1    **Fitting occupancy models with E-SURGE:**

2    **Hidden Markov modelling of presence-absence data**

3

4    Olivier Gimenez[1], Laetitia Blanc[1], Aurélien Besnard[1], Roger Pradel[1], Paul F. Doherty Jr[2], and

5                            Rémi Choquet[1]

6

7    [1]Centre d'Ecologie Fonctionnelle et Evolutive, UMR 5175, campus CNRS, 1919 Route de

8    Mende, 34293 Montpellier Cedex 5, France.

9    [2]Department of Fish, Wildlife and Conservation Biology, Colorado State University, Fort

10   Collins, CO 80523-1474, USA

11

Running title: occupancy in E-SURGE

1    *Abstract.*

2    1. Occupancy – the proportion of area occupied by a species – is a key notion for addressing

3    important questions in ecology, biogeography and conservation biology. Occupancy models

4    allow estimating and inferring about species occurrence while accounting for false absences

5    (or imperfect species detection).

6    2. Most occupancy models can be formulated as hidden Markov models (HMM) in which the

7    state process captures the Markovian dynamic of the actual but latent states while the

8    observation process consists of observations that are made from these underlying states.

9    3. We show how occupancy models can be implemented in program E-SURGE, which was

10   initially developed to analyse capture-recapture data in the HMM framework. Replacing

11   individuals by sites provides the user with access to several features of E-SURGE that are *not*

12   *available altogether* or *just not available* in standard occupancy software: i) user-friendly

13   model specification through a SAS/R-like syntax without having to write custom code, ii)

14   decomposition of the observation and state processes in several steps to provide flexible

15   parameterisation, iii) up-to-date diagnostics of model identifiability and iv) advanced

16   numerical algorithms to produce fast and reliable results (including site random effects).

17   4. To illustrate E-SURGE features, we provide simulated data and the details of the

18   implementation on the analysis of several occupancy models. These detailed examples are

19   gathered in a companion wiki platform http://occupancyinesurge.wikidot.com/.

20

21   *Key words*: capture-recapture; detectability; detection-nondetection; E-SURGE; hidden

22   Markov models; presence-absence; species occurrence.

23

2

1 **INTRODUCTION**

2  Occupancy models allow estimating and inferring about species occurrence while

3 accounting for false absences or imperfect species detection (MacKenzie *et al.* 2006). These

4 models have been extensively used to address important questions in fields as diverse as

5 conservation biology, biogeography, wildlife epidemiology, metapopulation dynamics and

6 community ecology (review in (Bailey, MacKenzie, & Nichols 2013)).

7  Following the seminal work of MacKenzie and colleagues (MacKenzie *et al.* 2006), it

8 was soon realized that occupancy models could be formulated as hidden Markov models

9 (HMMs) in which two time series run in parallel: the state process captures the Markovian

10 dynamic of the actual but latent states (e.g., site occupied vs. unoccupied) while the

11 observation process consists of observations that are made from these underlying states (e.g.,

12 species detected vs. undetected) (e.g., Royle & Kéry 2007).

13  There is an intimate connection between occupancy and capture-recapture models that

14 can be realized by interchanging individuals and sites. Interestingly, the formulation of

15 capture-recapture models as HMMs was also witnessed in the capture-recapture literature

16 (Pradel 2005; Gimenez *et al.* 2012).

17  Several software are available to fit occupancy models, either in the Frequentist

18 framework with programs PRESENCE (Hines 2013), MARK (White & Burnham 1999) and

19 the R package Unmarked (Fiske & Chandler 2011), or in the Bayesian framework using

20 WinBUGS (Kéry & Schaub 2011). WinBUGS requires writing custom code and knowledge

21 about the Bayesian theory. Programs PRESENCE and MARK often require the construction

22 of design matrices to specify models, a process that can be error-prone. PRESENCE, MARK

23 and Unmarked do not incorporate random effects.

