

Ecological inference with R: the *ecoreg* package

Version 0.1.1

Christopher Jackson
Department of Epidemiology and Public Health
Imperial College, London
`chris.jackson@imperial.ac.uk`

March 1, 2006

Abstract

In typical small-area studies of health and environment we wish to make inference on the relationship between individual-level quantities using aggregate, or ecological, data. Such ecological inference is often subject to bias and imprecision, due to the lack of individual-level information in the data. Simple regressions of area-level mean outcomes on area-level mean exposures are usually biased. Well-specified models for the group-level outcomes account for the within-area distribution of exposures. The *ecoreg* package can be used to fit this class of models for ecological inference from aggregate data. In addition, full outcome and covariate information from a survey of individuals within the areas can be used to improve bias and precision. *ecoreg* can be used to analyse ecological and individual data simultaneously using *hierarchical related regression*.

1 Ecological inference

Ecological studies analyse data defined at a group level, but aim to make inferences about the individuals within the groups. To make reliable individual-level inferences from these studies, a number of problems must be overcome. One crucial difficulty is that the group-level exposure-response relationship may not reflect the individual-level relationship, a problem known as *ecological bias*, or the *ecological fallacy*. See, for example, [1], [2], [3], [4] for discussion of these issues.

Denote the outcome count in area i , with population N_i , by y_i . To model y_i in terms of exposures measured as aggregate-level summaries, the usual model is a simple binomial regression on the area-level covariate means \bar{z}_{ir} . With binary covariates, \bar{z}_{ir} is the proportion exposed over the area. However, this only models the relationship between the *aggregate* exposures and outcomes. It is only justified as a method of estimating individual-level relationships if all individuals in the area have the same covariate value, or there is the same exposure-response relationship at the individual and aggregate levels, which is generally only true for linear models. With non-linear models, such as Binomial or Poisson models, using the same model form at both levels will lead to ecological bias [1].

Despite these problems, aggregate data *can* provide information about individual-level relationships. As the ratio of the between-area to the within-area variability of the exposure increases, the

aggregate data summarise the true distribution of the exposure more accurately, and contain more information about the true individual-level exposure-outcome relationship. With sufficient exposure information, and with *correctly-specified* models for the mean outcome ([1] [5] [6]), ecological bias can be reduced to negligible levels.

In general, successful ecological inference requires samples of individual-level data within areas. Individual exposure data are usually required to reduce ecological bias by accounting for the within-area variability of exposures. But further improvements can be made by using samples of exposures and *outcomes* for selected individuals, as discussed by Wakefield [7] for 2×2 tables, and more generally by Jackson *et al.*[8].

2 Models for aggregate and individual data

The models we describe here are also described in the papers by Jackson *et al.*[8] [9].

2.1 Individual data

We begin by specifying the form of the relationship between the individual-level risk of the binary outcome and the covariates. If individual-level exposure and outcome data are available, this is used to model them. It is also used as the basis for an equivalent model for the aggregate data. The risk p_{ij} of the individual-level outcome y_{ij} for the j th individual in area i is assumed to be a logit-linear function of the covariates. The most general model we consider is

$$\text{logit}(p_{ij}) = \mu_i + \sum_r \alpha_r x_{ir} + \sum_r \beta_r z_{ijr} + \gamma_{s_{ij}} \quad (1)$$

where x_{ir} are group-level covariates, and z_{ijr} are individual-level covariates. The group-level covariates may include descriptions of the socio-economic status of the area, or the health service provision in the area. Individual-level covariates might comprise individual behaviours such as smoking, demographics such as ethnicity, or individual indicators of wealth and social class. Individuals may be influenced by the overall average exposure in the area, in addition to their own, so that the group-level variables may include the means of certain individual-level variables. γ_s represents an additional contribution to the baseline risk for an individual occupying one of several strata s , usually defined by age and sex.

The baseline risk μ_i may be fixed at μ or considered as a random effect with some distribution across areas. This can account for any remaining overdispersion and heterogeneity between areas, after adjusting for observed area-level variables. A random effect also allows the borrowing of information across areas and can stabilise estimation from areas with small populations [10].

2.2 Aggregate data

2.2.1 Marginal model

Suppose the area-level exposures have been estimated from a survey. For example, in the UK census, aggregate data on social class and education are calculated using a 10% sample to maintain confidentiality. The proportion of smokers in the area might also be estimated from sales figures instead of

a census. Then, individual outcomes can be assumed to be independent and identically distributed, with risk equal to some marginal “group-level risk” p_i . We assume

$$y_i \sim \text{Bin}(N_i, p_i), \quad (2)$$

where p_i is determined by integrating the individual-level model over the joint within-area distribution of covariates [1, 5]. Thus p_i is the average risk for an individual in group i .

$$p_i = \int p_{ij}(\mathbf{x}) f_i(\mathbf{x}) dx = E_{\mathbf{x}}(p_{ij}(\mathbf{x})|i) \quad (3)$$

For a *single binary covariate*, observing the proportion exposed gives us enough information to estimate a binomial within-area distribution. For *continuous covariates* the mean is not sufficient to estimate the within-area distribution, and we generally need samples of individual covariate data to be able to estimate the within-area variability. For *multiple binary covariates* the joint distribution is estimated by the cross-classification of individuals between covariates. Typical census data do not usually have this cross-classification, and we need individual data to estimate it.

