

Methods

Matt Williamson

6/24/2019

Methods

Our interest was in evaluating the bias that arises from (potentially unmodeled) variation in the probability of reporting an easement (given that it is there), the probability that a spatial location is available for an easement, and spatial autocorrelation in both occupancy and reporting probability.

The models

The model used to generate the data was:

$$\begin{aligned} z_i &\sim \text{Bern}(\psi_i) \\ \text{logit}(\psi_i) &= \beta_0 + \beta \mathbf{X} + \phi_i \\ v_i | z_i &\sim \text{Bin}(n_{bg}, \alpha * z_i) \\ y_{i,j} | v_i z_i &\sim \text{Bern}(p_{i,j} * z_i * v_i) \\ \text{logit}(p_{i,j}) &= \gamma_0 + \gamma \mathbf{X} + \phi_{i,j} \end{aligned}$$

where $v_i=1$ if available and 0 if not and n_{bg} is the number of block groups within a tract.

The models that I actually fit were as follows: Note that I did not explicitly model α because this would be difficult to do in-practice and my interest was in whether the approach would be robust to leaving that component out of the model (despite its role in the generating process)

Occupancy with CAR on both components

$$\begin{aligned} z_i &\sim \text{Bern}(\psi_i) \\ \text{logit}(\psi_i) &= \beta_0 + \beta \mathbf{X} + \phi_i \\ y_{i,j} | z_i &\sim \text{Bern}(p_{i,j} * z_i) \\ \text{logit}(p_{i,j}) &= \gamma_0 + \gamma \mathbf{X} + \phi_{i,j} \end{aligned}$$

Occupancy with CAR on occupancy only

$$\begin{aligned} z_i &\sim \text{Bern}(\psi_i) \\ \text{logit}(\psi_i) &= \beta_0 + \beta \mathbf{X} + \phi_i \\ y_{i,j} | z_i &\sim \text{Bern}(p_{i,j} * z_i) \\ \text{logit}(p_{i,j}) &= \gamma_0 + \gamma \mathbf{X} \end{aligned}$$

Occupancy with CAR on detection only

$$\begin{aligned} z_i &\sim \text{Bern}(\psi_i) \\ \text{logit}(\psi_i) &= \beta_0 + \beta \mathbf{X} \\ y_{i,j} | z_i &\sim \text{Bern}(p_{i,j} * z_i) \\ \text{logit}(p_{i,j}) &= \gamma_0 + \gamma \mathbf{X} + \phi_{i,j} \end{aligned}$$

Occupancy with no CAR component

$$\begin{aligned}
z_i &\sim \text{Bern}(\psi_i) \\
\text{logit}(\psi_i) &= \beta_0 + \beta \mathbf{X} \\
y_{i,j}|z_i &\sim \text{Bern}(p_{i,j} * z_i) \\
\text{logit}(p_{i,j}) &= \gamma_0 + \gamma \mathbf{X}
\end{aligned}$$

Binomial fit to the tract (i) level

$$\begin{aligned}
y_i &\sim \text{Bern}(p_i) \\
\text{logit}(p_i) &= \beta_0 + \beta \mathbf{X} + \phi_i
\end{aligned}$$

Generating covariate observations

We used the tracts and block groups for Iowa as the basis for developing our simulated datasets. We generated random data values for each of three predictors at the tract-level and two predictors at the block group-level. Because predictors may also be spatially autocorrelated, we simulated these data from $x_i \sim \mathcal{N}_{||}(\boldsymbol{\mu} = 0, \boldsymbol{\Sigma} = \phi)$ where $\phi \sim \mathcal{N}(0, [\tau(D - \rho W)]^{-1})$ and τ , the spatially varying precision parameter is 1; ρ , the strength of spatial dependence is 0.3, D is a diagonal matrix containing the number of neighbors for a given location, and W is an adjacency matrix. We defined adjacency by determining the minimum distance necessary to ensure that all locations had at least one neighbor and considering to locations to be adjacent if they were within that distance from each other.

Latin Hypercube Sampling

We used a Latin hypercube design to evaluate model performance across a range of plausible parameter values. For each simulation replicate, we drew 300 uniformly distributed samples across the multi-dimensional space described by the limits in Table 1.

