

A simulation of nonresponse and imputation (or, a framework for narrative modeling)

Ben Klemens

Center for Statistical Research and Methodology

U.S. Census Bureau

`ben.klemens@census.gov`

20 August 2013



Disclaimer

THIS PAPER IS RELEASED TO INFORM INTERESTED PARTIES OF ONGOING RESEARCH AND TO ENCOURAGE DISCUSSION. THE VIEWS EXPRESSED ARE THOSE OF THE AUTHORS AND NOT NECESSARILY THOSE OF THE U.S. CENSUS BUREAU.

The Outline

- The theory part, defining a statistical model.
 - ▶ Using a random number generator as a statistical model.
- The simulation
 - ▶ Interviewer agents seek out respondent agents.
 - ▶ This is a work in progress.

My storyline

- I had agent-based models that needed statistical analysis. My options:
 - ▶ Write the ABM in the language with all the stats functions.
 - ▶ It was too slow.

My storyline

- I had agent-based models that needed statistical analysis. My options:
 - ▶ Write the ABM in the language with all the stats functions.
 - ▶ It was too slow.
 - ▶ Write stats functions in a language fast enough for ABM.
 - ▶ I wound up actually doing this

My storyline

- I had agent-based models that needed statistical analysis. My options:
 - ▶ Write the ABM in the language with all the stats functions.
 - ▶ It was too slow.
 - ▶ Write stats functions in a language fast enough for ABM.
 - ▶ I wound up actually doing this
- In natural programmer form, I started writing up a model object.
 - ▶ Slots for parameters, a pointer to the input data, functions describing the estimation, the score, &c.
 - ▶ Functions that take in a model object (see below)

My storyline

- I had agent-based models that needed statistical analysis. My options:
 - ▶ Write the ABM in the language with all the stats functions.
 - ▶ It was too slow.
 - ▶ Write stats functions in a language fast enough for ABM.
 - ▶ I wound up actually doing this
- In natural programmer form, I started writing up a model object.
 - ▶ Slots for parameters, a pointer to the input data, functions describing the estimation, the score, &c.
 - ▶ Functions that take in a model object (see below)
- The epiphany: The same model object that describes regressions and distributions also works for ABMs.

Apophenia

- A library of stats functions written around a model object as described here.

Apophenia

- A library of stats functions written around a model object as described here.
- Pretty mature:
 - ▶ ~ 250 functions for data shunting and modeling
 - ▶ Documentation is 70,000 words long (plus a book)
 - ▶ Testing suite about as large as the code base
 - ▶ ~ 13,500 nontrivial lines of code
 - ▶ Slipped into Census production here and there

Apophenia

- A library of stats functions written around a model object as described here.
- Pretty mature:
 - ▶ ~ 250 functions for data shunting and modeling
 - ▶ Documentation is 70,000 words long (plus a book)
 - ▶ Testing suite about as large as the code base
 - ▶ ~ 13,500 nontrivial lines of code
 - ▶ Slipped into Census production here and there
 - ▶ Raking (IPF) for sparse tables

Apophenia

- A library of stats functions written around a model object as described here.
- Pretty mature:
 - ▶ ~ 250 functions for data shunting and modeling
 - ▶ Documentation is 70,000 words long (plus a book)
 - ▶ Testing suite about as large as the code base
 - ▶ ~ 13,500 nontrivial lines of code
 - ▶ Slipped into Census production here and there
 - ▶ Raking (IPF) for sparse tables
- v1 coming soon
- Interested? Contribute!

Apophenia

- A library of stats functions written around a model object as described here.
- Pretty mature:
 - ▶ ~ 250 functions for data shunting and modeling
 - ▶ Documentation is 70,000 words long (plus a book)
 - ▶ Testing suite about as large as the code base
 - ▶ ~ 13,500 nontrivial lines of code
 - ▶ Slipped into Census production here and there
 - ▶ Raking (IPF) for sparse tables
- v1 coming soon
- Interested? Contribute!
- Not interested? This talk will be language independent—implement it in your favorite computing platform!

