

Developing a daily time step individual level demographic simulation model

Andy Turner

<http://bit.ly/TStpJP>

Outline

- Why?
- What?
- How?
- Results
- Plans and Next Steps
- Feedback

Why?

- Generally...
 - Demographic data is used in a wide range of applications
 - Epidemiology for estimating prevalence and incidence rates
 - Service planning
 - Risk management
 - Commercial
 - Census data tend to be years old by the time outputs are made available
 - Contemporary populations are assumed to be highly and increasingly mobile and fertility in terms of live births is perhaps becoming more variable

- Demographic forecasting is important
 - Planning is a key to sustainability
 - Dependency ratios are increasing in many countries with increasing aged populations
 - Pensions
 - Welfare
 - Services and infrastructure
- Many countries (including the UK) do not have official residential registration data that tracks the location of people (between censuses)
- Our ability to track where everyone has lived improves continually, but we need the models and the data to forecast and provide the best estimates

- Why daily time steps?
 - People tend to be born, die and move residences on specific days
 - It intrinsically makes sense to model at this resolution
 - Modelling for multiple days either misses important events or becomes much more complicated
 - Consider
 - migrations of individuals within a time period
 - births (same mother) at different times within a year
 - Allows for linkage with models of daily activity that work on sub-daily time steps

- Allows for variation in mortality, fertility and migration rates over the year to be modelled
 - Mortality, Fertility, Miscarriage and Migration are seasonal
 - Student migration
 - Holidays and fertility
 - Winter mortality
 - Power cuts/flood events and birth spikes
- Allows for new and exciting aggregate data/statistics to be produced
 - Distribution of the total number of
 - births per month in a region
 - moves per person in a year
 - maximum, minimum, average, variance

- Why individual level?
 - Data can be linked with other individual level data
 - e.g. Disease data
 - It has the possibility that other individual data can be augmented or linked
 - Linkage and substitution with "real data"
 - Everybody is different
 - Individuals have their own mortality, fertility and migration probabilities and history
 - There is scope for specifying these in the model
 - In such a way as to keep overall counts of births and deaths at observed levels
 - Return migration

What?

- Stages of development
- The nature of the model
- Initialisation
- Daily Simulation
 - Death
 - Birth
 - Migration
- Results

Stages of Development

- Natural change Simulation Model
 - [ESRC](#) funded [GENESIS Project](#)
 - Leeds Output Area level results produced
- e-Infrastructure
 - [JISC](#) funded [NeISS Project](#)
 - Web Portal based User Interface
 - Simulation Models configured, run and results stored on e-Science resources
- Migration model component
 - Not externally funded
 - Developed since July 2012

The nature of the model

- Open Source

- Development repositories

- Sourceforge

- <https://sourceforge.net/p/neiss/code/328/tree/genesis/>

- University of St Andrews

- <https://e-research.cs.st-andrews.ac.uk/repos/sim/projects/genesis/>
 - Thanks to Alex Voss

- Java

- [http://en.wikipedia.org/wiki/Java_\(software_platform\)](http://en.wikipedia.org/wiki/Java_(software_platform))

- Dependencies
 - [Generic Library](#)
 - [MoSeS Code](#)
 - For loading 2001 UK Population Census Data
- Grid enabled
 - Thanks to NeISS collaboration with Tom Doherty based at University of Glasgow
- Run for multiples of a year
- Individual representation
 - Males
 - Females

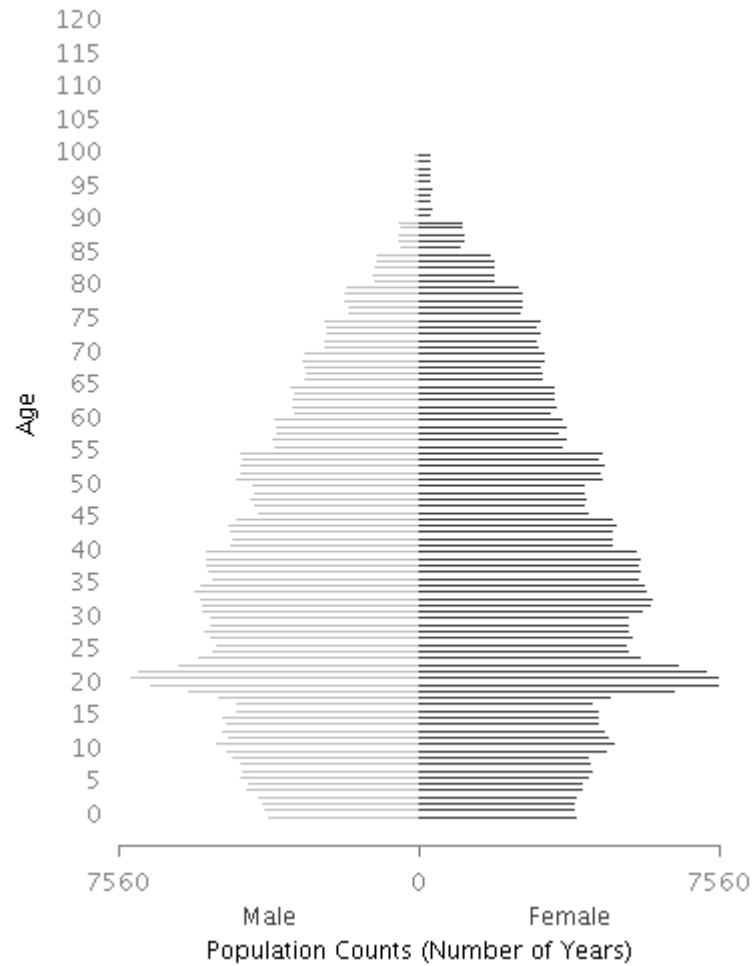
- **Stochastic yet deterministic**
 - Based on pseudo-random sequences
 - Results replicable
- **Study Region**
 - Comprised of regions and subregions
- **2 stages to modelling**
 - Initialisation
 - Simulation
- **Simulation proceeds for each subregion in turn, and for each individual in turn**
- **Synchronisation needed for each daily step**

- **There are many simplifying assumptions**
 - Many things are assumed to be evenly distributed
 - Some things are not explicitly modelled
- **There are interesting model details**
 - Pregnancy and miscarriage
 - Multiple births
- **Input data**
 - Population count data by age and gender
 - Either birth and death counts or fertility and mortality probabilities
 - Migration data

