



Institute  
and Faculty  
of Actuaries



# Obtaining public data for public purposes

Prof Elena Kulinskaya (UEA)



# Content

- What are the sources of public data
- What other data would be useful
- Can private companies help public research?



# Data Sources: Randomised Clinical Trials and Meta-analyses of RCTs

- Randomised clinical trials are the ‘gold standard’ for evidence of effectiveness of health interventions
  - Confounders randomized equally to both groups (in theory)
- Generalisability from trial participants to general population
  - Exclusion on grounds of age, comorbidity, intolerance to intervention
- Short follow-up (Maximum 5 years)
- Commercial trial data not available for individual scrutiny
  - Lack of transparency
- Large observational datasets can fill these gaps with robust statistical analyses



# Data Sources: Observational Data

## Administrative Data

- Primary care CPRD and THIN
- Hospital Data HES
- Social care data: care and nursing homes, older people in the community, etc.



## Specialised registries

- Cancer Registries
- National Joint Registry



# Survey and Census data

- Health Surveys for England, Wales, Scotland and Northern Ireland (yearly) A total of 8,795 adults and 2,185 children were interviewed in 2013.
- Census Data (every 10 years from 1971, 1% of the population of England and Wales, together with records for other people in their households.
- The ONS Longitudinal Study (LS) is a linked set of individual census responses. Results from the censuses of 1971, 1981, 1991, 2001 and 2011 and life event information. Linked to health data, Cancer Registry, <http://calls.ac.uk/>
- The Health and Social Care Information Centre (HSCIC) commissions and holds many data collections, and is supposed to provide linkages.



# Longitudinal Cohort Studies

- An estimated 3.5 % of the population have taken part in one or more UK cohort studies (MRC, 2014).
- **Birth cohorts**
- MRC National Survey of Health and Development, a representative sample (N=5362) of men and women born in England, Scotland or Wales in March 1946
- Aberdeen Children Of the Nineteen Fifties (N=12,150) [all Aberdeen primary school children born between 1950 and 1956]
- The 1970 British Cohort Study (BCS70) (N=17,000) follows the lives of more than 17,000 people born in England, Scotland and Wales in a single week of 1970



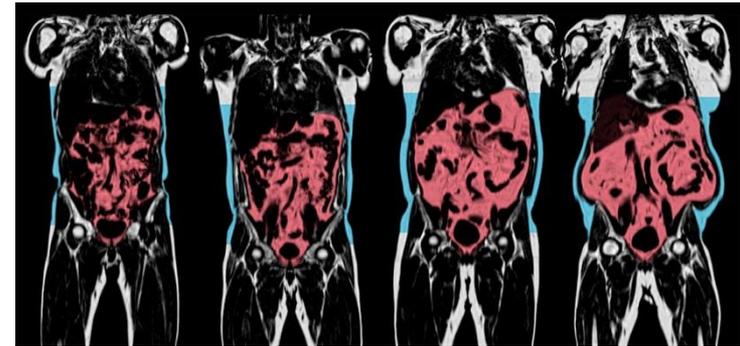
# Representative Population Studies

- English Longitudinal Study of Ageing (ELSA): Waves 0-7, 1998-2015 (over 50+ in 2002, N=12,000, adding extra samples every 2 years)
- Cognitive Function and Ageing Studies (CFAS(N=7635), CFAS II(N=7762))
- European Prospective Investigation of Cancer (EPIC) (N=30,000 in Norfolk, recruitment 1993-97, 40-79 yo, longitudinal information on diet+health)
- Dementias Platform combines 30 cohort studies



# UK Biobank

- 500,000 volunteers, aged 40-69, recruited 2006-2010
- biological samples
- detailed demographic and lifestyle
- MRI and DEXA scanning, n=10,000
- Diet questionnaire longitudinally
- Genetic information
- linked to ONS, Cancer Register, HES data, diagnostic imaging and primary care data by 2017.
- At the moment, there are about 10,000 deaths among the participants
- Self-selected sample, so may be biased.



