

'Medium' Data and Socio-Economic Mortality

Andrew J.G. Cairns

Heriot-Watt University, Edinburgh

Director, Actuarial Research Centre, IFoA

Joint work with J.Wen, T.Kleinow and G.Peters

Insurance Data Science Conference
ETH Zurich, 14 June 2019



Actuarial
Research Centre
Institute and Faculty
of Actuaries



The Actuarial Research Centre (ARC)

A gateway to global actuarial research

The Actuarial Research Centre (ARC) is the Institute and Faculty of Actuaries' (IFoA) network of actuarial researchers around the world. The ARC seeks to deliver cutting-edge research programmes that address some of the significant, global challenges in actuarial science, through a partnership of the actuarial profession, the academic community and practitioners.

The 'Modelling, Measurement and Management of Longevity and Morbidity Risk' research programme is being funded by the ARC, the SoA and the CIA.

www.actuaries.org.uk/arc

Outline

- Context and background
- Data
- Methodology
- Results

Context and background

- Life insurance and pensions
- Mortality: traditional \rightarrow big data
 \Rightarrow improved pricing and reserving
- Considering here:
male mortality in England
(results for females similar and consistent)
- Stylised facts:
 - Mortality varies by socio-economic group
 - Mortality varies by region

Background: A Map of England!



North East
North West
Yorkshire & Humber
East Midlands
West Midlands
East of England
London
South East
South West

Not in dataset:
Scotland, Wales,
Northern Ireland

Background: Relative mortality by region

England: Variation by region (males 60-69)

North East	118%
North West	116%
Yorkshire and The Humber	107%
East Midlands	98%
West Midlands	105%
East	88%
London	105%
South East	89%
South West	87%

Values show standardised mortality (ages 60-69) by region as a percentage of national standardised mortality

Regional variation < variation by income deprivation

Background: Regional Variation

- *How much of this can be explained by underlying **socio-economic** differences?*
- *And how much variation is **geographical**?*



Data: LSOA's

- England only
- $L = 32,844$ small geographical areas (LSOA's)
- Socio-economically homogeneous
- Average size ≈ 1600 persons
- LSOA's $i = 1, \dots, L$,
single years ($t = 2001-2016$), single ages, x :
 - Deaths: $D(i, t, x)$
 - Exposures: $E(i, t, x)$ (population)
- Plus many *static* predictive variables for each LSOA

Predictive variables by LSOA

- **Indices of deprivation (2015)** (single scores per LSOA)
 - **income deprivation** (benefits)
 - **employment deprivation** (unemployment)
 - education deprivation
 - crime
 - barriers to housing and services
 - geographical barriers (distance to services)
 - **wider barriers** (overcrowding; homelessness)
 - **living environment** (housing quality; unmodernised; air quality)
- Educational attainment (levels \times age groups)
- Occupation groups (types \times age groups)
- Average weekly income
- **Average number of bedrooms**
- **# people in care homes with/without nursing**
- **Urban/rural classification** (categorical)
-

- $D(i, t, x)$, $E(i, t, x)$ deaths and exposures by LSOA
- National death rates (all t and x)

$$m(t, x) = \frac{\sum_{i=1}^L D(i, t, x)}{\sum_{i=1}^L E(i, t, x)}$$

- LSOA's ($i = 1, \dots, L$) local death rates: $m(i, t, x)$
General Model: $D(i, t, x) \sim \text{Poisson}\left(m(i, t, x)E(i, t, x)\right)$

Methodology (cont.)

General approach:

- Over a limited age range (e.g. 60-69); and
- over a limited range of years:

$$m(i, t, x) = m(t, x)F_1(i)F_2(i)$$

- $F_1(i)$ = relative risk due to socio-economic characteristics
 - local (weighted) linear regression
- $F_2(i)$ = additional relative risk capturing spatial effects
 - kernel smoothing

Methodology (cont.)

- Years: $t = t_0, \dots, t_1$
- Ages: $x = x_0, \dots, x_1$
- Actual deaths by LSOA

$$D(i) = \sum_{t=t_0}^{t_1} \sum_{x=x_0}^{x_1} D(i, t, x)$$

- Expected deaths by LSOA (no modelled effects)

$$\hat{D}_0(i) = \sum_{t=t_0}^{t_1} \sum_{x=x_0}^{x_1} m(t, x)E(i, t, x)$$

- Actual-over-expected by LSOA

$$R_0(i) = D(i)/\hat{D}_0(i)$$

Stage 1: Introduce Predictive Variables

- LSOA's: $i = 1, \dots, L$
- Predictive variables (PV): $j = 1, \dots, n_P$
- Standardised: PV type j , LSOA i

$$X(i, j) \sim N(0, 1)$$

- Purpose of standardisation:
Simplifies the system of weighting later in Stage 1
- Vector: $X(i) = (X(i, 1), \dots, X(i, n_P))'$