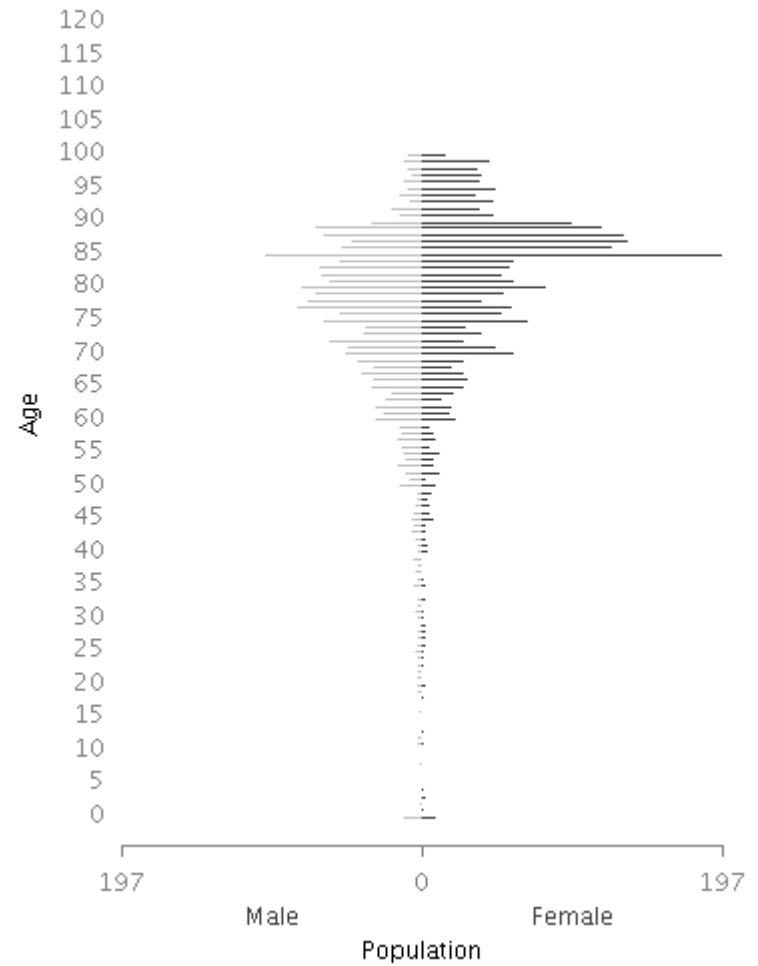
- **Output data**

- Produced annually for study region, regions, subregions and aggregates of subregions
- Includes raw ASCII data ([XML](#),[CSV](#)), binary serialised Java object data, and images ([PNG](#))
 - Population count estimates
 - Mortality and fertility estimates
 - Migration estimates
 - Comparisons with an annual time step model
 - Which uses mid year population estimation
 - An individual level population data set

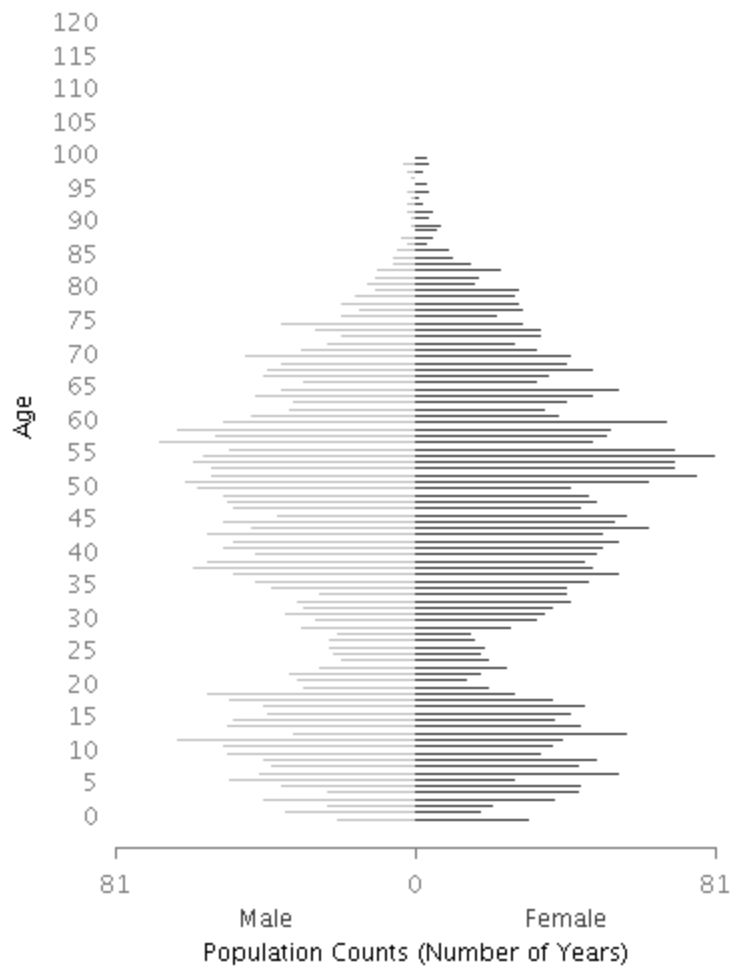
OODA Population _Simulated Living_ Year 2002



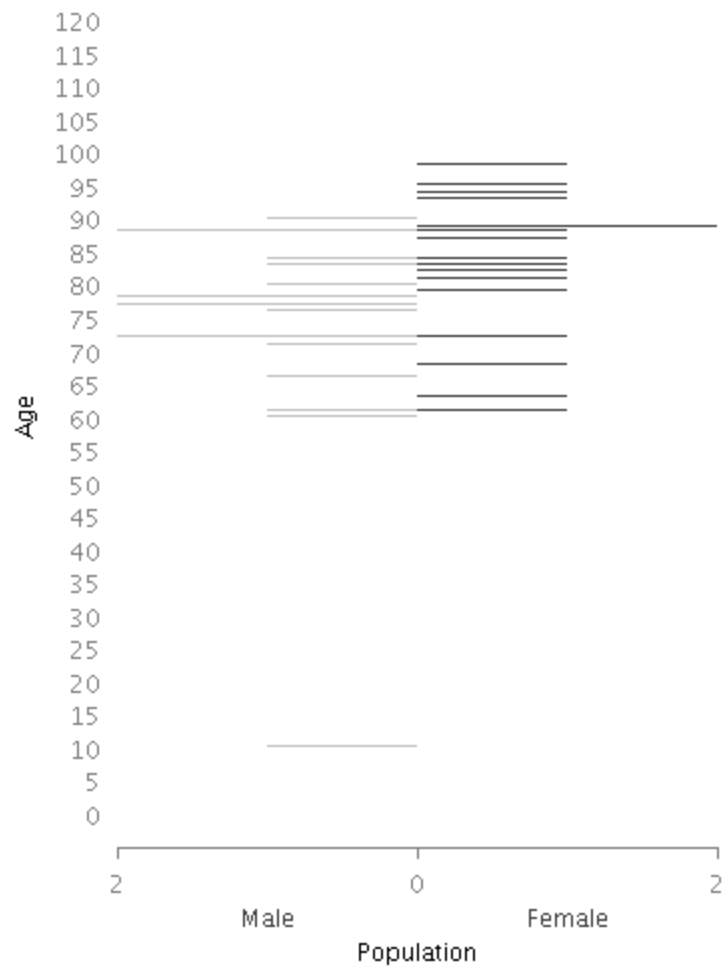
OODA Population _Simulated Dead_ Year 2002



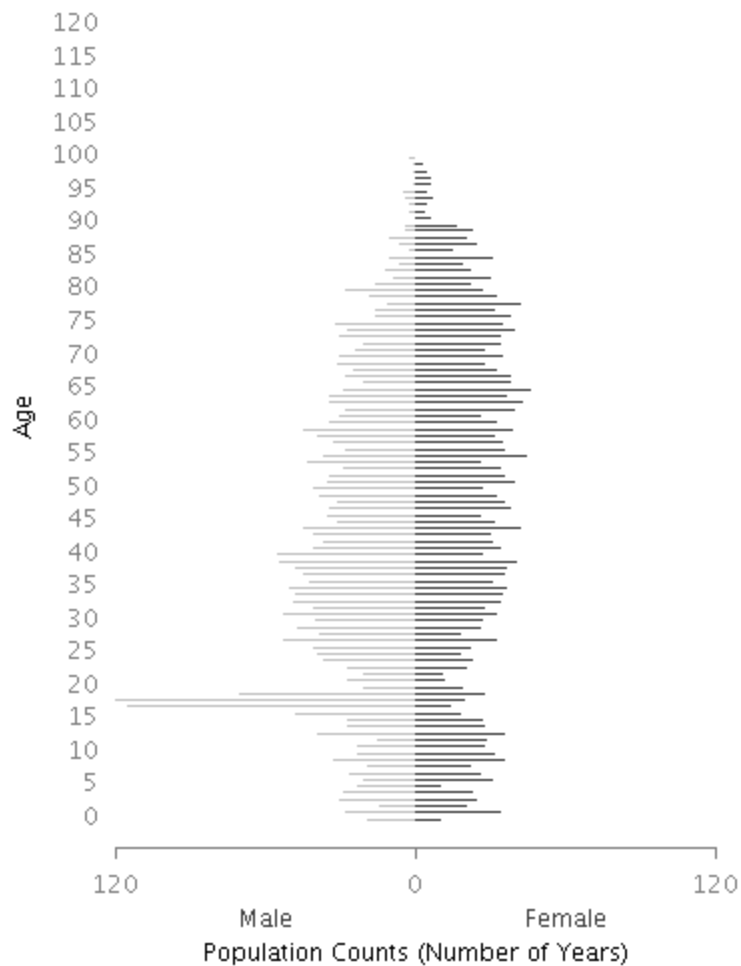
E02002330 Population _Simulated Living_ Year 2002