A bundle of functions

Notation

- \mathbb{D} : Data space. Anything required by the model; 'private' to the model unless otherwise noted. \leq is defined. [sample space]
- \mathbb{P} : Parameter space. Similarly model-specific. [state space]
- \mathbb{M} : Model space. The set of categories of $\{\mathbb{D}, \mathbb{P}, \text{ML-consistent functions as per the next slide}\}$.

A bundle of functions

- Likelihood: $(\mathbb{D}, \mathbb{P}) \rightarrow \mathbb{R}^+$.
 - ▶ Integrates to a finite value; always nonnegative.
 - ▶ In some cases, better described as an 'objective function'.
Doesn't have to integrate to one.
- Estimation: $\mathbb{D} \rightarrow \mathbb{P}$
 - ▶ ML-consistency: $L(\mathbf{d}, \mathbf{p})$ is maximized by $\mathbf{p} = \text{EST}(\mathbf{d})$.
- RNG: \mathbb{P} (and uniform prng) $\rightarrow \mathbb{D}$.
 - ▶ Likelihood of draw $\mathbf{d} = \text{RNG}(\mathbf{p}) \propto L(\mathbf{d}, \mathbf{p})$.
- CDF: $(\mathbb{D}, \mathbb{P}) \rightarrow [0, 1]$.
 - ▶ Proportion of random draws $\text{RNG}(\mathbf{p}) \leq \mathbf{d} \rightarrow \text{CDF}(\mathbf{d}, \mathbf{p})$.

The Normal example

- Likelihood: $\mathcal{N}(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-(x - \mu)^2/2\sigma^2)$ or
- Estimation: $\hat{\mu} = \text{mean of } D$; $\hat{\sigma} = \sqrt{\sum(d - \hat{\mu})^2/n}$.
- RNG: See Devroye (1986).
- CDF: `gsl_cdf_gaussian_P(d-mu, sd)` (or see `erf`).

Just a likelihood

Given only a likelihood function, $P(\mathbf{D}, \mathbf{P})$.

- Score (dlog likelihood): numeric deltas.
- Estimation: Use Maximum likelihood estimation.
 - ▶ All MLE algorithms repeatedly sample from the likelihood.
Some use the score.
- RNG: ARMS (Gilks 1995)
- CDF: make random draws, count the percent up to a given point

Just an RNG

Given only $\text{RNG}(\mathbf{P})$.

- CDF: make random draws, count the percent up to a given point
- Likelihood: make a million draws, write down a PMF using those draws. Optional: smooth the PMF.
 - ▶ The problem has been reduced to the prior slide

The category of $\{\mathbb{M}, \text{functions}\}$

- Transformations are of the form $t : \mathbb{M} \rightarrow \mathbb{M}$
 - ▶ Truncation
 - ▶ Jacobian: e.g., map $d \rightarrow \sqrt{d}$
 - ▶ Constrain the parameters.
- Transformations that join models: $j : (\mathbb{M}, \mathbb{M}) \rightarrow \mathbb{M}$
 - ▶ Stack it
 - ▶ Mix it
 - ▶ Bayesian-update it ($\mathbb{D}_{\text{out}} = \mathbb{P}_{\text{in}}$)
 - ▶ Data-compose it ($\mathbb{D}_{\text{out}} = \mathbb{D}_{\text{in}}$)

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution
 - ▶ $M_{\text{pri}} = \text{Stack}(\mathcal{N}, \text{InvWishart})$

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution
 - ▶ $M_{\text{pri}} = \text{Stack}(\mathcal{N}, \text{InvWishart})$
 - ▶ $M_{\text{L}} = \text{Truncate}_{x>0}(\mathcal{N})$

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution
 - ▶ $M_{\text{pri}} = \text{Stack}(\mathcal{N}, \text{InvWishart})$
 - ▶ $M_{\text{L}} = \text{Truncate}_{x>0}(\mathcal{N})$
 - ▶ $M_{\text{sq}} = \text{Update}(M_{\text{pri}}, M_{\text{L}})$

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution
 - ▶ $M_{\text{pri}} = \text{Stack}(\mathcal{N}, \text{InvWishart})$
 - ▶ $M_{\text{L}} = \text{Truncate}_{x>0}(\mathcal{N})$
 - ▶ $M_{\text{sq}} = \text{Update}(M_{\text{pri}}, M_{\text{L}})$
 - ▶ $M = \text{Jacobian}_{\sqrt{x}}(M_{\text{sq}})$

