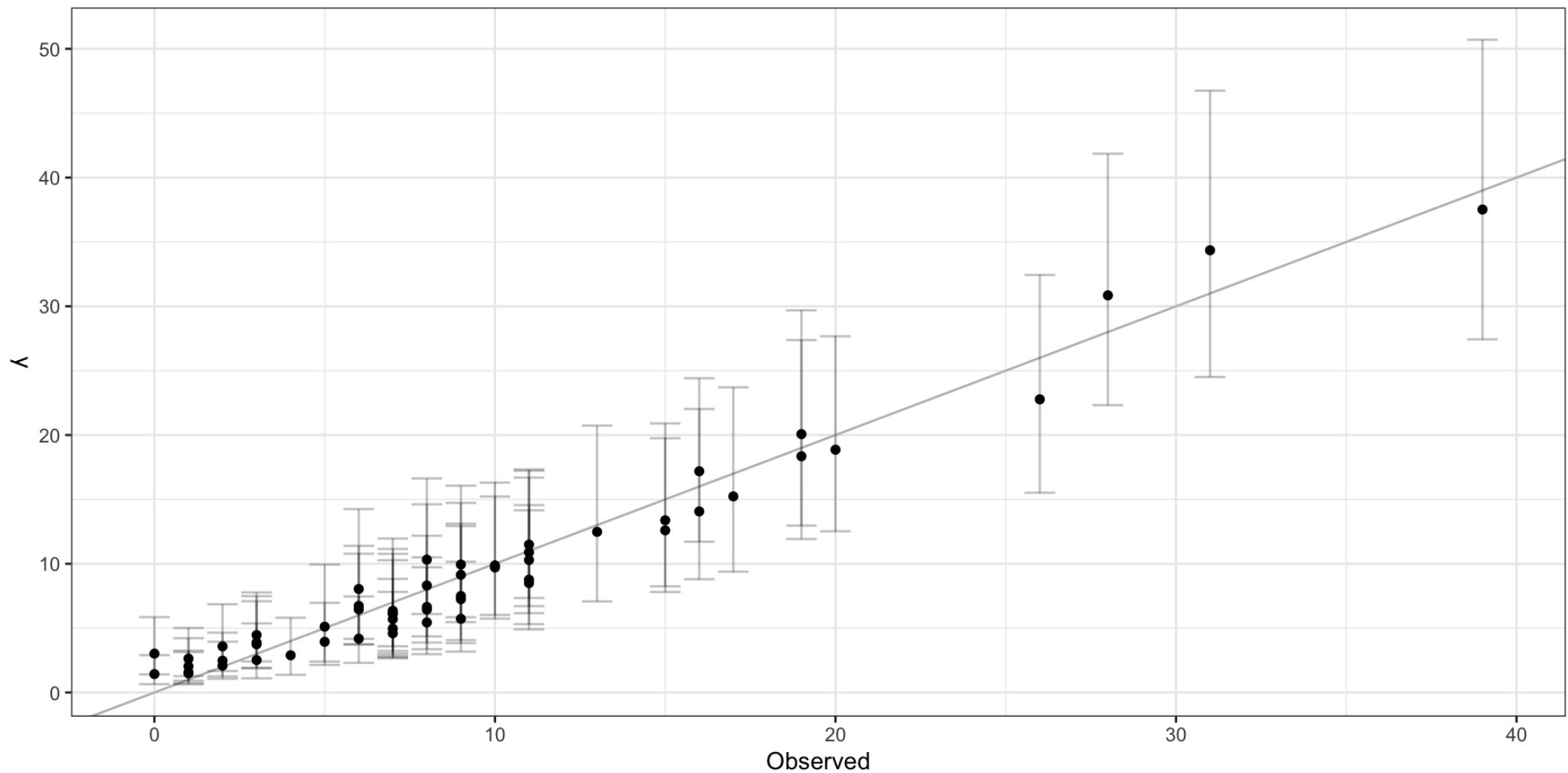
Observed Cases



Predicted Cases

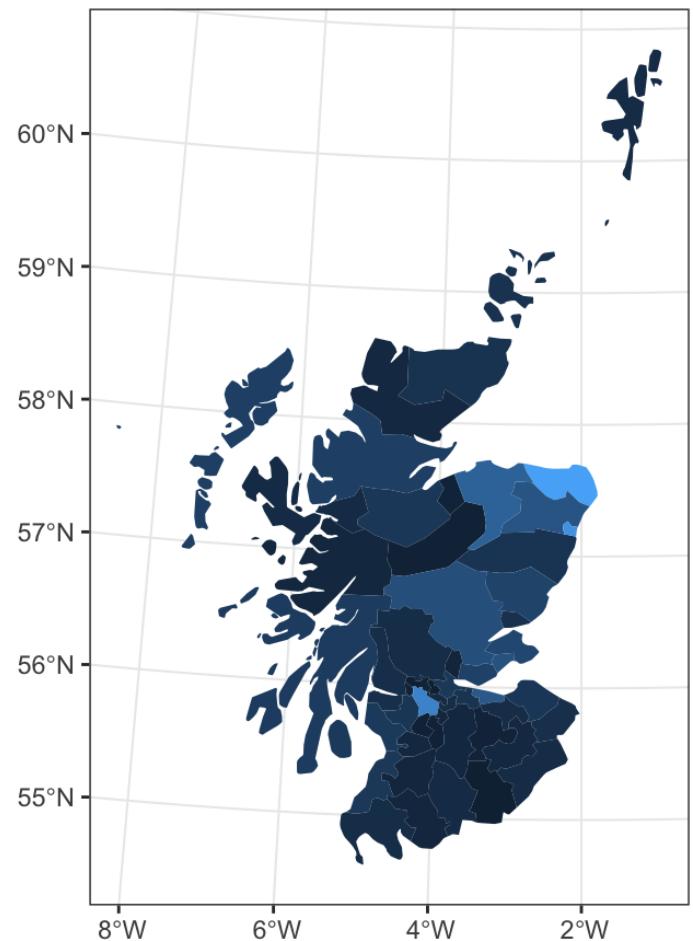


CAR Predictions

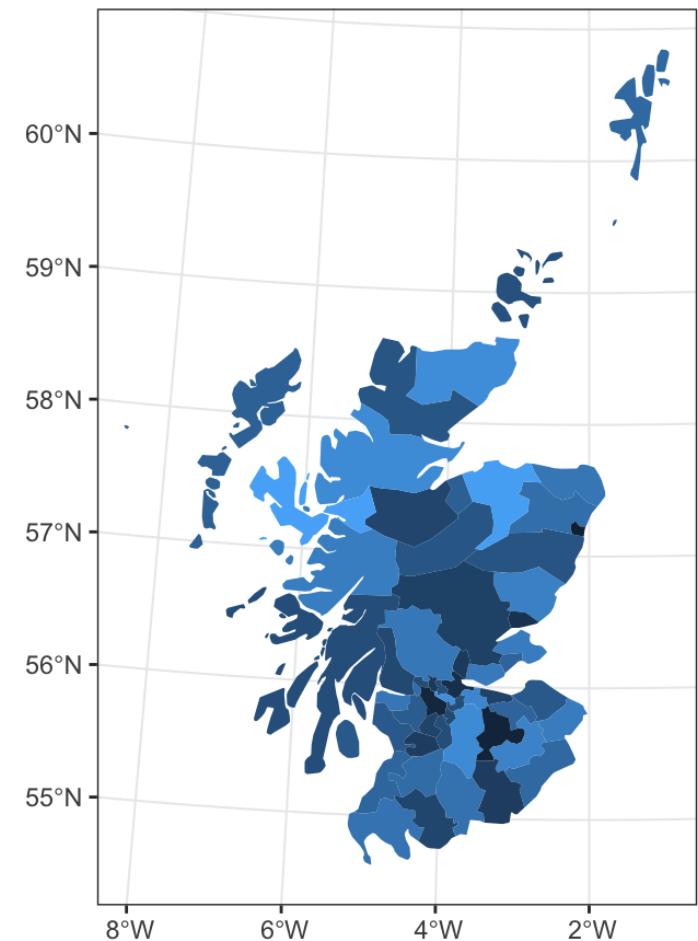


CAR Residuals

Predicted Cases



Residuals



CAR Results

```
1 #RMSE  
2 yardstick::rmse_vec(car_lip_cancer$Observed, car_lip_cancer$y_pred)
```

```
[1] 1.599187
```