24  Here, by exploiting the equivalence between occupancy and capture-recapture models,

25 we illustrate how occupancy models can be implemented in program E-SURGE (Choquet,

1   Rouan, & Pradel 2009) which was initially developed to analyse capture-recapture data in the

2   HMM framework. We aim at providing the user with access to features of E-SURGE that are

3   *not implemented altogether* in available occupancy software, namely i) user-friendly model

4   specification through a SAS/R-like syntax without having to write custom code, ii) advanced

5   numerical algorithms to produce fast and reliable results, including the incorporation of site

6   random effects, and several other features that are simply *not implemented* in these software,

7   namely iii) decomposition of the observation and state processes in several steps to provide

8   flexible parameterisation and iv) up-to-date diagnostics of model identifiability, in other

9   words a reliable way of counting the number of parameters entering the calculation of the

10  Akaike Information Criterion.

11

12  **HIDDEN MARKOV MODELLING OF OCCUPANCY DATA**

13          In Figure 1, we provide the HMM formulation of the general dynamic occupancy

14  models to carry out inference about occurrence and how extinction and colonization drive

15  changes in occurrence.

16

17                              [FIGURE 1 AROUND HERE]

18

19          The parameters of interest are the probability of local extinction $\varepsilon$ and of colonization

20  $\gamma$ as well as the detection probability $p$ and the probability of initial occupancy $\psi_1$ where we

21  have assumed all parameters constant across periods and sites for simplicity. A HMM is built

22  around three pieces of information: the vector of initial state probabilities, the matrix of

23  transition probabilities linking states in successive sampling occasions and the matrix of

24  observation probabilities linking observations and states at a given occasions. At the first

1    sampling occasion $t = 1$, with the first state being 'unoccupied' and the second 'occupied', the

2    vector of initial state probabilities is:

3

$$\begin{bmatrix} 1-\psi_1 & \psi_1 \end{bmatrix} \tag{1}$$

4

5    Then, the states are distributed as a first-order Markov chain governed by the transition matrix

6    with states unoccupied and occupied at $t$ in rows and states unoccupied and occupied at $t + 1$

7    in columns:

8

$$\begin{bmatrix} 1-\gamma & \gamma \\ \varepsilon & 1-\varepsilon \end{bmatrix} \tag{2}$$

9

10    The observation process conditional on underlying occupancy states is summarized by a

11    matrix with unoccupied and occupied states at $t$ in rows and undetected and detected

12    observations at visits $j$ in columns:

13

$$\begin{bmatrix} 1 & 0 \\ 1-p & p \end{bmatrix} \tag{3}$$

14

15        Single-season occupancy models can be reformulated as HMMs and fitted in E-

16    SURGE by imposing no extinction ( $\varepsilon = 0$ ) and no colonisation ( $\gamma = 0$ ) in the dynamic

17    model. The extension to multiple states with uncertainty (Nichols *et al.* 2007) is illustrated

18    with breeding states. We consider the states a site is unoccupied, occupied by non-breeders

19    and occupied by breeders, while the observations are species undetected (coded 0), species

20    detected without young (coded 1) and species detected with young (coded 2). We use $\psi^1$

1    (resp. $\psi^2$) the probability that the site is occupied by non-breeders (resp. by breeders), $p^1$

2    (resp. $p^2$) the detection probability of non-breeders (resp. of breeders). There is also a

3    possibility to accommodate uncertainty on a state, here for example on the breeder state to

4    acknowledge that even though reproduction occurs on a site, young might be missed. We

5    introduce $\delta$ the probability of detecting evidence of reproduction, given the site is occupied

6    with young. Then, the vector of initial state probabilities is:

7

$$\begin{bmatrix} 1-\psi^1-\psi^2 & \psi^1 & \psi^2 \end{bmatrix} \tag{4}$$

8

9    while the transition matrix is the identity matrix. The main modifications are in the

10    observation matrix which can be written as a product of two matrices, highlighting the

11    successive processes of detection and breeding state ascertainment:

12

$$\begin{bmatrix} 1 & 0 & 0 \\ 1-p^1 & p^1 & 0 \\ 1-p^2 & 0 & p^2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1-\delta & \delta \end{bmatrix} \tag{5}$$

13

14    **MAIN FEATURES OF E-SURGE**

15        To illustrate the main features of E-SURGE, we go through the workflow provided in

16    Figure 2. We refer to the E-SURGE manual (Choquet & Nogué 2013) as well as (Choquet

17    2008) and (Choquet, Rouan, & Pradel 2009) for more details.