So we can build complex models via transformations

- The square of x is Normally distributed; has uncertain μ, σ .
 - ▶ \mathcal{N} = Normal distribution
 - ▶ InvWish = Inverse Wishart distribution
 - ▶ $M_{\text{pri}} = \text{Stack}(\mathcal{N}, \text{InvWishart})$
 - ▶ $M_{\text{L}} = \text{Truncate}_{x>0}(\mathcal{N})$
 - ▶ $M_{\text{sq}} = \text{Update}(M_{\text{pri}}, M_{\text{L}})$
 - ▶ $M = \text{Jacobian}_{\sqrt{x}}(M_{\text{sq}})$

$$M = \text{Jacobian}_{\sqrt{x}}(\text{Update}(\text{Stack}(\mathcal{N}, \text{InvWishart}), \text{Trunc}_{x>0}(\mathcal{N})))$$

$$\mathbb{D}_{\text{out}} = \mathbb{D}_{\text{in}}$$

- Given \mathbb{P}_1 , M_1 produces d via RNG.
- Given \mathbb{P}_2 , M_2 gives likelihood of d .
- $M_3 = Dcompose(M_1, M_2)$ has params $\mathbb{P}_1 \times \mathbb{P}_2$.

$$\mathbb{D}_{\text{out}} = \mathbb{D}_{\text{in}}$$

- Given \mathbb{P}_1 , M_1 produces via RNG.
- Given \mathbb{P}_2 , M_2 gives likelihood of d .
- $M_3 = Dcompose(M_1, M_2)$ has params $\mathbb{P}_1 \times \mathbb{P}_2$.
- Find the most likely value of λ for a Poisson or Exponential based on data generated via a simulation.
- Likelihood for M_2 is the parameterless distance function between d and fixed $d_g \Rightarrow$ a traditional model calibration.

$$\mathbb{D}_{\text{out}} = \mathbb{D}_{\text{in}}$$

- Given \mathbb{P}_1 , M_1 produces via RNG.
- Given \mathbb{P}_2 , M_2 gives likelihood of d .
- $M_3 = Dcompose(M_1, M_2)$ has params $\mathbb{P}_1 \times \mathbb{P}_2$.
- Find the most likely value of λ for a Poisson or Exponential based on data generated via a simulation.
- Likelihood for M_2 is the parameterless distance function between d and fixed $d_g \Rightarrow$ a traditional model calibration.

More computational epistemology

- Frequentist linear models: $g(Y) = f(X\beta) + \epsilon$
 - ▶ Deterministic narrative, with uncertainty added at the end

More computational epistemology

- Frequentist linear models: $g(Y) = f(X\beta) + \epsilon$
 - ▶ Deterministic narrative, with uncertainty added at the end
- Bayesian hierarchies:

$$\mu_L \sim \mathcal{N}(\mu_p, \sigma_p)$$

$$\sigma_L \sim \text{InvWish}(\Sigma)$$

$$x \sim \mathcal{N}(\mu_L, \sigma_L)$$

- ▶ Uncertainty starts at the outset; goes throughout the model
- ▶ Textbook distributions expressing a narrative of how the data was generated

More computational epistemology

- Agent-based model:

$$\eta \sim \mathcal{N}(\mu_\rho, \sigma_\rho)$$

$$\psi \sim \text{InvWish}(\Sigma)$$

$$x \sim \img alt="A drawing of a four-sided die (tetrahedron) with a pointer on top, enclosed in a square box. The die has faces with numbers 1, 2, 3, and 4." data-bbox="502 465 583 571"/>$$

- ▶ Black box model expressing a narrative of how the data was generated
- ▶ Uncertainty also expressed throughout the model

Simulating a survey

The survey and impute process

- \mathcal{D} is the true population distribution
- In N dimensions, where N is large: \mathcal{D} (age, sex, race, block, income, commute time, gov't attitude, work hours, ...).