# Use of Biobank Data for longevity research



## ARTICLE

Received 28 Sep 2015 | Accepted 29 Feb 2016 | Published 31 Mar 2016

DOI: 10.1038/ncomms11174

OPEN

## Variants near *CHRNA3/5* and *APOE* have age- and sex-related effects on human lifespan

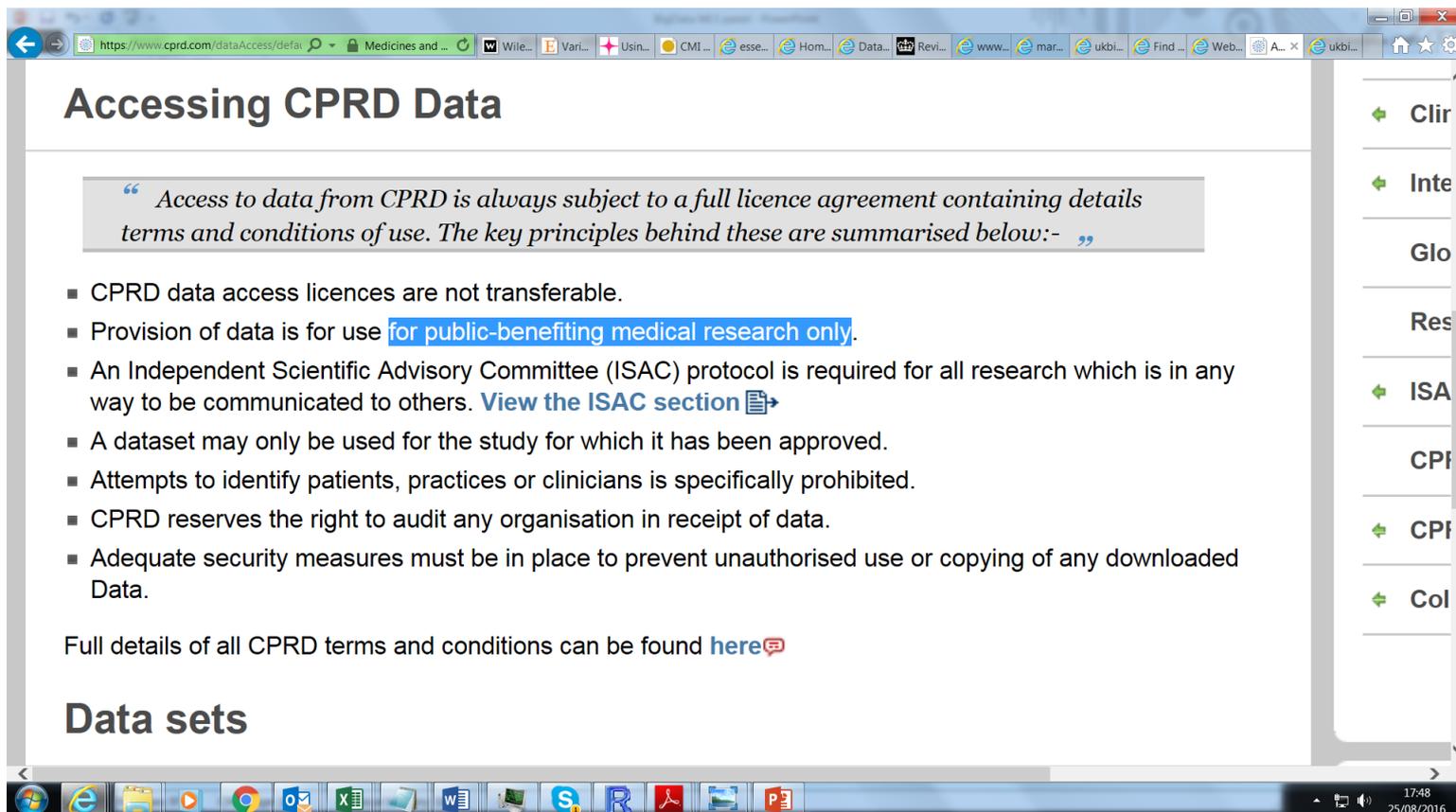
Peter K. Joshi<sup>1</sup>, Krista Fischer<sup>2</sup>, Katharina E. Schraut<sup>1,3</sup>, Harry Campbell<sup>1</sup>, Tõnu Esko<sup>2,4,5,6</sup> & James F. Wilson<sup>1,7</sup>

Lifespan is a trait of enormous personal interest. Research into the biological basis of human lifespan, however, is hampered by the long time to death. Using a novel approach of regressing (272,081) parental lifespans beyond age 40 years on participant genotype in a new large data set (UK Biobank), we here show that common variants near the apolipoprotein E and nicotinic acetylcholine receptor subunit alpha 5 genes are associated with lifespan. The effects are strongly sex and age dependent, with *APOE*  $\epsilon 4$  differentially influencing maternal lifespan ( $P = 4.2 \times 10^{-15}$ , effect  $-1.24$  years of maternal life per imputed risk allele in parent; sex difference,  $P = 0.011$ ), and a locus near *CHRNA3/5* differentially affecting paternal lifespan ( $P = 4.8 \times 10^{-11}$ , effect  $-0.86$  years per allele; sex difference  $P = 0.075$ ). Rare homozygous carriers of the risk alleles at both loci are predicted to have 3.3-3.7 years shorter lives.