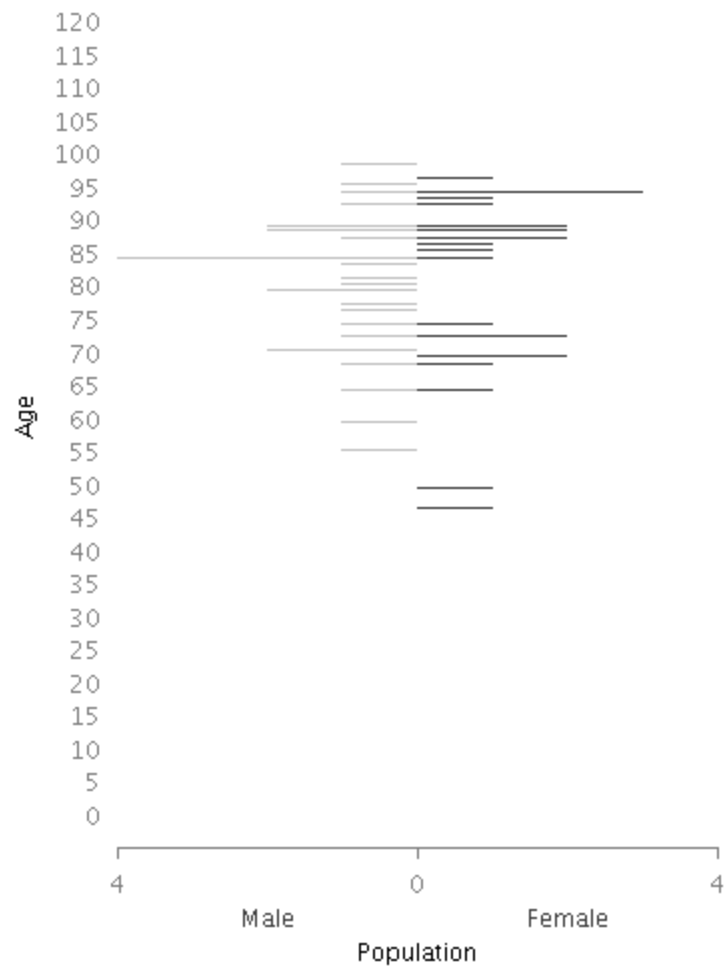
E02002330 Population _Simulated Dead_ Year 2002

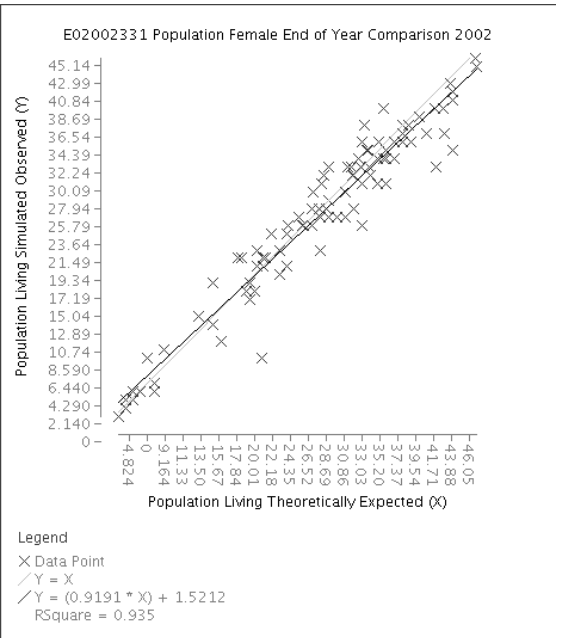
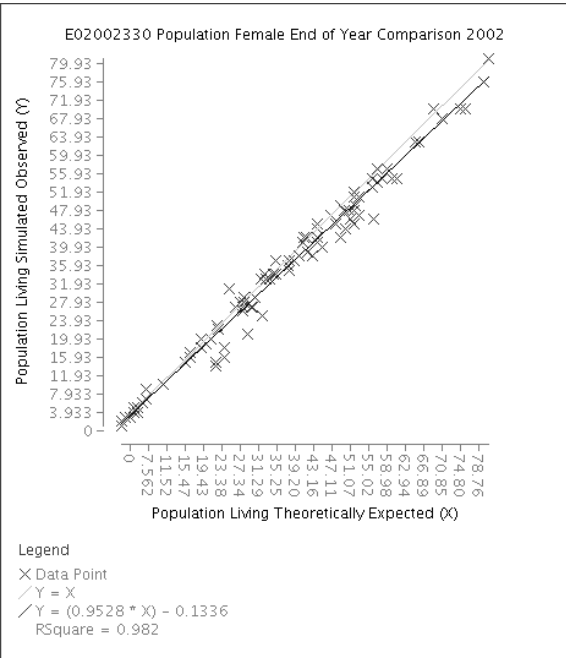
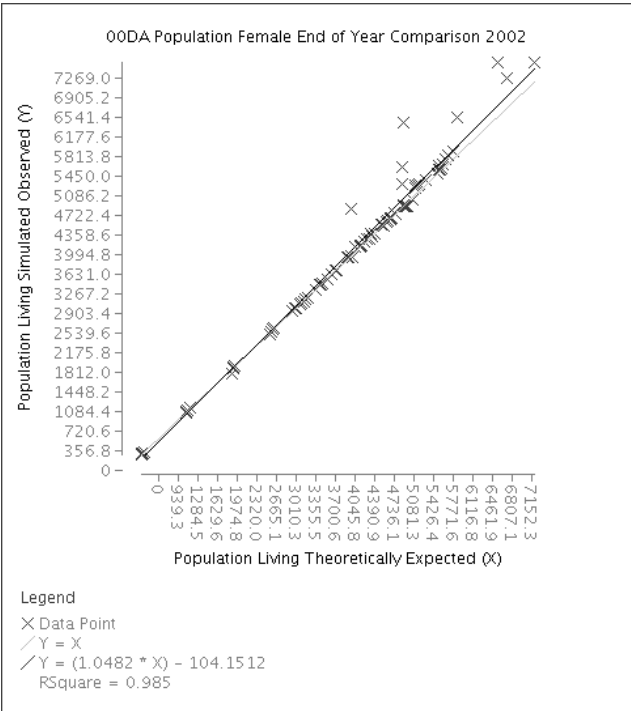