Generating observations

We simulated true easement occurrence for each sample generating a probability of easement occurrence for each tract by multiplying ψ (the mean occupancy probability for a given Latin hypercube sample [i.e., the intercept]) and three randomly generated regression coefficients by the tract-level design matrix and adding ϕ_Ψ where ϕ is described above and τ and ρ_Ψ were taken from the appropriate Latin hypercube sample. This resulted in true easement occurrence (z_i) for each tract which we then converted to an observed occupancy dataset by generating an estimate of p (the block group reporting probability) by multiplying p (the mean reporting probability for a given Latin hypercube sample [i.e., the intercept]) and two randomly generated regression coefficients by the block group-level design matrix and adding ϕ_p with τ and ρ_p as above. Finally, because the probability of an easement being reported (conditional on its existence) depends on both the probability of reporting (p) and the probability that a block group is available for an easement (α), we simulated easement observations following $y_{i,j} \sim \text{Bern}(\alpha p)$.

Table 1: Range of values used in Latin Hypercube sampler where Ψ is the probability of easement occurrence, p is detection probability, $\rho_{...}$ is the strength of spatial autocorrelation for Ψ or p , τ is the precision for the conditional autoregressive term, and α is the probability that a location is available for an easement.

Variables	Lower	Upper
Ψ	0.20	0.80
p	0.30	0.98
ρ_Ψ	0.50	1.00
ρ_p	0.50	1.00
τ	0.10	1.00
α	0.20	0.80

Fitting models

We generated 50 replicates [only showing results from 4 here] of each Latin hypercube sample (1500 simulated datasets total) and fit each of five models (a standard occupancy model, a binomial model with a conditional autoregressive [CAR] term for spatial autocorrelation, an occupancy model with a CAR term for occurrence only, an occupancy model with a CAR term for reporting only, and an occupancy model with a CAR terms for both components) to the simulated dataset. All models were fit in R using rstan, a wrapper to the Stan language (stan models are at end of document). Models were fit using an adaptation parameter of 0.98 and a maximum treedepth of 16 with four chains each run for 3700 iterations (with a warmup period of 3200 iterations). We calculated the relative bias for the intercepts and regression coefficients as:

$$RelBias = \frac{\hat{\beta} - \beta}{|\beta|}$$

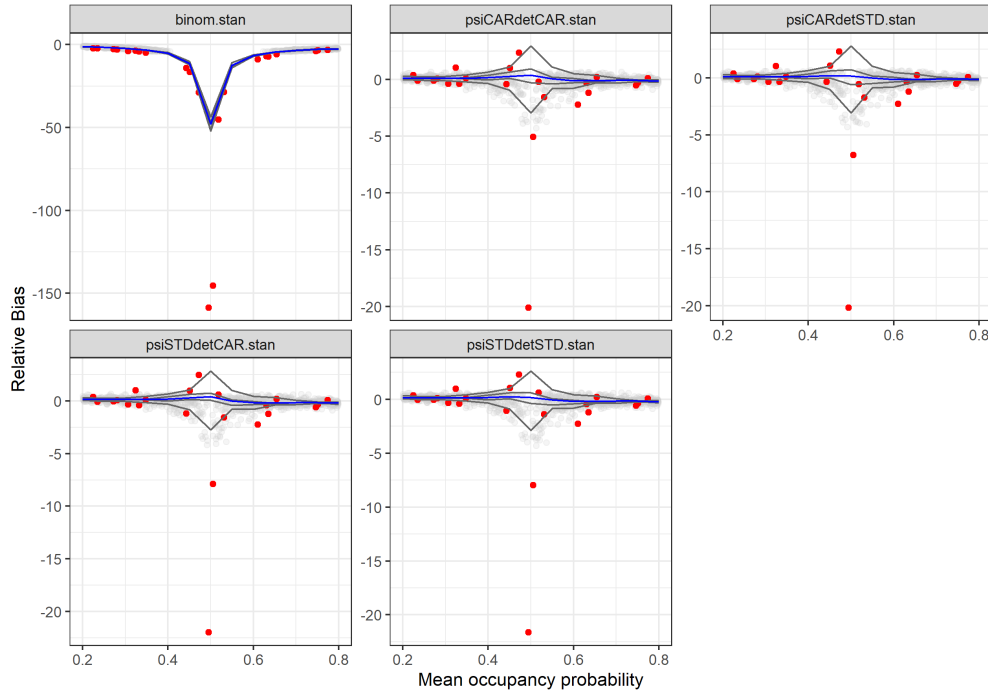


Figure 1: Relative bias of the posterior estimates of the intercept in the binomial model and the intercept in the occupancy component of the occupancy modes for each simulated dataset. Gray dots are median estimates of relative bias from each model run, gray lines are replicate-specific the relationship between median values of relative bias and the simulated mean occupancy probability, blue lines depict the relationship between median values of relative bias and occupancy probability across all simulation replicates, red dots indicate simulation runs where both the probability of detection and the probability of availability were in the lower decile.

