# Rigour not Rigidity
*The Industrialisation and Professionalisation of Data Science*

Dr Andrew Garrett

IFoA Data Science Summit

13th September 2017

(All views are my own)

# Background

 Board member / Approvals Panel

 Chair, Data Science Section (previously VP Professional Affairs)

 Independent review of methodology function

 Drug development – manage Biostatistics, Imaging etc.

# Big Data Landscape 2016

## Infrastructure

### Hadoop On-Premise
cloudera, Hortonworks, MAPR, Pivotal, IBM InfoSphere, splice MACHINE, bluedata, jethro

### Hadoop in the Cloud
amazon web services, Microsoft Azure, Google Cloud Platform, IBM InfoSphere, CAZENA, TREASURE DATA, altiscale, Qubole, xplenty

### Spark
databricks, GridGain, TACHYON NEXUS

### Cluster Services
amazon web services, kubernetes, HPCC SYSTEMS, docker, MESOSPHERE, CoreOS, pepperdata, StackIQ

### NoSQL Databases
amazon DynamoDB, Google Cloud Platform, ORACLE, Microsoft Azure, MarkLogic, mongoDB, DATASTAX, AEROSPIKE, Couchbase, SequoiaDB, redislabs, influxdata

### NewSQL Databases
SAP HANA, Clustrix, Pivotal, paradigm4, memsql, nuoDB, MariaDB, VOLTDB, citusdata, deepdb, Trafodion, Cockroach LABS

### Graph Databases
neo4j, GIRAPH, OrientDB, InfiniteGraph

### MPP Databases
TERADATA, VERTICA, NETEZZA, kognitio, dremio

### Cloud EDW
amazon web services, Microsoft Azure, Pivotal, snowflake, WATERLINE DATA, Infoworks

### Data Transformation
alteryx, TRIFACTA, tamr, Paxata, StreamSets

### Data Integration
informatica, MuleSoft, snapLogic, BedrockData

### Management / Monitoring
New Relic, APPDYNAMICS, amazon web services, actifio, Numerify, splunk, DATADOG, Rocana, Anodot

### Security
TANIUM, illumio, CODE42, DataGravity, CipherCloud, VECTRA, sqrrl, BlueTalon

### Storage
amazon web services, Google Cloud Platform, Microsoft Azure, panasas, nimblestorage, Qumulo

### App Dev
apigee, CASK, Keen IO, Typesafe, CONCURRENT

### Crowd-sourcing
amazon mechanical turk, CrowdFlower, WorkFusion

## Analytics

### Analyst Platforms
Palantir, AYASDI, Quid, enigma, Digital Reasoning, ORBITAL INSIGHT

### Analytics Platforms
Microsoft, guavus, Datameer, inter|ana

### Data Science Platforms — context relevant
CONTINUUM, DataRobot, Alpine, MODE, plotly, ADATAO, data.ai, nutonian, DOMINO, sense, yhat, ALGORITHMIA

### Visualization
tableau, Google Cloud Platform, Roambi, GOODDATA, Qlik, CHARTIO

### BI Platforms
Power BI, amazon web services, DOMO, salesforce, Wave Analytics, birst, GoodData, platfora, looker, atscale, ARCADIA, SISENSE

### Statistical Computing
SAS, SPSS, MATLAB

### Log Analytics
splunk, sumologic, kibana, CLOUD PHYSICS, loggly

### Social Analytics
NETBASE, DATASIFT, tracx, bitly, synthesio, bottlenose, simplereach

### Real-Time
amazon web services, METAMARKETS, confluent, DATATORRENT, dataArtisans

### Machine Learning
Azure Machine Learning, amazon web services, H2O.ai, Dato, SKYTREE, rapidminer, DATARPM, deepsense.io, ViSENZE, PredictionIO, glowfish

### Speech & NLP
NarrativeScience, api.ai, NUANCE, Gridspace, semantic machines, cortical.io, MindMeld, idibon, yseop, clarifai

### Horizontal AI
IBM Watson, Cortana, sentient, VIV, nervana, vicarious, nara, Numenta, HyperScience, MetaMind, Geometric Intelligence

### Search
HP, Autonomy, ORACLE, ENDECA, EXALEAD, Lucidworks, elastic, ThoughtSpot, MAANA, swiftype, kaggle, datascope, Algolia, SINEQUA

### Data Services
LIO, OPERA, Mu Sigma, DO THE MATH, DATASCIENCE, SILICON VALLEY DATA SCIENCE, kaggle, DataKind

### For Business Analysts
OrigamiLogic, ClearStory, CIRRO, import.io

### SMB / Commerce
Google Analytics, AMPLITUDE, RJMetrics, BLUECORE, sumall, granify, Airtable, retention SCIENCE, custora

## Applications

### Sales & Marketing
RADIUS, Gainsight, bloomreach, Zeta, liveyre, blue yonder, kahuna, Lattice, SAILTHRU, persado, infer, sense, AVISO, ACTIONIQ, QUANTIFIND, ENGAGIO, DigitalGenius, appuri, fuse machines