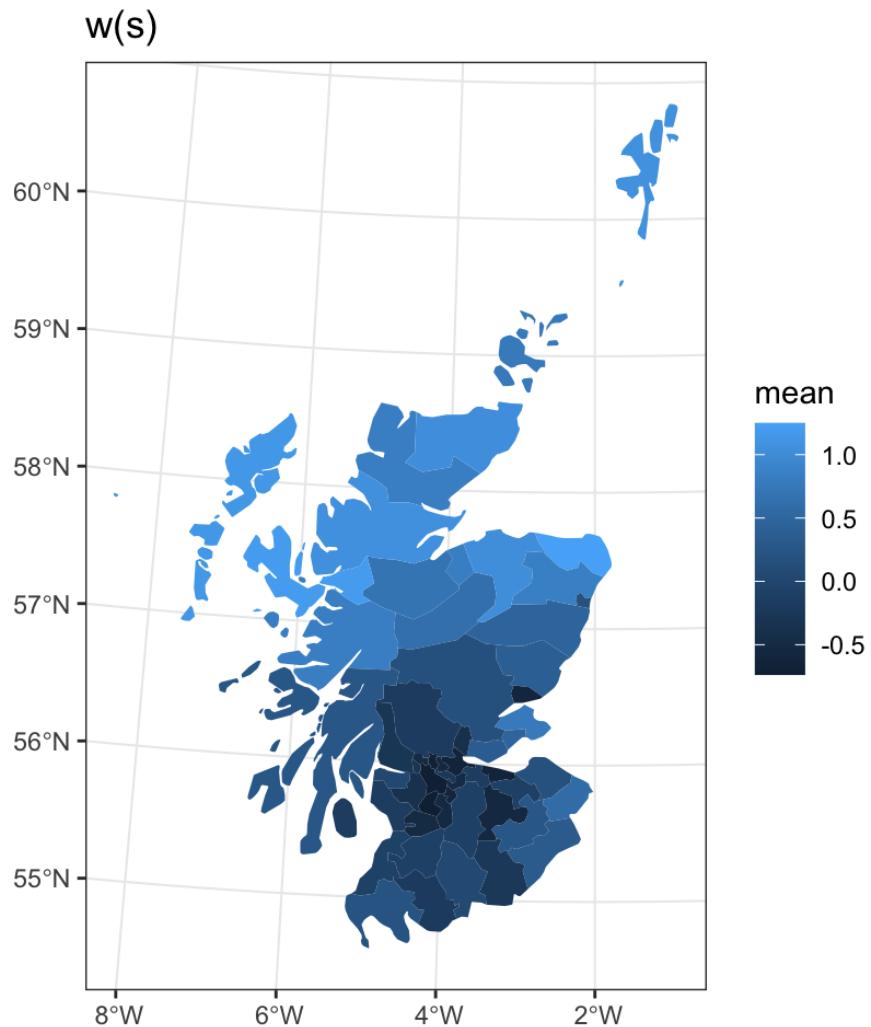
```
1 #Moran's I  
2 spdep::moran.test(car_lip_cancer$Observed - car_lip_cancer$y_pred, listw)
```

Moran I test under randomisation

```
data: car_lip_cancer$Observed - car_lip_cancer$y_pred  
weights: listw
```

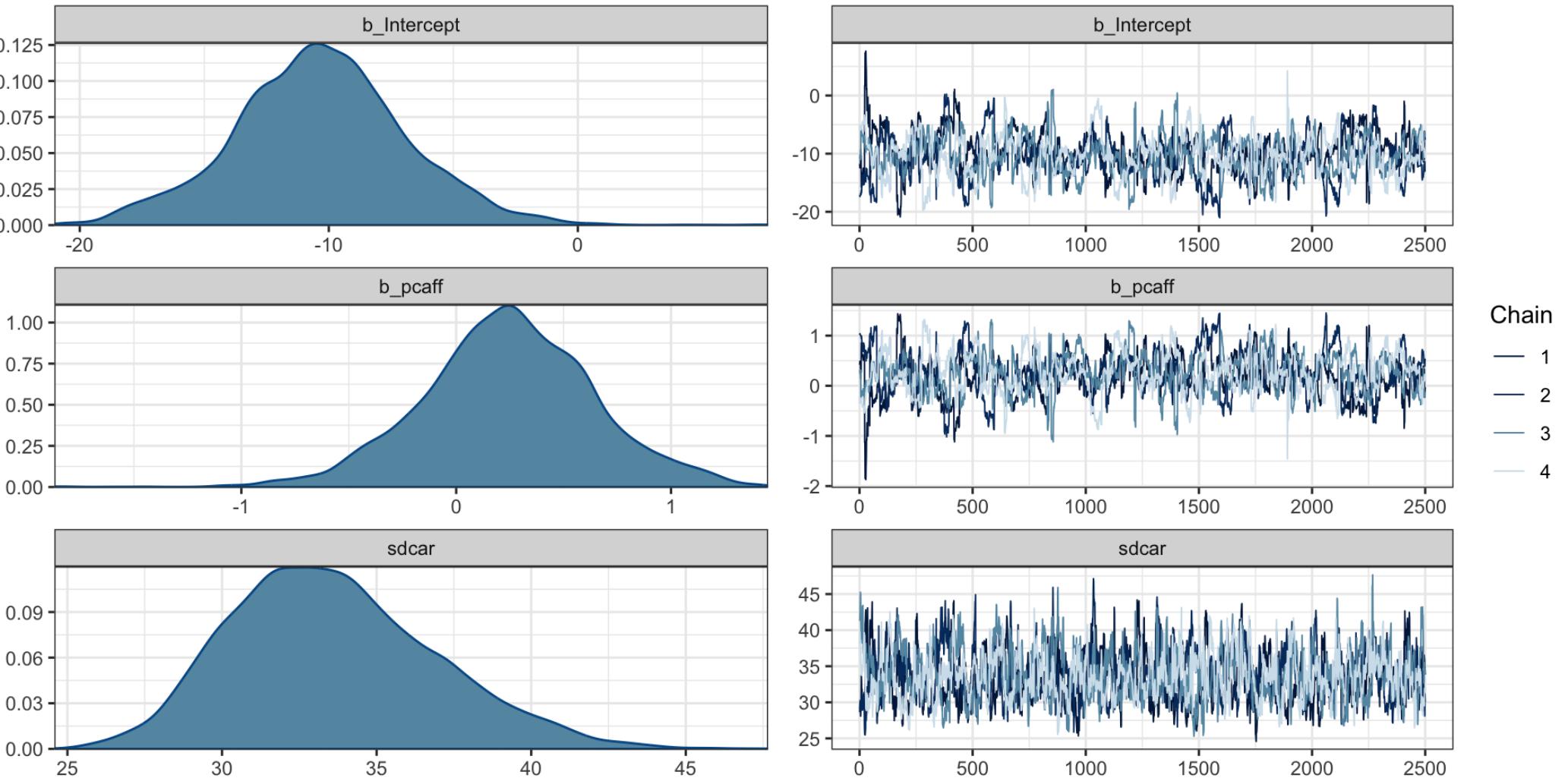
```
Moran I statistic standard deviate =  
0.73353, p-value = 0.2316  
alternative hypothesis: greater  
sample estimates:  
Moran I statistic      Expectation  
0.036963635      -0.018181818  
Variance  
0.005651802
```