2.2.2 Conditional model

An alternative to the marginal model was proposed by Wakefield [7]. The binomial model (2) is based on the assumption that each individual in group i has an identical marginal probability p_i of outcome, integrating over their unknown exposures. If the binary exposure is known for *all* the individuals in the area from a full population census, but not necessarily coupled with exposures for the same individuals, then these should be conditioned on, leading to a likelihood based on a convolution of binomial distributions. The exact likelihood is related to the extended (or non-central) hypergeometric distribution, however Wakefield describes a more computationally convenient normal approximation.

2.2.3 Binary covariates

For clarity we demonstrate the marginal model for 3 binary covariates, however the framework extends immediately to any number of binary covariates. The integral to obtain the group-level risk is equivalent to a sum. Each individual falls into one of $S \times 2^3$ categories, defined by the distinct combinations of the 3 covariates and the S age-sex strata, and indexed by k . Let ϕ_{ik} be the probability that an individual occupies category k . Let q_{ik} be the probability of the individual outcome conditionally on occupying category k . Rewriting the index k as $\{k_1, k_2, k_3, s\}$, where k_r (0 or 1) indicates the presence or absence of covariate $r = 1, \dots, 3$.

$$p_i = \sum_k \phi_{ik} q_{ik} = \sum_{k_1, k_2, k_3, s} \phi_{\{i, k_1, k_2, k_3, s\}} q_{\{i, k_1, k_2, k_3, s\}}. \quad (4)$$

The outcome model conditionally on the unobserved category is

$$\text{logit}(q_{\{i, k_1, k_2, k_3, s\}}) = \mu_i + \text{logit}(e_s) + \sum_r \alpha_r x_{ir} + \sum_r k_r \beta_r \quad (5)$$

The effect β_r of the binary individual-level covariate r only enters into this equation when the covariate r is present. This is a generalisation of the model presented for two covariates by Lasserre *et*

al.[11], except that the outcome is assumed to be non-rare and binomial instead of Poisson. $\text{logit}(e_s)$ is a fixed offset, where e_s is the risk of the outcome in stratum s , estimated from national population data. This is similar in spirit to “indirect standardisation”, see, for example, [12]. If population and outcome totals are available for strata within areas, then we could generalise this to a separate binomial model for each stratum within area, with coefficients for the other covariates shared between the models, assuming no stratum-covariate interactions. It is also important to check whether *exposures* are correlated with strata. If not accounted for, this can lead to “mutual standardisation bias” [13].

We normally replace ϕ_{ik} in (4) by an estimate $\hat{\phi}_{ik}$. Ideally this would be the proportion of individuals in area i occupying category k . But typical census data are not sufficient to know the complete cross-classification of individuals between all of these covariate and strata categories. Usually, our only information is the marginal proportions of single covariates, for example, the proportion $\phi_i^{(1)}$ of individuals who are economically inactive. If we assume that the covariates are independent, then ϕ_{ik} can be estimated by the product of the proportions of individuals occupying each marginal category defining k . But generally, socioeconomic indicators, such as unemployment and social class, are highly correlated. [11] demonstrated that, in a typical case, bias is negligible when the joint covariate distribution is estimated by the product of the two marginal distributions, even when the covariates are correlated. However we may wish to study more than two covariates.

To estimate ϕ_{ik} we can often use a combination of marginal proportions \bar{z}_{ir} and individual covariate data z_{ijr} . In the context of the UK census, for example, these correspond to district-level aggregate data and the Samples of Anonymised Records. Let C_{ik} be the number of individuals in area i in this individual-level dataset occupying category k , computed from the z_{ijr} . We use the following two principles. Firstly, the estimated $\hat{\phi}_{ik}$ corresponding to categories k in which binary covariate r is 1 must sum to z_{ir} . Secondly, the ratio of estimates $\hat{\phi}_{ik}/\hat{\phi}_{il}$ must be the same as the ratio C_{ik}/C_{il} . This gives, for example, where R is the set of categories in which covariate r is 1,

$$\hat{\phi}_{ik} = C_{ik}\bar{z}_{ir} / \sum_{l \in R} C_{il} \quad (6)$$

2.2.4 Continuous covariates

Suppose now the individual-level model (1) depends only on an intercept and one continuous covariate x_{ij} . The ecological data consist of the within-area mean m_i of x_{ij} . In some cases, as well as the within-area mean, we may also have an estimate \hat{s}_i^2 of the within-area variance of x_{ij} , for example, from geographical modelling of an environmental exposure surface (e.g. Best *et al.*[14]). Then we suppose that these exposures are normally distributed, with $x_{ij} \sim N(m_i, s_i^2)$. If an exposure is not naturally normally distributed, it can often be transformed to normality. We can then calculate the area-specific risks (3) by integrating over f_i , here the density function of the normal distribution.

$$p_i = \int \text{expit}(\mu_i + \beta x) f_i(x) dx \quad (7)$$

Our assumed underlying model for $p_{ij}(x)$ is a logit-linear model on the exposures (10). In this case, the integral is not available in closed form. However, if we approximate the logit by a probit link function, then (7) evaluates to

$$p_i = \text{expit} \left\{ (1 + c^2 \beta^2 s_i^2)^{-1/2} (\mu_i + \beta m_i) \right\} \quad (8)$$

where $c = 16\sqrt{3}/(15\pi)$ (Salway and Wakefield [15]).

If, instead, we were using a Poisson model for a rare outcome, $y_i \sim \text{Pois}(N_i p_i)$, and a log-linear individual-level model $\log(p_{ij}) = \mu_i + \beta x_{ij}$, then the integrals can be evaluated explicitly without an approximation. Instead of (8), we would have (see, for example, [1])

$$p_i = \exp \left\{ \mu_i + \beta m_i + \frac{\beta^2 s_i^2}{2} \right\} \quad (9)$$

These methods generalise to multiple jointly normally-distributed covariates.

