

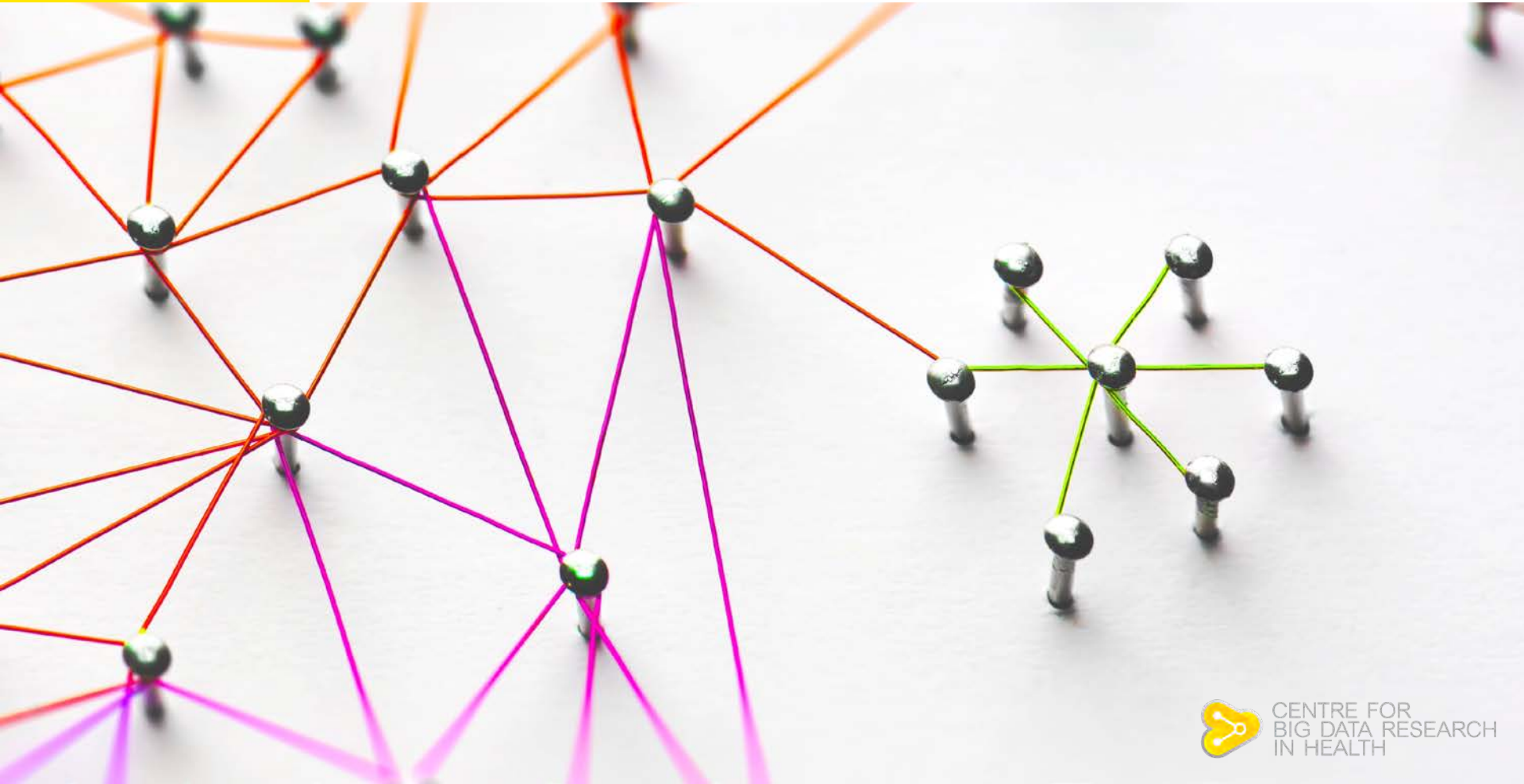
Big data in health and medicine: issues and challenges

Louisa Jorm

Health Research Symposium, Hong Kong, 16 June 2017

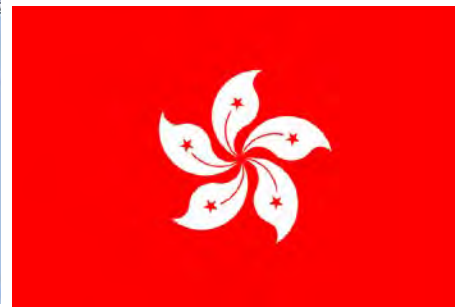


Australia's
Global
University

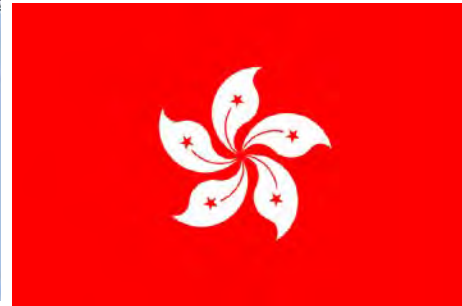


Outline

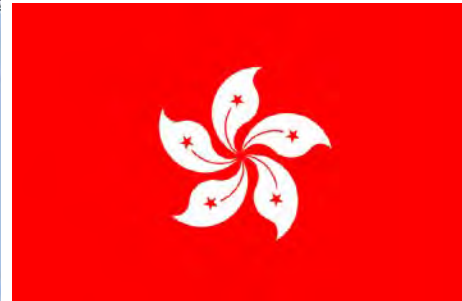
- What are big data?
- How are they being used in health and medicine?
- Issues and challenges:
 1. Upscaling technology
 2. New analytic paradigms
 3. Workforce shortages
 4. “Open data” vs. privacy protection



	Hong Kong SAR	Australia
Population ('000)	7,242	23,470
Total fertility rate	1.2	1.8
% population aged 65+	19.2	15.0
Life expectancy at birth M/F	81.2/86.9	80.4/84.5
Total health expenditure,% of GDP	5.2	9.0
Human development index	0.90	0.93
Per capita GDP (US\$)	38,074	44,820



	Hong Kong SAR	Australia
Population ('000)	7,242	23,470
Total fertility rate	1.2	1.8
% population aged 65+	19.2	15.0
Life expectancy at birth M/F	81.2/86.9	80.4/84.5
Total health expenditure,% of GDP	5.2	9.0
Human development index	0.90	0.93
Per capita GDP (US\$)	38,074	44,820



	Hong Kong SAR	Australia
Public hospitals	42	736
Hospital beds	25,000	60,300
Doctors per 10,000 population	18.3	35
Hospital beds per 10,000 population	34.5	25.7
Health ranking (Legatum prosperity index*)	7	8

*Combines indices of physical and mental health, health infrastructure, preventive care

Outline

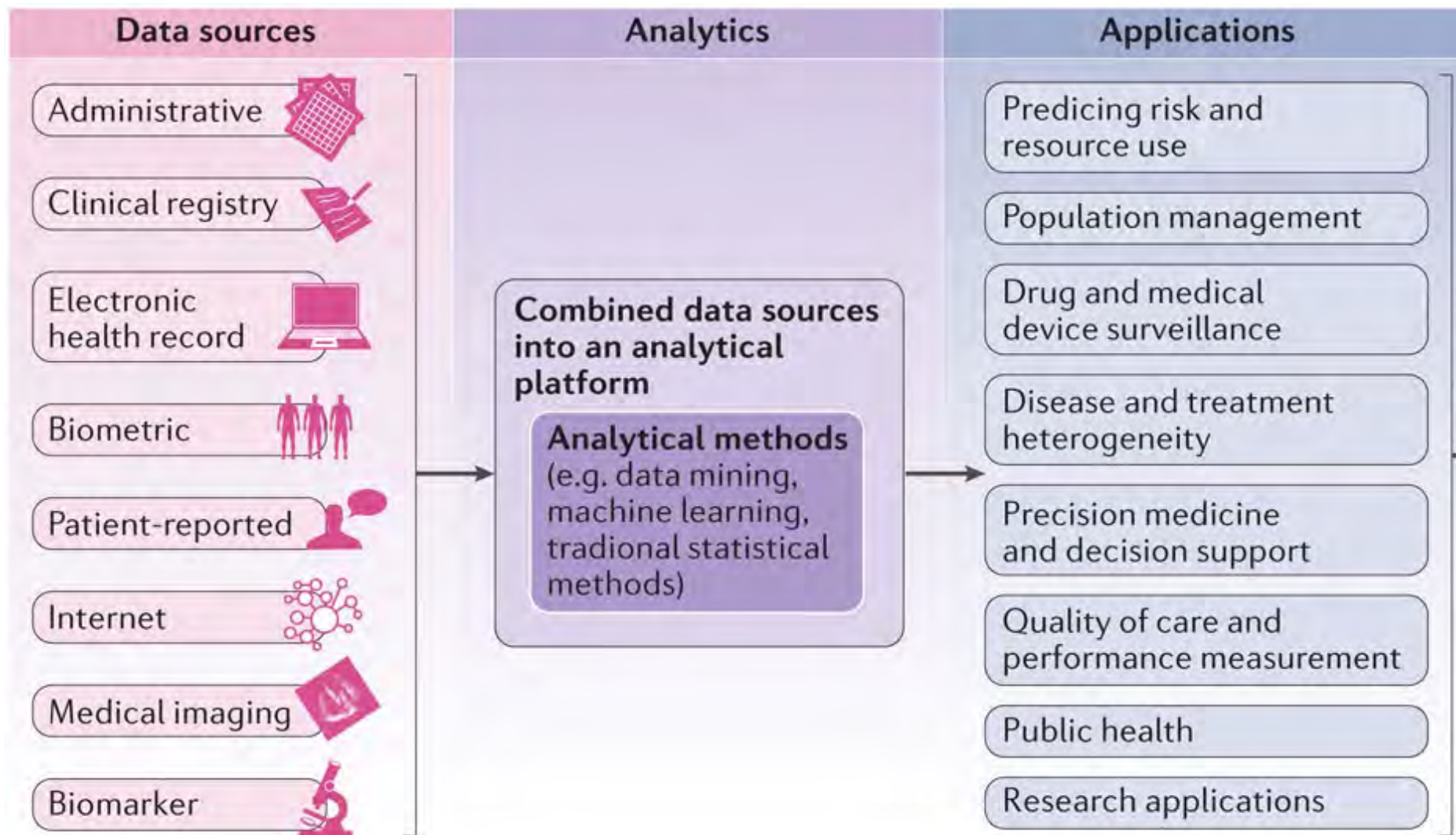
- What are big data?
- How are they being used in health and medicine?
- Issues and challenges:
 1. Upscaling technology
 2. New analytic paradigms
 3. Workforce shortages
 4. “Open data” vs. privacy protection

What are big data?

- Volume: Large scale of data (terabytes or petabytes)
- Variety: Variable format of data (structured, semi structured and unstructured)
- Velocity: Speed at which data are produced, processed, and analysed
- Value: Worth of information to stakeholders and decision makers

Outline

- What are big data?
- **How are they being used in health and medicine?**
- Issues and challenges:
 1. Upscaling technology
 2. New analytic paradigms
 3. Workforce shortages
 4. “Open data” vs. privacy protection



Real-Time Risk Prediction on the Wards: A Feasibility Study

Critical Care Medicine. 44(8):1468–1473, AUG 2016

Michael A. Kang; Matthew M. Churpek; [and 4 more](#)

Objective:

Failure to detect clinical deterioration in the hospital is common and associated with poor patient outcomes and increased healthcare costs. Our objective was to evaluate the feasibility and accuracy of real-time risk stratification using the electronic Cardiac Arrest Risk Triage score, an electronic health record-based early warning score.