Latent spatial process



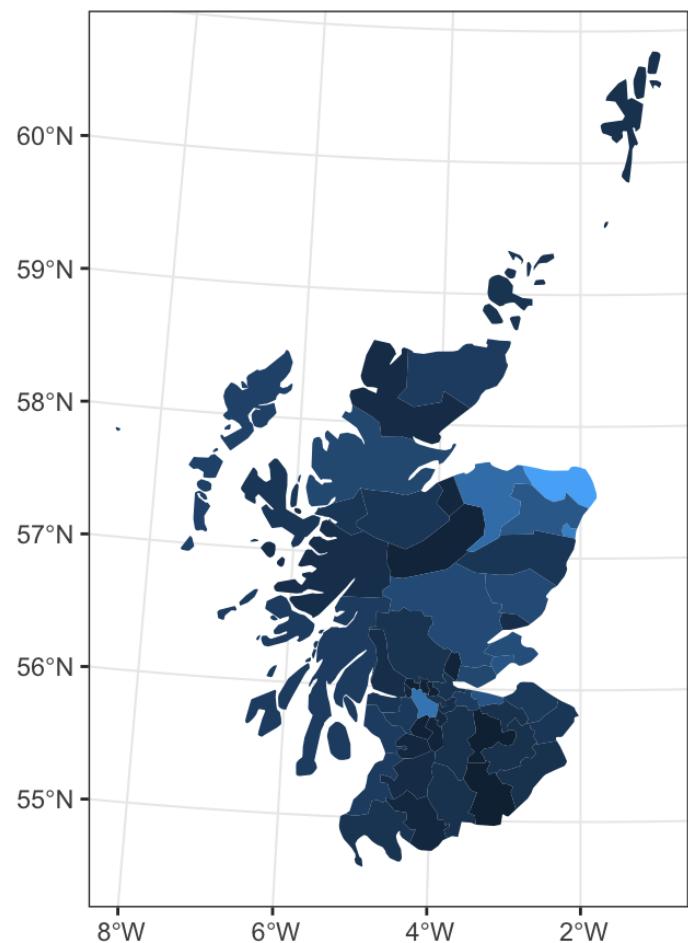
Intrinsic Autoregressive Model (IAR)

```
1 rownames(A) = lip_cancer$District  
2  
3 iar_m = brms::brm(  
4   Observed~offset(Expected)+pcaff+car(A, gr=District, type="icar"),  
5   data=lip_cancer, data2=list(A=A),  
6   family = poisson, cores=4, iter=5000  
7 )
```