18

19        [FIGURE 2 AROUND HERE]

20

21    *Preliminary steps*

1      We assume that the user has started a new session, loaded a data file and selected the

2      Occupancy option. In the 'Data status' section of the main window, the 'Modify' button

3      allows to specify the characteristics of your model. The number of age classes should be fixed

4      to 1 here. Age is the time elapsed since first detection, which is equivalent to time if all sites

5      start being monitored at $t = 1$. If site-specific covariates are to be used, then the number of

6      individual covariates should be amended accordingly (recall that the sites are the equivalent

7      of individuals in capture-recapture analyses). The number of events is the number of

8      observations (e.g., undetected and detected), while the number of states should always be the

9      number of states to be used (e.g., unoccupied, occupied by non-breeders, occupied by

10     breeders) plus one. Indeed, because E-SURGE was initially developed to estimate

11     demographic parameters, it always considers the absorbing state 'dead' which is useful in a

12     capture-recapture context to have individuals die but of little interest in an occupancy

13     analysis.

14     In the 'Advanced numerical options' section, the 'Compute C-I' option is deactivated

15     by default to save time by avoiding calculating the parameters' confidence intervals of each

16     model. Just tick this option to get standard errors and confidence intervals, e.g., for the best

17     fitting model. Yet in the same section, in the pull-down menu 'Initial values', E-SURGE

18     offers the possibility to use several sets of initial values randomly chosen or use estimates of a

19     simpler model (from last model) as initial values, which might be useful to avoid picking up

20     local minima in the deviance, an issue often encountered in HMMs.

21

22     *Model building*

23     The model construction is served in E-SURGE by two modules called GEPAT and

24     GEMACO (Choquet 2008). In GEPAT, you specify the structure of the vector of initial state

25     probabilities, the matrix of transitions governing the state process and the matrix of

1   observations (conditional on the states) governing the observation process. To specify a

2   parameter that will be estimated (i.e., that will be assigned an effect in GEMACO, see below),

3   you can use any letter. The minus sign '-' means that the parameter corresponding to this cell

4   will always be set to 0 while the star '*' means the complementary of the sum of all the other

5   parameters on the same row. This GEPAT step is very useful as it deactivates the relevant

6   matrix elements once and for all without having to fix values every time a model is fitted.

7        Note that in the vector of initial state probabilities, which corresponds in the capture-

8   recapture context to the state of individual at first encounter, the 'dead' state is always

9   removed as individuals are all alive when marked. In the transition and observation matrices,

10  this state is present and needs to be accounted for.

11       Of practical interest, the three elements of a HMM can be specified through a

12  multistep process that proves very useful in accommodating state uncertainty for example

13  (eqn. 5). After entering the size of the matrix, the default matrix options (diagonal, full or

14  empty matrix) can be used as a starting point to specify a matrix.

15       In GEMACO, we specify the effects (*sensu* the design matrix in programs MARK and

16  PRESENCE) using a R-like syntax: for example, a season effect will be specified by 't' for

17  time, a group effect by 'g' while a constant effect will be 'i' for intercept. If the effect of a site

18  covariate needs to be investigated, we use 'i + xind' where xind specifies the slope of the

19  relationship. The matrices defined at the GEPAT step can be manipulated using the syntax

20  'from' for rows and 'to' for columns. For example, f(1).to(2) will pick the element in row 1

21  and column 2 of the corresponding matrix, and if a time effect is required on this element,

22  then the syntax will be f(1).to(2).t. If the entire first row of a matrix with five columns needs

23  to be selected, then we use f(1).to(1,2,3,4,5). A colon ':' is useful to lump categories together

24  while the ampersand symbol '&' aggregates parameters corresponding to levels of different

1   factors. Additive and interactive effects can be specified with the plus sign '+' and dot '.'

2   respectively.

3       Shortcuts can be defined to assign a name to a given syntax, hence simplifying the

4   formulas in GEMACO.