2.2.5 Semi-parametric approach

We have described a fully parametric model for ecological inference. Prentice and Sheppard [6] described an alternative semi-parametric approach based on estimating functions. This is often simply called the *aggregate data* method. Suppose a sample of covariates (binary or continuous) are available from a subset of n_i individuals, but not the corresponding outcomes on the same individuals. Broadly, the mean and variance of the total disease count y_i are calculated in terms of an *aggregate* risk $\frac{1}{n_i} \sum_j p_{ij}$. p_{ij} is the risk for individual j in area i , conditionally on their covariate values. This method is not currently implemented in *ecoreg*. This approach does not require a within-area distribution to be specified for the covariates. It requires samples of covariate data as an explicit part of the model. On the other hand, the parametric approaches described above, and implemented in *ecoreg*, require individual covariate data implicitly to estimate an appropriate within-area distribution.

2.3 Combining aggregate and individual data

To summarise the model for the ecological data, we have a binomial model (2) for the area-level outcome y_i . The corresponding area-level risk p_i is calculated explicitly in terms of the transformed group baseline risk μ_i , the individual-level covariate effects α and β , and the within-area distributions of the covariates.

It is easy to extend this model to include information from a small sample of individuals in each area whose outcome and exposures are known. We simply assume the binary outcome y_{ij} for such an individual j in area i is Bernoulli(p_{ij}), with a logit linear model for p_{ij} (1). Then the covariate effects α and β and the intercept μ_i are shared by the models for both the aggregate and individual-level data. Thus, we can fit a joint model which combines the information between the two sources of data. This is termed *hierarchical related regression*.

Note that it is not necessary to have individual data within all of the areas i . In practice, sample survey data will be available from varying numbers of individuals between areas.

3 Using *ecoreg*

The *ecoreg* package implements a fairly general case of the models described in Section 2. We assume you have already downloaded and installed the *ecoreg* package. To apply these methods, you should have one or both of

- An aggregate dataset with one record for each aggregate group, for example a geographical area, or a stratum within area, for example from a population census.

- An individual-level dataset, for example from a sample survey study. There need not be the same number of individuals per area, and there may be some areas in the aggregate dataset with no individuals.

These contain

- a disease outcome available on individuals as a binary response, or from areas as a number or proportion of disease cases.
- any number of binary explanatory variables, or exposures, available directly from individuals or as proportions over areas.
- any number of continuous exposures with a multivariate normal distribution within areas. These are available directly from individuals, or from the aggregate data as area-level means and, optionally within-area covariances.
- any number of contextual explanatory variables, that is, characteristics specific to areas, constant within areas.
- optionally, the joint within-area distribution of the n binary exposures, available as a matrix with 2^n columns, and the same number of rows as the aggregate data, containing the number of individuals in each area with each of the distinct combinations of the exposures.

3.1 Simulated example

These models are illustrated here using a simulated dataset of the form studied in Jackson *et al.* [8]. There are 100 groups of 1,000 individuals each. The binary disease status of individual j in group i is simulated from the following model

$$\text{logit}(p_{ij}) = \mu_i + \alpha x_{ij}^{(1)} + \beta x_{ij}^{(2)}. \quad (10)$$

$x_{ij}^{(1)}$ represents smoking status (a binary covariate) and $x_{ij}^{(2)}$ represents pollution exposure (a continuous covariate). The full details of how the covariate data were constructed are described by Jackson *et al.* [8]. The covariate effects are $\alpha = \log(2)$, $\beta = \log(2.3)$. μ_i are chosen as 10 equally-spaced quantiles of a normal distribution, with mean $\text{logit}(0.1)$ and standard deviation 0.2, giving a 95% sampling interval for the baseline disease risk of (0.07, 0.14).

The complete individual-level data are aggregated over areas to form an aggregate dataset, and samples of 5% are taken from each within-area population to form an individual-level dataset. These hypothetical aggregate and individual datasets are provided with the *ecoreg* package, under the names `agg` and `indiv`.

3.2 Format of datasets in more detail

Format of aggregate dataset This should be a data frame with a row corresponding to an area or group. It should have at minimum an outcome variable, with the number of events, such as disease cases, in the area. In addition, any number of aggregate covariates can be specified. Often these will be binary covariates, expressed as proportions over the area. For covariates that are continuous at the individual-level, these should be specified as within-area means, and also within-area variances, after

transformation to an approximate normal distribution. For example, we print the first ten rows of the `agg` data.

```
> library(ecoreg)
> data(agg)
> options(digits = 2)
> agg[1:10, ]
```

	disease	n.smoke	n	p.smoke	poll.mean	poll.sd
1	218	17	990	0.017	1.41	0.24
2	198	18	990	0.018	1.04	0.53
3	169	18	990	0.018	0.98	0.43
4	199	18	990	0.018	1.01	0.35
5	258	17	990	0.017	1.46	0.32
6	256	18	990	0.018	1.25	0.33
7	253	18	990	0.018	1.23	0.36
8	310	18	990	0.018	1.39	0.36
9	297	18	990	0.018	1.35	0.39
10	303	18	990	0.018	1.16	0.27

The number of individuals with the disease, the number of smokers, the population of the area, the proportion of smokers, and the mean and standard deviation of the pollution exposure are labelled `disease`, `n.smoke`, `n`, `p.smoke`, `poll.mean` and `poll.sd`, respectively.

