

1 **Health By Wealth: Optical character recognition and LLMs with**
2 **population-level probate and administrative data uncover**
3 **substantial inequalities in health over the very long run**

4 Naomi Muggleton^{1,2}, Aaron Reeves^{3,4}, **Charles Rahal**^{5,6}, Paul Moore⁴, Linda Li⁷, and
5 Alexandra Rottenkolber⁸

6 ¹Warwick Business School, University of Warwick

7 ²Brasenose College, University of Oxford

8 ³Department of Sociology, London School of Economics and Political Science

9 ⁴Department of Sociology and Intervention, University of Oxford

10 ⁵Oxford Population Health, University of Oxford

11 ⁶Nuffield College, University of Oxford

12 ⁷Department of Methodology, London School of Economics

13 ⁸Department of Sociology, Linköping University

14 November 1, 2025

15 **This is a preliminary work in progress. Please do not cite or distribute.**

16 **Abstract**

17 Health inequalities in life expectancy between the rich and poor remain a persistent chal-
18 lenge. Yet, their long-run evolution has been poorly understood due to limited historical data
19 combining socio-economic status and mortality. Existing evidence is often fragmented across
20 time and place, leaving unresolved when and why disparities changed. We address this by
21 constructing a novel dataset spanning 130 years, linking 16 million digitized probate records to
22 66 million death registrations, providing unprecedented detail on wealth and mortality in Eng-
23 land and Wales. This enables macro-level analysis from individual microdata. Between 1860
24 and 1990, life expectancy inequalities narrowed markedly, particularly after 1940. The gap
25 at age 20 between the wealthiest 1% and poorest groups declined from 17.3 years for women
26 and 15.8 years for men to just 2–3 years in recent decades, coinciding with the emergence
27 of the welfare state. We further link genealogical data to the same probate records to assess
28 intergenerational effects. Contrary to earlier debate, we find clear evidence that individuals
29 born into wealthier families lived longer before 1900. Our results show that health inequali-
30 ties, though enduring, can be reduced. The remarkable rise in life expectancy after 1850 was
31 driven largely by gains among the least wealthy, reflecting collective societal efforts—through
32 public health initiatives, poverty reduction, and universal healthcare—to extend healthy lives
33 to all.
34

35 **Keywords:** *Machine Learning, Social Inequality, Population Health, Mortality*
36

37 It is not contentious to hypothesize that health inequalities are both pervasive, and durable.
38 They exist in every country for which we have data (Murtin et al., 2017), and they have
39 been largely stable in recent years (Mackenbach et al., 2018), not least but especially in
40 countries which have completed their demographic transition. The degree to which health
41 inequalities are in fact durable over the long run is, however, unclear. In part this is because
42 the data used to analyse this question only covers a relatively short period of time (a few
43 decades at most), and health inequalities are produced over and reproduced across entire life
44 courses. Intergenerational analysis over the long run is sincerely and especially lacking. Other
45 historical work is almost exclusively focused on individual – often geographically unique –
46 parts of space and time, resulting in a piecemeal understanding of this critical demographic
47 and population health problem. Limited examples include, for example, Chadwick (1843)
48 which estimates that the average age at death in Liverpool in 1840 for professional persons
49 was 35, but just 15 for laborers, mechanics, and servants. Other, more recent examples come
50 from the influential ‘Whitehall Studies’ (such as Marmot and Brunner, 2005) which show that
51 each hierarchical level of civil servant job seniority experienced a higher mortality rate than
52 the group one step higher. This results in academics having to rely on snapshots of information
53 to draw conclusions, despite a multitude of theories which explicate the link between things
54 like socioeconomic status, wealth, and health outcomes (Phelan et al., 2010). Indeed, data on
55 health inequalities over the very long-run is extraordinarily uncommon, and is usually only
56 stitched together into snapshots based on non-comparable projects which are based on diverse
57 and disparate sources of information. We simply do not know when, if at all, and – if so –
58 why inequalities in life expectancy have changed.