Institute  
and Faculty  
of Actuaries

# Difficulties: Health Research only



**Accessing CPRD Data**

*“ Access to data from CPRD is always subject to a full licence agreement containing details terms and conditions of use. The key principles behind these are summarised below:- ”*

- CPRD data access licences are not transferable.
- Provision of data is for use **for public-benefiting medical research only**.
- An Independent Scientific Advisory Committee (ISAC) protocol is required for all research which is in any way to be communicated to others. [View the ISAC section](#)
- A dataset may only be used for the study for which it has been approved.
- Attempts to identify patients, practices or clinicians is specifically prohibited.
- CPRD reserves the right to audit any organisation in receipt of data.
- Adequate security measures must be in place to prevent unauthorised use or copying of any downloaded Data.

Full details of all CPRD terms and conditions can be found [here](#)

## Data sets



# Biobank: not available for actuarial research?

The screenshot shows a web browser window displaying the UK Biobank website. The address bar shows the URL <http://www.ukbiobank.ac.uk/register-apply/>. The page features a dark teal header with the 'biobank<sup>uk</sup>' logo on the left and the contact information 'Call Us On: 0800 0 276 276' on the right, with the tagline 'Your feedback is important to us, tell us what you think'. Below the header is a navigation menu with links for 'About', 'Participants', 'Resources', 'Scientists', 'Data Showcase', 'Register & Apply', 'Approved Research', and 'Publications'. The 'Register & Apply' link is highlighted. A sidebar on the left contains a menu with 'Home', 'About UK Biobank', 'Participants', and 'Scientists' (all with right-pointing arrows), and a list of links: 'How to register', 'Register to use data: direct link', 'Data Showcase: direct link', 'Getting started: helpful info', 'Useful resources', 'Genetic data', and 'Imaging data'. The main content area is titled 'Register & Apply' and contains the text: 'All bona fide researchers can apply to use the UK Biobank resource [for health related research that is in the public interest](#). Watch our short video animation below to find out about the steps from registration through to approval and the information you will need to supply to us. A guide to access is also provided in the document below.' Below this text is a video player with the title 'How to register & apply to UK Biobank' and a play button. The video thumbnail shows the text 'DOES THE RESEARCH MEET OUR CRITERIA?' and a checklist. The Windows taskbar at the bottom shows various application icons and the system clock indicating 17:44 on 25/08/2016.



# Difficulties: confidentiality, data sharing, data linkage

The screenshot shows the BBC News website. The main article is titled "NHS Care.data information scheme 'mishandled'" by Chris Vallance, dated 18 April 2014. The article features an image of a hand holding a pen over a document. To the right, there is a "Top Stories" section with three items: "British boy, 14, killed in Italy quake" (11 minutes ago), "Optometrist sentenced over boy's death" (1 hour ago), and "French court suspends 'burkini' ban" (3 minutes ago). Below this is a "Features" section. At the bottom of the page, a "BREAKING" banner reads "French court suspends controversial 'burkini ban' in coastal town of Villeneuve-Loubet". The Windows taskbar at the bottom shows the time as 14:58 on 26/08/2016.

“The National Data Guardian proposes a New consent / opt-out model for consultation to enable people to opt out from their personal confidential data being used for purposes beyond their direct care.”



Institute  
and Faculty  
of Actuaries

# ESRC Business and Local Government Data Research Centre (BLG DRC)

- Funded under the ESRC's Big Data Network, £5m, 2014-19
  - An Eastern ARC (Essex, UEA & Kent) partnership
- Exploitation of data to benefit researchers, data owners and society
- Highest ethical standards, anonymised data used non-disclosively
- Data can be safely accommodated at the UK Data Archive at Essex, and safe access is possible in different formats, including through the safe rooms at UEA and Essex.
- Discussions with potential data providers to build trust and refine research questions of common interest
- Some nervousness about data sharing although route seems to be through well-specified research projects, subject to normal Research Governance arrangements
- Different LAs hold different data, in different ways
- Lack of a common approach to linking health and social care data



ESRC Business and Local Government  
Data Research Centre



University of Essex



Institute  
and Faculty  
of Actuaries

# Other sources of data



## Fitbit data is now being used in COURT: Wearable technology is set to revolutionise personal injury and accident claims

- Law firm in Calgary, Alberta is using a Fitbit to show effects of an accident
- Device collects steps taken, sleep patterns and other data about the victim
- Analytics firm Vivametrica has launched service to compare the victims' activity with the national average to determine the strength of a claim
- Data will be used alongside other evidence, and evaluated by a doctor

By SARAH GRIFFITHS FOR MAILONLINE

PUBLISHED: 16:56, 17 November 2014 | UPDATED: 16:57, 17 November 2014

Facebook Share | Twitter | Pinterest | Google+ | Email | RSS | 149 shares | View comments

In personal injury lawsuits, it can be notoriously difficult to establish whether a



# Google Flu example

“Big data hubris” is the often implicit assumption that big data are a substitute for, rather than a supplement to, traditional data collection and analysis.