### Customer Service
MEDALLIA, ATTENSITY, CLARABRIDGE, STELLAService, Preact, NGDATA, Wise.io, textio, entelo, hiQ

### Human Capital
gild, Connectifier

### Legal
RAVEL, JUDICATA, Everlaw, Brevia, PREMONITION

### Ad Optimization
MediaMath, Integral Ad Science, rocketfuel, OpenX, theTradeDesk, Algorithms, dstillery, LiveIntent, Data.XU, Appier, TAPAD

### Security
CYLANCE, CounterTack, cybereason, ThreatMetrix, AREA 1 SECURITY, SentinelOne, Guardian Analytics, Recorded Future, FORTSCALE, sift science, Keybase, feedzai, SIGNIFYD

### Vertical AI Applications
X, facebook, Clara, KASISTO, lumiata

### Publisher Tools
outbrain, mixpanel, Chartbeat, yieldbot, Yieldmo

### Govt/ Regulation
Socrata, OPENGOV, FiscalNote, enigma, PREDPOL, mark43, OpenDataSoft

### Finance
Affirm, LendingClub, OnDeck, Kreditech, zestfinance, LendUp, Kabbage, tidemark, payoff, INSIKT, Zuora, Dataminr, Lenddo, KENSHO, AIDYIA, iSENTIUM, Quantopian, sentient

### Education/ Learning
KNEWTON, Clever, declara, PANORAMA, knewre

### Life Sciences
23andMe, Counsyl, PATHWAY GENOMICS, Recombine, zymergen, HealthTap, KYRUUS, FLATIRON, METABIOTA, ZEPHYR HEALTH, ovia, Ginger.io, transcriptic, Glow, enlitic, AiCure, Atomwise

### Industries
OPOWER, eHarmony, RetailNext, duetto, STITCH FIX, WorkFusion, BLUE RIVER, TACHYUS, SwiftKey, Seeq, FarmLogs, HowGood, celect, statmuse, BOXEVER

## Cross-Infrastructure/Analytics
amazon web services, Google, Microsoft, IBM, SAP, SAS, Autonomy, vmware, talend, TIBCO, TERADATA, ORACLE, NetApp

## Open Source

### Framework
hadoop HDFS, hadoop MapReduce, YARN, MESOS, Spark, TEZ, Flink, CDAP

### Query / Data Flow
SLAMDATA, APACHE DRILL, APACHE HIVE, SciDB, OPENTSDB, Google Cloud Dataflow

### Data Access
accumulo, cassandra, HBASE, mongoDB, kafka, CouchDB, riak, nifi

### Coordination
talend, Apache Zookeeper, Apache Ambari

### Real-Time
STORM, Spark, APEX, Flink, TACHYON, druid

### Stat Tools
R, Scala, NumPy, SciPy

### Machine Learning
mlib, Aerossuch, Apache SINGA, MLlib, mahout, Caffe, torch, CNTK, TensorFlow, leetin, FeatureFu, WEKA, Jupyter, DL4J, VELES, DIMSUM

### Search
elasticsearch, Solr, Lucene

### Security
Apache Ranger

### Visualization
Zeppelin

## Data Sources & APIs

### Health
JAWBONE, GARMIN, practicefusion, fitbit, netatmo, Withings, VALIDIC, kinsa, Human API

### IOT
UPTAKE, ThingWorx, helium, samsara, AUGURY, estimote

### Financial & Economic Data
Bloomberg, DOW JONES, YODLEE, PREMISE, S&P CAPITAL IQ, quandl, xignite, CB INSIGHTS, mattermark, estimize, PLAID

### Air / Space / Sea
PLANET LABS, spire, WINDWARD, SKYCATCH, CRUISE, Airware, DroneDeploy

### Location/People/Entities
GARMIN, foursquare, InsideView, esri, STREETLINE, CARTODB, factual, Place IQ, Crimson Hexagon, placemeter, BASIS, Sense

### Other
qualtrics, panjiva, DATA.GOV

### Incubators & Schools
GA, DataCamp, INSIGHT, DataElite, METIS, The Data Incubator

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

# Content

- RSS DSS Remit
- Traditional Industrialised Statistics
- Big Data Challenges
- Big Data Opportunities
- Professionalisation (and Regulation)
- The 12 Questions
- Our Challenge

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS

# Data Science Section Remit

To be a professional body that represents data scientists in the UK.  The section will organise meetings for a broad range of attendees and generate outputs that are aimed at:

- Promoting good practice by addressing what good Data Science looks like (with exemplars) and what it does not look like.
- Promoting the statistical aspects of Data Science / re-enforcing the statistical framework
- Being a trusted voice on Data Science for employers, including inputting to consultation exercises
- Supporting the Data Science community throughout the UK
- Supporting the pipeline and career development of data scientists and statisticians by elevating skill sets to work in the modern world
- Supporting important emerging topics such as ethics, privacy, algorithmic responsibility and personalization - lifting the quality of the conversation
- Fostering multi-disciplinary connections and the exchanging of ideas

