

Estimating Population and Health Quantities and their Uncertainty from Data of Limited Quality

Adrian E. Raftery

University of Washington
<http://www.stat.washington.edu/raftery>

Joint work with Leontine Alkema, Samuel Clark, Patrick Gerland and Mark Wheldon
In collaboration with the UN Population Division

Supported by NICHD

UN EGM on Strengthening the Demographic Evidence Base for the
Post-2015 Development Agenda
New York
October 6, 2015

Estimating Population and Health Quantities

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex
 - HIV prevalence

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex
 - HIV prevalence
- Data:

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex
 - HIV prevalence
- Data:
 - High quality vital registration and health surveillance data for less than half of countries

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex
 - HIV prevalence
- Data:
 - High quality vital registration and health surveillance data for less than half of countries
 - In majority of countries, surveys and censuses only

Estimating Population and Health Quantities

- Goal: Estimate current and past demographic and health quantities and their uncertainty, e.g.
 - Vital rates (fertility, mortality, migration): summary and age-specific
 - Population by age and sex
 - HIV prevalence
- Data:
 - High quality vital registration and health surveillance data for less than half of countries
 - In majority of countries, surveys and censuses only
 - Multiple data sources, each with their own issues

Issues

Issues

- Systematic biases:

Issues

- Systematic biases:
 - Non-representative sampling

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?
 - general assessment of accuracy: now routine (e.g. UNAIDS, opinion polls, DHS, PMA2020, ACS)

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?
 - general assessment of accuracy: now routine (e.g. UNAIDS, opinion polls, DHS, PMA2020, ACS)
 - assessing changes and differences between outcomes and expectations

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?
 - general assessment of accuracy: now routine (e.g. UNAIDS, opinion polls, DHS, PMA2020, ACS)
 - assessing changes and differences between outcomes and expectations
 - making decisions that avoid *risks* (e.g. national finance ministries for pension planning, school closures)

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?
 - general assessment of accuracy: now routine (e.g. UNAIDS, opinion polls, DHS, PMA2020, ACS)
 - assessing changes and differences between outcomes and expectations
 - making decisions that avoid *risks* (e.g. national finance ministries for pension planning, school closures)
 - **Statistics NZ a leader: positive experience (Dunstan, Bryant)**

Issues

- Systematic biases:
 - Non-representative sampling
 - Poor geographic coverage
 - Recall bias
 - Undercount (censuses)
- Sampling variation: between individuals, between strata.
- Why is uncertainty assessment needed?
 - general assessment of accuracy: now routine (e.g. UNAIDS, opinion polls, DHS, PMA2020, ACS)
 - assessing changes and differences between outcomes and expectations
 - making decisions that avoid *risks* (e.g. national finance ministries for pension planning, school closures)
 - Statistics NZ a leader: positive experience (Dunstan, Bryant)
 - Raftery (2014, arXiv): experiences and types of user

Bayesian Statistical Modeling

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q),$$

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q),$$

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- **Advantages:**

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior
 - Multiple data sources can be included: If there are m data sources (e.g. different surveys) ($\text{Data}_1, \dots, \text{Data}_m$), the likelihoods are multiplied:

$$p(\text{Data}|Q) = p(\text{Data}_1|Q) \times \dots \times p(\text{Data}_m|Q).$$

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior
 - Multiple data sources can be included: If there are m data sources (e.g. different surveys) ($\text{Data}_1, \dots, \text{Data}_m$), the likelihoods are multiplied:

$$p(\text{Data}|Q) = p(\text{Data}_1|Q) \times \dots \times p(\text{Data}_m|Q).$$

- Estimates can be made for multiple countries at once, using multinational patterns, by a Bayesian hierarchical model.