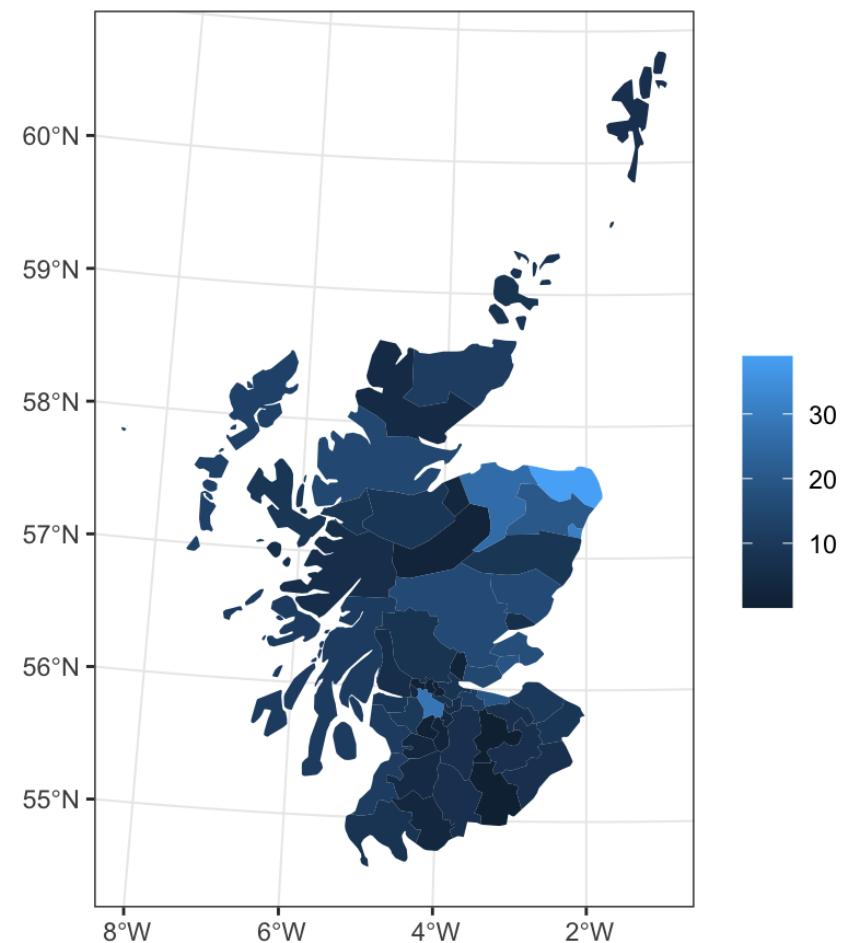


Predictions

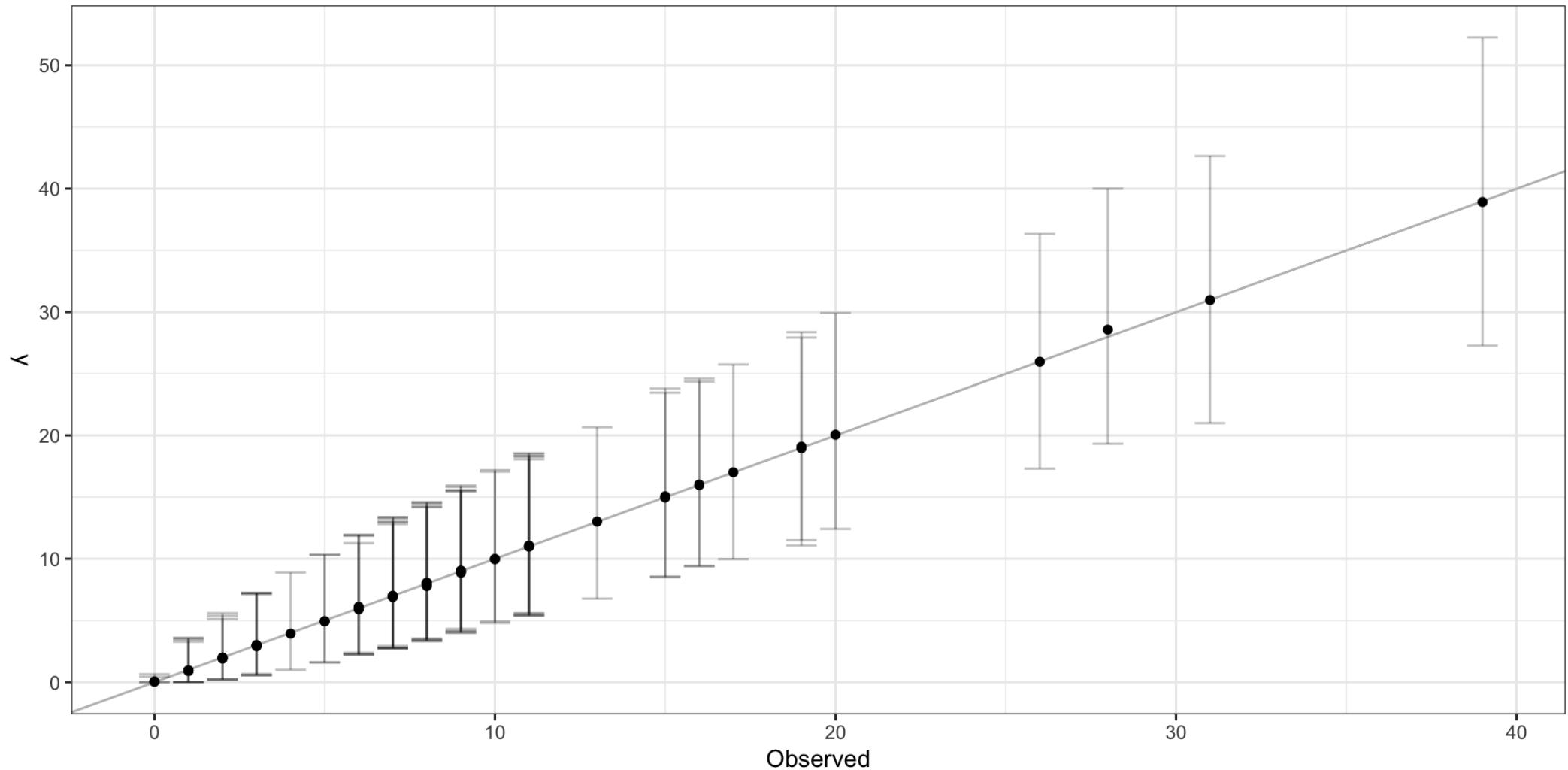
Observed Cases



Predicted Cases

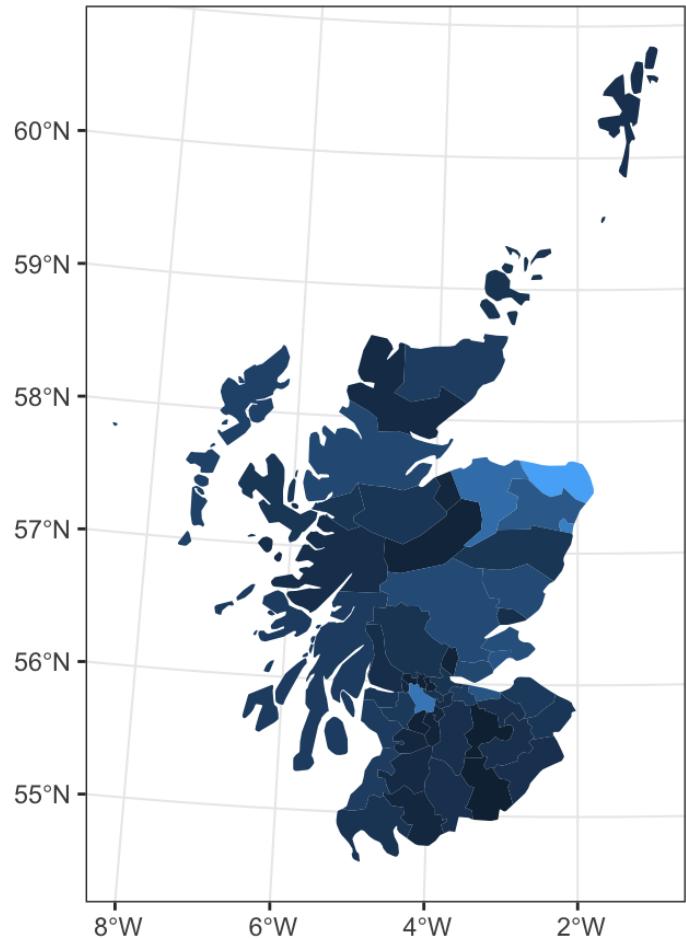


Observed vs predicted

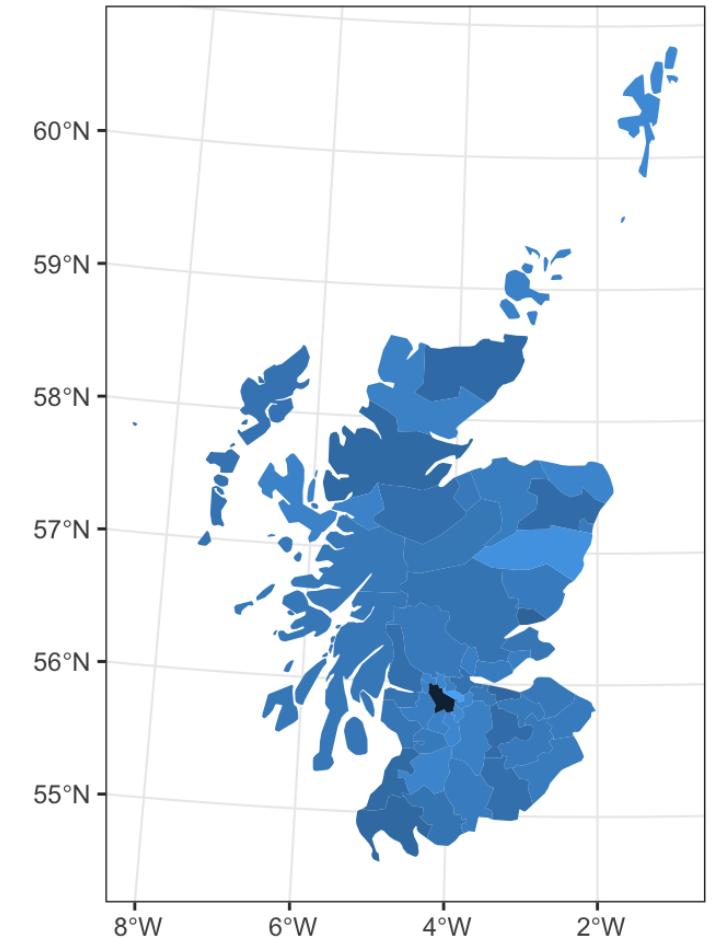


Residuals

Predicted Cases



Residuals



IAR Results

```
1 #RMSE  
2 yardstick::rmse_vec(iar_pred$Observed, iar_pred$y_pred)
```

```
[1] 0.09762396
```

```
1 #Moran's I  
2 spdep::moran.test(iar_pred$resid, listw)
```