5

6   *Final steps*

7   The last steps consist of specifying initial values, using the values by default, or fixing

8   parameters to some values if needed (IVFV step), and then running E-SURGE to fit the

9   current model (RUN step). Standard numerical results (e.g., maximum likelihood estimates,

10  AIC, confidence intervals) can be obtained in a text file or an Excel file. In particular, E-

11  SURGE provides a reliable number of parameters via an algorithm described in (Choquet &

12  Cole 2012), which is crucial in particular with parameter-redundant models, and is one of the

13  key steps for correct model selection using the AIC.

14

15  **CASE STUDIES IN E-SURGE**

16  To illustrate the use of E-SURGE for fitting occupancy models, we simulated data and used

17  existing simulated datasets. Results from fitting the three models described above to these

18  data are provided in Table 1. Estimates for all parameters were close to the true values and the

19  95% confidence intervals covered the true values in all cases.

20

21                          [TABLE 1 AROUND HERE]

22

23  We go through the most important steps of the implementation in E-SURGE and we refer to

24  the companion wiki website http://occupancyinesurge.wikidot.com/ for full details.

25

1 *Dynamic models*

2 We start with the dynamic occupancy models described in (1), (2) and (3) above. This

3 mathematical formulation of the model can be translated for E-SURGE as follows. In

4 GEPAT, the vector of initial state probabilities is specified as

5

6 $\quad * \ \Psi$

7

8 which corresponds to (1). Recall that the dead state is not displayed at this step. Then, the

9 transition matrix in (2) is specified as

10

11 $\quad * \ \gamma \ -$

12 $\quad \varepsilon \ * \ -$

13 $\quad - \ - \ *$

14

15 where the last row and column are for the dead state. Finally, the observation matrix in (3) is

16 coded as:

17

18 $\quad * \ -$

19 $\quad * \ p$

20 $\quad * \ -$

21

22 where the last row corresponds to the dead state. In GEMACO, we write `i` at the Initial state

23 and Event steps to impose a constant parameter. Regarding the Transition step, the data were

24 simulated with 3 seasons (primary sessions) and 3 visits within each season (secondary

25 session), therefore we use `to.t(1 2 4 5 7 8)+to.t(3 6)` (or `to.t(1 2 4 5 7 8,3 6)`).

26 The `to` makes the columns of the matrix different, resulting here in distinguishing the

27 colonization and extinction parameters (note that `from` would produce exactly the same result

1 by differentiating the rows). To handle the robust design, `t(1 2 4 5 7 8)` puts the intervals

2 between secondary occasions together, and the corresponding parameters will be fixed to 0 at

3 the IVFV step to impose closure within primary session (neither extinction nor colonization).

4 The term `t(3 6)` puts together the intervals between the last visit in a season and the first

5 visit of the next one (between primary sessions).

6

7 *Single-season models*

8 We start with the simplest model as described above. In GEPAT, the specification for the

9 Initial state and the Event steps is exactly the same as for the dynamic model. To satisfy the

10 closure assumption, we impose neither colonization nor extinction between sampling

11 occasions by specifying at the Transition step:

12

13        \* – –

14        – \* –

15        – – \*

16

17 In GEMACO, we use `i` for all steps to obtain constant parameters.

18     The extension of this model to multiple states with uncertainty is obtained as follows

19 in E-SURGE. In GEPAT, the initial state probabilities in (4) are:

20

21        \* Ψ Ψ

22

23 while the Transition step is the same as before. The originality lies in the specification of the

24 observation process that is accomplished in two steps to match (5) in GEPAT: step 1 and the

25 detection matrix is:

26

1
```
*  -  -  -
```

2
```
*  p  -  -
```

3
```
*  -  p  -
```

4
```
*  -  -  -
```

5  while step 2 gives the assignment matrix and is:

6

7
```
*  -  -
```

8
```
-  *  -
```

9
```
-  *  δ
```

10
```
*  -  -
```

11

12  In GEMACO, we use `to` at the Initial state step to distinguish the occupancy probabilities

13  according to states, `from` for the step 1 of the Event step to distinguish the detection

14  probabilities according to states and `i` for step 2 to have a constant assignment probability.