E02002331 Population _Simulated Living_ Year 2002

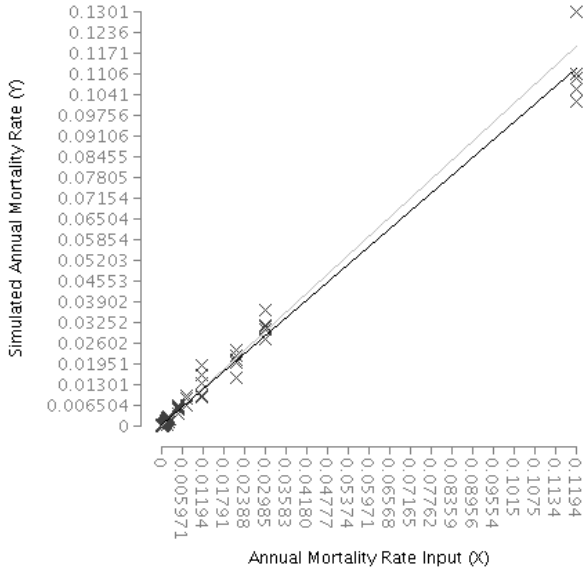


E02002331 Population _Simulated Dead_ Year 2002



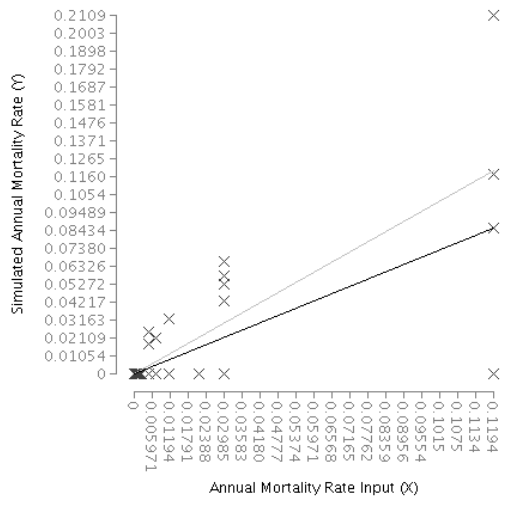


00DA Female Mortality Rate Comparison



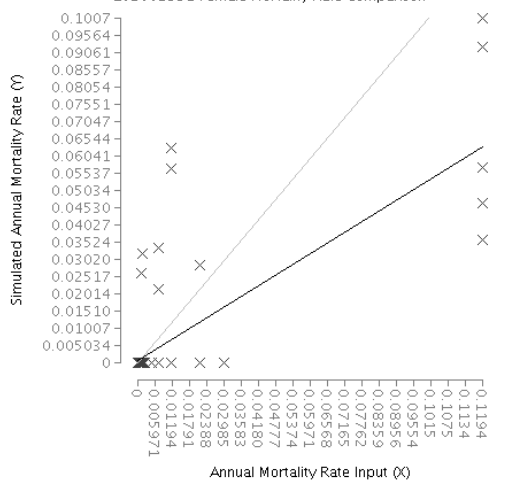
Legend
 X Data Point
 / Y = X
 / Y = (0.9373 * X) + 0.0005
 RSquare = 0.988

E02002330 Female Mortality Rate Comparison

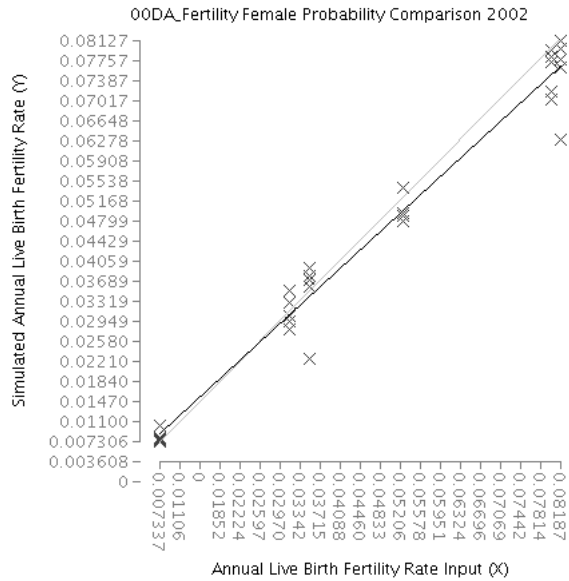


Legend
 X Data Point
 / Y = X
 / Y = (0.7158 * X) + 0.0001
 RSquare = 0.467

E02002331 Female Mortality Rate Comparison

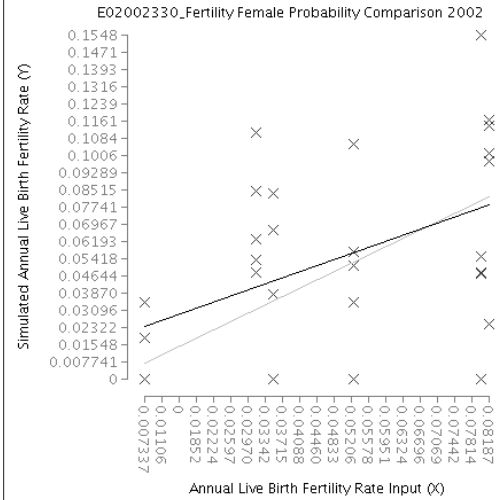


Legend
 X Data Point
 / Y = X
 / Y = (0.5227 * X) + 0.0008
 RSquare = 0.567



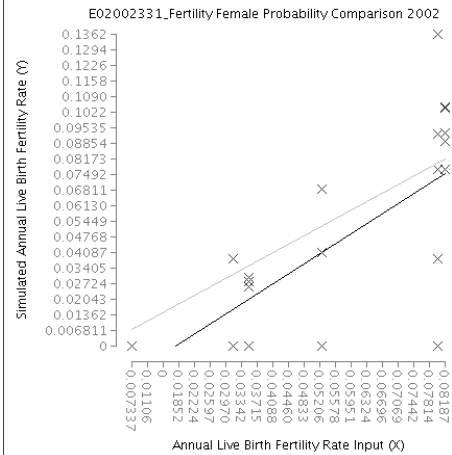
Legend

- X Data Point
- Y = X
- $Y = (0.9087 * X) + 0.0021$
- RSquare = 0.971



Legend

- X Data Point
- Y = X
- $Y = (0.731 * X) + 0.0186$
- RSquare = 0.214



Legend

- X Data Point
- Y = X
- $Y = (1.1771 * X) - 0.0209$
- RSquare = 0.577

Initialisation

- For each region
 - Daily survival probabilities are calculated for each age and gender
 - Death rate assumed to be even throughout the year
 - Daily pregnancy probabilities are calculated for each age of potential mother
 - Annual Live Birth Fertility Rates are factored for multiple births, miscarriage and death of mother
 - Pregnancy rate and miscarriage rate assumed to be even throughout the year

- Daily migration
 - Assumes migration evenly distributed throughout the year
 - General migration probability calculated
 - Internal migration rates are calculated for migration within the region
 - In migration rates are calculated for people moving from all regions not in the study region
- Cumulative sums of migration are calculated to help determine
 - The region destination for each out migration
 - The subregion destination for each in migration

- Each person is initialised
 - Assigned a date of birth
 - Assigned to a subregion as usually resident
 - Females are assigned pregnancies and due dates

Daily Simulation

- For each person
 - Do they die?
 - Is it their birthday?
 - If so update population statistics
 - Do they migrate?
 - If yes, find out where they move to
- For each female
 - If pregnant do they have a miscarriage
 - If due give birth
 - If not pregnant, determine if they become pregnant

- Having gone through the population for all regions in the study region
 - Migrate those migrating out of the study region
 - Migrate those migrating within the study region
 - For migration into the study region from outside of the study region
 - Create individuals
 - Assigning date of birth
 - Record migration origin location
 - Assign subregion usual resident location

How?

- Designed for (scalability) simulating large populations with large numbers of regions and subregions
 - Individual level data stored in collections which are swapped to and from slower access storage as required
 - Numerical indexes are stored in mapped collections
- Computational demands are considerable
 - Consider simulating a single region, population ~1 million, with ~10 thousand subregions
 - Can all the data be stored in the available fast access memory?