Moran I test under randomisation

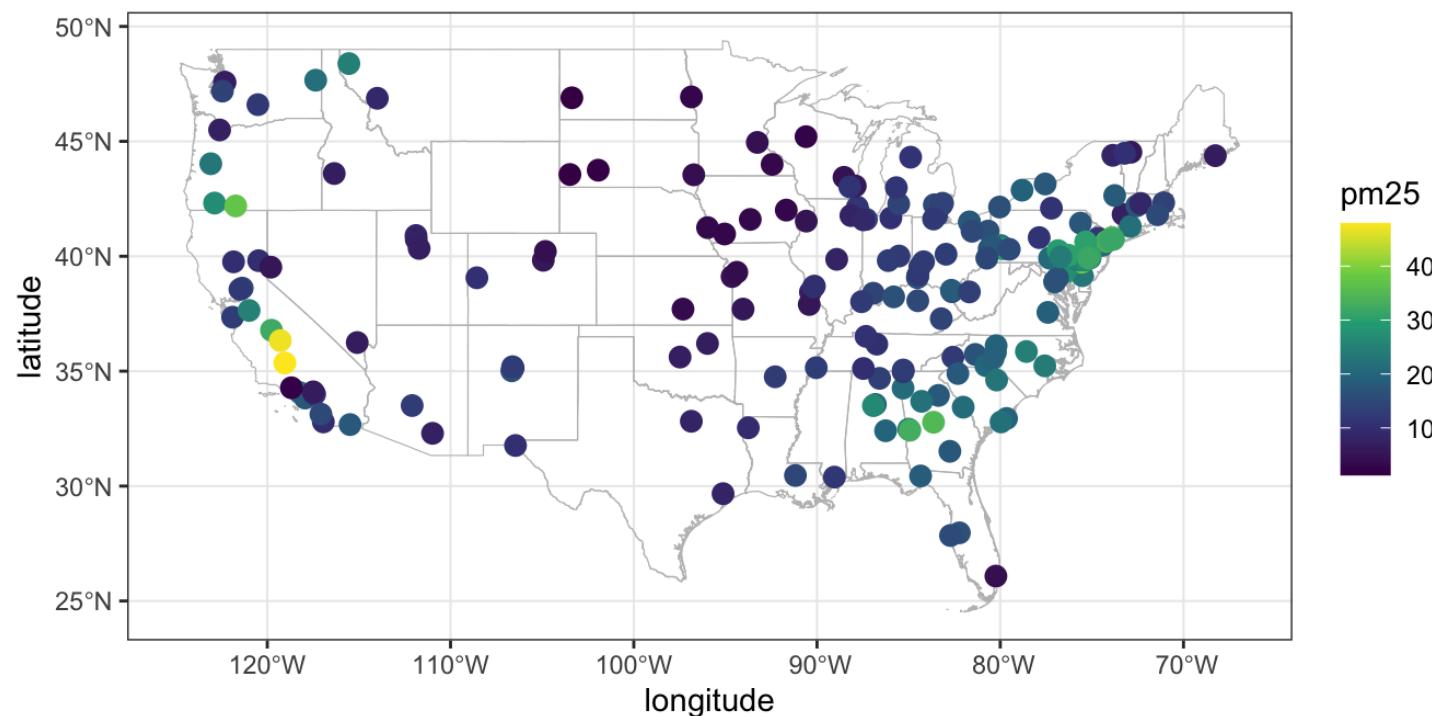
```
data: iar_pred$resid  
weights: listw
```

```
Moran I statistic standard deviate = 2.5724,  
p-value = 0.005051  
alternative hypothesis: greater  
sample estimates:  
Moran I statistic      Expectation  
0.131306700      -0.018181818  
Variance  
0.003377189
```

Point Referenced Data

Example - PM2.5 from CSN

The Chemical Speciation Network are a series of air quality monitors run by EPA (221 locations in 2007). We'll look at a subset of the data from Nov 11th, 2007 (n=191) for just PM2.5.

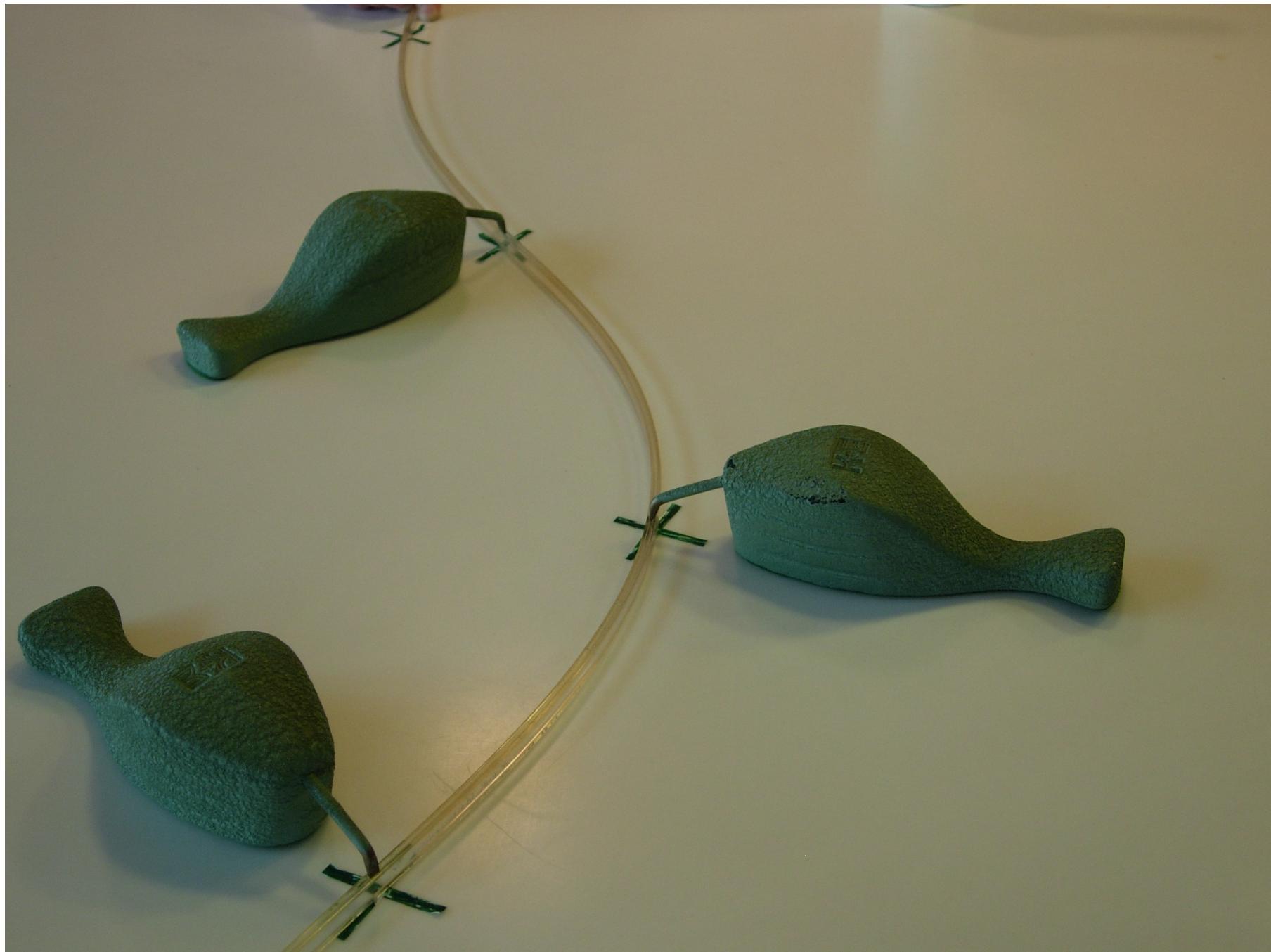


```
1 csn
```

```
# A tibble: 191 × 5
  site longitude latitude date
  <int>     <dbl>    <dbl> <dttm>
1 10730023     -86.8     33.6 2007-11-14 00:00:00
2 10732003     -86.9     33.5 2007-11-14 00:00:00
3 10890014     -86.6     34.7 2007-11-14 00:00:00
4 11011002     -86.3     32.4 2007-11-14 00:00:00
5 11130001     -85.0     32.5 2007-11-14 00:00:00
6 40139997    -112.      33.5 2007-11-14 00:00:00
7 40191028    -111.      32.3 2007-11-14 00:00:00
8 51190007    -92.3     34.8 2007-11-14 00:00:00
9 60070002   -122.      39.8 2007-11-14 00:00:00
```

Aside - Splines





Sta 344 - Fall 2022

Splines in 1d - Smoothing Splines

These are a mathematical analogue to the drafting splines represented using a penalized regression model.