Format of individual dataset The individual dataset should be a data frame with each row corresponding to an individual. Variables may include a binary outcome and any number of covariates. For example, the first 15 rows of the `indiv` data are illustrated.

```
> data(indiv)
> indiv[1:15, ]
```

	disease	smoke	poll	area
1	0	0	1.52	1
2	0	0	1.11	1
3	0	0	1.52	1
4	0	0	1.00	1
5	1	1	1.65	1
6	0	0	1.54	1
7	0	0	1.54	1
8	0	0	1.68	1
9	0	0	1.60	1
10	0	0	1.08	1
11	0	0	0.99	2
12	0	0	2.09	2
13	0	0	1.41	2
14	0	0	0.91	2
15	0	0	1.00	2

The disease status of the individual (0 or 1 corresponding to no disease and disease) , whether the individual smoked (0 or 1 corresponding to no and yes), the pollution exposure and the area indicator are labelled `disease`, `smoke`, `poll` and `area`, respectively. The area indicator is only necessary when using models with random area effects.

3.3 Calling `eco`

The main R function in *ecoreg* is called `eco`. This is used to fit a model to one of these datasets, or a combination of the two. The R help page for *eco* fully describes each of the function's arguments. Now we give simple examples of models that might be used in practice.

3.3.1 Models for aggregate data

Contextual model Firstly we fit a model which considers the proportion of smokers as a *contextual*, or area-level variable. This is a simple binomial regression of the area disease count in terms of the proportion of smokers. This should give a biased model for the individual-level association, because we have not accounted for the within-area distribution of smoking behaviour.

The first argument of `eco` (called `formula`, but it is not necessary to label the argument) is a formula, as used in most statistical modelling functions in R such as `lm` and `glm`. It specifies the *aggregate* component of the model, that is, the names of any covariates included in x_{ir} (equation (1)). The argument `data` specifies a data frame which should contain all aggregate-level variables specified in the call to `eco`.

```
> eco(cbind(disease, n) ~ p.smoke, data = agg)
```

Call:

```
eco(formula = cbind(disease, n) ~ p.smoke, data = agg)
```

Aggregate-level odds ratios:

	OR	195	u95
(Intercept)	0.31	0.3	0.32
p.smoke	3.61	2.7	4.74

No individual-level covariates

```
-2 x log-likelihood: 2141
```

The `eco` function returns objects of class `ecoreg`. Printing an object of this class displays the estimated odds ratios $\exp(\alpha)$ associated with aggregate-level covariates, and odds ratios $\exp(\beta)$ associated with individual covariates (equation (1), along with their 95% confidence intervals, and $-2 \times$ the maximised log-likelihood. In this example, the maximum likelihood estimate of the odds ratio associated with the aggregate proportion of smokers is 3.6, a substantial overestimate of the true individual-level odds ratio of 2.

Marginal model, binary covariate To use the aggregate proportion of smokers to determine the relative risk associated with *individual* smoking, using the marginal model (2–3), use the binary

argument to `eco`, as follows. Notice that the first argument contains a `1`, as there are no aggregate-level effects in this model.

```
> eco(cbind(disease, n) ~ 1, binary = ~p.smoke,
+     data = agg)
```

Call:

```
eco(formula = cbind(disease, n) ~ 1, binary = ~p.smoke, data = agg)
```

Aggregate-level odds ratios:

```
OR 195 u95
(Intercept) 0.31 0.3 0.32
```

Individual-level odds ratios:

```
OR 195 u95
p.smoke 3 2.4 3.8
```

```
-2 x log-likelihood: 2141
```

This gives a more accurate estimate of the true relative risk of 2.

Now we construct a regression model including the effect of both smoking and pollution. Firstly, we include area mean pollution exposure as a contextual variable, and smoking as an individual variable.

```
> eco(cbind(disease, n) ~ poll.mean, binary = ~p.smoke,
+     data = agg)
```

Call:

```
eco(formula = cbind(disease, n) ~ poll.mean, binary = ~p.smoke,
     data = agg)
```

Aggregate-level odds ratios:

```
OR 195 u95
(Intercept) 0.12 0.11 0.13
poll.mean 2.29 2.17 2.41
```

Individual-level odds ratios:

```
OR 195 u95
p.smoke 1.8 1.4 2.3
```

```
-2 x log-likelihood: 1246
```

As pollution exposure varies within the areas, this is a biased model for the individual-level association between pollution exposure and disease.

Marginal model, continuous covariate Here we account for the individual-level pollution exposure using a well-specified model which accounts for the within-area pollution variance. Pollution

exposure and smoking are considered as individual predictors of disease risk. The `normal` argument to `eco` is a formula whose right-hand side should contain variables denoting the group-level means of the normally-distributed covariates. These covariates will then be fitted as individual-level effects, using a model of the form of equation (8) by default, or (9) if `model = "poisson"` is specified. The `norm.var` is used to supply the corresponding group-level variances.

```
> eco(cbind(disease, n) ~ 1, binary = ~p.smoke,
+     normal = ~poll.mean, norm.var = poll.sd, data = agg)
```

Call:

```
eco(formula = cbind(disease, n) ~ 1, binary = ~p.smoke, normal = ~poll.mean,
    data = agg, norm.var = poll.sd)
```

Aggregate-level odds ratios:

```
          OR  195  u95
(Intercept) 0.12 0.11 0.13
```

Individual-level odds ratios:

```
          OR  195  u95
p.smoke    1.9  1.5  2.4
poll.mean  2.2  2.1  2.3
```

```
-2 x log-likelihood: 1266
```