- For a simple model, a 10 year simulation might take many days with only one CPU
 - Each individual in the population is updated ~3650 times
- The amount of persistent data produced and that we want to store is in the order of tens of GigaBytes
- For a UK Simulation there are in the order of 60 million individuals and 200 thousand subregions
- Grid enabled
- Parallelisation
- Numerical precision
 - [Java BigDecimal](#)

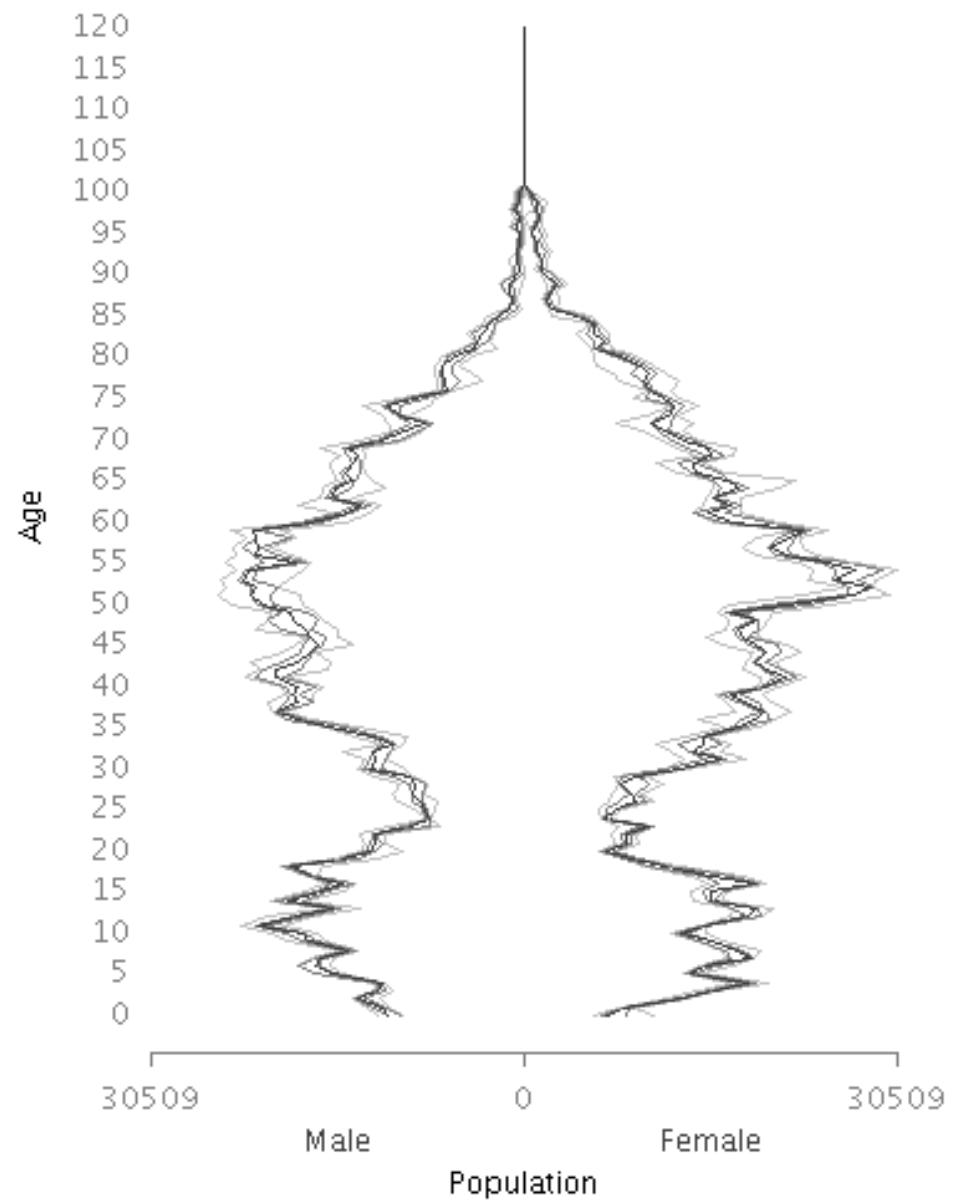
Results

- Results for simulations without migration
 - Provide confidence in daily probability calculations for natural processes
 - The expected amounts of deaths, pregnancies, miscarriages and births result at a regional level
 - Variation
 - At sub-regional level can be large
 - At regional level are generally small
 - At aggregated sub-regional level are intermediate
 - For less frequently occurring events is greater

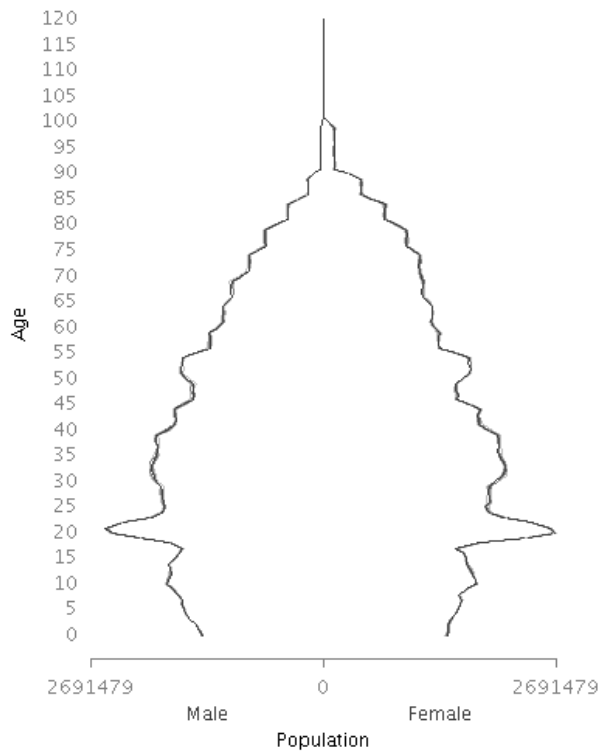
Variation in results

- 10 runs
 - Everything the same except the pseudo-random seed start point

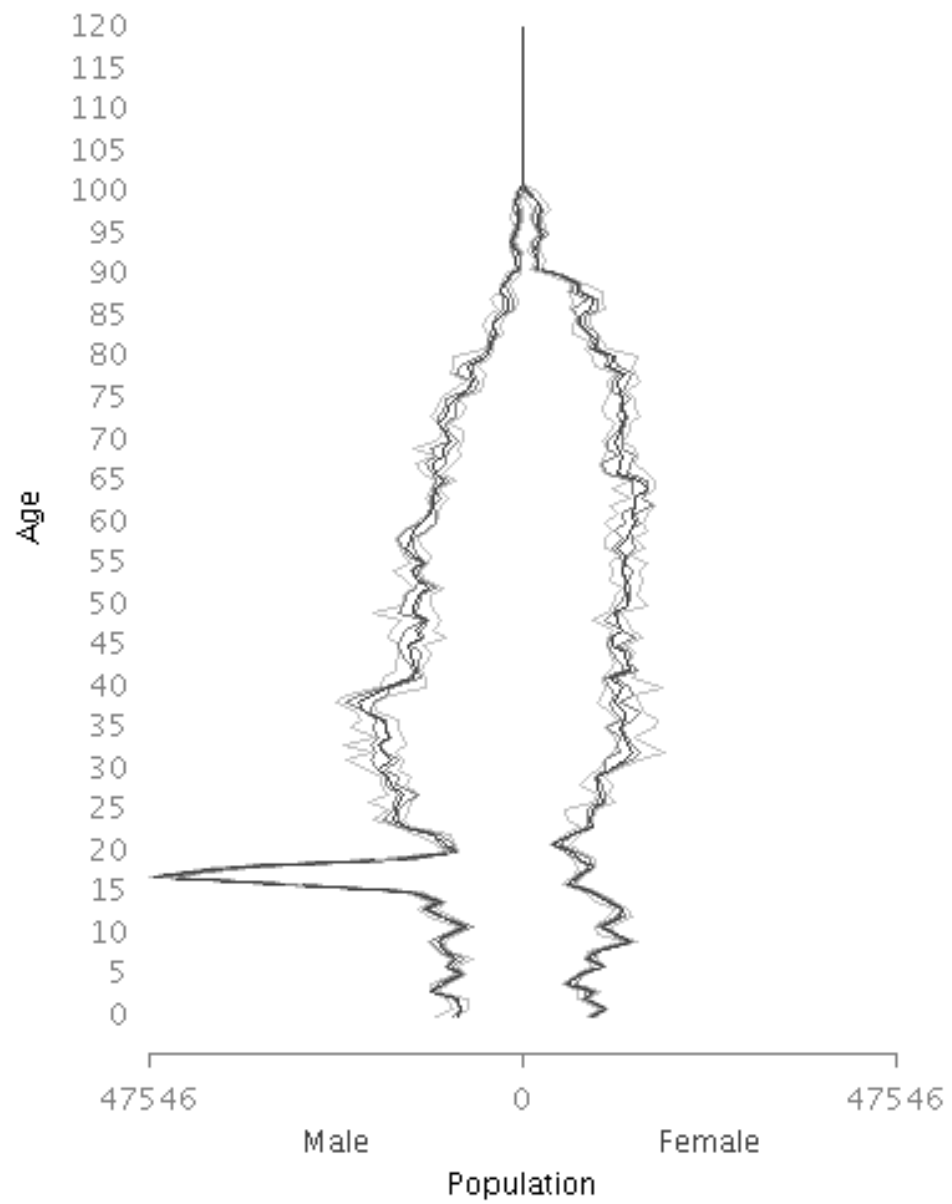
E02002330 Variation in Total Population Simulated Living Days 2001