Design:

We conducted a prospective black-box validation study. Data were transmitted via HL7 feed in real time to an integration engine and database server wherein the scores were calculated and stored without visualization for clinical providers. The high-risk threshold was set a priori. Timing and sensitivity of electronic Cardiac Arrest Risk Triage score activation were compared with standard-of-care Rapid Response Team activation for patients who experienced a ward cardiac arrest or ICU transfer.

Setting:

Three general care wards at an academic medical center.

Patients:

A total of 3,889 adult inpatients.

Conclusions:

Electronic Cardiac Arrest Risk Triage score identified significantly more cardiac arrests and ICU transfers than standard Rapid Response Team activation and did so many hours in advance.

ORIGINAL ARTICLE

Single Reading with Computer-Aided Detection for Screening Mammography

Fiona J. Gilbert, F.R.C.R., Susan M. Astley, Ph.D., Maureen G.C. Gillan, Ph.D., Olorunsola F. Agbaje, Ph.D., Matthew G. Wallis, F.R.C.R., Jonathan James, F.R.C.R., Caroline R.M. Boggis, F.R.C.R., and Stephen W. Duffy, M.Sc., for the CADET II Group[†]

BACKGROUND

The sensitivity of screening mammography for the detection of small breast cancers is higher when the mammogram is read by two readers rather than by a single reader. We conducted a trial to determine whether the performance of a single reader using a computer-aided detection system would match the performance achieved by two readers.

RESULTS

The proportion of cancers detected was 199 of 227 (87.7%) for double reading and 198 of 227 (87.2%) for single reading with computer-aided detection ($P=0.89$). The overall recall rates were 3.4% for double reading and 3.9% for single reading with computer-aided detection; the difference between the rates was small but significant ($P<0.001$). The estimated sensitivity, specificity, and positive predictive value for single reading with computer-aided detection were 87.2%, 96.9%, and 18.0%, respectively. The corresponding values for double reading were 87.7%, 97.4%, and 21.1%. There were no significant differences between the pathological attributes of tumors detected by single reading with computer-aided detection alone and those of tumors detected by double reading alone.

CONCLUSIONS

Single reading with computer-aided detection could be an alternative to double reading and could improve the rate of detection of cancer from screening mammograms read by a single reader. (ClinicalTrials.gov number, NCT00450359.)

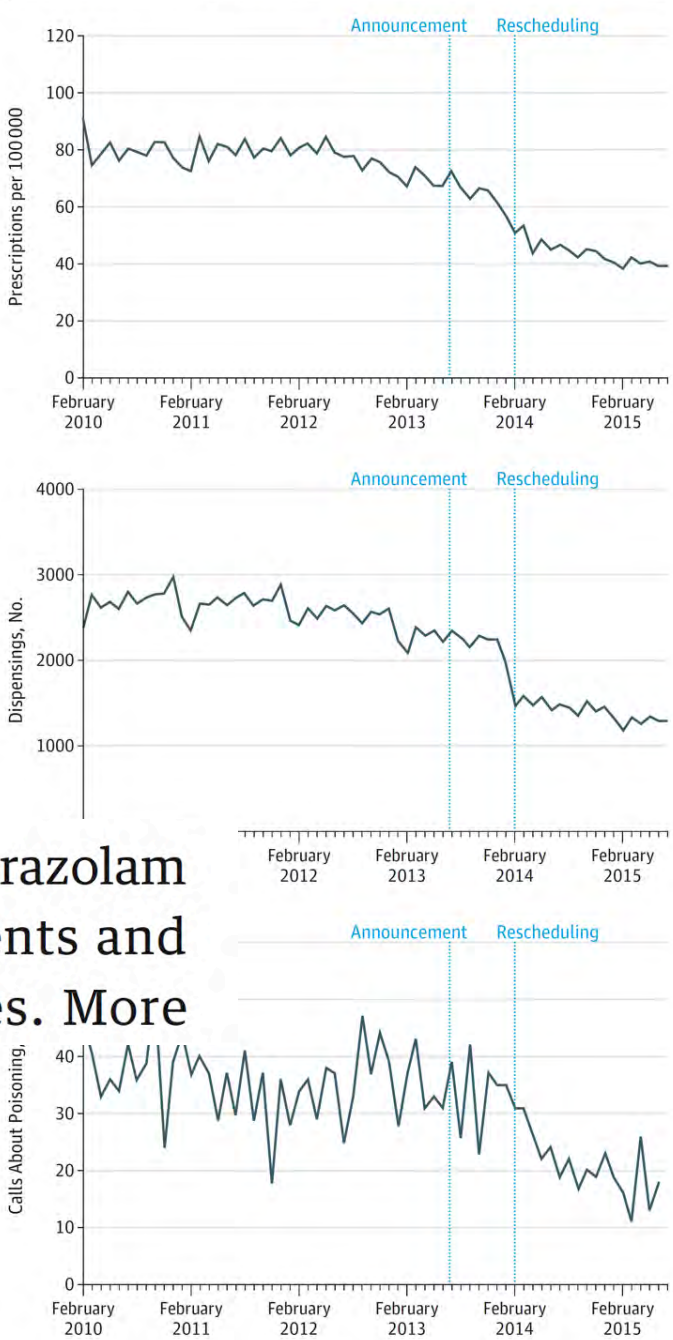
Interrupted Time Series Analysis of the Effect of Rescheduling Alprazolam in Australia: Taking Control of Prescription Drug Use

Alprazolam is significantly more toxic, has no additional therapeutic benefit, and is increasingly misused compared with other benzodiazepines.¹⁻³ Due to concerns about increasing use of alprazolam, in February 2014, the Australian Therapeutic Goods Administration selectively reclassified alprazolam from Schedule 4 (Prescription Only Medicine) of the Poisons Standard to Schedule 8 (Controlled Drug), equivalent to Schedule II in the United States.

Discussion | In Australia, selectively rescheduling alprazolam led to a reduction in overall use and adverse events and increased switching to less toxic benzodiazepines. **More**

Schaffer AL, Buckley NL, Cairns R, Pearson S. *JAMA Intern Med.* Published online July 05, 2016. doi:10.1001/jamainternmed.2016.2992

Figure. Monthly Time Series of Alprazolam Prescriptions per 100 000 Population, Dispensings, and Calls to Poisons Information Centre



Linkable de-identified 10% sample of Medicare Benefits Schedule (MBS) and Pharmaceutical Benefits Schedule (PBS)

Followers

5

Organisation



Department of Health

Department of Health [las mer](#)

Dataset

Grupper

Activity Stream

Use Cases

ISO19115/ISO19139 XML

RDF

JSON

Linkable de-identified 10% sample of Medicare Benefits Schedule (MBS) and Pharmaceutical Benefits Schedule (PBS)

This data is a collection of the current and historical use of Medicare and PBS services. This data release contains approximately 1 billion lines of data relating to approximately 3 million Australians. The data sets have been designed to enable other datasets to be linked in the future, for example hospital data, immunisation data. The addition of these data sets will greatly increase the amount of data and open new areas of analysis.

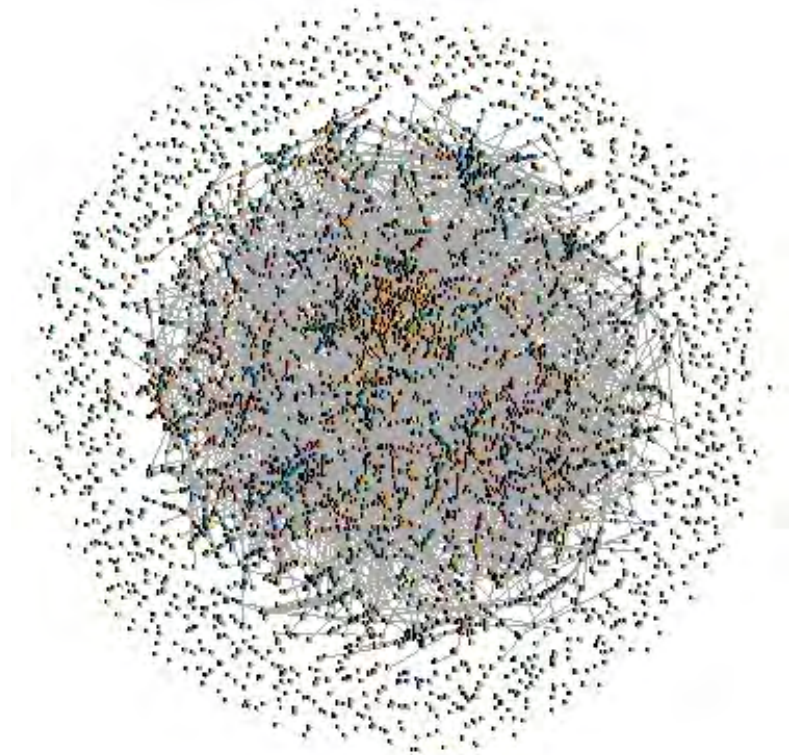
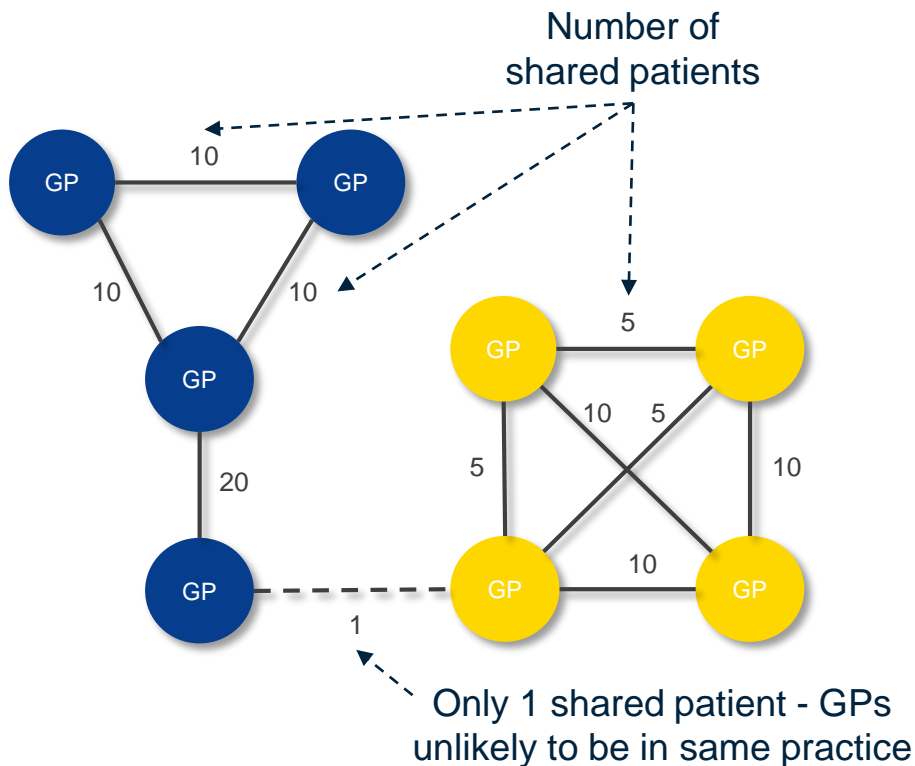
A suite of confidentiality measures including encryption, perturbation and exclusion of rare events has been applied to safeguard personal health information and ensure that patients and providers cannot be re-identified.