3.3.2 Combining aggregate and individual data

Next we attempt to improve these estimates further by using the information from samples of individuals, as described in Section 2.3. The individual-level regression model is given in the `iformula` argument. The name of the individual-level dataset, in which the variables in the individual-level model should appear, is given in the `idata` argument.

```
> eco(cbind(disease, n) ~ 1, binary = ~p.smoke,
+     normal = ~poll.mean, norm.var = poll.sd, data = agg,
+     iformula = disease ~ smoke + poll, idata = indiv)
```

Call:

```
eco(formula = cbind(disease, n) ~ 1, binary = ~p.smoke, normal = ~poll.mean,
    iformula = disease ~ smoke + poll, data = agg, idata = indiv,
    norm.var = poll.sd)
```

Aggregate-level odds ratios:

```
          OR  195  u95
(Intercept) 0.12 0.11 0.13
```

Individual-level odds ratios:

```
          OR  195  u95
smoke    1.9  1.6  2.4
```

```
poll 2.2 2.1 2.3
```

```
-2 x log-likelihood: 2351
```

The bias and precision of the estimates are much improved.

When combining individual and aggregate data, the models at both levels *should be the same*, with the same covariates. *All* contextual, binary, categorical and continuous covariates from the aggregate component should appear in the `iformula` formula.

Individual dataset alone `eco` can also be run using individual data alone.

```
> eco(iformula = disease ~ smoke + poll, idata = indiv)
```

Call:

```
eco(iformula = disease ~ smoke + poll, idata = indiv)
```

Aggregate-level odds ratios:

```
OR 195 u95  
(Intercept) 0.11 0.07 0.17
```

Individual-level odds ratios:

```
OR 195 u95  
smoke 2.0 1.3 3.1  
poll 2.2 1.6 3.0
```

```
-2 x log-likelihood: 1083
```

This is a logistic regression, which `eco` fits using maximum likelihood. Of course, this model can also be fitted using the `glm` function in R, which uses an iteratively reweighted least squares algorithm. Here we manipulate the output of `glm` from linear regression coefficients and standard errors on the logit scale to odds ratios with 95% confidence intervals, which are the same as produced by `eco`.

```
> indiv.glm <- glm(disease ~ smoke + poll, data = indiv,  
+ family = "binomial")  
> est <- coef(indiv.glm)  
> se <- coef(summary(indiv.glm))[, 2]  
> exp(cbind(est, est - qnorm(0.975) * se, est +  
+ qnorm(0.975) * se))
```

```
est  
(Intercept) 0.11 0.07 0.17  
smoke 2.04 1.33 3.12  
poll 2.21 1.63 3.00
```

The precision of these estimates is much less than those from combining the aggregate and individual data.

3.4 Importance of the between-area exposure contrasts

The amount of individual-level information in ecological data increases as the between-area to the within-area variability of the exposure increases. When there are low exposure contrasts between areas, inference may be improved by combining the ecological data with individual-level data, as described by Jackson *et al.* [8]. We demonstrate this with a simulated example, which also appears at the end of `help(eco)`.

Firstly, we simulate aggregate data consisting of 50 groups of 100 individuals. Two contextual covariates (labelled deprivation and mean income) are generated as standard normal variables. Two binary covariates, interpreted as the proportion of non-white individuals and smokers in each area, are generated from uniform distributions. The data frame `sim.df` contains the ecological covariate data.

```
> ng <- 50
> N <- rep(100, ng)
> set.seed(31412)
> ctx <- cbind(deprivation = rnorm(ng), mean.income = rnorm(ng))
> phi <- cbind(nonwhite = runif(ng), smoke = runif(ng))
> sim.df <- as.data.frame(cbind(ctx, phi))
> sim.df[1:5, ]
```

	deprivation	mean.income	nonwhite	smoke
1	0.813	-0.8614	0.64	0.80
2	-1.029	-1.6877	0.87	0.80
3	-0.045	-0.0035	0.18	0.12
4	0.089	1.1411	0.66	0.45
5	-1.788	-1.0944	0.36	0.35

A disease outcome with approximate 5% baseline prevalence, and odds ratios of 1.01, 1.02, 1.5 and 2 respectively for the four covariates, is now simulated. The function `sim.eco` is provided to simulate ecological outcome data and individual sample data, in terms of known covariates, baseline risks and odds ratios. Firstly we use `sim.eco` to generate an ecological outcome alone (`sim1`) and secondly an ecological outcome and individual sample data (`sim2`).

```
> mu <- qlogis(0.05)
> alpha.c <- log(c(1.01, 1.02))
> alpha <- log(c(1.5, 2))
> sim1 <- sim.eco(N, ctx = ~deprivation + mean.income,
+   binary = ~nonwhite + smoke, data = sim.df,
+   mu = mu, alpha.c = alpha.c, alpha = alpha)
> sim2 <- sim.eco(N, ctx = ~deprivation + mean.income,
+   binary = ~nonwhite + smoke, data = sim.df,
+   mu = mu, alpha.c = alpha.c, alpha = alpha,
+   isam = 7)
```

The return value of `sim.eco` has a component `y` containing the ecological outcome data (the number of individuals in each area with the outcome), and a component `idata` containing the individual sample data. Here we have specified `isam=7` in the call to `sim.eco`, producing an individual sample dataset with 7 individuals for each of the 50 areas.

```

> sim1$y[1:5]

[1] 14 17 7 8 6

> sim2$idata[1:15, ]

  group y deprivation mean.income nonwhite smoke
1     1 0      0.813     -0.8614        1     1
2     1 1      0.813     -0.8614        0     1
3     1 1      0.813     -0.8614        1     0
4     1 0      0.813     -0.8614        1     1
5     1 0      0.813     -0.8614        0     0
6     1 0      0.813     -0.8614        0     1
7     1 0      0.813     -0.8614        1     0
8     2 0     -1.029     -1.6877        1     1
9     2 1     -1.029     -1.6877        1     1
10    2 0     -1.029     -1.6877        1     1
11    2 0     -1.029     -1.6877        1     1
12    2 0     -1.029     -1.6877        1     0
13    2 0     -1.029     -1.6877        1     0
14    2 0     -1.029     -1.6877        1     1
15    3 0     -0.045     -0.0035        1     1