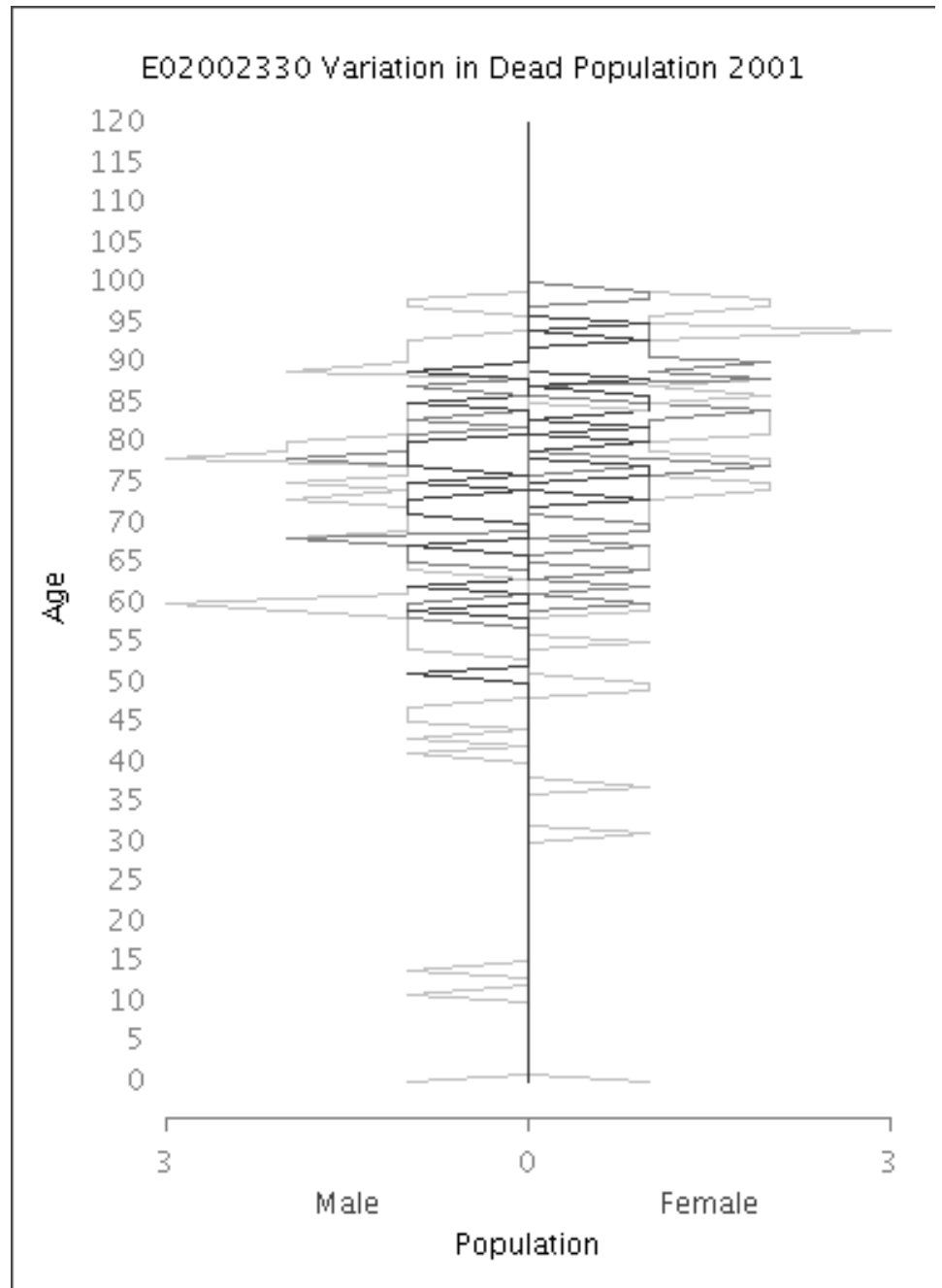
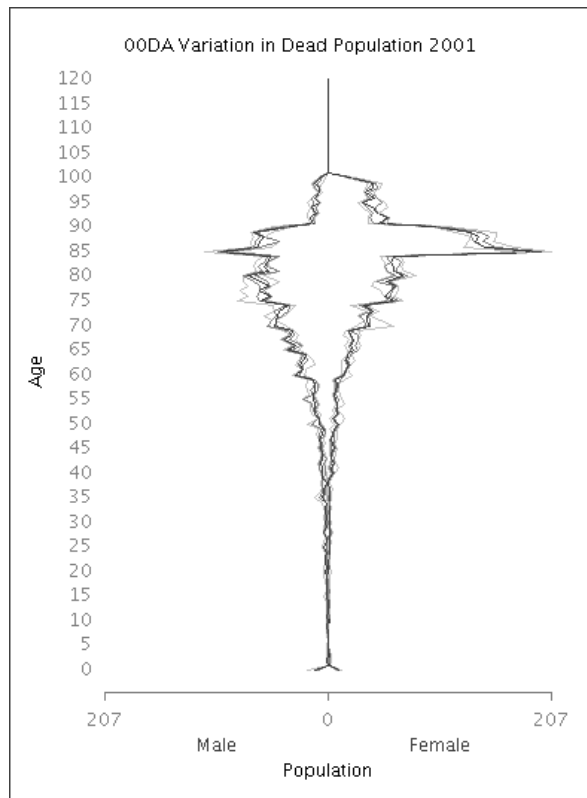