We want to find a function $f(x)$ that best fits our observed data $\mathbf{y} = y_1, \dots, y_n$ while being *smooth*.

$$\arg \min_{f(x)} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_{-\infty}^{\infty} f''(x)^2 dx$$

Interestingly, this minimization problem has an exact solution which is given by a mixture of weighted natural cubic splines (cubic splines that are linear in the tails) with knots at the observed data locations (x_s).

Splines in 2d - Thin Plate Splines

Now imagine we have observed data of the form (x_i, y_i, z_i) where we wish to predict z_i given x_i and y_i for all i . We can extend the smoothing spline model in two dimensions,

$$\arg \min_{f(x,y)} \sum_{i=1}^n (z_i - f(x_i, y_i))^2 + \lambda \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(\frac{\partial^2 f}{\partial x^2} + 2 \frac{\partial^2 f}{\partial x \partial y} + \frac{\partial^2 f}{\partial y^2} \right) dx dy$$

The solution to this equation has a natural representation using a weighted sum of *radial basis functions* with knots at the observed data locations (x_i)

$$f(x) = \sum_{i=1}^n w_i d(x, x_i)^2 \log d(x, x_i).$$

Prediction locations

```
1 r_usa = stars::st_rasterize(  
2   usa,  
3   stars::st_as_stars(st_bbox(usa),  
4     nx = 100, ny = 50, values=NA_real_  
5   )  
6 plot(r_usa)
```

ID



Fitting a TPS

```
1 coords = select(csn, long=longitude, lat=latitude) |>
2   as.matrix()
3 (tps = fields:::Tps(x = coords, Y=csn$pm25, lon.lat=TRUE))
```

Call:

```
fields:::Tps(x = coords, Y = csn$pm25, lon.lat = TRUE)
```

Number of Observations: 191
Number of parameters in the null space 3
Parameters for fixed spatial drift 3
Model degrees of freedom: 64
Residual degrees of freedom: 127
GCV estimate for tau: 4.461
MLE for tau: 4.286
MLE for sigma: 15.35
lambda 1.2
User supplied sigma NA
User supplied tau^2 NA
Summary of estimates:
lambda trA GCV tau1Hat