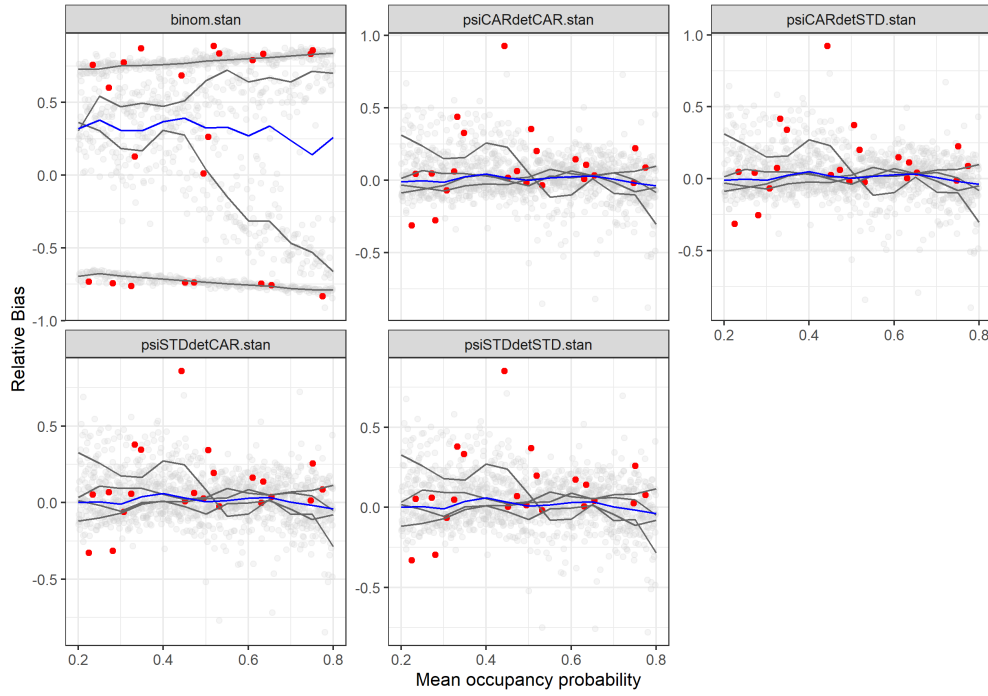


Figure 2: Relative bias of the posterior estimates of the regression coefficients in the binomial model and the intercept in the occupancy component of the occupancy modes for each simulated dataset. Gray dots are median estimates of relative bias from each model run, gray lines are replicate-specific the relationship between median values of relative bias and the simulated mean occupancy probability, blue lines depict the relationship between median values of relative bias and occupancy probability across all simulation replicates, red dots indicate simulation runs where both the probability of detection and the probability of availability were in the lower decile.