00DA Variation in Total Population Simulated Living Days 2001

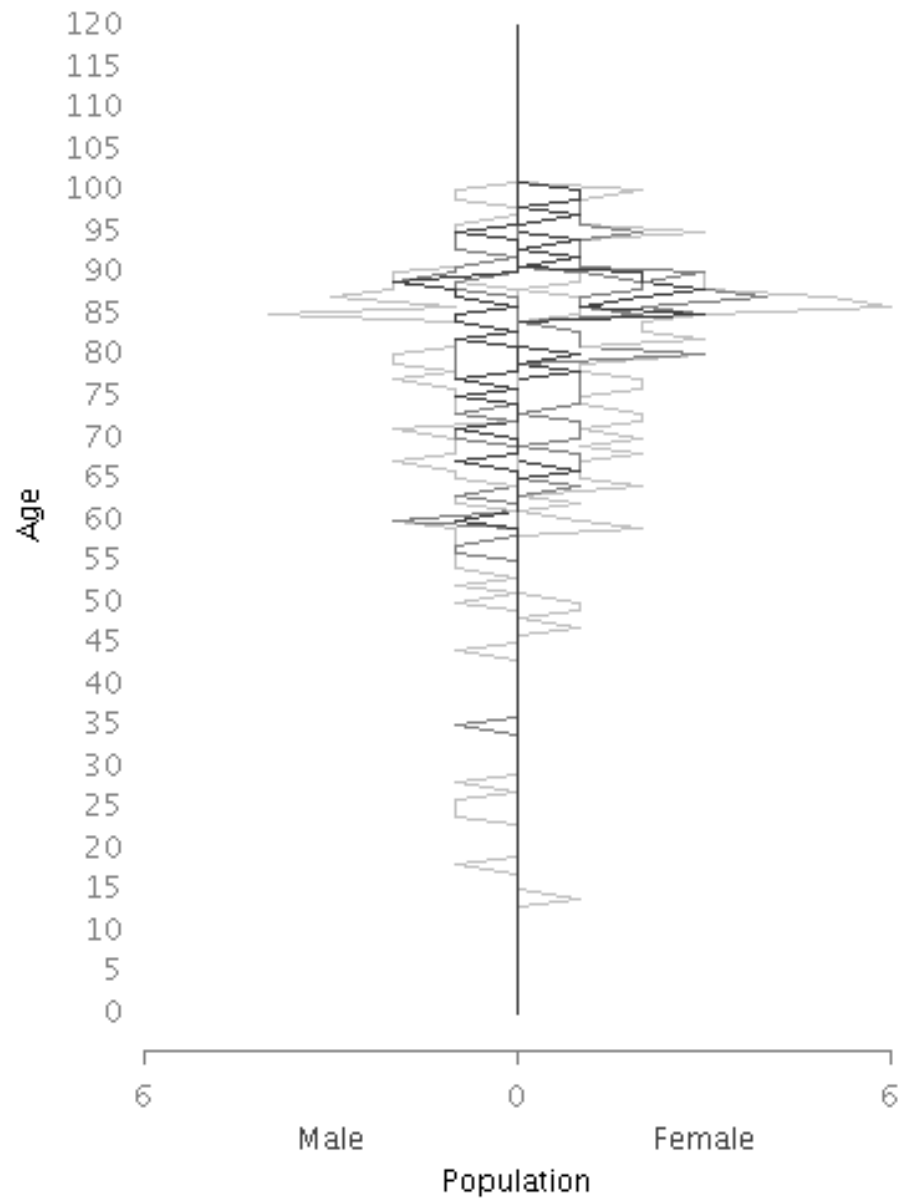


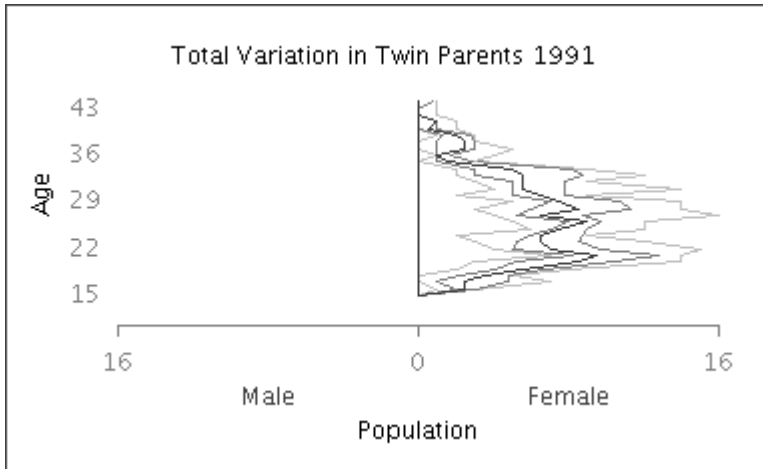
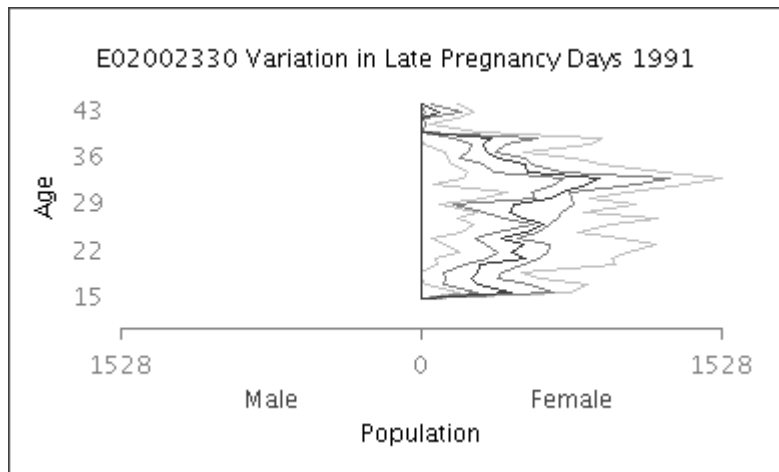
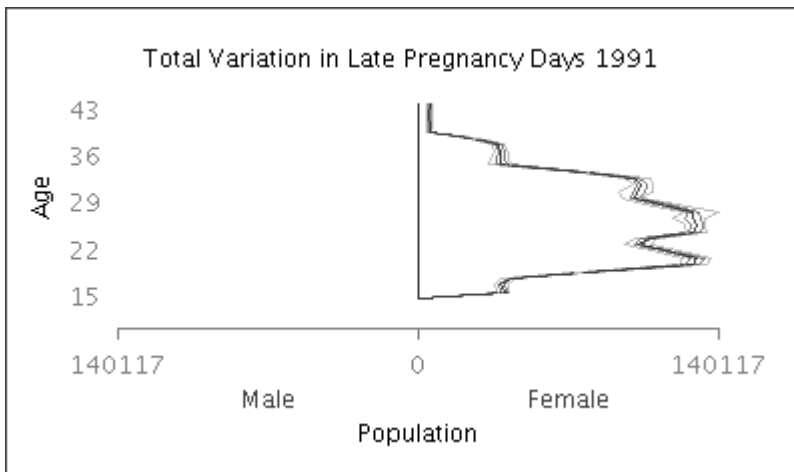
E02002331 Variation in Total Population Simulated Living Days 2001

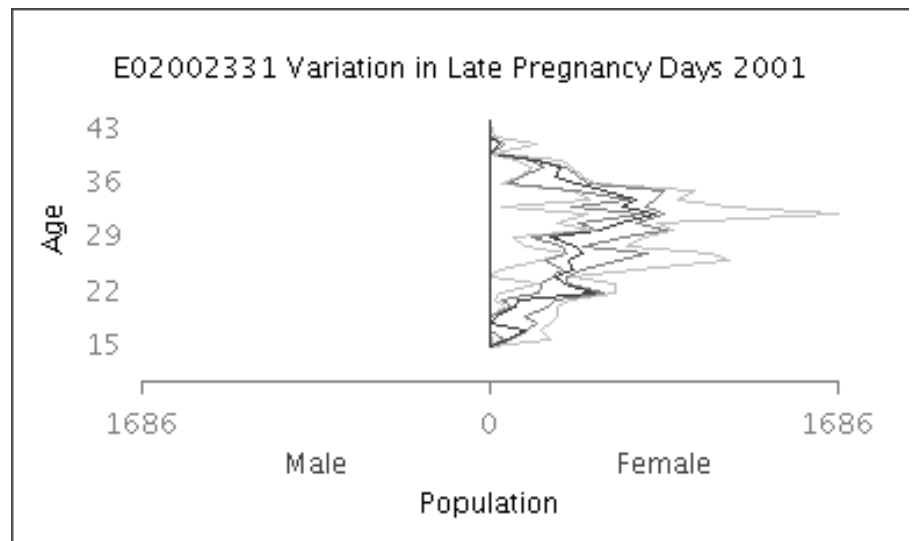
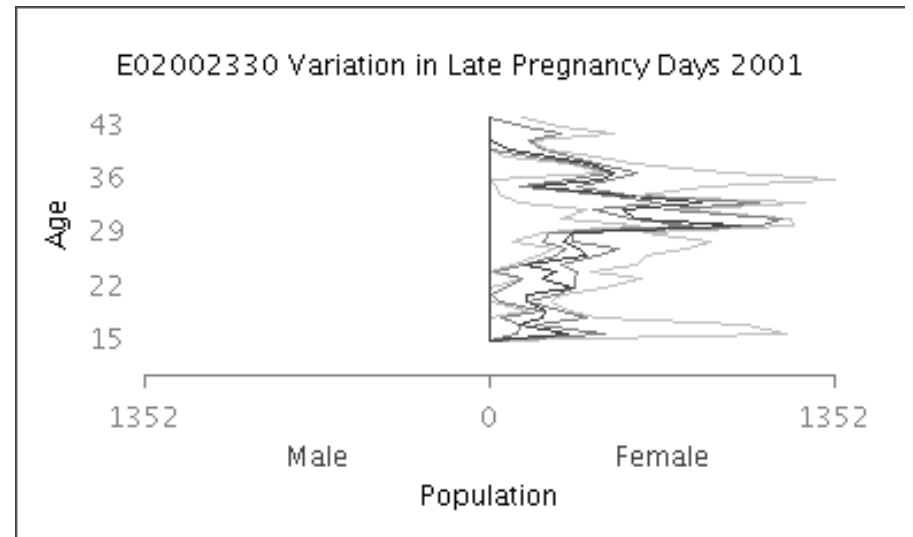
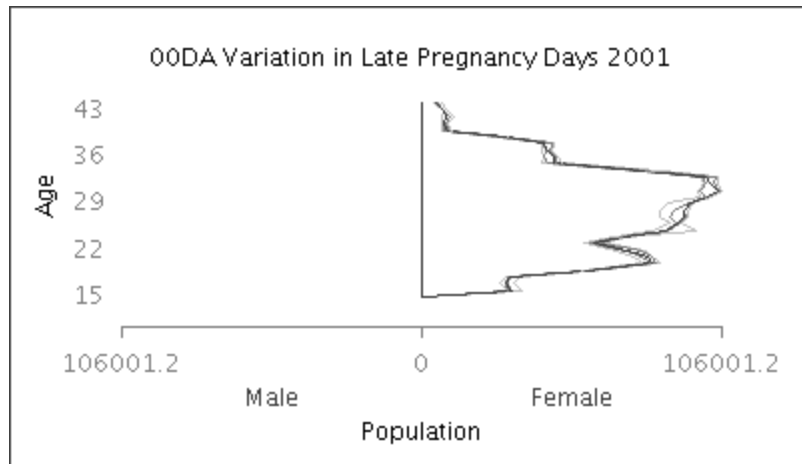