15

16  **FURTHER E-SURGE CAPABILITIES**

17       E-SURGE offers the possibility to include heterogeneity in the detection using finite

18  mixtures, to fit multiple species models can also be fitted in E-SURGE, incorporate covariates

19  measured at the site or season level as well as site random effects. The companion wiki

20  website http://occupancyinesurge.wikidot.com/ presents such examples.

21

22  **CONCLUSIONS**

23  Although initially developed for capture-recapture data, E-SURGE can be efficiently used to

24  build and analyse a variety of occupancy models via the HMM framework. E-SURGE

25  includes a user-friendly syntax for specifying models without having to write custom code

26  and used advanced numerical algorithms to produce fast and reliable results. By making the

27  link between the two fast growing user communities of capture-recapture and occupancy, E-

1    SURGE has the potential to provide a unified framework for the construction and analysis of

2    hidden-Markov models in ecology.

1  **REFERENCES**

2  Bailey, L.L., MacKenzie, D.I. & Nichols, J.D. (2013) Advances and applications of

3     occupancy models. *Methods in Ecology and Evolution*.

4  Choquet, R. (2008) Automatic generation of multistate capture recapture models. *The*

5     *Canadian Journal of Statistics*, **36**, 43–57.

6  Choquet, R. & Cole, D.J. (2012) A Hybrid Symbolic-Numerical Method for Determining

7     Model Structure. *Mathematical Biosciences*, **236**, 117–125.

8  Choquet, R. & Nogué, E. (2013) *E-SURGE 1.8 User's Manual*. Montpellier.

9  Choquet, R., Rouan, L. & Pradel, R. (2009) Program E - SURGE : A Software Application

10    for Fitting Multievent Models. *Environmental and Ecological Statistics* (eds D.L.

11    Thomson, E.G. Cooch & M.J. Conroy), pp. 845–865. Springer US.

12  Fiske, I.J. & Chandler, R.B. (2011) unmarked : An R Package for Fitting Hierarchical Models

13    of Wildlife Occurrence and Abundance. *Journal Of Statistical Software*, **43**, 1–23.

14  Gimenez, O., Lebreton, J.-D., Gaillard, J.-M., Choquet, R. & Pradel, R. (2012) Estimating

15    demographic parameters using hidden process dynamic models. *Theoretical Population*

16    *Biology*, **82**, 307–316.

17  Hines, J.E. (2013) PRESENCE 5.9 – Software to estimate patch occupancy and related

18    parameters.

19  Kéry, M. & Schaub, M. (2011) *Bayesian Population Analysis Using WinBUGS: A*

20    *Hierarchical Perspective*. Academic Press.

21  MacKenzie, D.I., Nichols, J.D., Royle, J.A., Pollock, K.H., Bailey, L.L. & Hines, J.E. (2006)

22    *Occupancy Estimation and Modeling: Inferring Patterns and Dynamics of Species*

23    *Occurrence*. Academic Press.

1  Nichols, J.D., Hines, A.J.E., Mackenzie, D.I., Seamans, M.E. & Gutiérrez, R.J. (2007)

2     Occupancy estimation and modeling with multiple states and state uncertainty. *Ecology*,

3     **88**, 1395–1400.

4  Pradel, R. (2005) Multievent: an extension of multistate capture-recapture models to uncertain

5     states. *Biometrics*, **61**, 442–7.

6  Royle, J.A. & Kéry, M. (2007) A Bayesian state-space formulation of dynamic occupancy

7     models. *Ecology*, **88**, 1813–23.

8  White, G.C. & Burnham, K.P. (1999) Program MARK: survival estimation from populations

9     of marked animals. *Bird Study*, **46**, 120–139.