Stan models

Binomial

```
functions{
/**
* Return log probability of a unit-scale proper conditional autoregressive
* (CAR) prior with a sparse representation for the adjacency matrix
*
* @param phi Vector containing the parameters with a CAR prior
* @param alpha Dependence (usually spatial) parameter for the CAR prior (real)
* @param W_sparse Sparse representation of adjacency matrix (int array)
* @param n Length of phi (int)
* @param W_n Number of adjacent pairs (int)
* @param D_sparse Number of neighbors for each location (vector)
* @param lambda Eigenvalues of  $D^{-1/2} * W * D^{-1/2}$  (vector)
*
* @return Log probability density of CAR prior up to additive constant
*/
real sparse_car_lpdf(vector phi, real alpha,
  int[, ] W_sparse, vector D_sparse, vector lambda, int nobs, int W_n) {
  row_vector[nobs] phit_D; //  $\phi' * D$ 
  row_vector[nobs] phit_W; //  $\phi' * W$ 
  vector[nobs] ldet_terms;

  phit_D = (phi .* D_sparse)';
  phit_W = rep_row_vector(0, nobs);
  for (i in 1:W_n) {
    phit_W[W_sparse[i, 1]] += phi[W_sparse[i, 2]];
    phit_W[W_sparse[i, 2]] += phi[W_sparse[i, 1]];
  }

  for (i in 1:nobs) ldet_terms[i] = loglm(alpha * lambda[i]);
  return 0.5 * (sum(ldet_terms)
    - (phit_D * phi - alpha * (phit_W * phi)));
}
}
data {
  // site-level occupancy covariates
  int<lower = 1> n_site;
  int<lower = 1> m_psi;
  matrix[n_site, m_psi] X_tct;

  // summary of whether species is known to be present at each site
  int<lower = 0, upper = 1> any_seen[n_site];

  // number of surveys at each site
  int<lower = 0> n_survey[n_site];

  matrix<lower = 0, upper = 1>[n_site, n_site] W_tct; //adjacency matrix tract
  int W_n_tct; //number of adjacent pairs
}
transformed data {
```

```

int W_sparse_occ[W_n_tct, 2]; // adjacency pairs
vector[n_site] D_sparse_occ; // diagonal of D (number of neighbors for each site)

vector[n_site] lambda_occ; // eigenvalues of invsqrtD * W * invsqrtD
{ // generate sparse representation for W
int counter_occ;
counter_occ = 1;
// loop over upper triangular part of W to identify neighbor pairs
for (i in 1:(n_site - 1)) {
for (j in (i + 1):n_site) {
if (W_tct[i, j] == 1) {
W_sparse_occ[counter_occ, 1] = i;
W_sparse_occ[counter_occ, 2] = j;
counter_occ += 1;
}
}
}
}
for (i in 1:n_site) D_sparse_occ[i] = sum(W_tct[i]);
{
vector[n_site] invsqrtD_occ;
for (i in 1:n_site) {
invsqrtD_occ[i] = 1 / sqrt(D_sparse_occ[i]);
}
lambda_occ = eigenvalues_sym(quad_form(W_tct, diag_matrix(invsqrtD_occ)));
}
}

parameters {
vector[n_site] phi_occ;
real<lower = 0, upper = 0.999> alpha_occ;
real<lower = 0> sigma_occ;
vector[m_psi] beta_psi;
}

transformed parameters {
vector[n_site] logit_psi = X_tct * beta_psi + phi_occ * sigma_occ;
}

model {
phi_occ ~ sparse_car(alpha_occ, W_sparse_occ, D_sparse_occ, lambda_occ, n_site, W_n_tct);
alpha_occ ~ beta(4,1);
sigma_occ ~ normal(0, 1.5);
beta_psi ~ student_t(7.763,0, 1.566);

any_seen ~ binomial_logit(n_survey, logit_psi);
}

