# Big data in health care
## – a traditionalist's words of caution

Gustaf Edgren, MD PhD
Associate professor of Epidemiology and Hematologist
Karolinska Institutet/Karolinska University Hospital

---

## Sweden: a nation of hamsters

- National registration number for all inhabitants since 1948

- Nationwide health registers starting in 1958 (cancer), 1960 (cause of death), 1964 (patient care), 1973 (medical birth), 2004 (drug prescriptions), 2003 (outpatient and primary care)

- Additional registers on: income, family structure, occupation, school grades, immigration/emigration, etc.
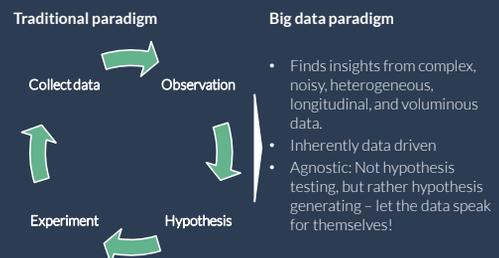
BIG DATA?

---

## Outline

- Background

- Introductory example(s) of big data in health

- An epidemiologist's words of caution:
  - Common problems
  - Common solutions

- More examples

---

## BACKGROUND

---

## So, what's big data?

- The traditional definition of 'Big data' is typically technocratic – 3V:
  - Volume
  - Velocity
  - Variation
- > 100GB?, >100TB?, or >100PB?
- Does it matter?
- A practical definition is perhaps that the 3 V's together poses a technological challenge that can only be resolved with non-traditional methods?

---

## The big data paradigm in health

**Traditional paradigm**

Collect data → Observation → Hypothesis → Experiment → (cycle)

**Big data paradigm**

- Finds insights from complex, noisy, heterogeneous, longitudinal, and voluminous data.
- Inherently data driven
- Agnostic: Not hypothesis testing, but rather hypothesis generating – let the data speak for themselves!

## Similarities to 'omics' revolution

- One can easily see a parallel to the genomics revolution

- Transition from small scale datasets with one-dimensional data – all hypothesis driven

- Now: all massive-scale data with multidimensional data – all using agnostic methods

- Key lesson: replication is key

## New paradigm? Perhaps, but…

- Health care 'big data' applications are typically based on routine administrative data (e.g. insurance claims data, discharge data, prescription data, etc)
- To allow any inference form such data, one must understand the data context:
  - How the data was collected
  - Why the data was collected
  - In what situations the data was collected
  - What incentive structures might have affected the data collection

## Example 1



*Lancet* 2012; 379: 244–49

**Risk of pulmonary embolism in patients with autoimmune disorders: a nationwide follow-up study from Sweden**

## Example 1: What?

- An analysis of whether autoimmune disease increases risk of pulmonary embolism

- Important implications: common diseases, feared complication, effective prophylactic treatment

- Nationwide cohort study including entire Swedish population 1964-2008 (~15 million individuals, 500,000 with autoimmune disease)

## Example 1: Results



## Example 1: Interpretation

- The risk of VTE is very high for individuals with a new diagnosis of autoimmune disease;
- the risk then decreases gradually, ostensibly with successful treatment?
- Perhaps patients with autoimmune disease should routinely be given thrombosis prophylaxis?
- Important and powerful illustration of the power of big(ish) data?

# No!
# Why?

## Example 2



ORIGINAL CONTRIBUTION

JAMA 2012; 308 (13)

### Risk of Venous Thromboembolism in Patients With Rheumatoid Arthritis and Association With Disease Duration and Hospitalization

Marie E. Holmqvist, MD, PhD
Martin Neovius, PhD
Jonas Eriksson, MSc
Ängla Mantel
Solveig Wållberg-Jonsson, MD, PhD
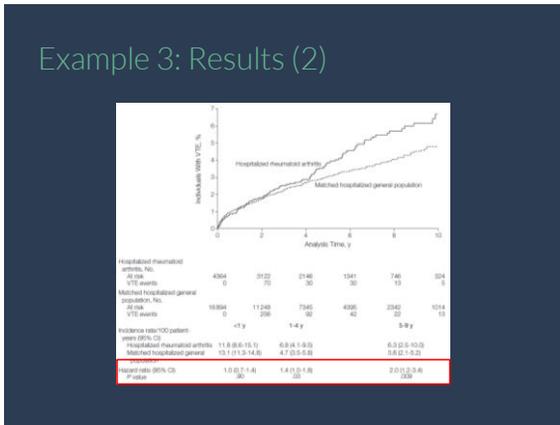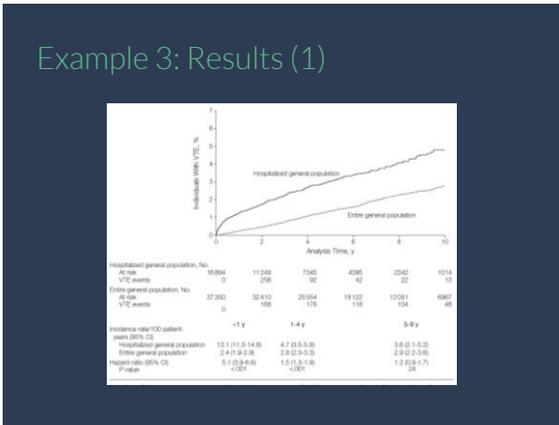Lennart T. H. Jacobsson, MD, PhD
Johan Askling, MD, PhD

**Context** Recent reports suggest that rheumatoid arthritis (RA) may be a risk factor for venous thromboembolism (VTE), particularly in conjunction with hospitalization. Using hospitalization data to identify RA and VTE may identify patients when they are at elevated risk for other reasons, obscuring the incompletely understood underlying association between RA and VTE and leading to inappropriate institution or timing of interventions.