Confidentialisation Methodology

All Medicare and PBS claims for a random 10% sample of patients are included in the release. To be clear, it is a 10% sample of patients, not a 10% sample of Medicare or PBS claiming activity for the selected patients. Although the data held by the Department does not contain identifiers such as individual patient names, a number of steps have been taken to further protect the confidentiality of the released data.

Creating networks of health practitioners

- Connect GPs based on patterns of shared patients
 - GPs who share patients are more likely to work in the same practice
- e.g. National Medicare data for 10% of all Australians
 - 2,923 GP communities created from 25,338 GPs



We can explore the impact of different practice patterns on patient care



7 alerts for all diseases, current location, in the past week

Outbreaks in Current Location

Zika outbreak



4 Respiratory Alerts

Avian Influenza H5N1 (1), Avian Influenza H5N6 (1), Avian Influenza H7N9 (2)



1 Animal Alerts

Avian Influenza (1)



1 Skin/Rash Alerts

Hand, Foot and Mouth Disease (1)

Potentially Preventable Hospitalisations

Potentially prevented by timely and effective provision of primary and preventative care

Vaccine preventable

- Influenza and pneumonia
- Other vaccine-preventable conditions

Acute

- Dehydration & gastroenteritis
- Pyelonephritis
- Perforated/bleeding ulcer
- Cellulitis
- Pelvic inflammatory disease
- Ear, nose & throat infections
- Dental conditions
- Appendicitis with generalised peritonitis
- Convulsions & epilepsy
- Gangrene

Chronic

- Asthma
- Congestive heart failure
- Diabetes complications
- Chronic obstructive pulmonary disease
- Angina
- Iron deficiency anaemia
- Hypertension
- Nutritional deficiencies
- Rheumatic heart disease

Impact Of Socioeconomic Status On Hospital Use In New York City

by John Billings, Lisa Zeitel, Joanne Lukomnik, Timothy S. Carey,
Arthur E. Blank, and Laurie Newman

Abstract: This DataWatch examines the potential impact of socioeconomic differences on rates of hospitalization, based on patterns of hospital use in New York City in 1988. The research suggests that lack of timely and effective outpatient care may lead to higher hospitalization rates in low-income areas. For certain conditions identified as ambulatory care sensitive, hospitalization rates were higher in low-income areas than they were in higher-income areas where appropriate outpatient care was more readily available. Further study is needed to determine the relative impact of various economic, structural, and cultural factors that affect access to care.

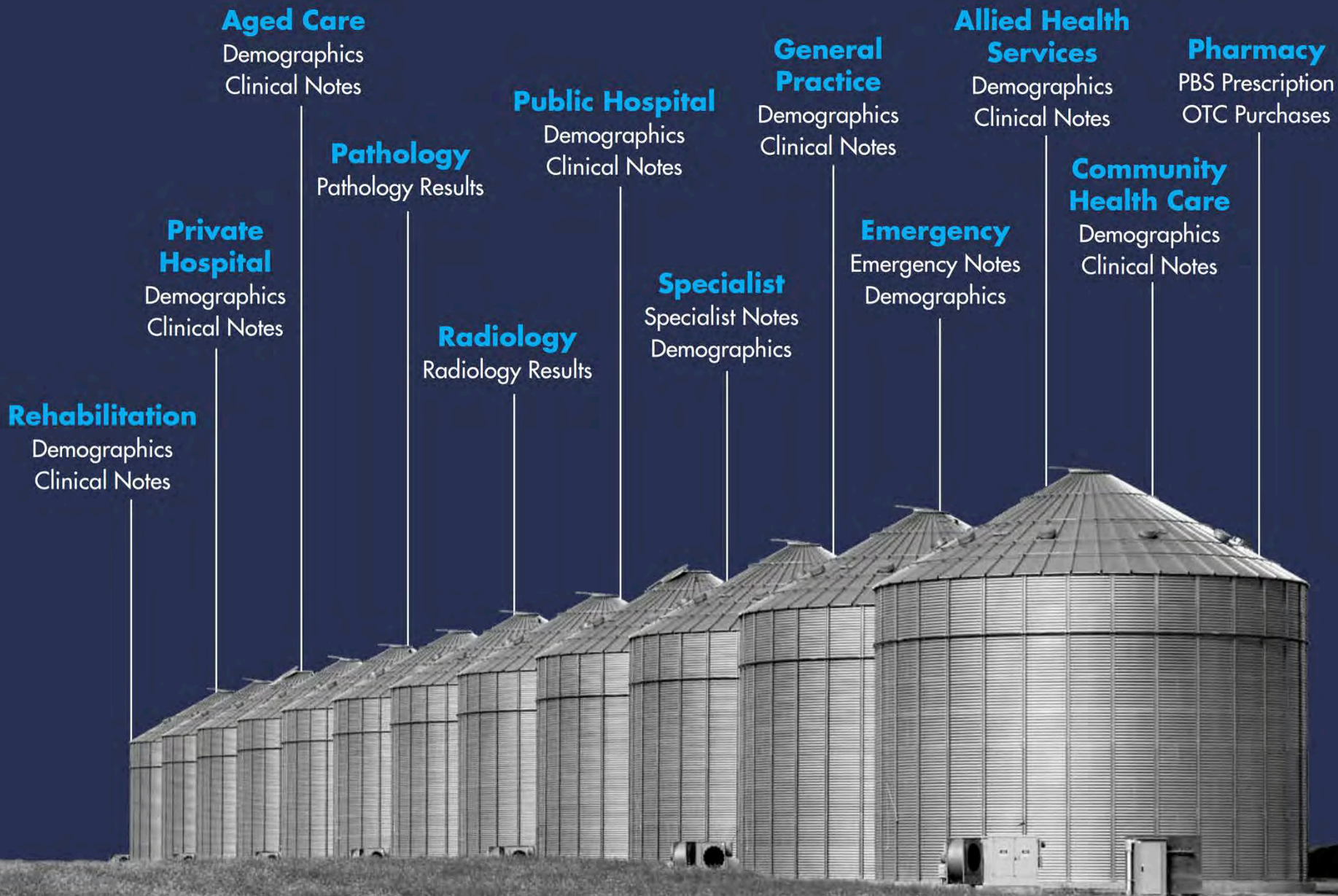
Cite this article as:

J Billings, L Zeitel, J Lukomnik, T S Carey, A E
Blank and L Newman

Impact of socioeconomic status on hospital use in
New York City

Health Affairs, 12, no.1 (1993):162-173

doi: 10.1377/hlthaff.12.1.162



The APHID Study

45 & Up Study

- Prospective cohort of 267,091 people aged over 45 in NSW.
- Study entry 2006-2008
- Questionnaire
 - Demographics
 - Health status
 - Risk factors

NSW Admitted Patient Data Collection

- Census of all hospital separations in NSW public and private hospitals and day procedure centres.
- Linked data, 2000-2013
- N=1,761,178 records

MBS

- Claims for subsidised medical and diagnostic services in Australia
- Linked data, 2004-2011
- N=45,754,339 records

PBS

- Claims for subsidised pharmaceuticals in Australia
- Linked data, 2004-2011
- N= 34,978,006 records

NSW Emergency Department Data Collection

- Presentations to 80 EDs (75% of NSW presentations)
- Linked data, 2006-2013
- N= 586,131 records

+ Fact of death to 2013

Data linkage

45 & Up
Study

NSW
Admitted
Patient Data
Collection

MBS

PBS

Emergency
Department
Data
Collection

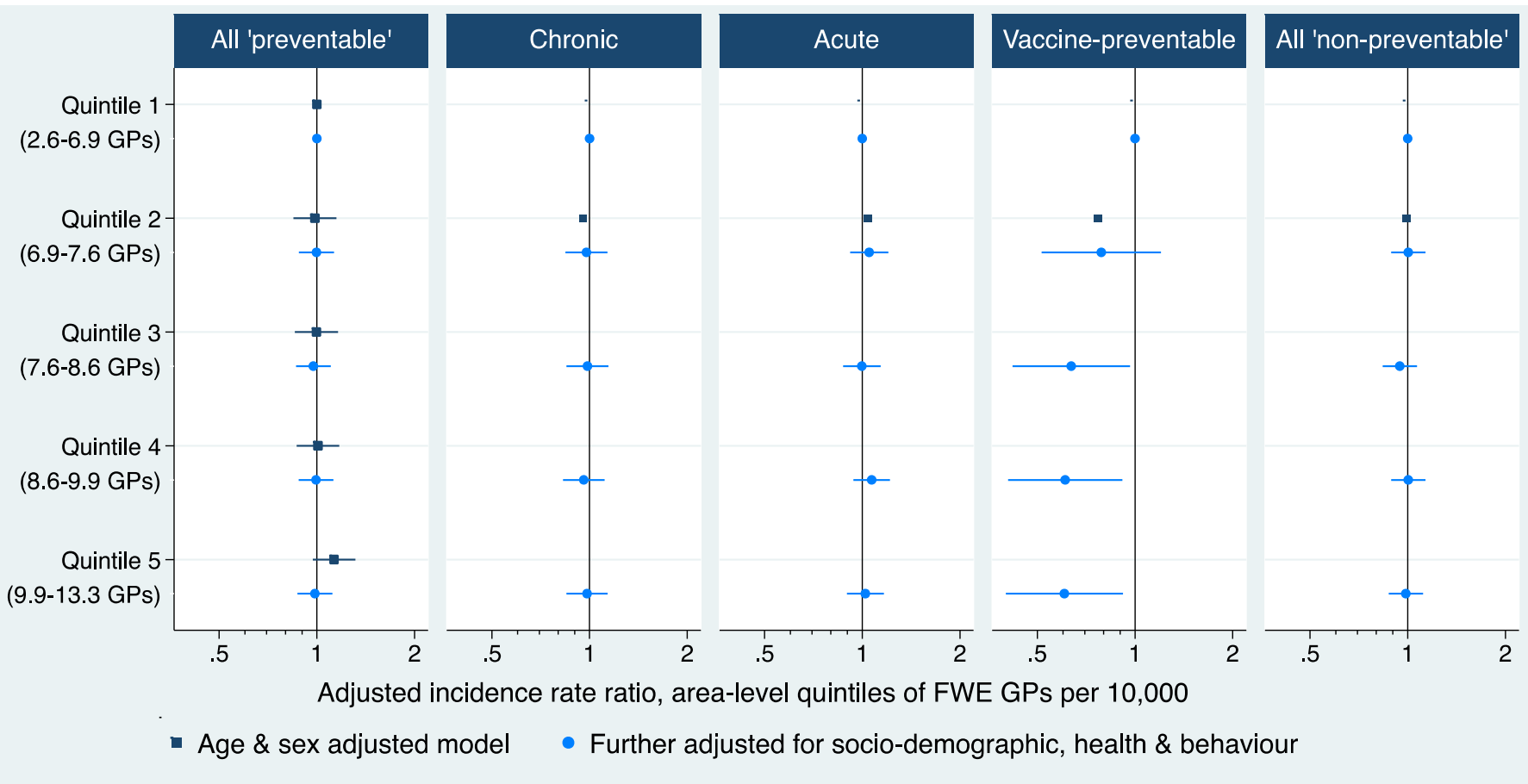
Detailed data for 267,153 people...