```

Occupancy without CAR

Note that although this code has the functions to estimate ϕ , the actual model does not do so.

```
functions{
/**
* Return log probability of a unit-scale proper conditional autoregressive
* (CAR) prior with a sparse representation for the adjacency matrix
*
* @param phi Vector containing the parameters with a CAR prior
* @param alpha Dependence (usually spatial) parameter for the CAR prior (real)
* @param W_sparse Sparse representation of adjacency matrix (int array)
* @param n Length of phi (int)
* @param W_n Number of adjacent pairs (int)
* @param D_sparse Number of neighbors for each location (vector)
* @param lambda Eigenvalues of  $D^{-1/2} * W * D^{-1/2}$  (vector)
*
* @return Log probability density of CAR prior up to additive constant
*/
real sparse_car_lpdf(vector phi, real alpha,
  int[, ] W_sparse, vector D_sparse, vector lambda, int nobs, int W_n) {
  row_vector[nobs] phit_D; // phi' * D
  row_vector[nobs] phit_W; // phi' * W
  vector[nobs] ldet_terms;

  phit_D = (phi .* D_sparse)';
  phit_W = rep_row_vector(0, nobs);
  for (i in 1:W_n) {
    phit_W[W_sparse[i, 1]] += phi[W_sparse[i, 2]];
    phit_W[W_sparse[i, 2]] += phi[W_sparse[i, 1]];
  }

  for (i in 1:nobs) ldet_terms[i] = loglm(alpha * lambda[i]);
  return 0.5 * (sum(ldet_terms)
    - (phit_D * phi - alpha * (phit_W * phi)));
}
}
data {
  // site-level occupancy covariates
  int<lower = 1> n_site;
  int<lower = 1> m_psi;
  matrix[n_site, m_psi] X_tct;

  // survey-level detection covariates
  int<lower = 1> total_surveys;
  int<lower = 1> m_p;
  matrix[total_surveys, m_p] X_bg;

  // survey level information
  int<lower = 1, upper = n_site> site[total_surveys];
  int<lower = 0, upper = 1> y[total_surveys];
  int<lower = 0, upper = total_surveys> start_idx[n_site];
  int<lower = 0, upper = total_surveys> end_idx[n_site];

  // summary of whether species is known to be present at each site
```

```

int<lower = 0, upper = 1> any_seen[n_site];

// number of surveys at each site
int<lower = 0> n_survey[n_site];

}
parameters {
  vector[m_psi] beta_psi;
  vector[m_p] beta_p;
}
transformed parameters {
  vector[total_surveys] logit_p = X_bg * beta_p;
  vector[n_site] logit_psi = X_tct * beta_psi;
}
model {
  vector[n_site] log_psi = log_inv_logit(logit_psi);
  vector[n_site] log1m_psi = log1m_inv_logit(logit_psi);

  beta_psi ~ student_t(7.763,0, 1.566);
  beta_p ~ student_t(7.763,0, 1.566);

  for (i in 1:n_site) {
    if (n_survey[i] > 0) {
      if (any_seen[i]) {
        // site is occupied
        target += log_psi[i]
                  + bernoulli_logit_lpmf(y[start_idx[i]:end_idx[i]] |
                                          logit_p[start_idx[i]:end_idx[i]]);
      } else {
        // site may or may not be occupied
        target += log_sum_exp(
          log_psi[i] + bernoulli_logit_lpmf(y[start_idx[i]:end_idx[i]] |
                                          logit_p[start_idx[i]:end_idx[i]]),
          log1m_psi[i]
        );
      }
    }
  }
}

```


Occupancy with CAR on both components

```

functions{
/**
* Return log probability of a unit-scale proper conditional autoregressive
* (CAR) prior with a sparse representation for the adjacency matrix
*
* @param phi Vector containing the parameters with a CAR prior
* @param alpha Dependence (usually spatial) parameter for the CAR prior (real)
* @param W_sparse Sparse representation of adjacency matrix (int array)
* @param n Length of phi (int)
* @param W_n Number of adjacent pairs (int)
* @param D_sparse Number of neighbors for each location (vector)
* @param lambda Eigenvalues of  $D^{-1/2} * W * D^{-1/2}$  (vector)
*
* @return Log probability density of CAR prior up to additive constant
*/
real sparse_car_lpdf(vector phi, real alpha,
  int[, ] W_sparse, vector D_sparse, vector lambda, int nobs, int W_n) {
  row_vector[nobs] phit_D; // phi' * D
  row_vector[nobs] phit_W; // phi' * W
  vector[nobs] ldet_terms;

  phit_D = (phi .* D_sparse)';
  phit_W = rep_row_vector(0, nobs);
  for (i in 1:W_n) {
    phit_W[W_sparse[i, 1]] += phi[W_sparse[i, 2]];
    phit_W[W_sparse[i, 2]] += phi[W_sparse[i, 1]];
  }

  for (i in 1:nobs) ldet_terms[i] = loglm(alpha * lambda[i]);
  return 0.5 * (sum(ldet_terms)
    - (phit_D * phi - alpha * (phit_W * phi)));
}
}
data {
  // site-level occupancy covariates
  int<lower = 1> n_site;
  int<lower = 1> m_psi;
  matrix[n_site, m_psi] X_tct;

  // survey-level detection covariates
  int<lower = 1> total_surveys;
  int<lower = 1> m_p;
  matrix[total_surveys, m_p] X_bg;

  // survey level information
  int<lower = 1, upper = n_site> site[total_surveys];
  int<lower = 0, upper = 1> y[total_surveys];
  int<lower = 0, upper = total_surveys> start_idx[n_site];
  int<lower = 0, upper = total_surveys> end_idx[n_site];

  // summary of whether species is known to be present at each site
  int<lower = 0, upper = 1> any_seen[n_site];