Stage 1: Urban versus Rural

- Urban-rural classification
 - 1: Conurbation; London (4810 LSOA's)
 - 2: Conurbation: not London (7921)
 - 3: City or town (14515)
 - 4: Rural town (3056)
 - 5: Rural village and dispersed (2542)
- Preliminary experiments \Rightarrow
contribution and importance of specific predictive variables varies significantly between urban and rural LSOA's

Stage 1: Local linear regression

- LSOA i
- Estimate the socio-economic-specific Relative Risk, $F_1(i)$
- For each i , fit an n_P -dimensional sheet around $X(i)$

$$F(i, \mathbf{x}) = a(i) + \mathbf{b}(i)^T \mathbf{x}$$

- n_P predictive variables exclude urban-rural classification
urban-rural handled in the weights, $w_1(i, j)$
- Minimise

$$S(a(i), b(i)) = \sum_j w_1(i, j) (R_0(j) - a(i) - b(i)^T X(j))^2$$

over $a(i)$ and $b(i)$

Stage 1: Local linear regression (cont.)

- Then set

$$F_1(i) = a(i) + b(i)^T X(i)$$

⇒ relative risk accounting for socio-economic factors

- Update estimated deaths:

$$\hat{D}_1(i) = \hat{D}_0(i)F_1(i)$$

Stage 1: Local linear regression (cont.)

How to calculate the weights?

- $w(i, j)$ depends on the “distance” between predictive variables $X(i)$ and $X(j)$
- $w(i, j) \rightarrow 0$ as the distance gets larger
- $w(i, i) = 0$ (facilitates cross validation)
- $w(i, j) = 0$ if LSOA's i and j are in different urban-rural groups

Stage 1 → Stage 2

$D(i)$ = LSOA actual deaths

$\hat{D}_0(i)$ = LSOA expected deaths with no predictive variables

$\hat{D}_1(i)$ = LSOA expected deaths with predictive variables

$R_1(i) = \frac{D(i)}{\hat{D}_1(i)}$ = updated actual-over-expected

Stage 2: Add location data:

$Y(i)$ = LSOA location co-ordinates
= (latitude, longitude)

Kernel smooth the $R_1(i)$ using location data.

Stage 2: Smooth A/E by Location

Estimate the *additional* location-specific relative risk

$$F_2(i) = \frac{\sum_j w_2(i, j) R_1(i)}{\sum_j w_2(i, j)}$$

Then the fitted expected deaths are

$$\hat{D}_2(i) = \hat{D}_0(i) F_1(i) F_2(i)$$

How to calculate the weights, $w_2(i, j)$?

Geographical distance \rightarrow ranks \rightarrow exponentially decaying weights

Data and Results So Far

- 2001-2015; 2001-2008; 2009-2016
- Ages: 40-49, 50-59, 60-69, 70-79, 80-89
- Predictive variables:
 - income deprivation (**elderly**; receiving government benefits)
 - employment deprivation (unemployment)
 - average number of bedrooms
 - living environment deprivation (housing quality and air quality)
 - wider barriers (overcrowding)

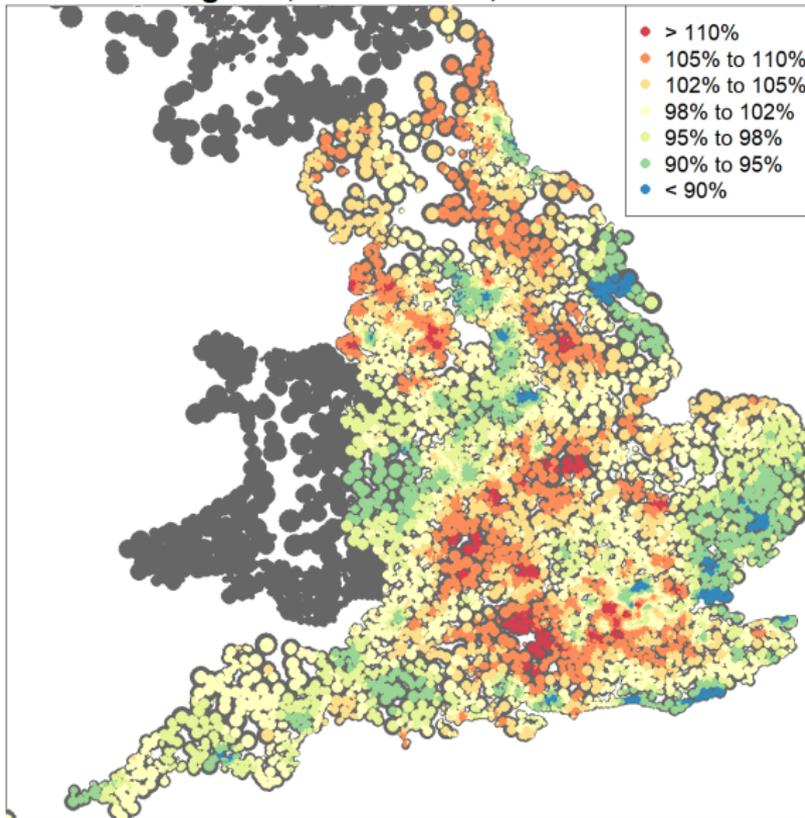
 - % in care home (age 60+ with nursing)
 - % in care home (age 60+ without nursing)