... WHO they are
... HOW they have interacted with
the primary health system
... WHETHER they were admitted to
hospital

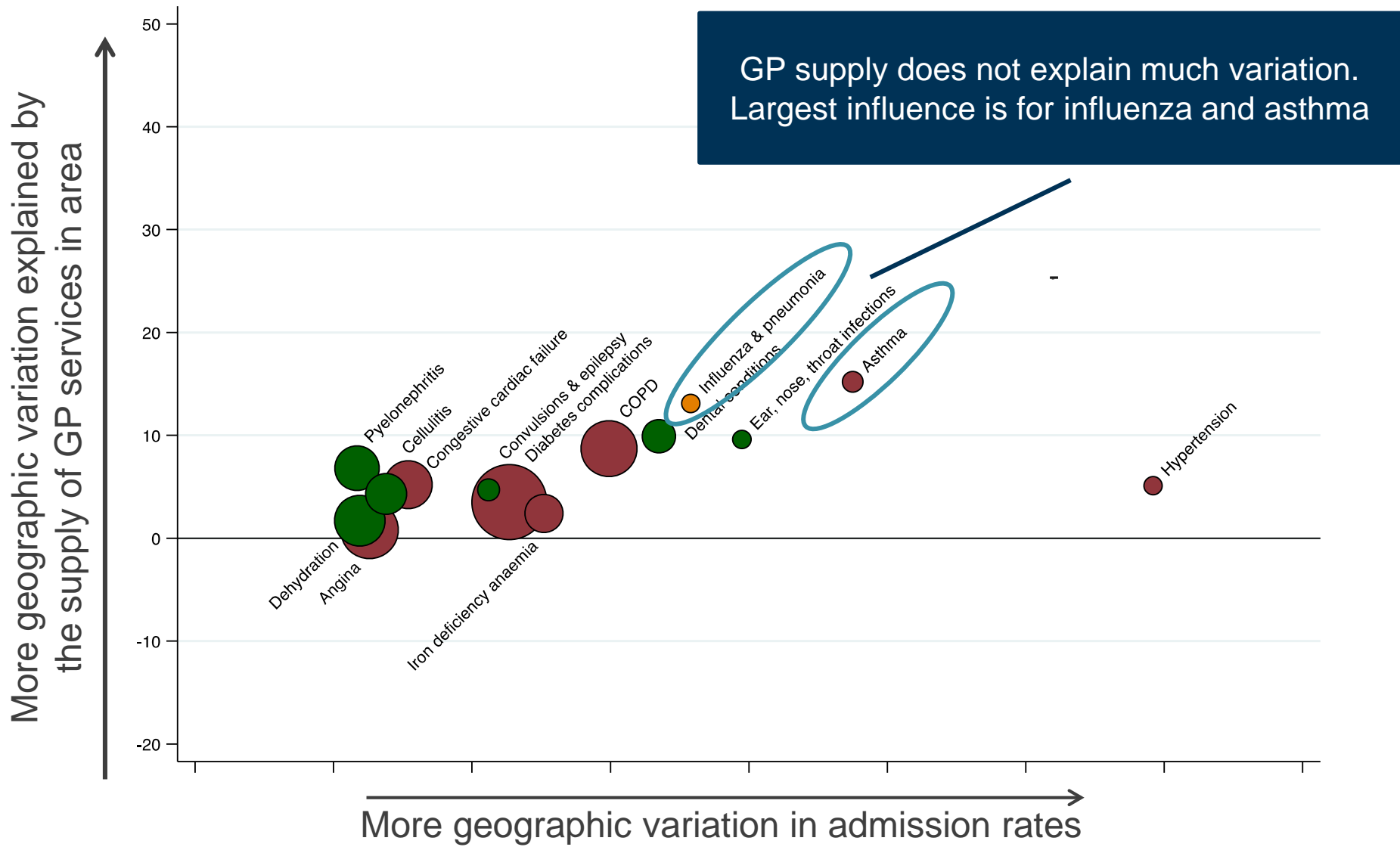
What factors explain geographic variation in admission?

Primary care supply	Socio-demographic factors	Factors amenable to behaviour change & disease management
<ul style="list-style-type: none">▪ Number of full time workload equivalent (FWE) GPs, per 10,000 residents in area	<ul style="list-style-type: none">▪ Education▪ Language spoken at home▪ Marital status▪ Aboriginal status▪ Income▪ Employment▪ Health insurance status▪ Number of people can depend on	<ul style="list-style-type: none">▪ Healthy behaviours▪ Body Mass Index▪ Self-reported health▪ Number of co-morbidities▪ Functional status▪ Psychological distress

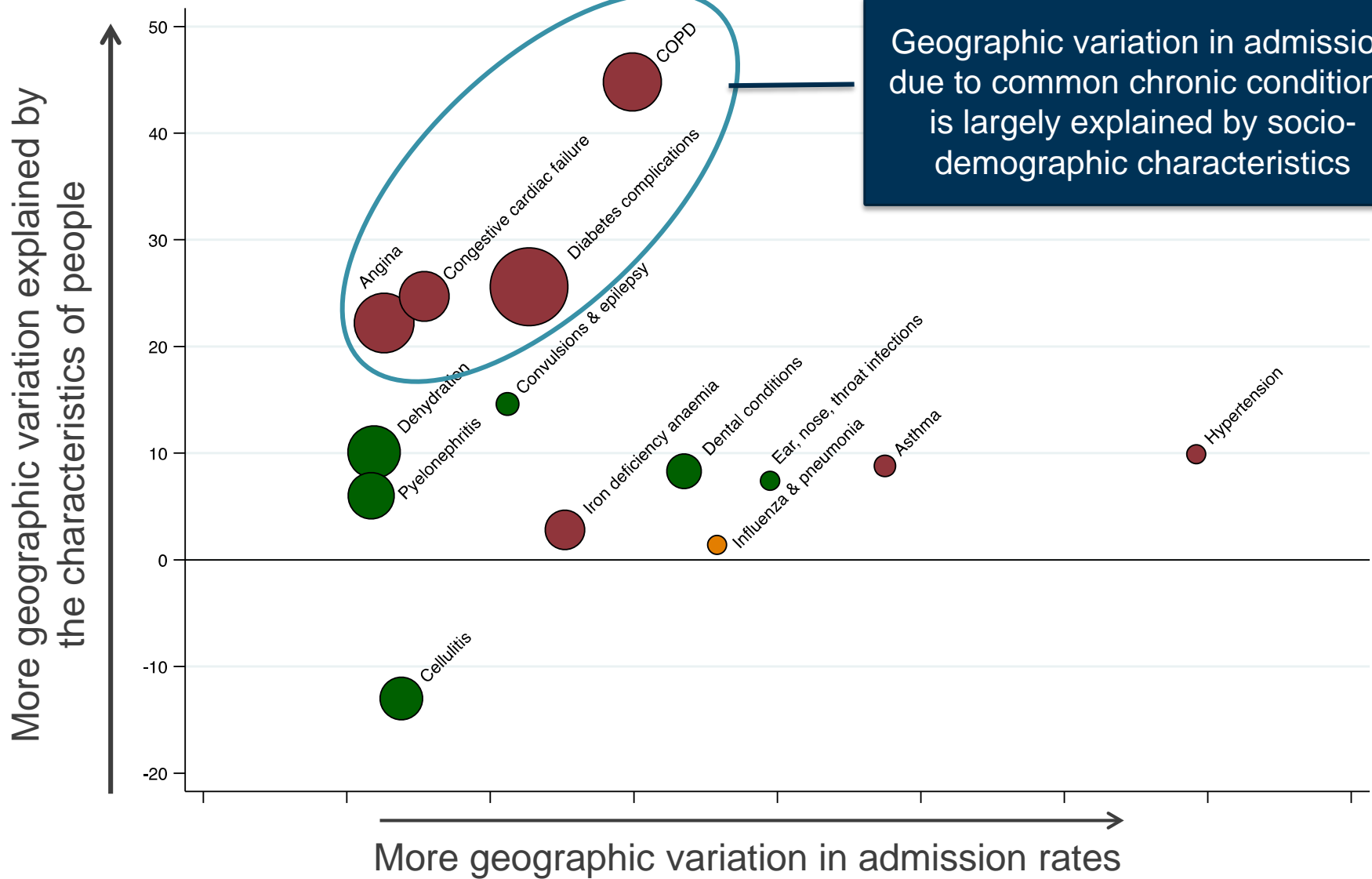
Area-level supply of full time workload equivalent GPs and PPH admissions



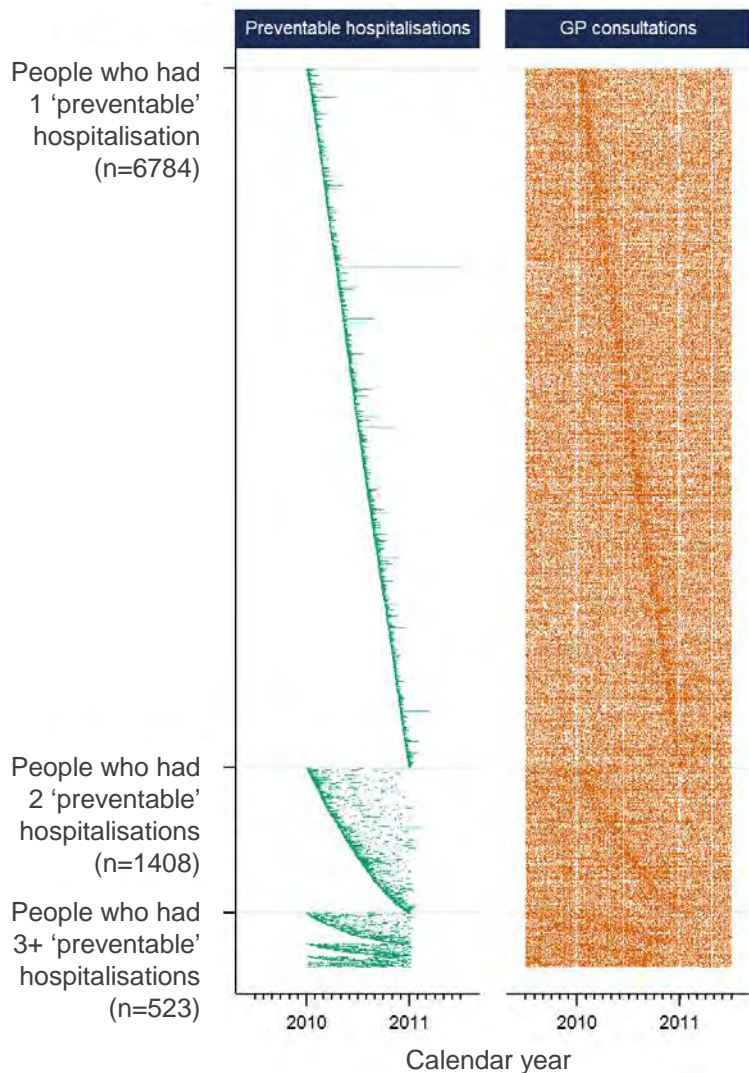
For which conditions does GP supply explain most variation?



For which conditions do personal factors explain most variation?