Lazer et al, *Science*, 2014

## Why Google Flu Is A Failure



**Steven Salzberg**, CONTRIBUTOR

Fighting Pseudoscience **FULL BIO** ✓

Opinions expressed by Forbes Contributors are their own.

It seemed like such a good idea at the time.

People with the flu (the influenza virus, that is) will probably go online to find out how to treat it, or to search for other information about the flu. So **Google** **GOOG +0.58%** decided to track such behavior, hoping it might be able to predict flu outbreaks even faster than traditional health authorities such as the Centers for Disease Control (CDC).

Instead, as the authors of [a new article in Science](#) explain, we got “big data hubris.” David Lazer and colleagues explain that:

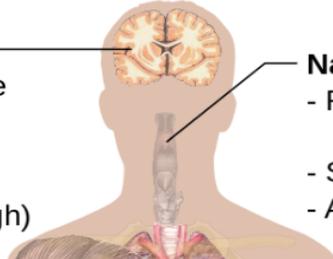
“Big data hubris” is the often implicit assumption that big data are a substitute for, rather than a supplement to, traditional data collection and analysis.

The folks at Google figured that, with all their massive data, they could outsmart anyone.

### Symptoms of Influenza

**Central**  
- Headache

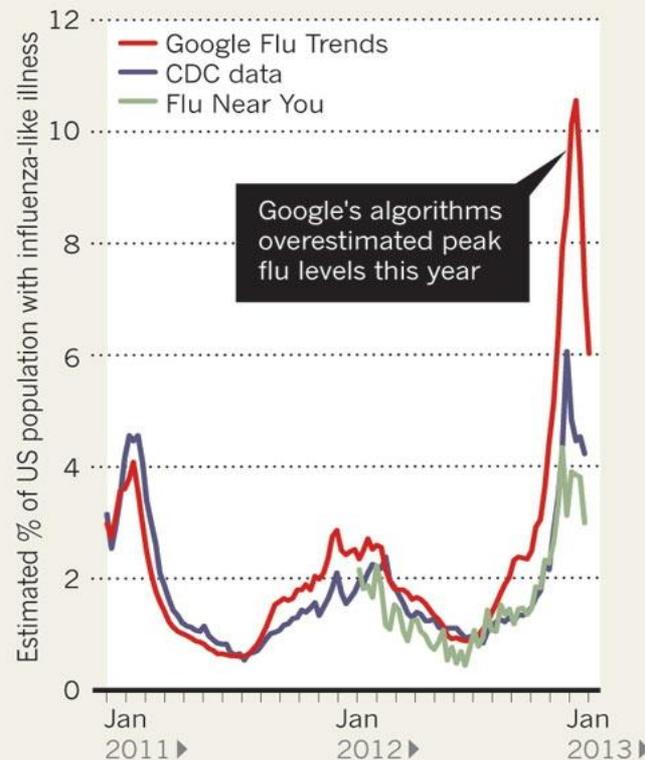
**Systemic**  
- Fever  
(usually high)



**Nasopharynx**  
- Runny or stuffy nose  
- Sore throat  
- Aches

### FEVER PEAKS

A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.



Institute and Faculty of Actuaries

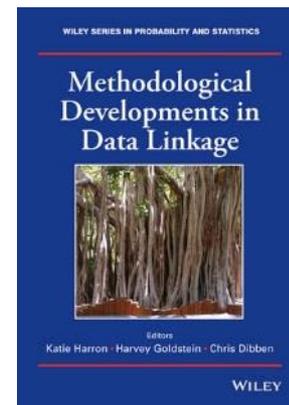
# Statistical Problems in Big Data

- Big Data is observational data!
- Methodology of observational data analysis/ meta-analysis:
  - False positives arising from multiple exploratory analyses
  - Biases due to peculiarities of units, outcomes or settings
  - Missing data
  - Inadequate linkage strategies
  - Data gathered at varied levels such as transaction, person, organization, community, and state
  - Causal Inference from (mostly) correlational data
  - Modelling heterogeneity



# Need for linkages between private and public data

- **Basis risk and Longevity**
- Longevity and morbidity risks are evaluated on the general population data.
- CMI data, link to primary care data, ONS and to Mosaic codes!
- **Does private medical insurance increase longevity? HLE?**
- Health Insurance data, link to primary care data, ONS, Mosaic
- **Is a private care home better than an LA home?**
- LA care home census, BUPA care home census, link to primary care data, ONS
- ***Need for a trusted intermediary able to provide such linkages***



# Questions

# Comments

Expressions of individual views by members of the Institute and Faculty of Actuaries and its staff are encouraged.

The views expressed in this presentation are those of the presenter.



Institute  
and Faculty  
of Actuaries