```

Next we fit the correct model to the simulated data, with two contextual covariates and two individual binary covariates. We see that in this case, combining with the individual sample data does not improve the precision of the estimates.

```

> aggdata <- as.data.frame(cbind(y = sim1$y, sim.df))
> aggdata[1:5, ]

  y deprivation mean.income nonwhite smoke
1 14      0.813     -0.8614    0.64 0.80
2 17     -1.029     -1.6877    0.87 0.80
3 7      -0.045     -0.0035    0.18 0.12
4 8       0.089      1.1411    0.66 0.45
5 6      -1.788     -1.0944    0.36 0.35

> agg.eco <- eco(cbind(y, N) ~ deprivation + mean.income,
+               binary = ~nonwhite + smoke, data = aggdata)
> agg.eco

```

Call:

```

eco(formula = cbind(y, N) ~ deprivation + mean.income, binary = ~nonwhite +
    smoke, data = aggdata)

```

Aggregate-level odds ratios:

```

OR    195    u95

```

```
(Intercept) 0.054 0.039 0.075
deprivation 0.953 0.868 1.045
mean.income 1.007 0.922 1.100
```

Individual-level odds ratios:

```
      OR 195 u95
nonwhite 1.6 1.1 2.3
smoke    2.0 1.4 2.8
```

-2 x log-likelihood: 238

```
> agg.indiv.eco <- eco(cbind(y, N) ~ deprivation +
+   mean.income, binary = ~nonwhite + smoke, iformula = y ~
+   deprivation + mean.income + nonwhite + smoke,
+   data = aggdata, idata = sim2$idata)
> agg.indiv.eco
```

Call:

```
eco(formula = cbind(y, N) ~ deprivation + mean.income, binary = ~nonwhite +
  smoke, iformula = y ~ deprivation + mean.income + nonwhite +
  smoke, data = aggdata, idata = sim2$idata)
```

Aggregate-level odds ratios:

```
      OR 195 u95
(Intercept) 0.053 0.039 0.071
deprivation 0.962 0.880 1.052
mean.income 0.999 0.918 1.087
```

Individual-level odds ratios:

```
      OR 195 u95
nonwhite 1.6 1.2 2.2
smoke    2.0 1.5 2.8
```

-2 x log-likelihood: 464

However, suppose we simulate data with a much lower between-area exposure variance, such as a $\text{uniform}(0, 0.2)$ distribution for proportion of non-white ethnicity and $\text{uniform}(0.1, 0.3)$ for proportion of smokers. The aggregate data now contain little information about the individual-level effects, and we obtain highly imprecise estimates of the true individual model.

```
> phi <- cbind(nonwhite = runif(ng, 0, 0.2), smoke = runif(ng,
+   0.1, 0.3))
> sim.df <- as.data.frame(cbind(ctx, phi))
> sim1 <- sim.eco(N, ctx = ~deprivation + mean.income,
+   binary = ~nonwhite + smoke, data = sim.df,
+   mu = mu, alpha.c = alpha.c, alpha = alpha)
> sim2 <- sim.eco(N, ctx = ~deprivation + mean.income,
```

```

+   binary = ~nonwhite + smoke, data = sim.df,
+   mu = mu, alpha.c = alpha.c, alpha = alpha,
+   isam = 10)
> aggdata <- as.data.frame(cbind(y = sim1$y, sim.df))
> agg.eco <- eco(cbind(y, N) ~ deprivation + mean.income,
+   binary = ~nonwhite + smoke, data = aggdata)
> agg.eco

```

Call:

```

eco(formula = cbind(y, N) ~ deprivation + mean.income, binary = ~nonwhite +
  smoke, data = aggdata)

```

Aggregate-level odds ratios:

	OR	195	u95
(Intercept)	0.074	0.047	0.12
deprivation	1.018	0.912	1.14
mean.income	1.059	0.953	1.18

Individual-level odds ratios:

	OR	195	u95
nonwhite	0.18	0.00011	286
smoke	0.85	0.09077	8

-2 x log-likelihood: 233

To be able to estimate the individual-level effects most accurately, we combine the aggregate data with the individual-level sample data.

```

> agg.indiv.eco <- eco(cbind(y, N) ~ deprivation +
+   mean.income, binary = ~nonwhite + smoke, iformula = y ~
+   deprivation + mean.income + nonwhite + smoke,
+   data = aggdata, idata = sim2$idata)
> agg.indiv.eco

```

Call:

```

eco(formula = cbind(y, N) ~ deprivation + mean.income, binary = ~nonwhite +
  smoke, iformula = y ~ deprivation + mean.income + nonwhite +
  smoke, data = aggdata, idata = sim2$idata)

```

Aggregate-level odds ratios:

	OR	195	u95
(Intercept)	0.061	0.048	0.078
deprivation	1.048	0.946	1.162
mean.income	1.056	0.956	1.167

Individual-level odds ratios:

	OR	195	u95

```
nonwhite 0.82 0.28 2.5
smoke    1.37 0.60 3.1
```

```
-2 x log-likelihood: 443
```