```

```

// number of surveys at each site
int<lower = 0> n_survey[n_site];

//adjacency data
matrix<lower = 0, upper = 1>[n_site, n_site] W_tct; //adjacency matrix tract
int W_n_tct; //number of adjacent pairs
matrix<lower = 0, upper = 1>[total_surveys, total_surveys] W_bg; //adjacency matrix bg
int W_n_bg; //number of adjacent pairs bg
//real<lower = 0, upper =1> alpha_occ;
//real<lower = 0, upper =1> alpha_det;

}
transformed data {
  int W_sparse_occ[W_n_tct, 2]; // adjacency pairs
  int W_sparse_det[W_n_bg, 2];
  vector[n_site] D_sparse_occ; // diagonal of D (number of neighbors for each site)
  vector[total_surveys] D_sparse_det;

  vector[n_site] lambda_occ; // eigenvalues of invsqrtD * W * invsqrtD
  vector[total_surveys] lambda_det;
  { // generate sparse representation for W
    int counter_occ;
    counter_occ = 1;
    // loop over upper triangular part of W to identify neighbor pairs
    for (i in 1:(n_site - 1)) {
      for (j in (i + 1):n_site) {
        if (W_tct[i, j] == 1) {
          W_sparse_occ[counter_occ, 1] = i;
          W_sparse_occ[counter_occ, 2] = j;
          counter_occ += 1;
        }
      }
    }
  }
  for (i in 1:n_site) D_sparse_occ[i] = sum(W_tct[i]);
  {
    vector[n_site] invsqrtD_occ;
    for (i in 1:n_site) {
      invsqrtD_occ[i] = 1 / sqrt(D_sparse_occ[i]);
    }
    lambda_occ = eigenvalues_sym(quad_form(W_tct, diag_matrix(invsqrtD_occ)));
  }
  { // generate sparse representation for W_det
    int counter_det;
    counter_det = 1;
    // loop over upper triangular part of W to identify neighbor pairs
    for (i in 1:(total_surveys - 1)) {
      for (j in (i + 1):total_surveys) {
        if (W_bg[i, j] == 1) {
          W_sparse_det[counter_det, 1] = i;
          W_sparse_det[counter_det, 2] = j;
          counter_det += 1;
        }
      }
    }
  }
}

```

```

    }
  }
}
for (i in 1:total_surveys) D_sparse_det[i] = sum(W_bg[i]);
{
  vector[total_surveys] invsqrtD_det;
  for (i in 1:total_surveys) {
    invsqrtD_det[i] = 1 / sqrt(D_sparse_det[i]);
  }
  lambda_det = eigenvalues_sym(quad_form(W_bg, diag_matrix(invsqrtD_det)));
}
}

parameters {
  vector[n_site] phi_occ;
  vector[total_surveys] phi_det;
  real<lower = 0, upper =0.999> alpha_occ;
  real<lower = 0, upper =0.999> alpha_det;
  real<lower = 0> sigma_occ;
  real<lower = 0> sigma_det;

  vector[m_psi] beta_psi;
  vector[m_p] beta_p;
}
transformed parameters {
  vector[total_surveys] logit_p = X_bg * beta_p + phi_det * sigma_det;
  vector[n_site] logit_psi = X_tct * beta_psi + phi_occ * sigma_occ;
}
model {
  vector[n_site] log_psi = log_inv_logit(logit_psi);
  vector[n_site] log1m_psi = log1m_inv_logit(logit_psi);

  phi_occ ~ sparse_car(alpha_occ, W_sparse_occ, D_sparse_occ, lambda_occ, n_site, W_n_tct);
  alpha_occ ~ beta(4,1);
  sigma_occ ~ normal(0, 1.5);
  phi_det ~ sparse_car(alpha_det, W_sparse_det, D_sparse_det, lambda_det, total_surveys, W_n_bg);
  alpha_det ~ beta(4,1);
  sigma_det ~ normal(0, 1.5);

  beta_psi ~ student_t(7.763,0, 1.566);
  beta_p ~ student_t(7.763,0, 1.566);

  for (i in 1:n_site) {
    if (n_survey[i] > 0) {
      if (any_seen[i]) {
        // site is occupied
        target += log_psi[i]
                  + bernoulli_logit_lpmf(y[start_idx[i]:end_idx[i]] |
                                          logit_p[start_idx[i]:end_idx[i]]);
      } else {
        // site may or may not be occupied
        target += log_sum_exp(
          log_psi[i] + bernoulli_logit_lpmf(y[start_idx[i]:end_idx[i]] |
                                          logit_p[start_idx[i]:end_idx[i]]),

```

```
        log1m_psi[i]
    );
}
}
```