Visualising unit record data on preventable hospitalisations, GP consultations... and other health events



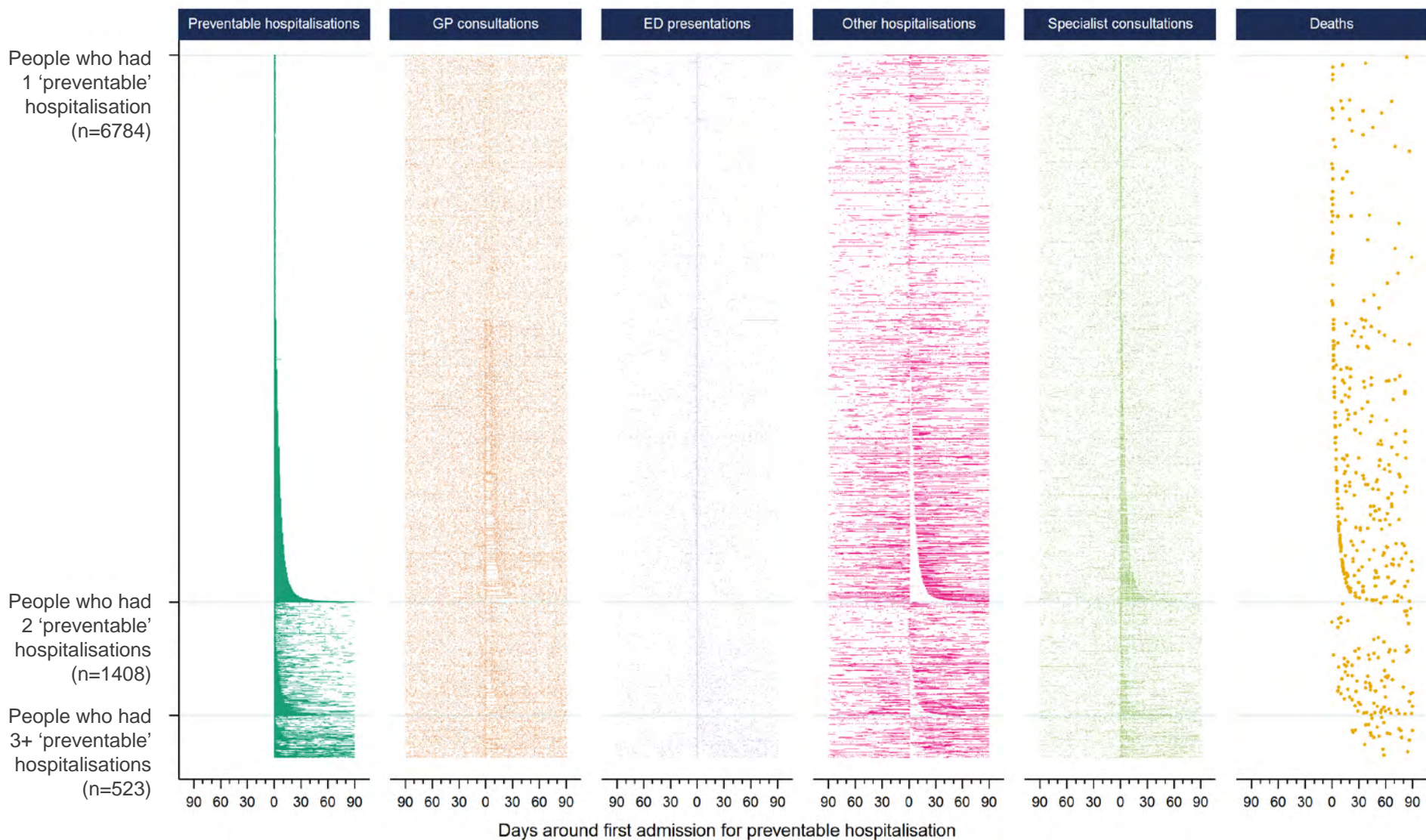
Each row is a person

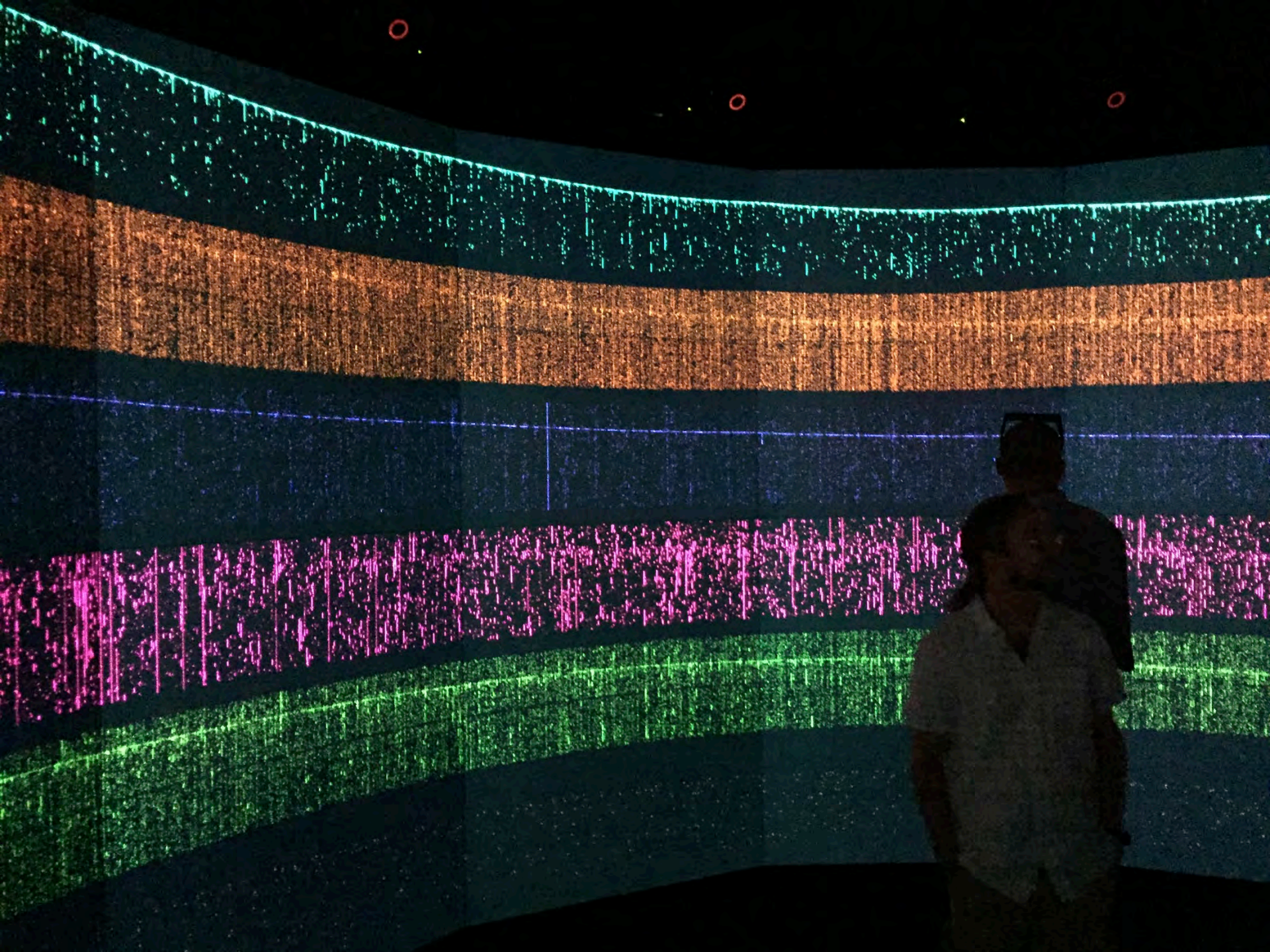
Each line/dot is a health event over time

Patterns become clear in the position and density of health events in the plot

e.g. Patients visiting GPs at the time of hospitalisation, not at Christmas

Zooming in on the time around preventable hospitalisation





Outline

- What are big data?
- How are they being used in health and medicine?
- **Issues and challenges:**
 1. **Upscaling technology**
 2. New analytic paradigms
 3. Workforce shortages
 4. “Open data” vs. privacy protection

Big Data Landscape 2016 (Version 2.0)

Infrastructure

Hadoop On-Premise
 cloudera, Hortonworks, Pivotal, IBM InfoSphere, Aludata, jethro

Hadoop in the Cloud
 Amazon, Microsoft Azure, Google Cloud Platform, IBM InfoSphere, A. V. Brink, Treasoft, Altiscale, Databricks

Spark
 databricks, GridGain, TACHYON

Cluster Services
 Amazon, Docker, Mesosphere, Core OS, Splunk, StackIQ

Analytics

Analyst Platforms
 Palantir, AYASDI, Quid, Digital Reasoning, ORBITALINSIGHT

Analytics Platforms
 Microsoft, GUVVUS, Datarameer, Bottlenose, interana

Data Science Platforms
 Continuum, DataRobot, Alpine, MODE, plury, ADATA, Qubol, DataSift, DOMINO, Sense, Aladdin

Visualization
 Tableau, Qlik, Looker, PowerBI, Rudderstack, CHARTIQ

Applications

Sales & Marketing
 RADIUS, Gainsight, bloomreach, Zeta, EVERSTRING, livefyre, blueyonder, Lattice, Kaltura, Liner, SALTIRU, persado, AVISO, sense, QUANTIFIND, ACTIONIQ

Customer Service
 MEDALLIA, ATTENTIFY, CLARABRIDGE, CLICKFOX, STELLASERVICE, NGDATA, Friend, DigitalGenius, epur, WISEO

Human Capital
 gild, Connective, Textile, enelo, hiQ

Legal
 RAVEL, JUDICIAL, Everlaw, Brevia, FRENCHMOUNTAIN

NoSQL Databases
 Amazon DynamoDB, Google, Microsoft Azure, Oracle, MarkLogic, MongoDB, DATASTAX, KEROFSPIKE, Couchbase, SequoiaDB, redislabs, InfluxDB

NewSQL Databases
 SAP, Clustrix, Pivotal, paradigm4, nuora, memsql, YOLDB, splice, MariaDB, VOLTDB, citusdata, deep db, Truvision, Cockroach Labs

BI Platforms
 Power BI, Amazon, D3.js, Qlik, Tableau, GoodData, Birst, platforma, Alacritas

Statistical Computing
 SAS, SPSS, MATLAB

Log Analytics
 Splunk, Sumologic, kibana, ELBUD, WINSYS, loggly

Social Analytics
 Hootsuite, Netbase, DataSift, track, bitty, syntaxis, reach

Ad Optimization
 AppNexus, Criteo, MediaMath, OpenX, RocketHub, Integral, theTradeDesk, Livestorm, dsillery, DataXu, Oppier, TAYO

Security
 Cylance, CounterTrack, ThreatMetrix, AREA1 SECURITY, SiftOne, Recorded Future, Guardian Analytics, FortScale, Sift Science, Feedzai, SCNIFY

Vertical AI Applications
 Facebook, Clara, KASIST, Lumina

Graph Databases
 Neo4j, OrientDB, InfluxDB

MPP Databases
 Teradata, Vertica, Nivezza, Conon, Kognitio, SQL, Greenplum

Cloud EDW
 Amazon Redshift, Microsoft Azure, Pivotal, Snowflake, InfoWorks

Data Transformation
 Alteryx, Talend, Trifacta, Tamr, StreamSets, Alation

Data Integration
 Informatica, Mulesoft, SnapLogic, BedrockData, Splinty

Real-Time
 Amazon, HEIMARKETS, Streamium, Confluent, DataArtisans

Machine Learning
 H2O, DataRobot, Skytree, Rapidminer, PANGEA, DataCamp, YIELDIQ, IBM Watson

Speech & NLP
 Narrative Science, ARRIA, Nuance, Semantic Machines, Capital, Lintell, IBM Watson

Horizontal AI
 IBM Watson, Corana, Viv, Nomo, Memento, InterScience, Darifai, Numenta

Publisher Tools
 Outbrain, Taboola, Quantcast, Chartbeat, Yieldbot, Yieldmo

Govt / Regulation
 Socrata, OpenGov, EN, FiscalNote, Fingma, Mark43, OpenDataSoft

Finance
 Affirm, LendingClub, OnDeck, Kreditech, Kabbage, Bidmark, INSIKT, ZUORO, Dataminr, Lenddo, KENSHO, AIDYIA, ISENTIUM, Quantopian, Sentient