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior
 - Multiple data sources can be included: If there are m data sources (e.g. different surveys) ($\text{Data}_1, \dots, \text{Data}_m$), the likelihoods are multiplied:

$$p(\text{Data}|Q) = p(\text{Data}_1|Q) \times \dots \times p(\text{Data}_m|Q).$$

- Estimates can be made for multiple countries at once, using multinational patterns, by a Bayesian hierarchical model.
 - Now basis for UN population projections

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior
 - Multiple data sources can be included: If there are m data sources (e.g. different surveys) ($\text{Data}_1, \dots, \text{Data}_m$), the likelihoods are multiplied:

$$p(\text{Data}|Q) = p(\text{Data}_1|Q) \times \dots \times p(\text{Data}_m|Q).$$

- Estimates can be made for multiple countries at once, using multinational patterns, by a Bayesian hierarchical model.
 - Now basis for UN population projections
- **Automatically gives uncertainty**

Bayesian Statistical Modeling

- Inference about a quantity of interest, Q , summarized by its *posterior distribution* given all data and evidence:

$$p(Q|\text{Data}) \propto p(\text{Data}|Q) \times p(Q), \quad \text{i.e.}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

- Unknown parameters (e.g. bias and measurement error variance of surveys), can be included and estimated.
- Advantages:
 - Information from other countries and expert knowledge can be included through the prior
 - Multiple data sources can be included: If there are m data sources (e.g. different surveys) ($\text{Data}_1, \dots, \text{Data}_m$), the likelihoods are multiplied:

$$p(\text{Data}|Q) = p(\text{Data}_1|Q) \times \dots \times p(\text{Data}_m|Q).$$

- Estimates can be made for multiple countries at once, using multinational patterns, by a Bayesian hierarchical model.
 - Now basis for UN population projections
- Automatically gives uncertainty
- **Complex models can be estimated by Monte Carlo methods**

Estimating HIV Prevalance in Generalized Epidemics

Estimating HIV Prevalance in Generalized Epidemics

- Data:

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative

Estimating HIV Prevalence in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:
 - **Standard SIR epidemic model**

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:
 - Standard SIR epidemic model
 - Bias in ANC data

Estimating HIV Prevalance in Generalized Epidemics

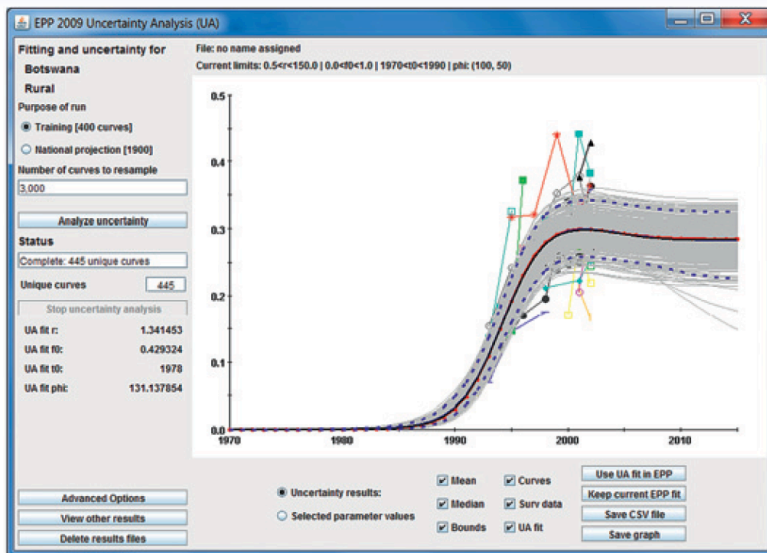
- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:
 - Standard SIR epidemic model
 - Bias in ANC data
 - Measurement error in ANC and DHS data

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:
 - Standard SIR epidemic model
 - Bias in ANC data
 - Measurement error in ANC and DHS data
- Evaluated by UNAIDS Reference Group

Estimating HIV Prevalance in Generalized Epidemics

- Data:
 - HIV prevalence at ante-natal clinics:
 - Frequent measurements \implies Good for trends
 - Unrepresentative
 - Poor geographic coverage
 - National household surveys (e.g. DHS):
 - (more) Representative
 - Infrequent (e.g. 0-2 DHS's in a country).
- Bayesian model includes:
 - Standard SIR epidemic model
 - Bias in ANC data
 - Measurement error in ANC and DHS data
- Evaluated by UNAIDS Reference Group
 - Now used for UNAIDS estimation and projection (EPP/Spectrum software)



Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- **Inputs:**

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.
- Outputs: Joint posterior distribution of all population quantities of interest.