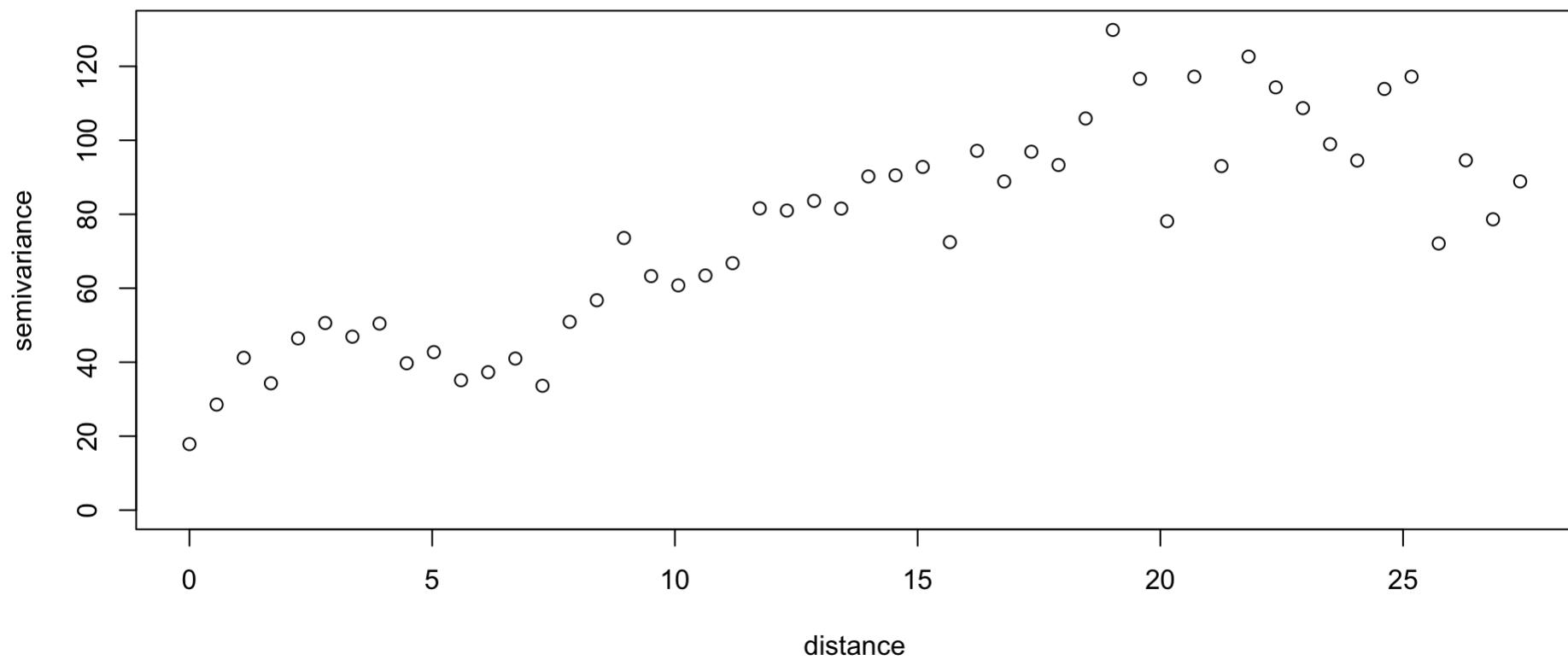
Predictions



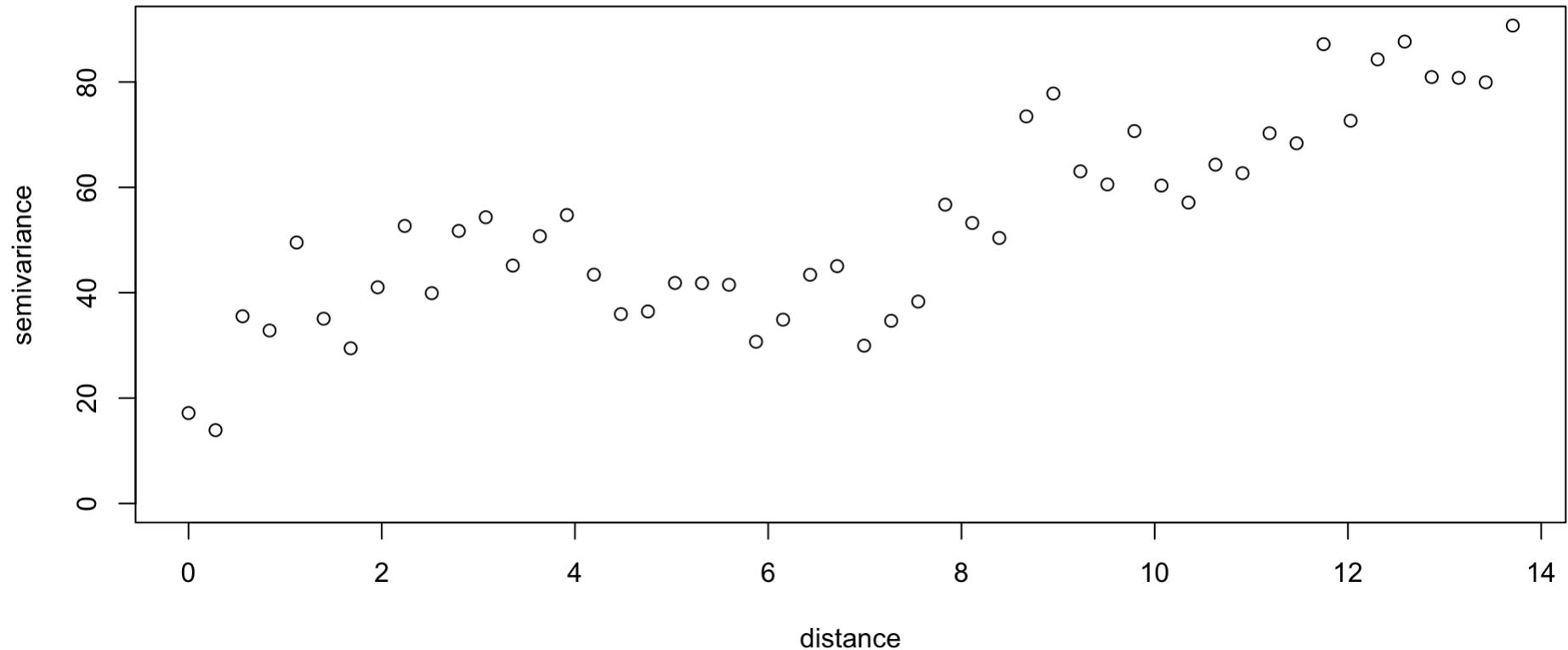
Gaussian Process Models / Kriging

Variogram

```
1 coords = csn %>% select(latitude, longitude) %>% as.matrix()
2 d = fields:::rdist(coords)
3
4 geoR:::variog(
5   coords = coords, data = csn$pm25, messages = FALSE,
6   uvec = seq(0, max(d)/2, length.out=50)
7 ) %>%
8 plot()
```

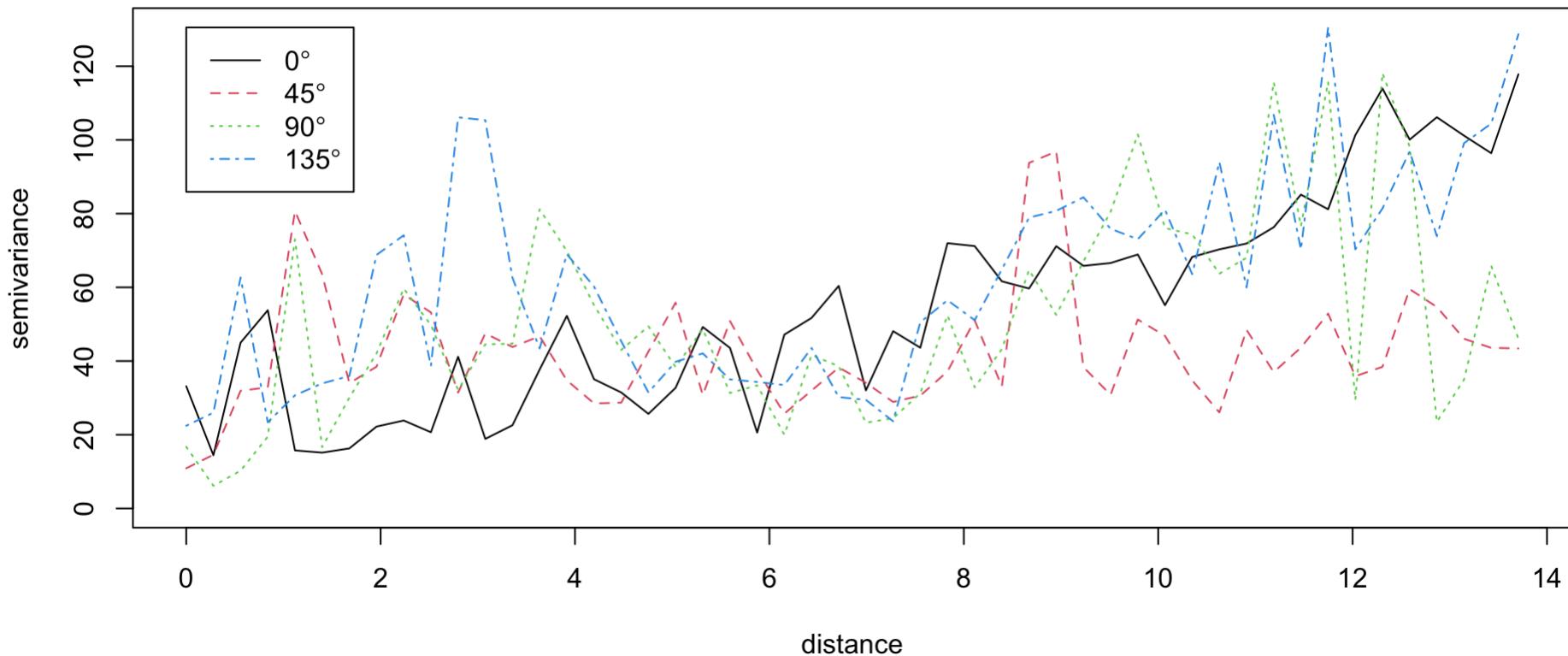


```
1 geoR::variog(  
2   coords = coords, data = csn$pm25, messages = FALSE,  
3   uvec = seq(0, max(d)/4, length.out=50)  
4 ) %>% plot()
```



Isotropy / Anisotropy

```
1 v4 = geoR:::variog4(  
2   coords = coords, data = csn$pm25, messages = FALSE,  
3   uvec = seq(0, max(d)/4, length.out = 50)  
4 )  
5 plot(v4)
```



GP Spatial Model

If we assume that our data is *stationary* and *isotropic* then we can use a Gaussian Process model to fit the data. We will assume an exponential covariance structure.

$$\mathbf{y} \sim (\boldsymbol{\mu}, \Sigma)$$

$$\{\Sigma\}_{ij} = \sigma^2 \exp(-1 \|s_i - s_j\|) + \sigma_n^2 \mathbf{1}_{i=j}$$

we can also view this as a spatial random effects model where

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

$$w(s) \sim (0, \Sigma')$$

$$\epsilon(s_i) \sim (0, \sigma_n^2)$$

$$\{\Sigma'\}_{ij} = \sigma^2 \exp(-r \|s_i - s_j\|)$$

Fitting with gplm() (spBayes)

```
1 max_range = max(dist(csn[,c("longitude", "latitude")])) / 4
2
3 m = gplm(
4   pm25~1, data = csn, coords=c("longitude", "latitude"),
5   cov_model = "exponential",
6   starting = list(phi = 3/3, sigma.sq = 33, tau.sq = 17),
7   tuning = list("phi"=0.1, "sigma.sq"=0.1, "tau.sq"=0.1),
8   priors = list(
9     phi.Unif = c(3/max_range, 3/(0.5)),
10    sigma.sq.IG = c(2, 2),
11    tau.sq.IG = c(2, 2)
12  ),
13  thin=10,
14  verbose=TRUE
15 )
```