**Objective** To estimate risks for VTE in patients with RA, including the relation of these risks to disease duration and hospitalization.

**Design, Setting, and Patients** Prospective, population-based cohort study of 1 prevalent RA cohort (n=37 856), 1 incident RA cohort (n=7904), and matched general population comparison cohorts, all from Sweden, with follow-up from 1997 through 2010.

**Main Outcome Measure** First-time VTE.

## Example 2: What?

- An analysis of whether rheumatoid arthritis (an autoimmune disease) increases risk of thromboembolism

- Important implications: common disease, feared complications, effective prophylactic treatment

- Large-scale cohort study of >40,000 patients with rheumatoid arthritis

## Example 3: Results (1)



## Example 3: Results (2)



## So, lessons learned?

- Even something seemingly trivial, such as assessing whether one patient group is at increased risk of a second disease is not all that trivial

- First, we must recall some key things that pertain to the analysis of any medical data...
  - Confounding (by indication)
  - Reverse causation
  - Surveillance bias
  - Misclassification of diagnosis
  - Will Roger's phenomenon
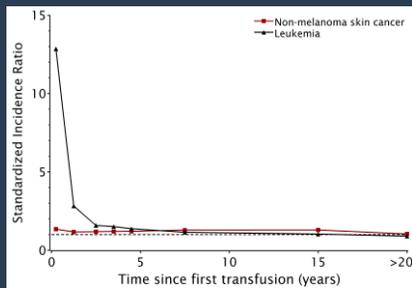
## 1. Indication is key (confounding)

- Confounding is a central element in epidemiology and medical research

- It results from the confusion of the effect of one variable on a particular outcome with a third variable:

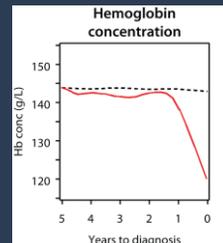Mass migration  ⟶  Cholera outbreak

Flooding

## 1. Indication is key (2)

- In any analysis of clinical data 'confounding by indication' is a an ever-present problem

- It results from the conscious choice of different treatments for patients with different prognosis

  – Sicker patients get more aggressive treatment
  – *Or, really sick patients can't get aggressive treatment because they are too sick*

## 2. Sequence isn't all that easy



## 2. Post hoc ergo propter hoc…



## 2. Reverse causation

- It will make it appear from your data that different events happen in one particular order

- However, due to incompleteness, that order will often be obfuscated, leading to incorrect characterization of sequence/causation
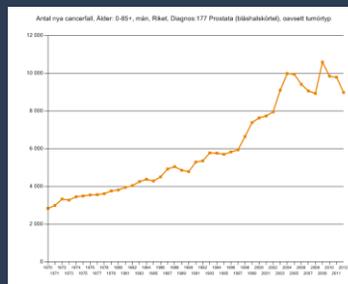
## 3. Surveillance bias

- As a rule of thumb: the more you look, the more you will find

- In any situation when patients undergo intense scrutiny (i.e. screening), disease rates automatically go up

- Implicitly, and disturbingly, complication rates therefore often go up with the quality of care
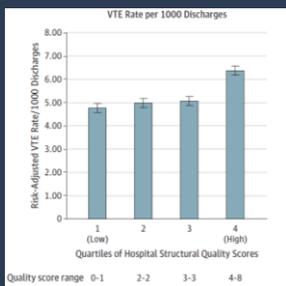
## 3. If you look, you shall find



Figure 1. Hospital Venous Thromboembolism (VTE) Prophylaxis Adherence Rates and Risk-Adjusted VTE Event Rates

## 3. If you look you shall find



## 3. Surveillance bias



## 4. Doctors are sloppy/lazy/greedy

- Most EHR data are based on diagnosis and procedure codes
- Such coding systems are generally inherently accurate and well structured
- However, the accuracy of the coding use is dependent on:
  1. Quality of care
  2. Misunderstanding of coding practices
  3. Financial incentives for coding use
  4. Diagnostic progress – i.e. misdiagnoses

## 5. Will Rogers' phenomenon

*When the Okies left Oklahoma and moved to California, they raised the average intelligence in both states*

## CASE STUDIES

## Application 1: ER frequent fliers

- In Sweden (as in most other parts of the world) a small part of the population account for a large fraction of healthcare spending

- It is possible to predict who will have large healthcare needs in the near future

- We hypothesized that interventions targeted at the small high-risk population may prevent unnecessary healthcare consumption

---

Hälso- och sjukvårdsförvaltningen
STOCKHOLMS LÄNS LANDSTING

**A telephone-based case-management intervention reduces healthcare utilization for frequent emergency department visitors**

Peter Reinius[a], Magnus Johansson[a], Ann Fjellner[b], Joachim Werr[c], Gunnar Öhlén[a] and Gustaf Edgren[d,e]

**Background** A small group of frequent visitors to emergency departments accounts for a disproportional large number of total emergency department visits. Previous interventions in this population have shown mixed results.