This raises the question of whether the aggregate data do contribute any information in this case, and whether we can do just as well by analysing the individual data alone. But we find that precision is lower when only the individual data are included.

```
> indiv.eco <- eco(iformula = y ~ deprivation +
+   mean.income + nonwhite + smoke, idata = sim2$idata)
> indiv.eco
```

Call:

```
eco(iformula = y ~ deprivation + mean.income + nonwhite + smoke,
    idata = sim2$idata)
```

Aggregate-level odds ratios:

```
OR 195 u95
(Intercept) 0.046 0.028 0.076
```

Individual-level odds ratios:

```
OR 195 u95
deprivation 1.3 0.95 1.9
mean.income 1.1 0.77 1.5
nonwhite    1.1 0.31 3.8
smoke       1.4 0.53 3.5
```

```
-2 x log-likelihood: 206
```

3.5 Other features

3.5.1 Within-area distribution of binary covariates

To build the model (3) with more than one binary covariate, by default `eco` assumes that the covariates are independent within areas. Often this assumption is not appropriate, especially when considering, for example, socio-economically related factors.

To account for the joint within-area distribution of a set of binary or categorical covariates, use the `cross` argument to `eco`. This should be a matrix containing the same number of rows as the aggregate data, and number of columns equal to the distinct number of covariate categories into which an individual can belong. For full details of how to specify `cross`, refer to the help page for the `eco` function. The `cross` needs to be calculated by the user before calling `eco`. Individual data may be required to estimate `cross`, as typical census data do not give detailed cross-classification tables.

This is now illustrated with a hypothetical example. We introduce another covariate called `p.soclass` into the `agg` data representing the proportion of individuals in a lower social class. Suppose this varies uniformly over areas between 0.1 and 0.6. This is likely to be correlated with

smoking at the individual level. Suppose that we have an information from an individual-level survey, suggesting that an individual is twice as likely to smoke if in a lower social class. We wish to use this information to construct a matrix with 100 rows and four columns, representing estimates of the proportion of individuals who are in each of four categories:

1. neither smoke nor are in the lower social class (p_{00})
2. smoke but are not in the lower social class (p_{10})
3. do not smoke but are in the lower social class (p_{01})
4. smoke and are in the lower social class. (p_{11})

Let p_A be the proportion of smokers, and p_B be the proportion of individuals in the lower social class, in an area. We know that $p_{10} + p_{11} = p_A$, $p_{01} + p_{11} = p_B$ and, from the survey, $(p_{11}/p_B)/(p_{10}/(1 - p_B)) = 2$. This leads, for example, to

$$p_{11} = 2p_A p_B / (1 + p_B)$$

This is enough information to construct all the estimated cross-classification probabilities.

```
> agg$p.socclass <- seq(0.1, 0.6, length = 100)
> pa <- agg$p.smoke
> pb <- agg$p.socclass
> p11 <- pa * pb * 2/(1 + pb)
> p10 <- pa - p11
> p01 <- pb - p11
> p00 <- 1 - (p01 + p10 + p11)
> cross <- cbind(p00, p10, p01, p11)
> cross[1:5, ]
```

```
      p00  p10  p01  p11
[1,] 0.89 0.014 0.097 0.0031
[2,] 0.88 0.015 0.102 0.0035
[3,] 0.88 0.015 0.106 0.0036
[4,] 0.87 0.014 0.111 0.0038
[5,] 0.87 0.013 0.117 0.0037
```

This is the cross-classification matrix that we can supply to `eco` if we wanted to construct an ecological model including both smoking and social class.

```
> eco(cbind(disease, n) ~ 1, binary = ~p.smoke +
+       p.socclass, cross = cross, data = agg)
```

To account for the joint within-area distribution of a set of continuous covariates, use the `norm.var` argument to `eco`, which specifies the joint covariance matrix of the covariates, assumed normally distributed, for each area. For further details of how to specify `norm.var`, refer to the help page for the `eco` function.

Currently, `ecoreg` does not support specifying the within-area covariance between binary and continuous covariates.

3.5.2 Stratification

To account for differing disease risks in strata defined by age and sex, using fixed offsets determined from the whole population (as in equations 1 and 5), use the `strata`, `pstrata` and `istrata` arguments to `eco`.

- `pstrata` should be a vector with one element for each stratum, giving the assumed baseline outcome probabilities for the strata.
- `strata` should be a matrix with the same number of rows as the aggregate data. Rows represent areas, and columns represent the strata occupancy *probabilities* for those areas (often estimated as observed occupancy *proportions*). Alternatively, to account for within-area correlation between strata membership and binary covariate status, the cross-classification between strata and covariates can be specified in the `cross` argument. See the help page to `eco`.
- If individual data are modelled, `istrata` should be a variable containing the individual-level variable indicating the stratum an individual occupies. This should be a factor, whose levels correspond to the columns of the matrix `strata`.

3.5.3 Categorical covariates

As well as binary covariates, categorical covariates can also be fitted as individual-level predictors. The aggregate data for categorical covariates must be supplied separately from the main aggregate dataset, in the `categorical` argument to `eco`. See the help page for `eco`. In practice, there is not likely to be enough information in ecological data for successful ecological inference on categorical variables with large numbers of categories.

3.5.4 Random effects models

By default, `eco` assumes the baseline risk μ_i (equation 1) is constant μ between areas i . Optionally, `eco` can also fit μ_i as a normally-distributed random effect, using adaptive Gauss-Hermite integration [16].