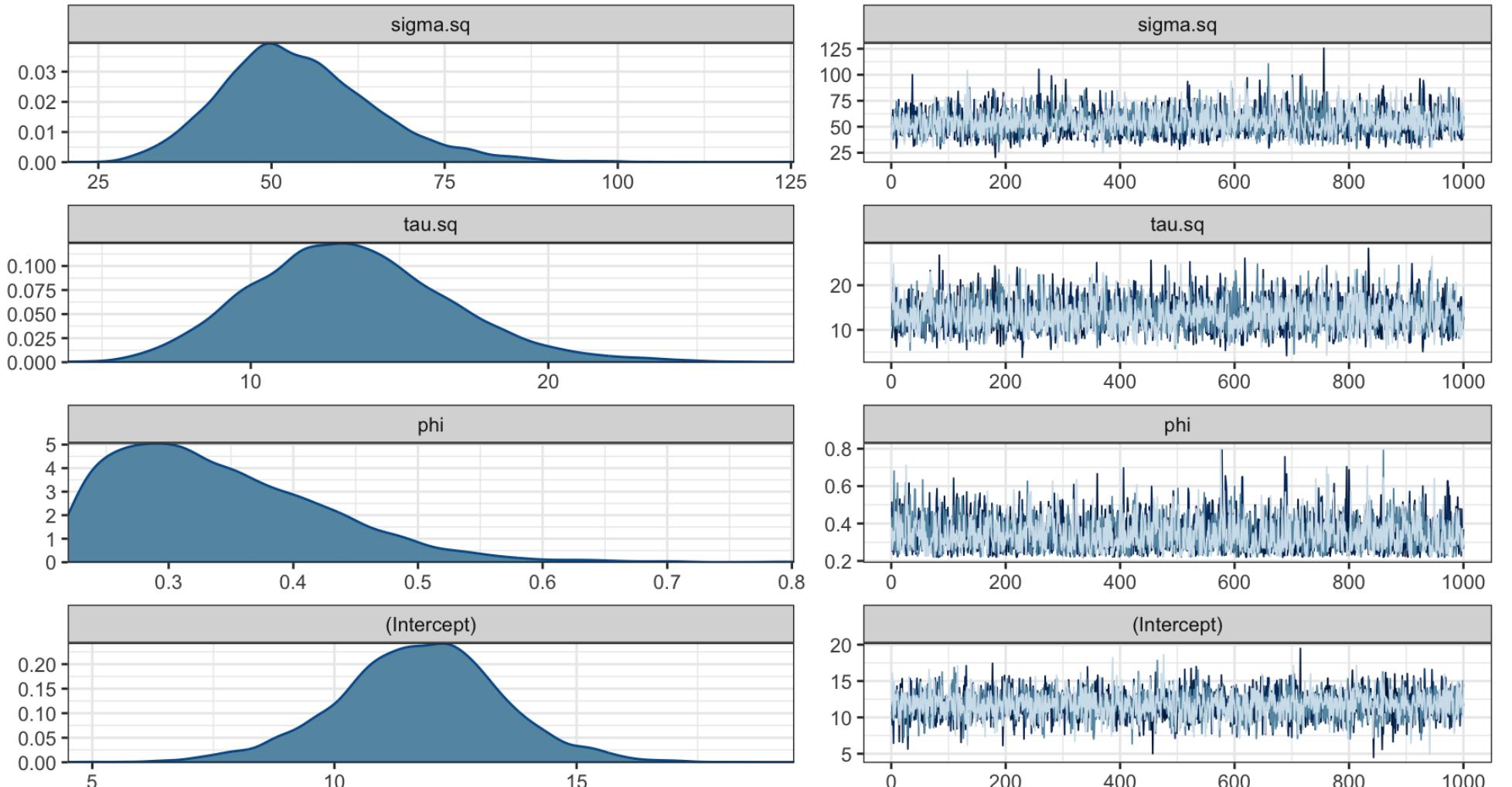
```
1 m
```

```
# A gplm model (spBayes spLM) with 4 chains, 4 variables, and 4000
iterations.

# A tibble: 4 × 10
  variable      mean   median       sd     mad     q5
  <chr>        <dbl>   <dbl>    <dbl>   <dbl>   <dbl>
1 sigma.sq     54.0    52.9    11.3    10.4    37.7
2 tau.sq       13.4    13.2    3.33    3.20    8.28
3 phi          0.341   0.325   0.0848  0.0835  0.232
4 (Intercept) 11.8    11.8    1.71    1.59    8.87
# ... with 4 more variables: q95 <dbl>, rhat <dbl>,
#   ess_bulk <dbl>, ess_tail <dbl>
```

Parameter values

```
1 plot(m)
```



Predictions

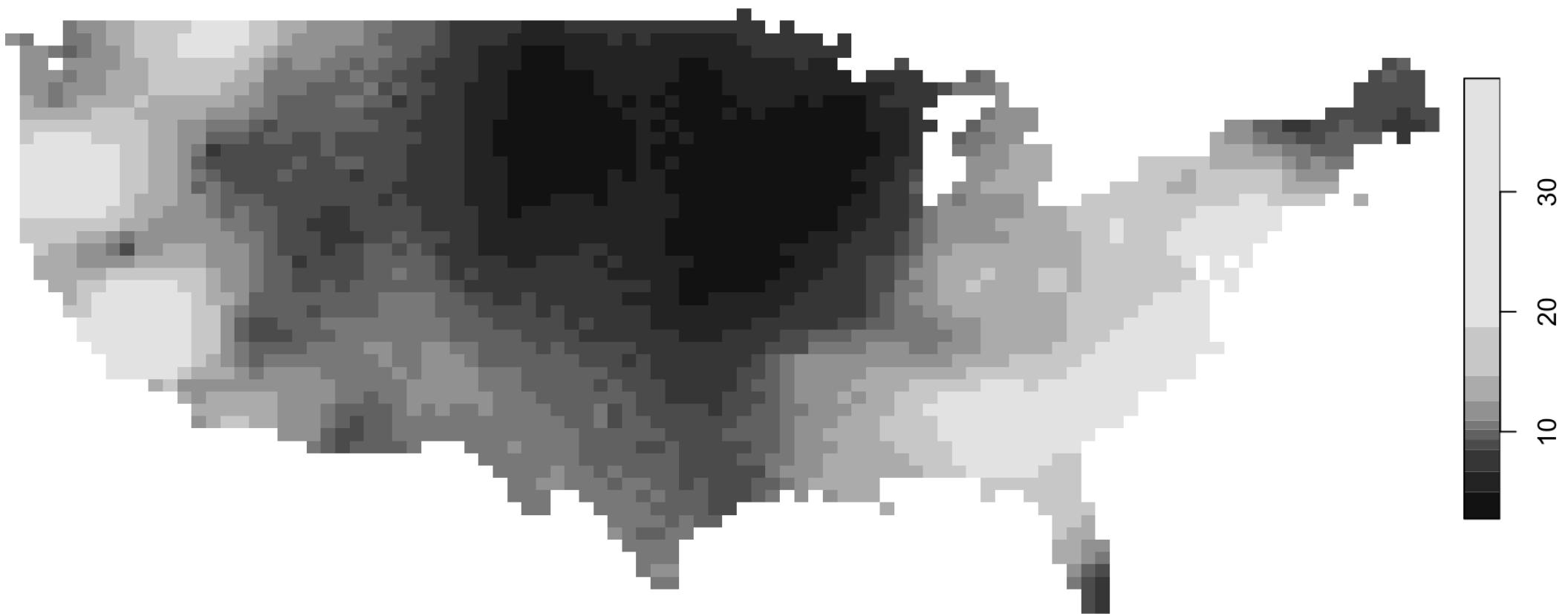
```
1 (p = predict(m, newdata=pred, coords=c("longitude", "latitude")))
```

```
# A draws_matrix: 1000 iterations, 4 chains, and 2828 variables
```

```
variable
```

draw	y[1]	y[2]	y[3]	y[4]	y[5]	y[6]	y[7]	y[8]
1	14.03	-4.073	15.0	4.8	-8.8	7.84	21	4.9
2	11.71	0.052	10.2	5.8	11.3	14.58	20	10.7
3	-3.37	17.307	18.4	20.2	23.7	28.46	9	20.4
4	7.31	2.500	4.6	7.3	23.7	14.63	15	11.8
5	0.47	10.014	10.4	17.2	14.6	11.17	10	10.0
6	7.57	11.004	10.6	9.2	10.6	14.56	23	10.0
7	7.16	6.791	12.8	5.0	22.4	0.88	16	20.1
8	16.54	9.611	1.8	23.9	23.9	19.23	38	10.0
9	16.03	3.135	23.7	1.1	12.4	13.10	34	20.1

mean



sd