10

11

1   **Table 1**. Estimates of the parameters in the occupancy models fitted to the simulated data set

2   in E-SURGE. For each parameter, the true value, the maximum likelihood estimate and the

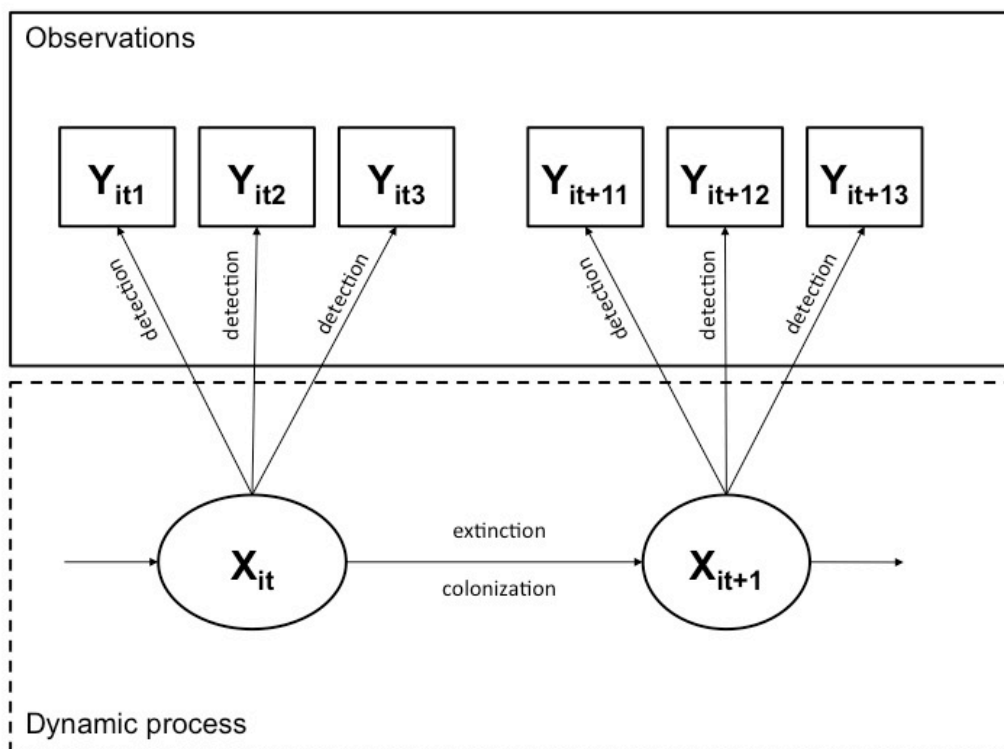3   95% confidence interval are provided. See text for details.

| Model | Parameter | True value | Estimate | 95% confidence interval |
|---|---|---|---|---|
| *Single-season* | $\psi$ | 0.8 | 0.8 | (0.70, 0.87) |
| | $p$ | 0.5 | 0.5 | (0.44, 0.56) |
| *Multistate with uncertainty* | $\psi^1$ | 0.3 | 0.30 | (0.23, 0.37) |
| | $\psi^2$ | 0.5 | 0.50 | (0.43, 0.57) |
| | $p^1$ | 0.5 | 0.52 | (0.43, 0.61) |
| | $p^2$ | 0.7 | 0.70 | (0.65, 0.76) |
| | $\delta$ | 0.8 | 0.79 | (0.73, 0.85) |
| *Dynamic model* | $\psi_1$ | 0.6 | 0.59 | (0.53, 0.65) |
| | $\gamma$ | 0.3 | 0.32 | (0.27, 0.39) |
| | $\varepsilon$ | 0.5 | 0.53 | (0.47, 0.59) |
| | $p$ | 0.7 | 0.70 | (0.67, 0.73) |

4

5

1    **Figure 1**. Schematic hidden Markov representation of a dynamic occupancy model with

2    imperfect detection. We consider site $i$ between seasons (primary periods) $t$ and $t + 1$ and

3    three visits (secondary periods) within each primary period. Each site is closed (i.e. no change

4    in its occupancy status) within primary periods but open (i.e., allowing for changes in

5    occupancy status) between primary periods. The first layer (circles) is a succession of hidden

6    states or latent states (0 for unoccupied vs. 1 for occupied) of site $i$ at season $t +1$ ( $X_{it+1}$ )

7    depending on its states at time $t$ ( $X_{it}$ ). The dynamic of the states is driven by transition

8    probabilities, here probabilities of colonization $\gamma = \Pr\left(X_{it+1} = 1 \middle| X_{it} = 0\right)$ and local extinction

9    $\varepsilon = \Pr\left(X_{it+1} = 0 \middle| X_{it} = 1\right)$. The second layer (squares) corresponds to the detection ( $Y_{itj} = 1$ ) or

10    not ( $Y_{itj} = 0$ )of the target species on site $i$ at visit $j = 1, 2$ or 3 conditional on site $i$ being in

11    state $X_{it}$ . These events are driven by the species detection probability $p = \Pr\left(Y_{itj} = 1 \middle| X_{it} = 1\right)$.