Management / Monitoring
 New Relic, Amazon, Dynatrace, Actifio, Numerify, Splunk, Datto, Streamio, Aconet

Security
 Tanium, Lumio, DDE 42, DataGravity, CyberArk, Vectra, BlueTalon

Storage
 Amazon, Microsoft Azure, Panasas, Cimble, COHO, Quimulo

App Dev
 Apigee, CPSK, Typesafe, Driven

Crowd-sourcing
 Amazon Mechanical Turk, CrowdFlower, Workfusion

Search
 HP, Oracle, Elastic, Ludicworks, MAANA, Swiftype, Algolia, Swiftype

Data Services
 IQ, OPERA, MeSema, Kaggle, Datastax, Datalend

For Business Analysts
 Origami, ClearStory, Cirro, Import IO

Web / Mobile / Commerce
 Google Analytics, Mixpanel, RJMetric, Bluecore, Amplitude, Granify, Sumall, Airtable, Retention, Custora

Education / Learning
 Knewton, Clever, Ceclara, Panorama, Knowit

Life Sciences
 Genentech, Genzyme, Recombinate, Xyrus, Flatiron, Zymogen, HealthTap, Metabota, Zephyr Health, OVI, Gingerio, Transcriptic, Glow, Enlinc, AICure, Axiom

Industries
 OPOWER, eHarmony, RetailNext, Stich Fix, WorkFusion, BlueRiver, Tachyus, Seeq, FarmLogs, Swifkey, HowGood, Select, StatMuse, Boxever

Cross-Infrastructure/Analytics

Amazon, Google, Microsoft, IBM, SAP, SAS, Hadoop, HP, Oracle, Vertica, VMware, TIBCO, Teradata, Oracle, NetApp

Open Source

Framework
 Hadoop, YARN, Spark, Mesos, TEZ, Flink, CDAP

Query / Data Flow
 SLAMDATA, DRILL, CouchDB, riak

Data Access
 Cassandra, HBase, MongoDB, Kafka, Scio.io, CouchDB, riak

Coordination
 Apache Zookeeper, Apache Ambari

Real-Time
 Storm, Spark, Flink, Tachyon, Druid

Stat Tools
 R, Scala, SciPy

Machine Learning
 MLlib, Aerosolve, Caffe, CNTK, FeatureFu, Dimsun, WEKA, Vektor, DL4J

Search
 ElasticSearch, Solr

Security
 Apache Ranger

Visualization
 Tableau, Qlik, Looker

Data Sources & APIs

Health
 Apple, Jawbone, Garmin, Withings, Fitbit, Validic, Human API

IOT
 Uptake, ThingWorx, Samsara

Financial & Economic Data
 Bloomberg, Dow Jones, Thomson Reuters, S&P, Capital IQ, Yodlee, Premise, CB Insights, Quandl, Xignite, PLAIN

Air / Space / Sea
 Planet Labs, Spire, Airware, DroneDeploy

Location / People / Entities
 Acxiom, Experian, Epsilon, Garmin, Foursquare, InsideView, Esri, Streetline, Clarifai, Factual, PlacIQ, Cision, Placemeter, Basis, Systems

Other
 Qualtrics, Panjiva, GATA.GOV

Incubators & Schools
 GA, DataCamp, Insight, DataElite, MetLife, The Data Incubator

Secure Unified Research Environment

- SURE hosts virtual project workspaces that are accessed remotely over encrypted internet connections
- Members of a research team have shared access to project workspaces
- Each project workspace has its own security perimeter so data from different projects cannot be combined
- All data are stored on dedicated servers housed in a highly secure data centre
- Movement of data in and out of SURE is carefully controlled via specially designed software called the Curated Gateway



Secure Unified
Research Environment

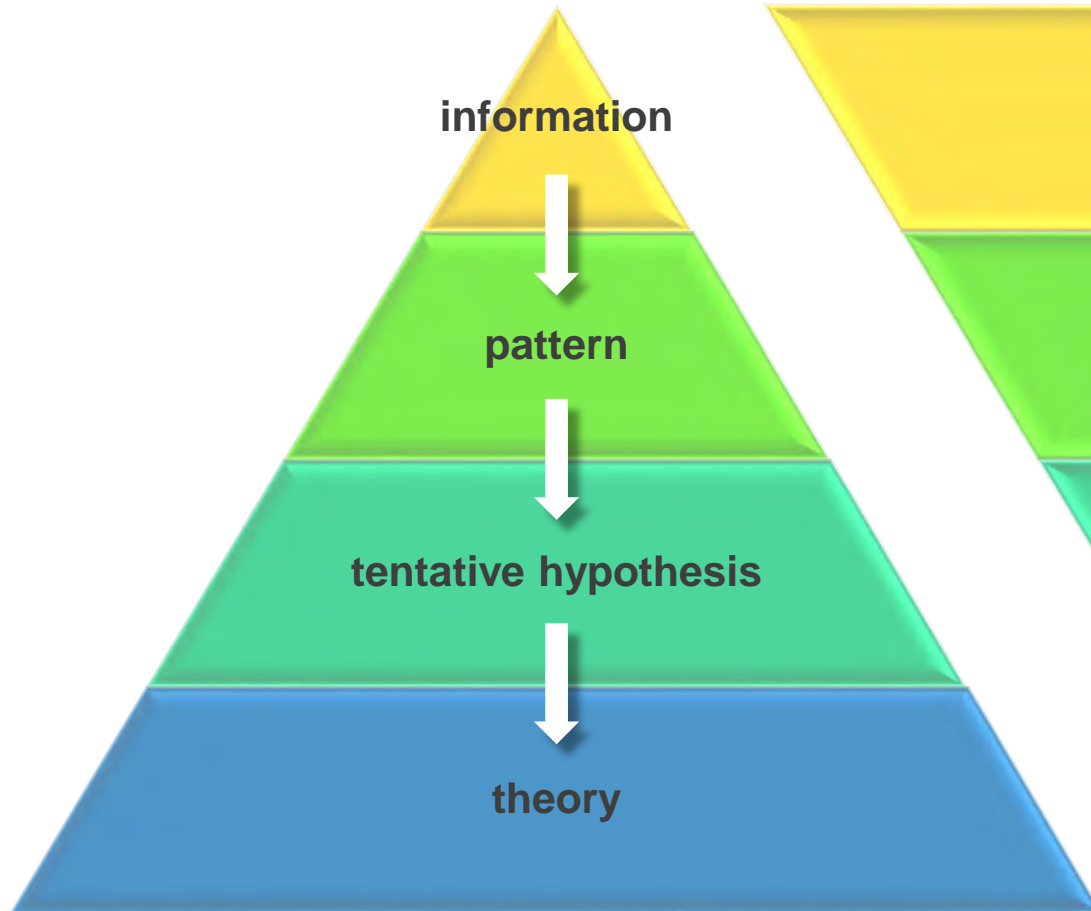
saxinstitute

Outline

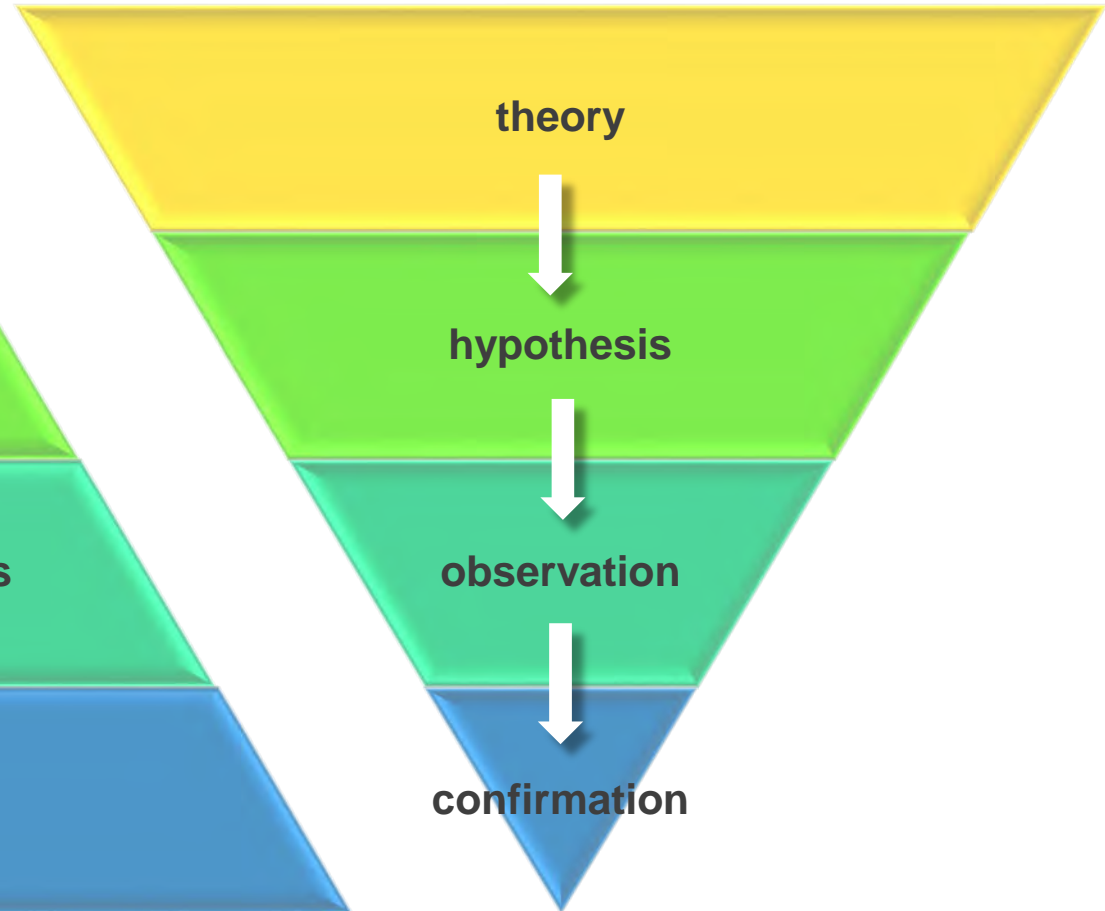
- What are big data?
- How are they being used in health and medicine?
- **Issues and challenges:**
 1. Upscaling technology
 2. **New analytic paradigms**
 3. Workforce shortages
 4. “Open data” vs. privacy protection

Deductive vs. inductive reasoning

Deductive



Inductive



Statistics vs. machine learning

	Statistics	Machine learning
Focus	Formal statistical inference (why)	Prediction (what)
Problem space	Low dimensional problems	High dimensional problems
Models	Models to explain and predict	Network/graphs to train and test
Assumptions	Explicit a priori assumptions	None (learn from the data)
Distribution	Defined <i>a priori</i>	Unknown <i>a priori</i>
Fit	Fit to the distribution	Best fit to learning models (generalisation)
Language	Estimation Data point Regression Classification Covariate	Learning Example/Instance Supervised Learning Unsupervised Learning Feature

Describing the Relationship between Cat Bites and Human Depression Using Data from an Electronic Health Record

Abstract Data mining techniques have been increasingly applied to the electronic health record and have led to the discovery of numerous clinical associations. In more detail we first used administrative diagnosis codes to identify patients with either a cat bite or a dog bite from a population of 1.3 million patients. We then conducted a manual chart review in the electronic health record to determine which were from cats or dogs. Overall there were 108 with dog bites, and approximately 117,000 patients with depression. Depression was found in 47.0% of those with cat bites, compared to 24.2% of those with dog bites. Furthermore, 85.5% of those with both cat bites and depression were women, compared to 64.5% of those with dog bites and depression. The probability of a woman being diagnosed with depression at some point in her life if she presented to our health system with a cat bite was 47.0%, compared to 24.2% of men presenting with a similar bite. The high proportion of depression in patients who had cat bites, especially among women, suggests that screening for depression could be appropriate in patients who present to a clinical provider with a cat bite. Additionally, while no causative link is known to explain this association, there is growing evidence to suggest that the relationship between cat bites and human mental illness, such as depression, warrants further investigation.