If `random=TRUE` is specified in the call to `eco`, an area-level random intercept is included in the model. In this case the data should indicate which area each row of the data corresponds to. In the individual data, `igroups` should give the name of a variable containing the group identifiers of the individual-level data. In the aggregate data, by default, the groups are the row numbers of the dataset. Alternatively, `groups` specifies a group-level variable containing the group identifiers to be matched with the groups given in `igroups`.

The Gauss-Hermite integration can be controlled by the arguments `gh.points` and `iter.adapt` to `eco`. `gh.points` gives the number of points to use for quadrature, while `iter.adapt` gives the number of iterations to use for the adaptive phase of the algorithm.

Random effects model fitting is relatively slow, and it may be useful to view the progress of the model fitting by specifying a `control` argument, such as `control=list(trace=1, REPORT=1)`. This is passed from `eco` to `optim`, the R function which performs optimisation of the likelihood. See `help(optim)` for further options to control optimisation.

4 Warnings and limitations

- It is easy to over-fit models, especially with several covariates. Often there is not enough information available in aggregate data.
- When fitting many covariates, it is essential to account for the within-area distribution.
- Continuous covariates must be normally-distributed or able to be transformed to normality.
- Only limited error-checking is performed. `eco` may fail with an incomprehensible error message if the model or data are specified wrongly or inconsistently.

5 *eco* reference guide

The R help page for `eco` gives details of all the allowed arguments and options to the `eco` function. To view this online in R, type:

```
> help(eco)
```

Similarly all the other functions in the package have help pages, which should always be consulted in case of doubt about how to call them. The web-browser based help interface may be convenient - type

```
> help.start()
```

and navigate to `Packages ... ecoreg`, which brings up a list of all the functions in the package with links to their documentation, and a link to this manual in PDF format.

6 Similar software

- WinBUGS (<http://www.mrc-bsu.cam.ac.uk/bugs>) can be used to fit these models from a Bayesian perspective using Markov Chain Monte Carlo simulation. This approach is amenable to extension to account for complexities such as random intercepts, random coefficients, spatial correlation and measurement error. WinBUGS files containing worked examples of a range of these models, using methods described by Jackson *et al.* [?] [9] are provided at <http://www.bias-project.org.uk/software>.
- The R package *MCMCpack*, available from CRAN, implements ecological inference for 2×2 tables using a Bayesian hierarchical model described by Wakefield [7].
- The R package *eco*, available from CRAN, implements ecological inference for 2×2 tables, using methods described by Imai and Lu [17].
- *EI* and *EzI* [18] by Kenneth Benoit and Gary King, implementing methods from King [19].

References

- [1] S. Richardson, I. Stucker, and D. Hémon. Comparison of relative risks obtained in ecological and individual studies: some methodological considerations. *16*(1):111–120, 1987.
- [2] S. Greenland and H. Morgenstern. Ecological bias, confounding and effect modification. *18*:269–284, 1989.
- [3] S. Greenland and J. Robins. Ecological studies — biases, misconceptions and counterexamples. *139*:747–760, 1994.
- [4] S. Richardson and C. Monfort. Ecological correlation studies. In *Spatial Epidemiology*, chapter 11, pages 205–220. Oxford University Press, Oxford, 2000.
- [5] J. Wakefield and R. Salway. A statistical framework for ecological and aggregate studies. *164*(1):119–137, 2001.
- [6] R. L. Prentice and L. Sheppard. Aggregate data studies of disease risk factors. *Biometrika*, *82*:113–125, 1995.
- [7] J. Wakefield. Ecological inference for 2×2 tables (with discussion). *167*(3):385–445, 2004.
- [8] C. H. Jackson, N. G. Best, and S. Richardson. Improving ecological inference using individual-level data. 2006. (in press).
- [9] C. H. Jackson, N. G. Best, and S. Richardson. Hierarchical related regression for combining aggregate and survey data in studies of socio-economic disease risk factors. *Technical report, Imperial College, London*, 2006. URL: <http://www.bias-project.org.uk/papers/hrr.pdf>.
- [10] S. Richardson and N. Best. Bayesian hierarchical models in ecological studies of health-environment effects. *Environmetrics*, *14*:129–147, 2003.
- [11] V. Lasserre, C. Guihenneuc-Jouyaux, and S. Richardson. Biases in ecological studies: utility of including within-area distribution of confounders. *19*:45–59, 2000.
- [12] D. Clayton and M. Hills. *Statistical Models in Epidemiology*. Oxford University Press, 1993.
- [13] P.R. Rosenbaum and D.B. Rubin. Difficulties with regression analyses of age-adjusted rates. *Biometrics*, *40*:437–443, 1984.
- [14] N. Best, S. Cockings, J. Bennett, J. Wakefield, and P. Elliott. Ecological regression analysis of environmental benzene exposure and childhood leukaemia: sensitivity to data inaccuracies, geographical scale and ecological bias. *164*(1):155–174, 2001.
- [15] R. Salway and J. Wakefield. Sources of bias in ecological studies of non-rare events. *Environmental and Ecological Statistics*, *12*(3):321–347, 2005.
- [16] Q. Liu and D. A. Pierce. A note on Gauss-Hermite quadrature. *Biometrika*, *81*:624–629, 1994.
- [17] K. Imai and Y. Lu. An incomplete data approach to the ecological inference problem. (*working paper, Princeton University*), 2006. (URL: <http://imai.princeton.edu/research/coarse.html>).

- [18] G. King. Ei: A program for ecological inference. *Journal of Statistical Software*, 11(7), 2004.
- [19] G. King. *A Solution to the Ecological Inference Problem: Reconstructing Individual Behavior from Aggregate Data*. Princeton University Press, 1997.