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.
- Outputs: Joint posterior distribution of all population quantities of interest.
- Improvements:

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.
- Outputs: Joint posterior distribution of all population quantities of interest.
- Improvements:
 - Uncertainty is assessed

Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.
- Outputs: Joint posterior distribution of all population quantities of interest.
- Improvements:
 - Uncertainty is assessed
 - All population quantities are estimated simultaneously: trends and uncertainty are estimated in a demographically consistent way.

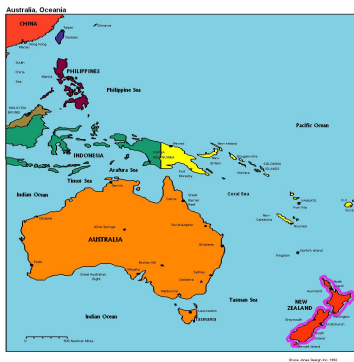
Bayesian Population Reconstruction

Wheldon et al., (2010, 2013, 2015)

- Bayesian hierarchical model for all population quantities and data
- Inputs:
 - Bias-corrected initial estimates of age-specific vital rates, net migration and population counts.
 - Expert knowledge about measurement error variances.
- Outputs: Joint posterior distribution of all population quantities of interest.
- Improvements:
 - Uncertainty is assessed
 - All population quantities are estimated simultaneously: trends and uncertainty are estimated in a demographically consistent way.
 - **Software: popReconstruct R package**

Laos and New Zealand

Laos and New Zealand



Laos and New Zealand Data

Laos and New Zealand Data

- We reconstruct female populations for

Laos and New Zealand Data

- We reconstruct female populations for
 - Laos, 1985–2005 (increased from 1.8 million to 3.0 million)

Laos and New Zealand Data

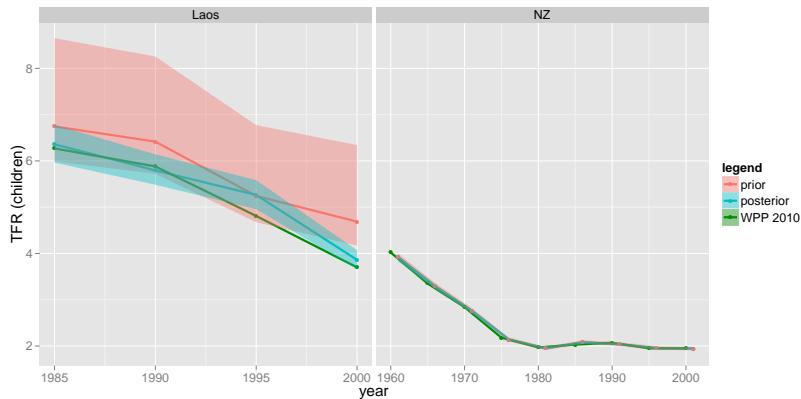
- We reconstruct female populations for
 - Laos, 1985–2005 (increased from 1.8 million to 3.0 million)
 - New Zealand, 1961–2006 (increased from 1.2 million to 2.1 million)

