# Statistical Problems in Big Data

Prof Elena Kulinskaya (UEA)
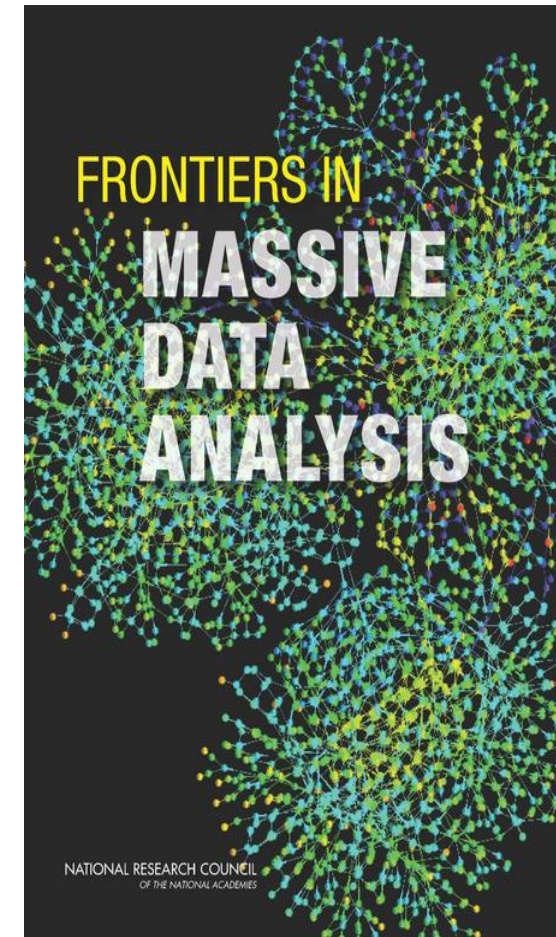
# Big Data:  Fields and Disciplines

Actuaries are the founders of the Big Data: the CMI!
Big Data also arises in such fields  as genomics, public health, environmental sciences, neuroscience, government and business.

Basic issues of management and storage have primarily implicated computer science, underpinning initiatives sometimes described as business analytics or data science.

Statistical science has not played a prominent role. Because practical problem solving has proceeded rapidly, the science has lagged behind and work to identify the statistical features associated with Big Data has been largely ad hoc.

The Alan Turing Institute:  Cambridge, Edinburgh, Oxford, Warwick and UCL.



FRONTIERS IN MASSIVE DATA ANALYSIS

NATIONAL RESEARCH COUNCIL
OF THE NATIONAL ACADEMIES

# 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

Evidence synthesis for data  at varied levels of aggregation
such as transaction, person, organization, community, and state

Causal Inference from (mostly) correlational data

Modelling heterogeneity

# Google Flu example

## Development of novel statistical and actuarial methods for:

modelling mortality

modelling trends in morbidity

assessing basis risk

evaluating longevity improvement based on Big Health and Actuarial Data

tools to forecast longevity risk of a book



Science

**Scientists and insurers develop 'death clock' to predict when customers will die**

A new computer algorithm will predict how long people will live  CREDIT: WALES NEWS SERVICE LTD.

# Data

The Health Improvement Network (THIN) data

> ➢ Medical records from primary care
> ➢ Representative of the UK when adjusted for deprivation

➢ All patients born before 1960 and followed to 01.01.2015, this includes 3.4 million patients

➢ Added various social economic status variables such as IMD and Mosaic

➢ The Continuing Mortality Investigation (CMI) data

# Design and methods

The most efficient way to analyse the data of variable quality  is to delete "bad" data.

For a particular condition we design a population-based prospective cohort study using an appropriate extract of the primary care data.

We intend to use a case-control design with cases matched with several  controls from the same GP practice.  This provides balanced and comparable cohorts of cases and controls and  simplifies the study of comparatively rare conditions without loss of efficiency.

To account for the interdependence of patients from the same GP practice, we use multilevel modelling and multiple imputation.