 - urban-rural classification

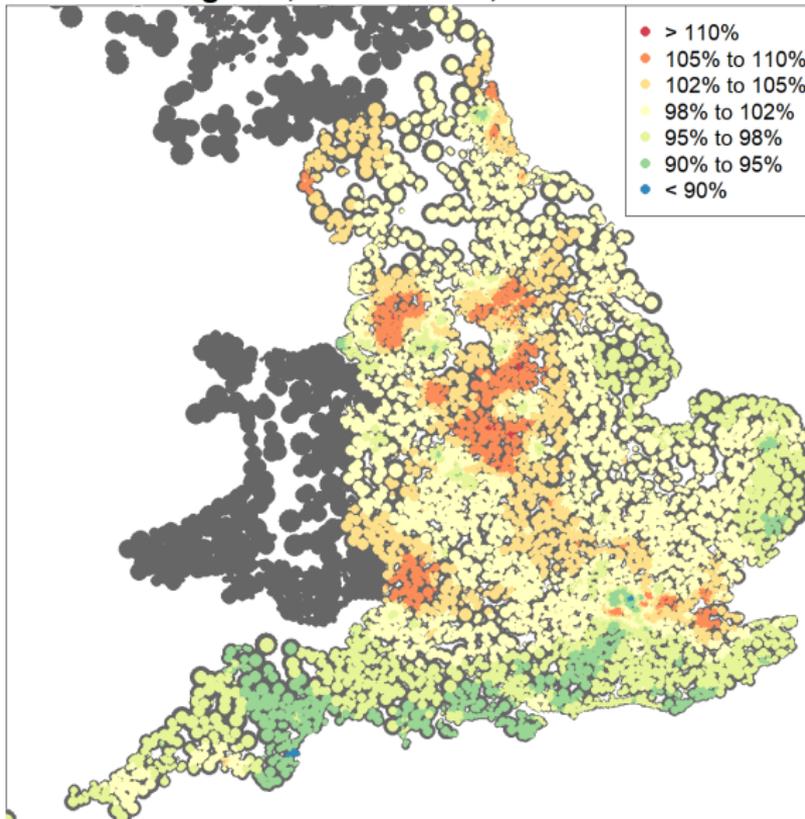
Role of Predictive Variables: Socio-Economic

- Income deprivation (elderly) and employment deprivation are the main drivers
- Employment deprivation is the main driver for younger age groups
- Income deprivation (elderly) is the main driver for older age groups
- Urban-rural classification is also an important driver
 - Rural areas: relative risk is less sensitive to variation in predictive variables
- Bedrooms, living environment and wider barriers are second order but significant
- Care homes:
 - “nuisance” variables when considering socio-economic effects
 - but including these predictive variables is very important
 - methodology allows us to filter out the impact of care homes on individual LSOA mortality
 - E.g. males 80-89 in a care home with nursing: mortality is 3x to 6x higher than not in a care home