Laos and New Zealand Data

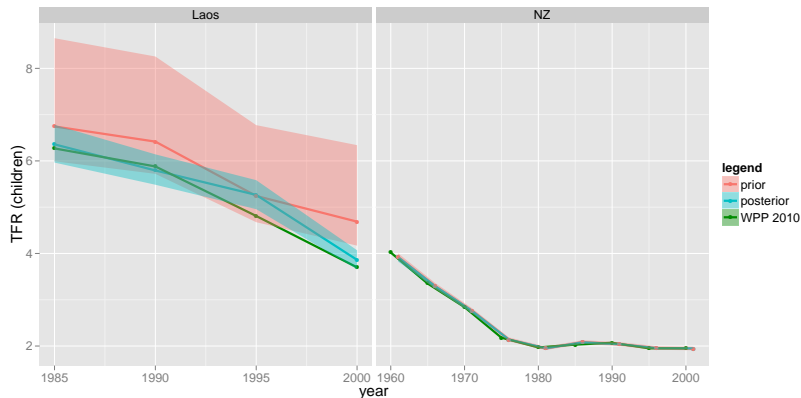
- We reconstruct female populations for
 - Laos, 1985–2005 (increased from 1.8 million to 3.0 million)
 - New Zealand, 1961–2006 (increased from 1.2 million to 2.1 million)
- Very different data qualities and demographics.