59 Here we present the first, original and yet unpublished results from an unusually and
60 extraordinarily detailed and uniquely constructed dataset. While our research design is ex-
61 traordinarily computational in nature, it addresses the dearth of studies as described above.
62 We parse digitized versions of handwritten probate records using the latest techniques in
63 Computer Vision and Optical Character Recognition. This data comes from the National
64 Probate Calendar (created by the Probate Registry); a register of proved wills and subse-
65 quent bequeathments available following the enactment of the Court of Probate Act 1857
66 across administrations in England and Wales since 1858. This data is at the almost entire
67 population level. Processing this data for demographic analysis involved a complex pipeline
68 to make the data ‘research ready’. We first set up 138 virtual machines; one for each year.
69 We then sent requests for probate images at a rate of – on average – one every 10 seconds. We
70 then used wildcard searches to iterate through all surnames between Aa* and Zz*, and saved
71 the images for downstream processing. This processing – on over one million .png images –
72 involved; inverting colours, binarizing grayscale, programatically removing ink blots, detect-
73 ing margins with Gaussian processes, programatically removing dust, removing handwritten
74 margin annotations, calculating pixel angulation and skew, page resizement (based on optimal
75 extraction), the artificial creation of bounding boxes, optical character recognition, the use of
76 large language models to detect extraction mistakes, and finally, fuzzy regular expressions to
77 create tabular data ready for analysis. The millions of images, when stacked on their edges,
78 would be over one hundred metres in length if they were printed in their physical version.
79 These data encompass all wills proved in the UK above certain administrative thresholds
80 (with rare exception, such as for members of the Royal Family), allowing us to capture the
81 full wealth distribution of almost all of a large population. This ranges – over 16,634,470
82 rows of cleaned data – from the upper echelons of British society (our dataset includes, for
83 example, Charles Darwin who left £145,911 in 1842, inflation adjusted to ~£18.75m in 2025),
84 to tradespeople and widows leaving a few pounds (such as John Croxton, a comparatively un-
85 known labourer who died in December 1854 and left £5 (inflation adjusted to £625 in 2025)
86 to a farmer. To note – and while we have full ‘ethical review board’ approval – our analysis
87 is exempt from GDPR compliance restrictions because it primarily involves individuals who
88 have since deceased. Based on a complex record linkage procedure which involves *extensive*
89 manual verification and hand labeling for model training, we reconcile the probate with admin-
90 istrative death records for the purpose of recovering dates of birth which allows us to compute
91 individual level longevity at scale for the first time. This comes in the form of digitised and

92 population wide death registration data (n=66,161,380) which involves thousands of manual
93 transcriptors, as well as with other auxiliary sources, such as tools of gender inference and
94 vast geneological data in order to examine the intergenerationla transmission of wealth and
95 its subsequent effect on life expectancy. Our computational pipeline is visualised in Figure 1.

96 Despite temporal variation, we can at times match the vast majority of death records
97 onto the probate due to our advanced record linkage techniques. This is done through a new,
98 novel, and custom-built Python 3.x library which also computes a range of life expectancy
99 based statistics, visualisations, and statistical tests. We aim to publish this library ‘Open
100 Source’ (under a GNU GPL 3.0 license) in advance of the conference, so that others may
101 gain utility from it. This library also computes an extremely large number of cohort based
102 life tables – tens of thousands – based on the entire feasible range of assumptions which
103 any reasonable demographer could have chosen to make. This includes not just different
104 conceptualisations of wealth, but also a multitude of different methodological assumptions
105 with regards to smoothing of mortality rates. While the foundations of this library are based
106 on [Preston \(2000\)](#), it also leans on ideas from [Caughley \(1966\)](#) and a multitude of other,
107 more recent methodological literature ([Riffe et al., 2019](#)). The existence of the ‘personal
108 effects’ variable which is programatically extracted from the raw images allows us to make
109 comparative statements about longevity by wealth at death at an unprecedented temporal
110 span and resolution. Figure 2, for example shows simple age at death calculations across each
111 of the fourth quartiles of the wealth at death distribution (by year) using observations on
112 individuals found in the probate records. However, by virtue of the fact that we have the
113 entire population register of deaths, we are able to decompose observations into those who
114 hold an amount of wealth which would put them in the ‘top 1%’ of population, those who
115 held some wealth, and those members of the population who held no wealth (or wealth below
116 the small threshold for inclusion). This is shown in Figure 3. For both male and females,
117 for all years in our sample, and for life expectancies at almost all ages, we observe enormous
118 stratification by personal effects at death. This gap narrows over time – and especially since
119 the 1940s – but in earlier parts of our sample, the gap in life expectancy at age 20 between
120 the ‘1%’ and those holding no or negligible wealth is as wide as 17.3 years for females, and
121 15.8 years for males. In addition to Figures 1-3, our analytical pipeline also outputs a large
122 number of other, illuminating examples. Further to this, we link individuals across time in
123 order to combat endogeneity, again showing that life expectancy is a function of wealth, but,
124 there, of wealth held by parents. This enables an analysis of the effects of transmission of
125 wealth before individuals have the opportunity to acruce wealth through other, non-hereditary
126 channels.