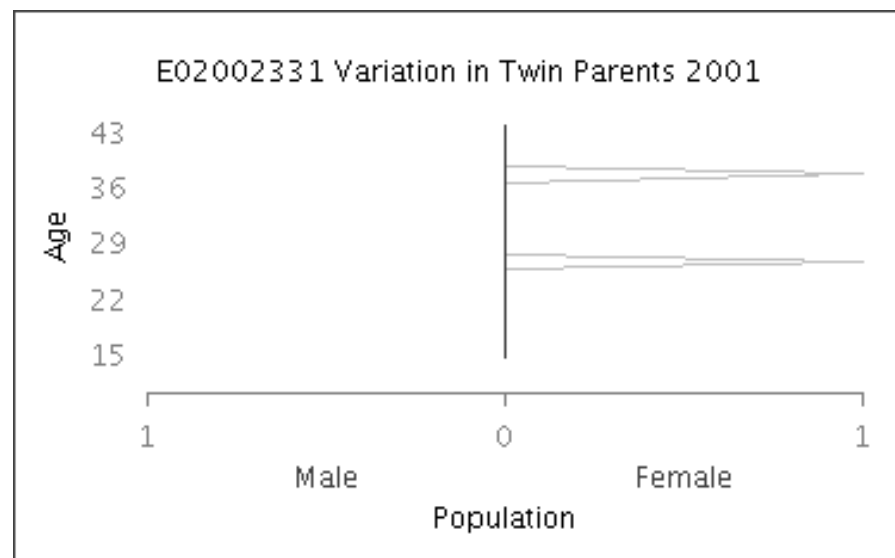
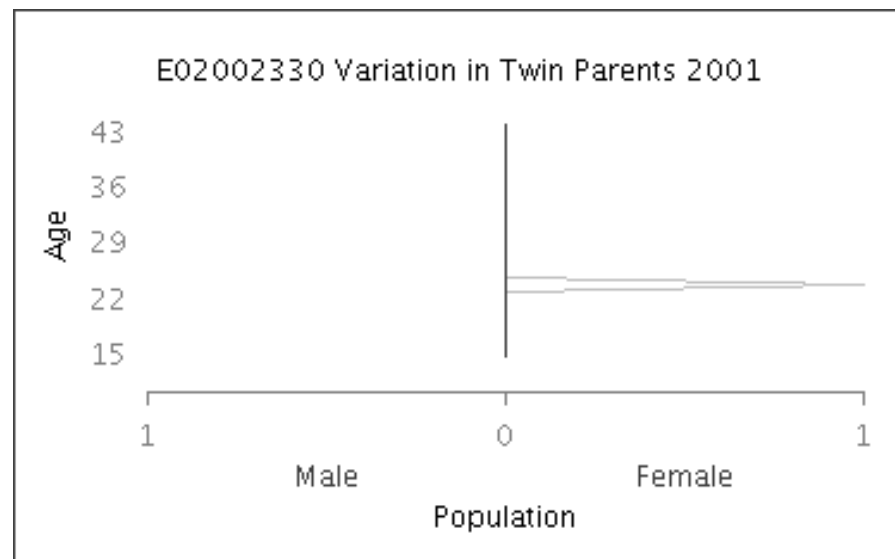
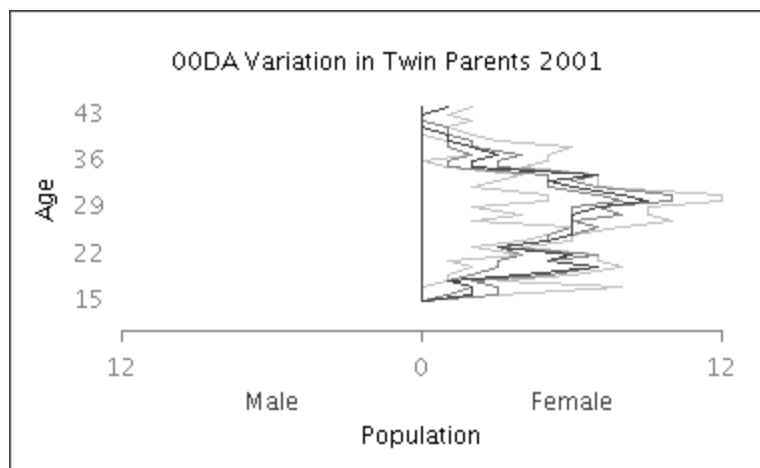


E02002331 Variation in Dead Population 2001









Migration

- Types of migration modelled
 - Immigration
 - In migration to Study Region
 - Out migration from Study Region
 - Internal migration within Study Region
- Input data
 - 2001 UK Population Special Migration Statistics
 - LAD to LAD flows by age and gender
 - OA to OA flows by age and gender
- Region (LAD to LAD) flows are primarily used

- Subregion (OA to OA) flows are used to assign individuals to subregions with each region
- A migration factor and a minimum flow

Plans and Next Steps

- Add emigration to the model
- Detailed results statistics for migration
- Simulate population change in West Yorkshire from 2001 to 2011
 - Vary migration factor and minimum flow
 - Present results at an appropriate event
 - Publish a paper on the demographic simulation model and the results for West Yorkshire
- Simulate population change for all of England from 2001 to 2011
 - Compare results with 2011 census data
 - More publication

- Further modelling

- Use [Nik Lomax](#)'s estimated migration flows for 2001 to 2011
- Constrain migration using subregion area classifications
- Allow for variations in mortality, pregnancy, miscarriage and migration rates over the year
 - Student migration
- Migrating groups (families/households)
- Fathers

- **Seek data for more detailed simulations**
 - Annual and regional miscarriage data
- **Seek collaboration with statistical offices**
- **Seek further funding**
 - Secondment to UK ONS funded by ESRC?

Feedback

- Much can be done to improve this work
- What has emerged is something like the simplest demographic model
 - There is much detail to add...
- Anyone interested in writing this up or collaborating in anyway?
- Any questions?

Thank You

<http://bit.ly/TStpJP>