The survey and impute process

- \mathcal{D} is the true population distribution
- In N dimensions, where N is large: \mathcal{D} (age, sex, race, block, income, commute time, gov't attitude, work hours, ...).
- Surveying process is a transformation: $\mathbb{M} \rightarrow \mathbb{M}$ that induces distortions.
- Imputation: $Imp : \mathbb{M} \rightarrow \mathbb{M}$
- Our job: select Imp such that $Imp(Survey(\mathcal{D})) \sim \mathcal{D}$.

Why an ABM

- We don't actually know $Survey(\cdot)$.
 - ▶ When are people home? What is their inclination toward gov't surveys? Where/when did the Census send interviewers?
- It's easier to write $Survey(\mathcal{D})$ as its own model, with an explicit description of agents.

Why an ABM

- We don't actually know $Survey(\cdot)$.
 - ▶ When are people home? What is their inclination toward gov't surveys? Where/when did the Census send interviewers?
- It's easier to write $Survey(\mathcal{D})$ as its own model, with an explicit description of agents.
- Information sponge: you give me a relevant fact about how respondents and interviewers interact, and I can fit it into the narrative.

The population

- Two types of agent: respondent households and interviewers.
- Draw households from DC PUMS (urban, square)
 - ▶ The geography in the simulation is terrible so far.
- Census Barriers, Attitudes, and Motivators Survey (CBAMS, Bates et al 2009):
 - ▶ Marketing survey about survey attitudes
 - ▶ Five categories (see next slide)
 - ▶ Mulry & Olson: a Logit model of type given race, income, urban/rural, U.S. born, Primary English.
 - ▶ Use their coefficients to calculate the odds that a person is of any given type; draw a type.

The mindsets

- Leading edge (26%): pro-Census, “say that they will inform family and friends about the Census”
- Cynical Fifth (19%), “anti-government and anti-institution”
- Head noddors (41%), largely positive, but not well informed. “might be vulnerable to negative publicity”
- Insulated (7%), have heard of the Census, but are unfamiliar
- Unacquainted (6%) Have never heard of the Census

The interviewers

- Based on a simulation by Bor-Chung Chen.
- Explicitly models drive time and drive costs.
- Uniform characteristics and abilities.
- Everything is one mode. With more data, one could model mail/telephone.

A day

- Interviewers are allocated
 - ▶ Current strategy: in proportion to demand
 - ▶ Uniform, no-reallocation, or other strategies are possible

A day

- Interviewers are allocated
 - ▶ Current strategy: in proportion to demand
 - ▶ Uniform, no-reallocation, or other strategies are possible
- Knock on a door:
 - ▶ Odds of respondent being at home is a three-hump density function.
 - ▶ If a householder is home, answers with probability determined by type.
- If respondent answers, all data is gathered with 100% accuracy.

A day

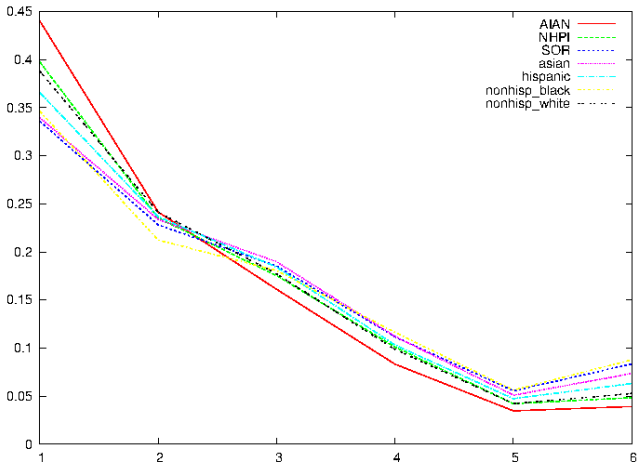
- Interviewers are allocated
 - ▶ Current strategy: in proportion to demand
 - ▶ Uniform, no-reallocation, or other strategies are possible
- Knock on a door:
 - ▶ Odds of respondent being at home is a three-hump density function.
 - ▶ If a householder is home, answers with probability determined by type.
- If respondent answers, all data is gathered with 100% accuracy.
- At end of day, interviewees tell their neighbors.