**Objective** To determine whether a nurse-managed telephone-based case-management intervention can reduce healthcare utilization and improve self-assessed health status in frequent emergency department users.

**Methods** We carried out a Zelen-design randomized-controlled trial among patients who were identified as frequent emergency department users (≥ 3 visits during

difference in mortality or other identified adverse outcomes between the intervention and the control groups. Patient self-assessed health status improved for the patients who received the case-management intervention.

**Conclusion** Our results indicate that the nurse-managed telephone-based case-management intervention represents a possible strategy to improve care for frequent emergency department users as well as decrease outpatient visits, admission days and healthcare costs. *European Journal of Emergency Medicine* 00:000–000 © 2012 Wolters Kluwer Health | Lippincott Williams & Wilkins.

---

Hälso- och sjukvårdsförvaltningen
STOCKHOLMS LÄNS LANDSTING

| | Intention to treat | |
|---|---|---|
| | Events/person years | Relative risk (95% CI) |
| All inpatient care | | |
| Intervention group | 444/238 | 0.90 (0.74–1.09) |
| Controls | 134/64 | 1.00 (reference) |
| Emergency inpatient care | | |
| Intervention group | 374/238 | 0.90 (0.73–1.12) |
| Controls | 112/64 | 1.00 (reference) |
| All outpatient care | | |
| Intervention group | 4794/238 | 0.80 (0.75–0.84) |
| Controls | 1629/64 | 1.00 (reference) |
| Emergency outpatient care | | |
| Intervention group | 1175/238 | 0.77 (0.69–0.86) |
| Controls | 413/64 | 1.00 (reference) |

CI, confidence interval.

33

---

## Application 2: New cancer syndromes?

- We all know there is a heritable component for many cancers
- For example, one estimate of the heritability of breast cancer is 28%
- This is an average – estimating how much the variability in risk is influenced by genetic factors
- However, rare highly penetrant genes should not have a strong effect on heritability measures

---

## Background (2)

- There are many well-characterized, monogenic cancer syndromes described in the literature, e.g.:
  - BRCA1/2: breast and ovarian cancer
  - Familial adenomatous polyposis (FAP) and lynch syndrome: colon polyps and/or cancer
  - Li-Fraumeni syndrome: sarcomas, breast, brain and adrenal gland cancers, as well as leukemia
  - MEN-1/2: endocrine cancers
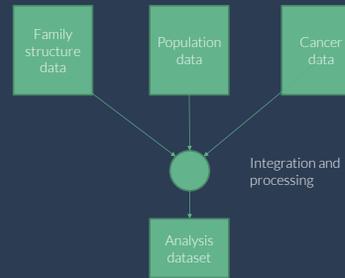  - Dicer1 syndrome: lung, kidney, ovaries, etc

---

## Background (3)

- Most of the known cancer syndromes are caused by highly penetrant single mutations
- Most have also been "detected" by careful individual researchers or clinicians with access to a particularly informative patient material
  - Selective in terms of cancer sites
  - Unlikely to pick up genes with low penetrance
- No (to me known) systematic attempts have been made to find new rare, monogenic cancer syndromes*
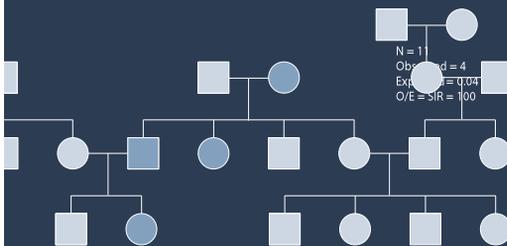
*GWAS:s don't count

## Approach (1)

- Very simple analytical approach:
  1. Use multi-generation register to find all "maximized" family trees (i.e. top people and their descendants)
  2. Use cancer register to identify all cancers in each family tree (i.e. observed)
  3. Use Poisson regression (or suitable alternative) to calculate expected number of cancers (given age, calendar year, sex, etc.)
  4. Reiterate step 2-3 for one cancer at a time, combinations of 2 cancers, 3 cancers, etc.

## Data process



## Example family



N = 11
Observed = 4
Expected = 0.04
O/E = SIR = 100

## Concluding thoughts

- Most important lesson: You have to get as close to the bottom of the data context as possible!

- Data scientists and data context experts have to work together
  – Who needs to be involved?

- Variation in data availability:
  – Behaviour, almost no data – 'hard outcomes,' a lot of data

  – How can we improve that?

gustaf.edgren@ki.se