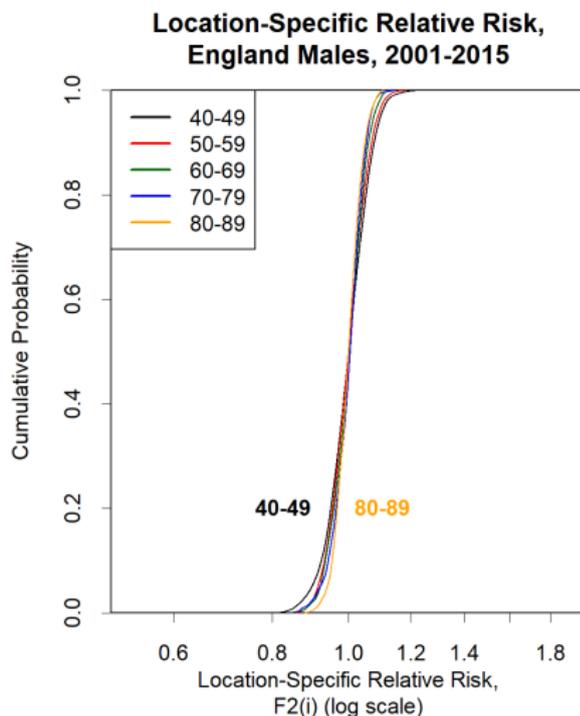
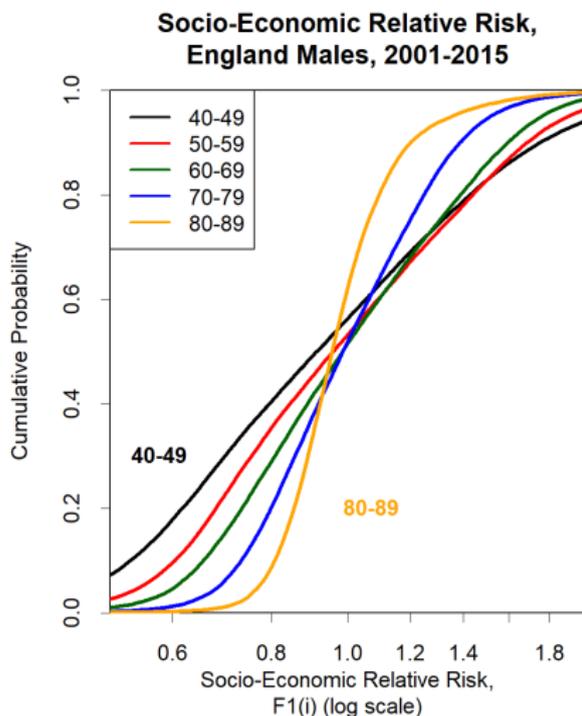
Location-Specific Relative Risk England, Males 40-49, 2001-2015



Location-Specific Relative Risk England, Males 80-89, 2001-2015



Socio-Economic vs Location-Specific Effects



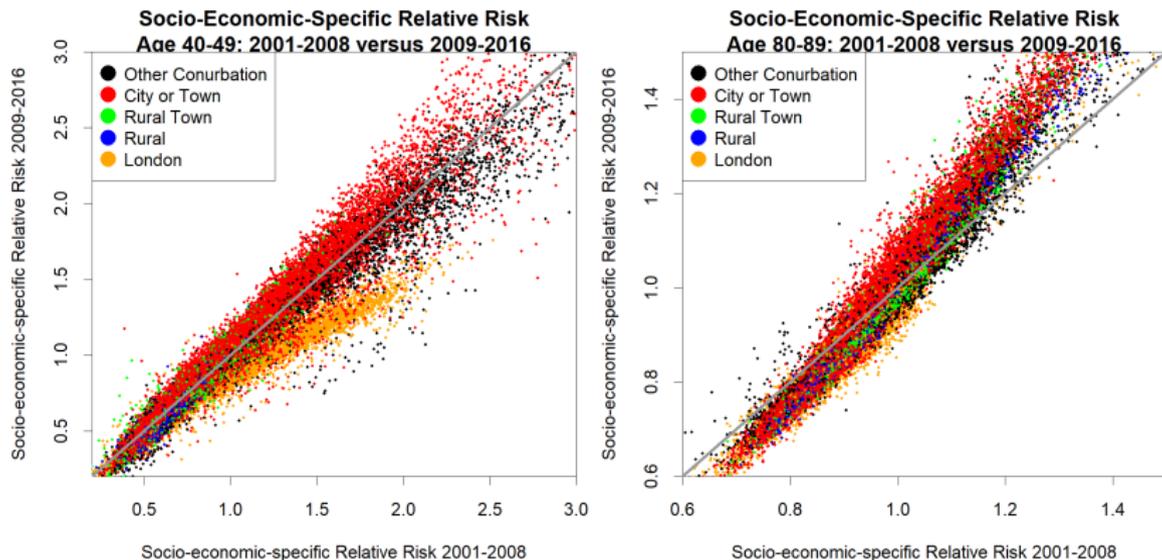
- Location contributes 1.3% to 3.5% of the variance in the relative risk

Actual-over-expected: Ages 60-69

Region	No effect	Socio-economic only	Full Model
North East	118	100	99
North West	116	102	100
Yorkshire and The Humber	107	100	100
East Midlands	98	100	99
West Midlands	105	99	100
East	88	96	98
London	105	100	99
South East	89	101	100
South West	87	94	99

- Similar patterns for other age groups and for females

2001-2008 versus 2009-2016: Ages 40-49 and 80-89



- Widening inequality gap at 80-89
- Stable gap at 40-49, except London: narrowing gap

Conclusions

- Spatial/regional effects are significant
- But much less important than socio-economic (non-regional) effects
- Both effects: can these be used to improve predictions of insurance and pensions mortality?
- Longer term objective:
Can we form e.g. 10 clusters of LSOA's with similar mortality experience over the period of observation?
- Work in progress

Thank You!

Questions?

E: A.J.G.Cairns@hw.ac.uk

W: www.macs.hw.ac.uk/~andrewc/ARCresources