Calibrating the model—parameters and settings

- Settings: details that are fixed and aren't worth optimizing:
 - ▶ population size
 - ▶ driving costs
 - ▶ for this run: staff size
- Parameters: the optimizer will search for a best value
 - ▶ $P(\text{response}|\text{unacquainted})$
 - ▶ $P(\text{response}|\text{insulated})$
 - ▶ $P(\text{response}|\text{head nodder})$
 - ▶ $P(\text{response}|\text{leading edge})$
 - ▶ $P(\text{response}|\text{cynical fifth})$

Response rate by race

Odds of responding on NRFU 1, 2, ..., 6+ broken down by race:



Calibrating the model

- We could do it nonparametrically by comparing K-L divergence between observed and simulated distributions; I parameterized it.
- Recall the dcompose transformation.

Calibrating the model

- We could do it nonparametrically by comparing K-L divergence between observed and simulated distributions; I parameterized it.
- Recall the dcompose transformation.
- Distributions are parameterized by λ_D^{hisp} , $\lambda_D^{\text{nonhisp black}}$, $\lambda_D^{\text{nonhisp white}}$...
- Data generates distributions of NRFUs with parameters λ_S^{hisp} , $\lambda_S^{\text{nonhisp black}}$, $\lambda_S^{\text{nonhisp white}}$...
- Loss function is

$$|\lambda_S^{\text{hisp}} - \lambda_D^{\text{hisp}}| + |\lambda_S^{\text{nonhisp black}} - \lambda_D^{\text{nonhisp black}}| + \dots$$

Optimization

- The optimization search happens here, via simulated annealing:
 - ▶ SimAn is ideal for optimizing a stochastic function.
 - ▶ Nelder-Mead Simplex also tends to work OK in practice.

```
estimated_model = apop_estimate(data, survey_sim)
```


Variation

Given the optimum from the last slide, I calculate variances using the local inverse Hessian.

- Inverse Hessian calculation easily takes a black-box model
- Comparative statics are valid: if σ_i is $3\sigma_j$, then in expectation whatever causes a 1-unit shift in β_i will cause, in expectation, a 3-unit shift in β_j .
- Parameter sensitivity is inversely proportional to Hessian^{-1} : as perturbances happen, does β_i fluctuate or hold steady?

Variation

Given the optimum from the last slide, I calculate variances using the local inverse Hessian.

- Inverse Hessian calculation easily takes a black-box model
- Comparative statics are valid: if σ_i is $3\sigma_j$, then in expectation whatever causes a 1-unit shift in β_i will cause, in expectation, a 3-unit shift in β_j .
- Parameter sensitivity is inversely proportional to Hessian^{-1} : as perturbances happen, does β_i fluctuate or hold steady?
- Asking whether zero or β_j is $2\sigma_i$ from β_i has limited utility—how much do we believe $\beta \sim \text{Multivariate Normal}$?

Odds of responding

head nodder	80% (7.2%)
leading edge	76% (1.3%)
unacquainted	74% (6.3%)
insulated	57% (8.5%)
cynical fifth	53% (1.7%)

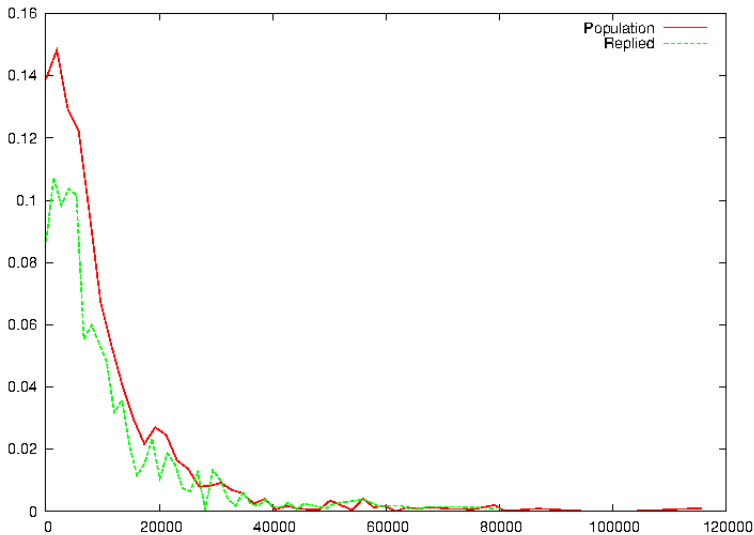
- This was a test (and vindication) of CBAMS: loss function improves when the ABM uses type labels; labels have meaningful interpretation.