**SHOULD I FOLLOW THE
DATA**

OR MY INSTINCTS?

memegenerator.net



Revealed: Google AI has access to huge haul of NHS patient data

A data-sharing agreement obtained by **New Scientist** shows that Google DeepMind's collaboration with the NHS goes far beyond what it has publicly announced



Gathering information
Oli Scarff/AFP/Getty Images

By Hal Hodson

It's no secret that [Google has broad ambitions in healthcare](#). But a document obtained by New Scientist reveals that the tech giant's collaboration with the UK's National Health Service goes far beyond what has been publicly announced.

The agreement gives DeepMind access to a wide range of healthcare data on the 1.6 million patients who pass through three London hospitals run by the Royal Free NHS Trust – Barnet, Chase Farm and the Royal Free – each year. The agreement also includes access to patient data from the last five years.

<https://www.newscientist.com/article/2086454-revealed-google-ai-has-access-to-huge-haul-of-nhs-patient-data>

Outline

- What are big data?
- How are they being used in health and medicine?
- **Issues and challenges:**
 1. Upscaling technology
 2. New analytic paradigms
 3. **Workforce shortages**
 4. “Open data” vs. privacy protection

At e-commerce site operator Etsy Inc., a biostatistics Ph.D. who spent years mining medical records for early signs of breast cancer now writes statistical models to figure out the terms people use when they search Etsy for a new fashion they saw on the street.

Another 28-year-old at Yelp, with a Ph.D. in applied mathematics, turned his dissertation research on genome mapping into a product used by the company's advertising team. The same genome-mapping algorithm is now used to measure the effect on consumers when multiple small changes are made to online ads.



DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

COMPARISON

Jobs in Data Science



Data Scientist

vs



Data Engineer

vs



Statistician

These people use their analytical and technical capabilities to extract meaning insights from data.

These people ensure uninterrupted flow of data between servers and applications. They are responsible for data architecture.

These people understand statistics theoretically and apply them to real life problems.

Responsibilities

Develop and plan required analytic projects in response to business needs.

Contribute to data mining architectures, modeling standards, reporting, and data analysis methodologies.

Collaborate with stakeholders to integrate data mining results with existing systems.

Monitor data mining system performance and implement efficiency improvements.

Design, construct, install, test and maintain highly scalable data management systems

Improve data foundational procedures, guidelines and standards

Integrate new data management technologies and software engineering tools into existing structures

Create custom software components (e.g. specialized UDFs) and analytics applications

Apply statistical theories and methods to solve practical problems of various industries

Determine methods for finding or collecting data
Design surveys or experiments or opinion polls to collect data

Analyze, interpret & undertake data analysis

Report conclusions from their analyses

Skills

Programming, Mathematics, Business Understanding, Statistics, Data Visualization, Machine Learning, Attention to detail

Database design, Production coding, Data collection, data warehousing, Data transformation, Work diligently with data

Technical and Analytics Skills, Mathematics, Operational Research, Writing skills, Ability to Analyze, Model and interpret data, Flair of explaining difficult concepts in simple manner

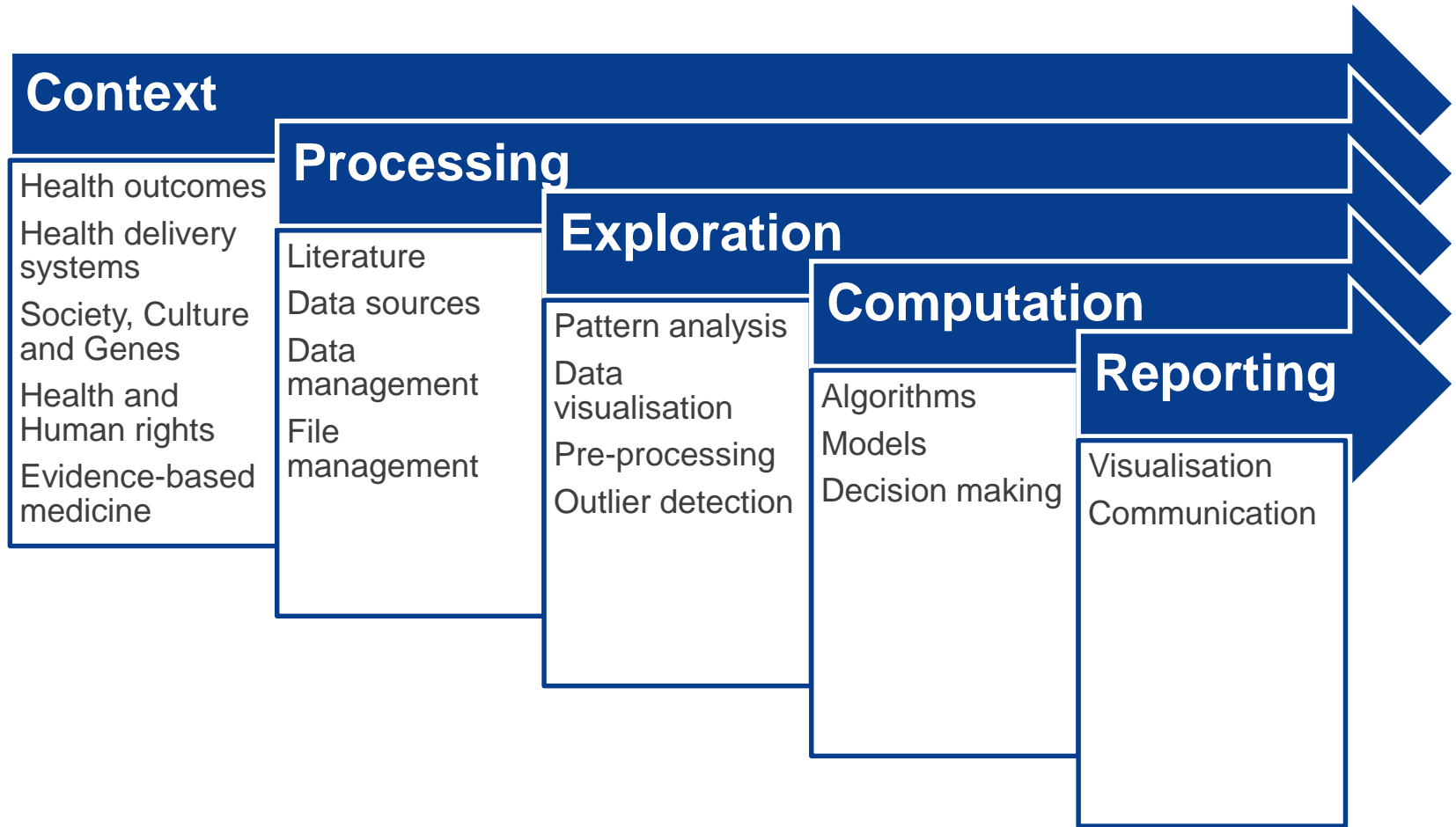
Salary (US)



Tools

<https://www.analyticsvidhya.com/wp-content/uploads/2015/10/infographic.jpg>

UNSW MSc in Health Data Science (commencing 2018)



UNSW MSc in Health Data Science: Overview

Who?

- Broad local and international target student base

What?

Grad Cert		Grad Dip		MSc	
Context of Health Data Science		Health Data Analytics: Machine Learning and Data Mining		Dissertation	Capstone
Statistical Foundations for Health Data Science		Health Data Analytics: Statistical Modelling I			Elective
Principles of Programming		Health Data Analytics: Statistical Modelling II			Elective
Management and Curation of Health Data		Visualisation and Communication of Health Data			Elective

When?

- Semester 1 2018

Outline

- What are big data?
- How are they being used in health and medicine?
- **Issues and challenges:**
 1. Upscaling technology
 2. New analytic paradigms
 3. Workforce shortages
 4. **“Open data” vs. privacy protection**

Benefits of data sharing

- Accelerates the pace of discovery
- Promotes open inquiry
- Supports diversity in analysis and interpretation
- Allows results to be replicated or alternative hypotheses to be tested
- Avoids unnecessary duplication of data collection

Barriers to data sharing

Political

- Lack of trust

Legal

- Ownership, IP
- Preserving privacy

Ethical

- Lack of proportionality
- Lack of reciprocity

Technical

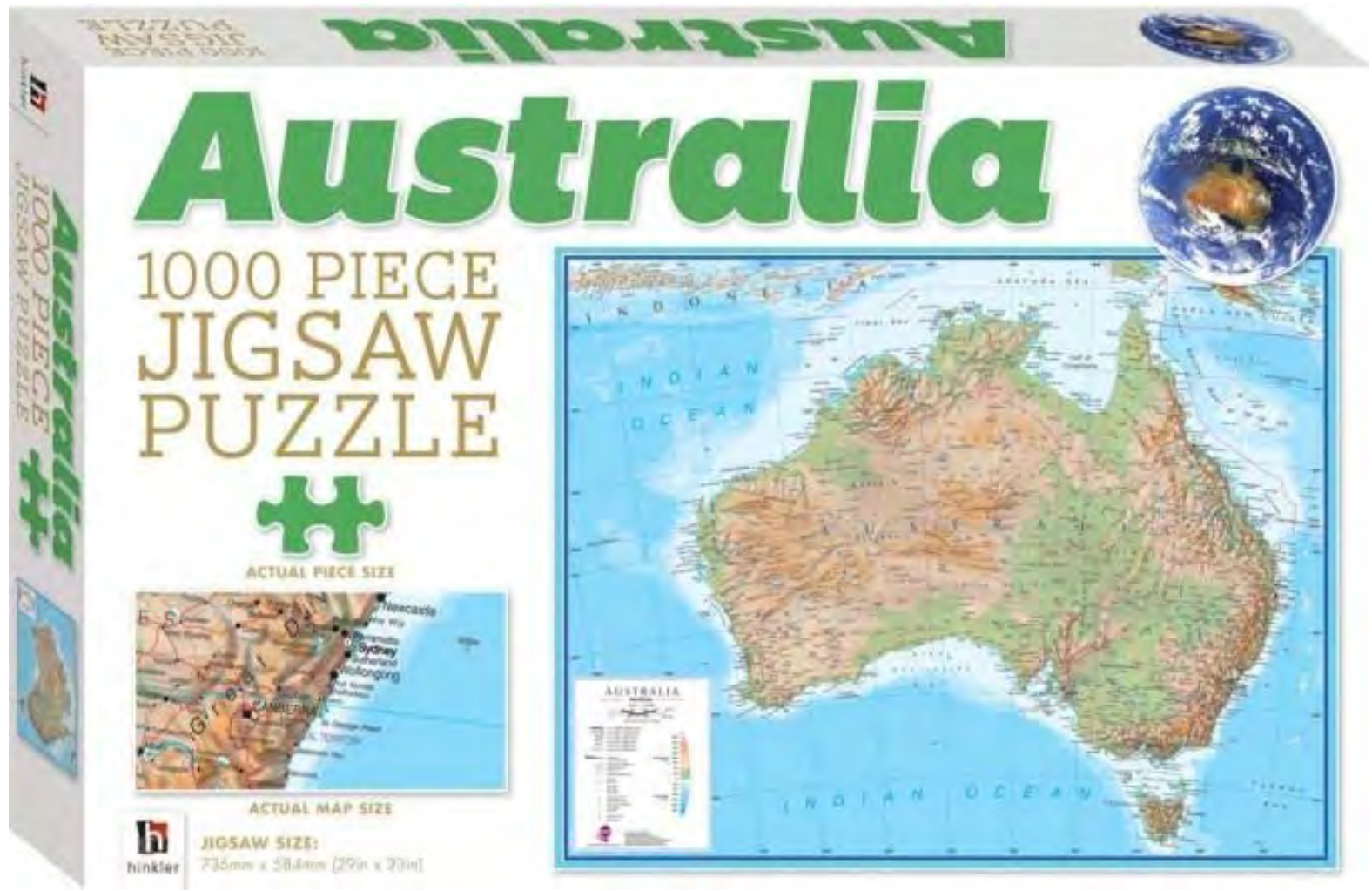
- Data not preserved or not found
- Lack of metadata and standards