12    The initial state probabilities $\psi_1 = \Pr\left(X_{i1} = 1\right)$ and $1 - \psi_1$ are not represented.
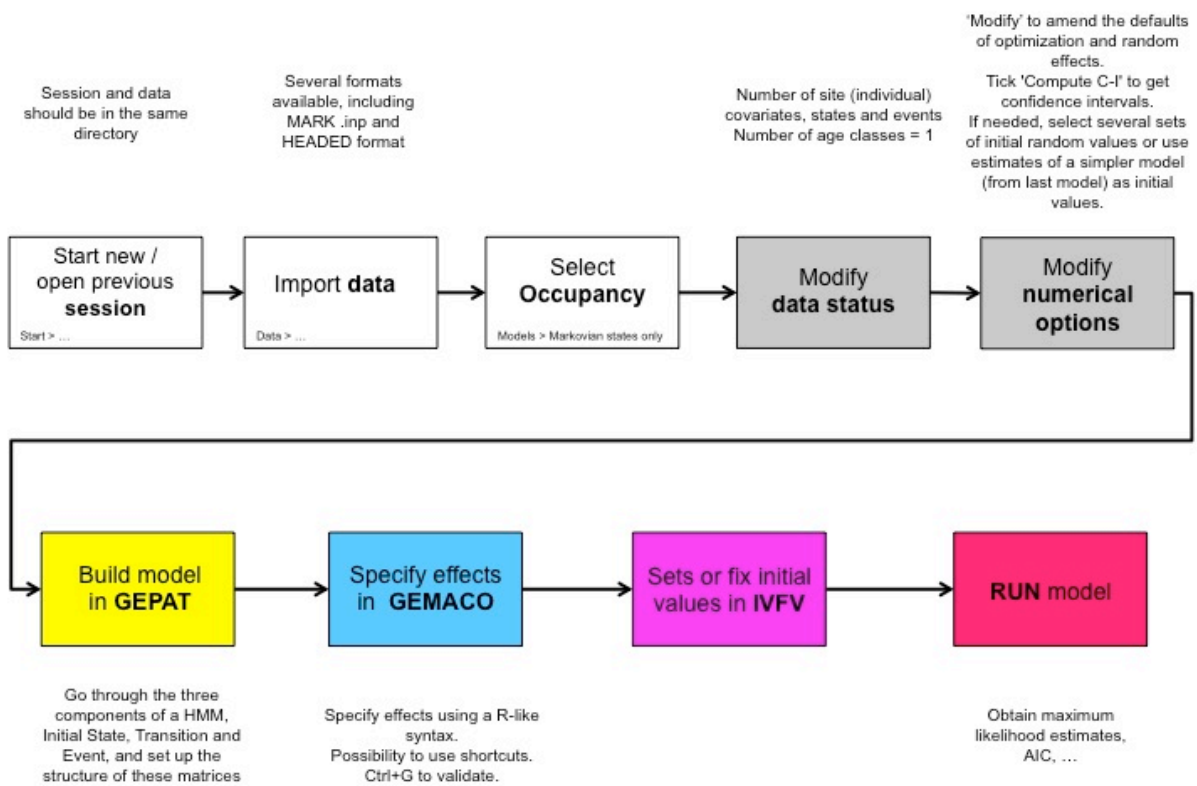


13

14

1 **Figure 2**. Workflow diagram for E-SURGE. We describe the successive steps of a typical

2 analysis in E-SURGE, from data input to model fitting through model building and effects

3 specification. Steps that need to be accomplished through pull-down menus are in white

4 boxes, the others can be done directly from the main interface. We provide details on key

5 steps in using E-SURGE in the text above or below the boxes.

6



7