Total Fertility Rate

Total Fertility Rate



Total Fertility Rate

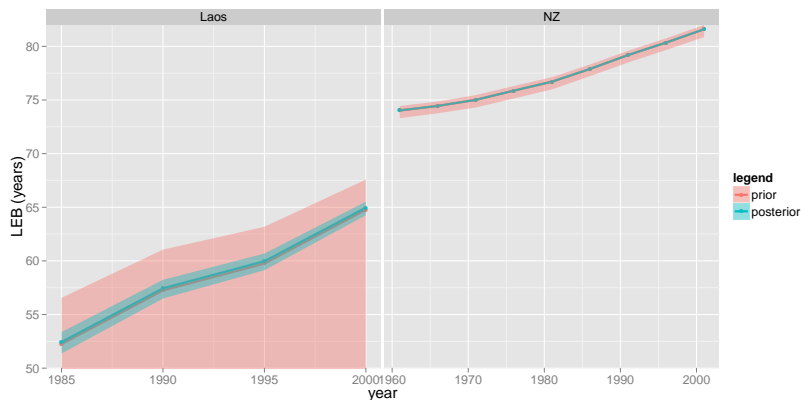


Average 95% Posterior Interval Half-widths

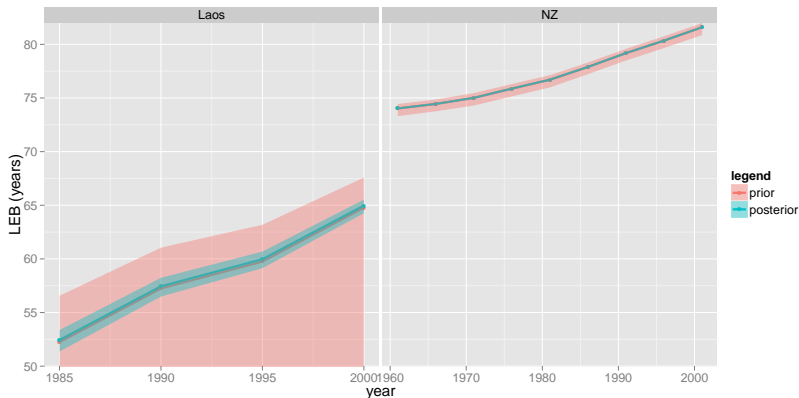
Laos	0.30
NZ	0.03

Life Expectancy at Birth

Life Expectancy at Birth



Life Expectancy at Birth

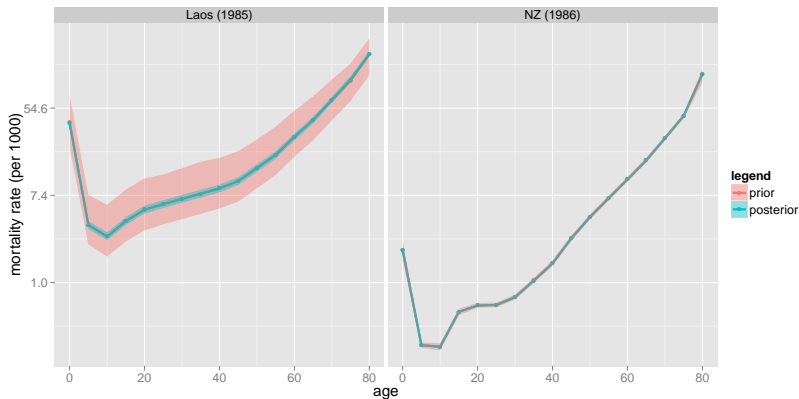


Average 95% Posterior Interval Half-widths

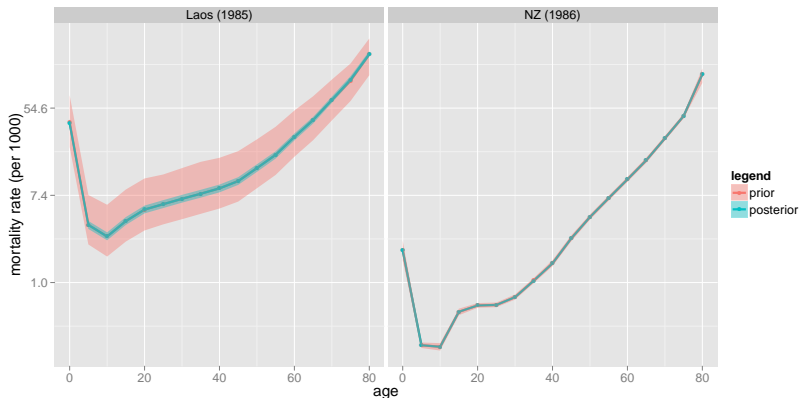
Laos	0.80
NZ	0.04

Age Specific Mortality Rate

Age Specific Mortality Rate



Age Specific Mortality Rate

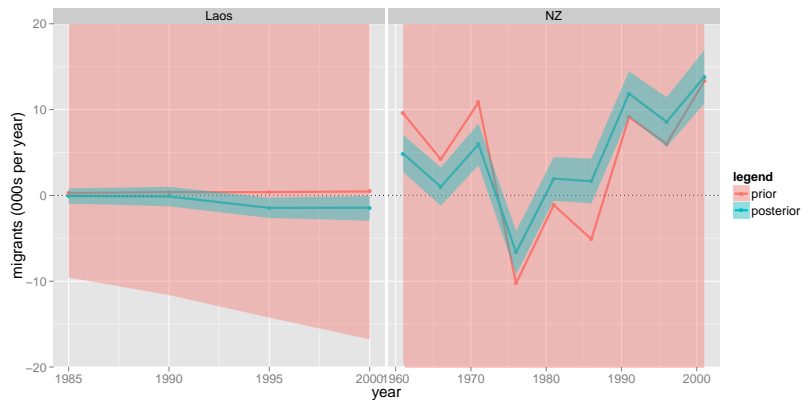


Average 95% Posterior Interval Half-widths

Laos	2.0
NZ	0.1

Total Net Migration

Total Net Migration



Total Net Migration



Average 95% Posterior Interval Half-widths

	% of Median
Laos	570
NZ	80

Summary

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:
 - Estimating and projecting HIV prevalence in generalized epidemics

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:
 - Estimating and projecting HIV prevalence in generalized epidemics
 - **Reconstructing past and current population from limited data**

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:
 - Estimating and projecting HIV prevalence in generalized epidemics
 - Reconstructing past and current population from limited data
- Require systematic consistent data for as long in the past as possible

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:
 - Estimating and projecting HIV prevalence in generalized epidemics
 - Reconstructing past and current population from limited data
- Require systematic consistent data for as long in the past as possible
- Papers available at <http://www.stat.washington.edu/raftery/Research/soc.html>

Summary

- Estimates for demography and health in the majority of countries are based on surveys and censuses from multiple sources with biases and measurement error
- Bayesian approaches can model and take account of all these issues
- Some success in:
 - Estimating and projecting HIV prevalence in generalized epidemics
 - Reconstructing past and current population from limited data
- Require systematic consistent data for as long in the past as possible
- Papers available at
<http://www.stat.washington.edu/raftery/Research/soc.html>
- Software: **popReconstruct R package**