A few imputations

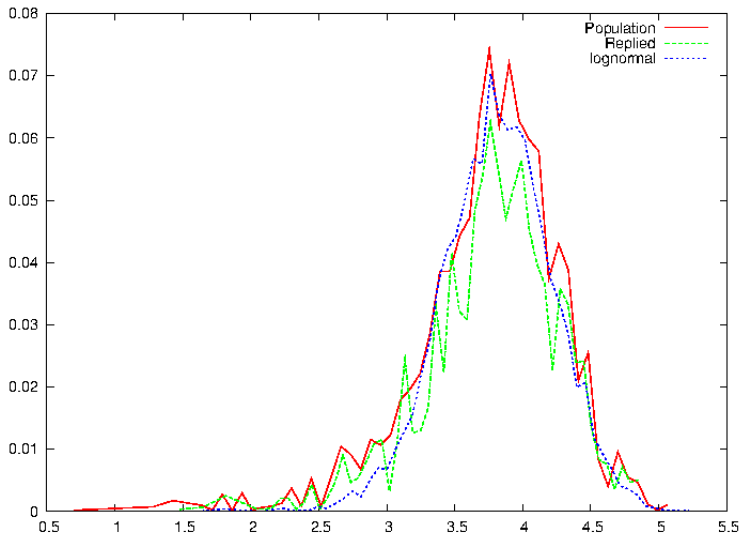
Imputing with the sim

- Distinct analysis from the above
- Takes parameters as fixed. Assigning a prior to them is left as an exercise for the reader

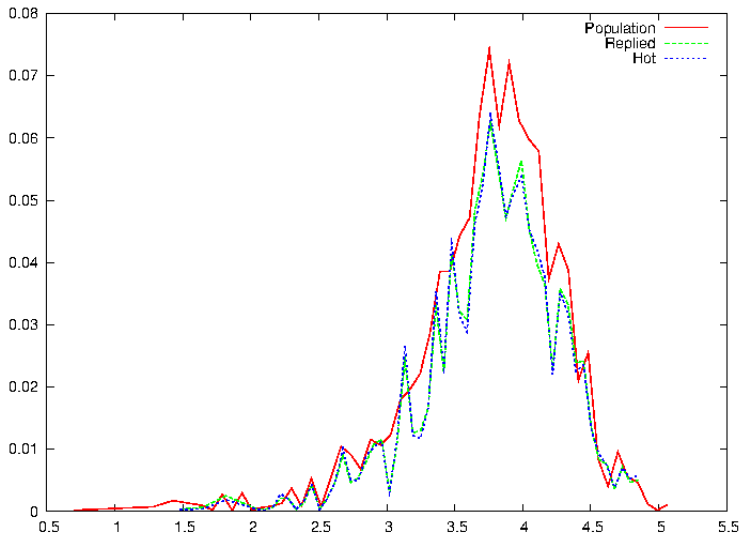
Our respondents



The imputations–lognormal



The imputations—hot deck



Log(income)–lognormal vs hot deck

	bias	MSE
lognormal distribution	0.035 (± 0.005)	.166 (± 0.002)
hot deck	-0.002 (± 0.002)	.237 (± 0.02)

Est (replication variance)

Conclusion

The conclusion slide

- A random number generator is sufficient to describe a statistical model.
- An ABM produces random outcomes from fixed input parameters.
- \therefore An ABM is sufficient to describe a statistical model.

The conclusion slide

- A random number generator is sufficient to describe a statistical model.
- An ABM produces random outcomes from fixed input parameters.
- \therefore An ABM is sufficient to describe a statistical model.
- Things we do with a statistical model:
 - ▶ Calibrate it, by searching for parameters that produce data as close as possible to real-world data.
 - ▶ Evaluate the robustness of those parameters.
 - ▶ Transform it, such as applying an imputation.