127 We provide the most granular picture on inequalities in life expectancy to date. This
128 analysis elucidates when, and among which parts of the wealth distribution these changes in
129 life expectancy occurred over more than a century of unprecedented data. Declines in health
130 inequalities are driven by economic growth and increased government support to low-income
131 households, and touch on important debates across various academic disciplines (demography,
132 economic history, sociology, and population health), contributing to theories of demographic
133 transition, the ‘fundamental causes’ of health inequalities, and the impact of industrialization
134 on society. Our advanced computational techniques extend and refine classical demographic
135 frameworks. Through the integration of probabilistic models of lifespan with detailed wealth
136 distribution analyses, our approach addresses persistent gaps in the literature concerning the
137 evolution of health disparities over the long run. Our custom-built Python library enables
138 extraction of demographic parameters and facilitates the computation of robust confidence
139 intervals, enhancing the statistical rigor of our findings. Moreover, the implementation of
140 both deterministic and probabilistic record linkage methods significantly improves the ac-
141 curacy of matching probate records with death registration. This framework solidifies our
142 historical findings and reveals subtle interactions between socioeconomic status and mortality
143 that have previously been overlooked. Consequently, our analysis contributes a comprehensive
144 understanding of the determinants of health inequalities through a new computational lens,
145 offering valuable insights for policymakers and researchers interested in mitigating the adverse
146 effects of social and economic disparities on population health.