Motivational

- No incentives
- Opportunity cost
- Potential criticism/embarrassment

Economic

- Lack of resources



Australia

1000 PIECE
JIGSAW
PUZZLE



ACTUAL PIECE SIZE



ACTUAL MAP SIZE



JIGSAW SIZE:
735mm x 584mm (29in x 23in)



1 General privacy legislation currently in place in Australia

Jurisdiction	Legislation	Regulator
Federal	<i>Privacy Act 1988 (Cwlth)</i>	Federal Privacy Commissioner
Australian Capital Territory	<i>Privacy Act 1988 (Cwlth)</i>	Federal Privacy Commissioner
New South Wales	<i>Privacy and Personal Information Protection Act 1998</i>	NSW Privacy Commissioner
Northern Territory	<i>Information Act 2002</i>	NT Information Commissioner
Queensland	<i>Information Privacy Act 2009</i>	QLD Information Commissioner
South Australia	<i>Cabinet Administrative Instruction 1/89 2009</i>	Privacy Committee of South Australia
Tasmania	<i>Personal Information Protection Act 2004</i>	Ombudsman Tasmania
Victoria	<i>Information Privacy Act 2000</i>	Victorian Privacy Commissioner
Western Australia	No laws	Not applicable

2 Health privacy legislation currently in place in Australia

Jurisdiction	Health privacy legislation	Regulator
Federal	<i>Privacy Act 1988 (Cwlth)</i>	Federal Privacy Commissioner
Australian Capital Territory	<i>Health Records (Privacy and Access) Act 1997</i>	Community and Health Services Complaints Commissioner
New South Wales	<i>Health Records and Information Privacy Act 2002</i>	Public sector: internal review Private sector: Privacy NSW
Northern Territory	None currently in place	Not applicable
Queensland	<i>Information Privacy Act 2009</i>	Health Quality and Complaints Commission
South Australia	<i>Code of Fair Information Practice</i>	Not applicable
Tasmania	None currently in place	Not applicable
Victoria	<i>Health Records Act 2001</i>	Health Services Commissioner
Western Australia	None currently in place	Not applicable

Cwlth = Commonwealth.

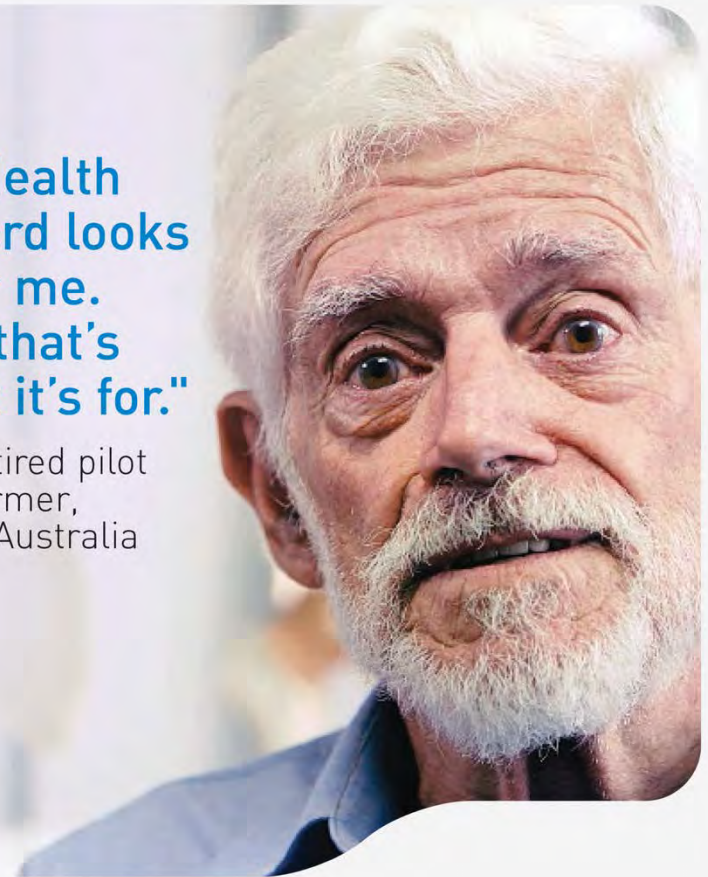
PUBLIC SECTOR DATA MANAGEMENT

July 2015

Principles for accessing and using publicly funded data for health research

"My Health
Record looks
after me.
And that's
what it's for."

Ian, retired pilot
and farmer,
South Australia



The National Health and Medical Research Council, the
Department of Health; the Australian Institute of Health and
Welfare, the Australian Bureau of Statistics, the Australian Government Department
of Health, the Australian Electoral Commission, the Australian
and Torres Strait Islander Studies, Universities Australia

Data Availability and Use

Productivity Commission
Inquiry Report

No. 82, 31 March 2017

My Health Record 'dumb and useless': Australian Privacy Foundation

Forget last week's Census debacle. Far more has been spent on an e-health system with little clinical value and fewer than 17 percent of Australians on board.



By Stilgherrian for [The Full Tilt](#) | August 19, 2016 -- 04:59 GMT (14:59 AEST) | Topic: [Security](#)

" [My Health Record](#) (MyHR) is not yet a f*** up because hardly anybody's using it, [but] it's a f*** up in terms of how much money the government has spent, and how little they've got for that expenditure," Dr Bernard Robertson-Dunn, who chairs the health committee of the Australian Privacy Foundation (APF), said.

"It's cost AU\$2 billion so far, it's costing over AU\$400 million every year, but the government has never told us how it has improved health care or reduced health costs. All it is doing is putting patient data at risk."

Blame all around as the Bureau of Statistics deflects criticism of Census 2016



Peter Martin



SHARE



TWEET



MORE

The Australian Bureau of Statistics has blamed the media for the failure of its census hotline and blamed an overseas denial-of-service attack for the failure of its census website.

In a [strongly worded submission to a Senate inquiry](#), the Bureau also attempts to deflect blame for the overwhelming of its website on to its contractor, IBM.



'Five safes' framework

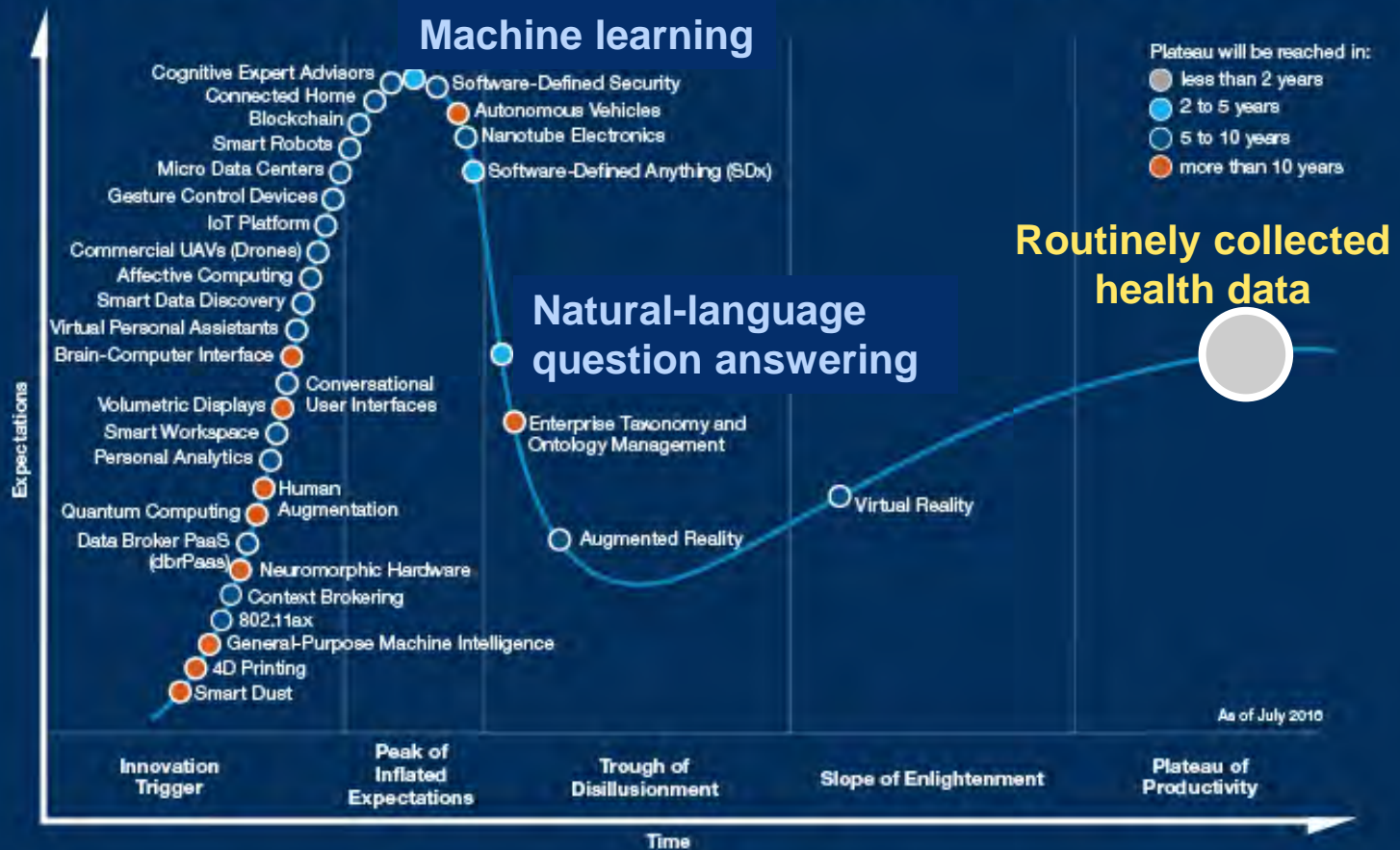
1. Safe Projects Is this use of the data appropriate?
2. Safe People Can the researchers be trusted to use it appropriately?
3. Safe Data Is there a disclosure risk in the data itself?
4. Safe Settings Does the access facility limit unauthorised use?
5. Safe Outputs Are the statistical results non-disclosive?



Outline

- What are big data?
- How are they being used in health and medicine?
- Issues and challenges:
 1. Upscaling technology
 2. New analytic paradigms
 3. Workforce shortages
 4. “Open data” vs. privacy protection

Gartner Hype Cycle for Emerging Technologies, 2016



gartner.com/SmarterWithGartner

Source: Gartner
© 2016 Gartner, Inc. and/or its affiliates. All rights reserved.

Gartner.

Thank you!
I.jorm@unsw.edu.au