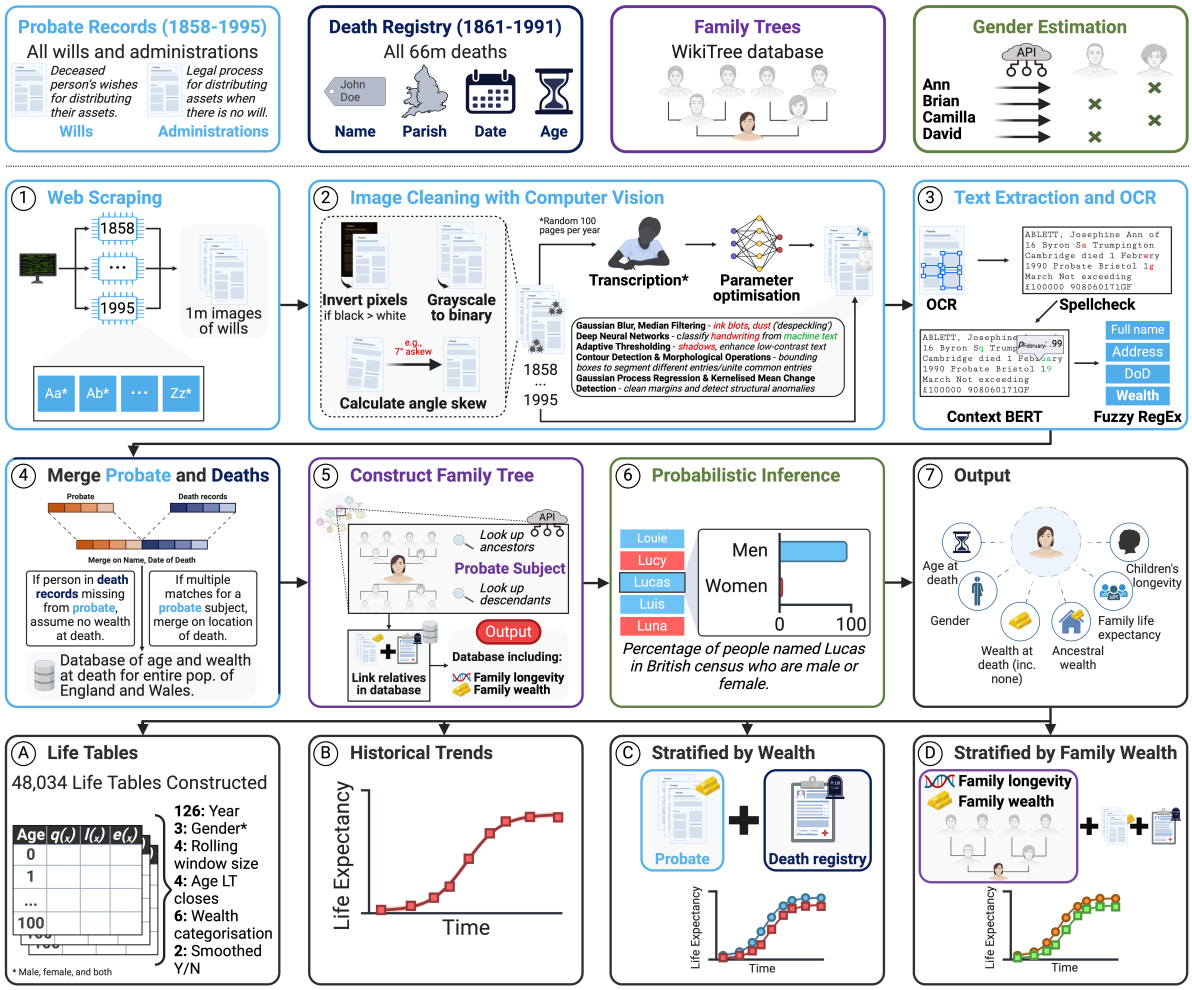


Figure 1: Our Computational Pipeline: A schematic which shows the different parts of our computational pipeline.

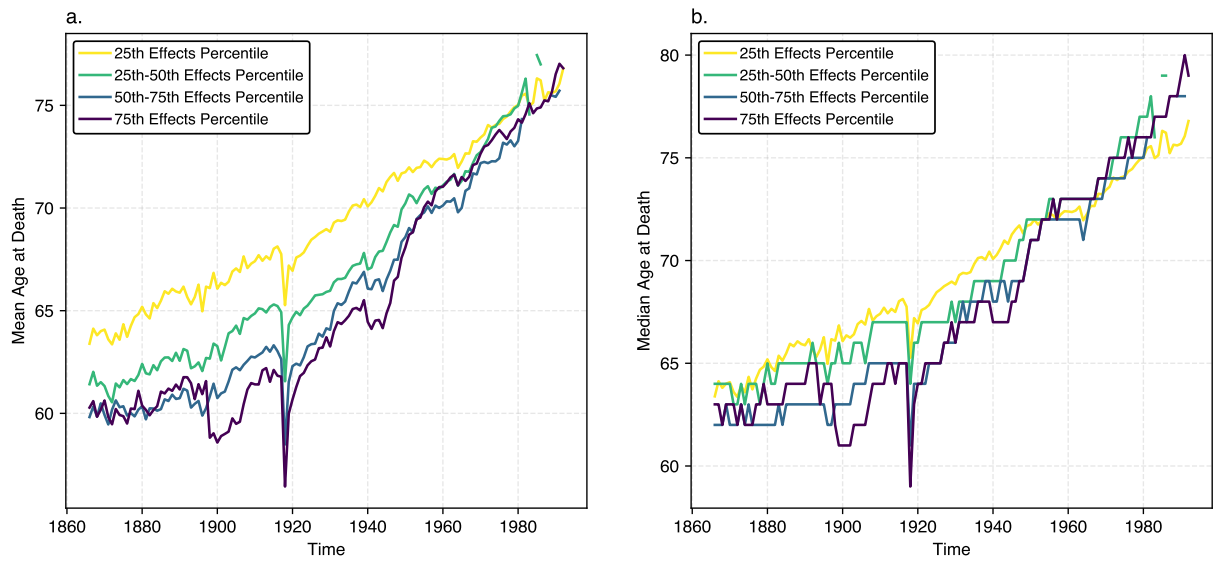


Figure 2: Age at Death by Effects: Subfigure 'a.' shows mean age at death across the four quartiles of the wealth distribution. Subfigure 'b.' shows the equivalent calculations, but instead using the median instead of the mean. Both figures are for male and females combined.

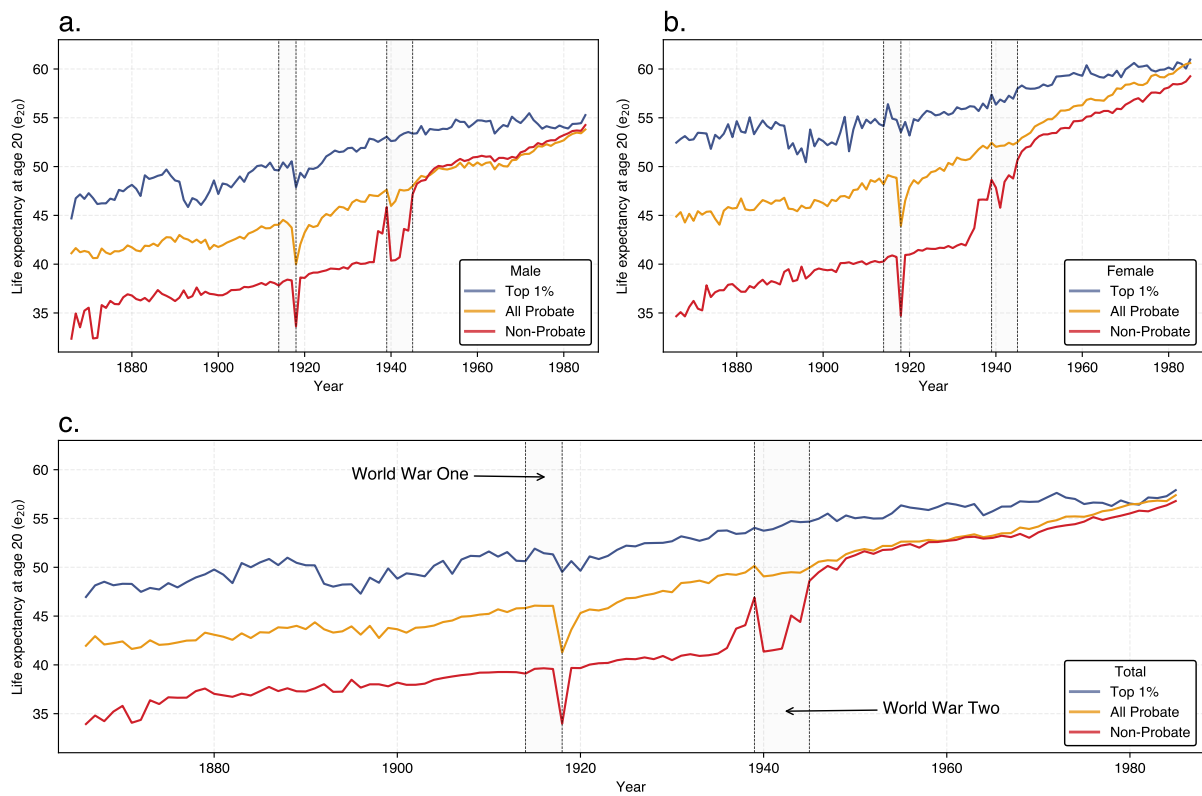


Figure 3: Life Expectancy at 20 years old: Subfigure 'a.' shows life expectancy at 20 years old for men, subfigure 'b.' for women, and both combined for subfigure 'c.'. Lines denote the wealth ('personal effects') of those estimated to be in the 'top 1%' of the population's wealth distribution, all people who have left *some* effects (on the probate register), and all other members of the general population.

References

Caughley, G. (1966). Mortality patterns in mammals. *Ecology* 47(6), 906–918.

Chadwick, E. (1843). *Report on the sanitary conditions of the labouring population of Great Britain: A supplementary report on the results of a special inquiry into the practice of interment in towns. Made at the request of Her Majesty's principal secretary of state for the Home department.* W. Clowes and sons.

Mackenbach, J. P., J. R. Valverde, B. Artnik, M. Bopp, H. Brønnum-Hansen, P. Deboosere, R. Kalediene, K. Kovács, M. Leinsalu, P. Martikainen, et al. (2018). Trends in health inequalities in 27 european countries. *Proceedings of the National Academy of Sciences* 115(25), 6440–6445.

Marmot, M. and E. Brunner (2005). Cohort profile: the whitehall ii study. *International journal of epidemiology* 34(2), 251–256.

Murtin, F., J. Mackenbach, D. Jasilionis, and M. M. d'Ercole (2017). Inequalities in longevity by education in oecd countries: Insights from new oecd estimates. *OECD Statistics*.

Phelan, J. C., B. G. Link, and P. Tehranifar (2010). Social conditions as fundamental causes of health inequalities: theory, evidence, and policy implications. *Journal of health and social behavior* 51(1_suppl), S28–S40.

Preston, S. (2000). *Demography: Measuring and modeling population processes.* (No Title).

Riffe, T., J. Aburto, M. Alexander, S. Fennell, I. Kashnitsky, M. Pascariu, and P. Gerland (2019). Demotools: An r package of tools for aggregate demographic analysis.