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Traditional industrialised statistics

- Data collected through a controlled process (experiment or survey)

- Drug development – Randomised Clinical Trial
  - e.g. Placebo Controlled Parallel Group Design

- Government – Surveys using Random sampling
  - e.g. British Crime Survey, International Passenger Survey

- Study population or sampling frame well defined

- Properties of "statistics" known
  - specifically in relation to bias and variability (precision)

- Focus on estimation, causation (?), decision-making (policy not patient)

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Big data challenges - Bias

- Data collected as part of service provision, not to produce "statistics"
    - Electronic Health (Medical) Records
    - Government Administrative Data
- Data collected passively / automatically (IoT)
- Representation (e.g. social media)?
- Meta data are key
    - How (order) data are collected / how they are checked (consistency, missing)
    - Challenge of adding variables to legacy systems
    - Merging and appending data sets
        - Consistency across data sets / providers (incl. within Government)
        - Identifiers (de-identifiers), temporal consistency

# Big data challenges – Precision

- 1/√n

- Precise in what context?

- Estimating a "constant"?

  - Data=ALL as one realisation in time

    - Super-population?

    - Prediction?

# Big data challenges - Ethics

- New ethical challenges due to:
  - Digitalisation of images/text/voice
  - Multiple data sources
  - Admin data = ALL

- Algorithmic transparency
  - Fairness
  - Understanding
  - Hacking and malevolence





ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Big Data Opportunities

- Disease trajectories
- Investigating co-morbidities
- Internet of Things





EHR: Treatment pathways of Lung cancer patients treated with erlotinib after diagnosis

Total Number by Gender

Female 6006 — 5060 Male

Export: Summary Statistics (complete cohort)

| Therapy | Sequence | Number of Patients | Average Time (Days) | Median Time (Days) |
|---|---|---|---|---|
| Carboplatin,Paclitaxel | 1 | 592 | 114.3 | 84 |
| Erlotinib | 1 | 3735 | 310.5 | 95 |
| Erlotinib | 2 | 1674 | 222.9 | 62 |

Median Age at First Diagnosis

Female 64 — 65 Male

Courtesy of J Weatherall, Head Advanced Analytics Centre, Astra Zeneca

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Big Data Opportunities

- Early estimates of economic statistics

- Lead indicators

- Private big data – useful for thinking about the consequences of new things (Bean review)

- Granularity – e.g. local inflation rates, sector specific data

- Supplementary and experimental statistics

- How to measure the digital economy

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS

# Professionalisation

- Actuaries (IFoA)
  - Certification (<u>to practice/sign</u>), pass exams/ gain exemptions /experience /CPD

- Statisticians (RSS)
  - Chartered Statistician – degree level (must have breadth/depth), 5 years experience, annual re-certification (previously award for life, if paid fee)
  - Code of Conduct
  - <u>Not</u> required to practice
  - Little traction in drug development (global), and government statisticians struggle to meet the "qualification" criterion

- Data Scientists – is something required here?

# Professionalised?

- Regulated (Drug development)
  - ICH E9 Section 1.2: actual responsibility for all statistical work… will lie with an appropriately <u>qualified</u> and <u>experienced</u> statistician… ensure that statistical principles are applied appropriately…… statistician should have a combination of <u>education/training</u> and <u>experience sufficient</u> to implement the principles….
  - In practice most will be Masters level or above
  - Many will be CStat (in UK)
- Regulated (Government)
  - National statistics are a subset of official statistics certified by the UKSA as compliant with its Code of Practice for Official Statistics

**ROYAL STATISTICAL SOCIETY**
DATA | EVIDENCE | DECISIONS

# Impact of regulation

- Quality and Reliability
- Consistency
- Trust
- Processes / procedures
- Documentation
- Audit
- Pre-specification?
- Qualified people?

- Slow
- Cumbersome
- Expensive
- Blinkered
- Risk-averse / conservative
- Always done it this way
- Uncritical
- Specialised

# Data Science – plus ca change?

Data Ingestion = data capture and storage

Data Wrangling = data management / preparation

- DS complain they spend too much time on it

Curation = managing data through its lifecycle

Munging = mapping data from one format to another

Parsing = processing text

Scraping = getting unstructured data from the Web

Data Lakes = Storing data in its native format

*Didn't Biostatistics start with unicorns?*

# But let's not forget…

Regulation is here for a reason – and it evolves in response to events



What might the future look like?

- Processes and procedures introduced in the face of material errors
- Data privacy and fairness concerns bring guidelines and regulation
- Data Scientists specialise and more routine work is done by sub-specialties
- Good Data Science Practice emerges

# The Industrialisation and Professionalisation of DS: 12 Questions

1. **What does great DS look like?**

2. **(*) What does a good DS workflow look like?**

3. **What kind of problems can be addressed by DS?**

4. **What are the characteristics of the ideal Data Scientist?**

5. **How should an organisation start a DS function**

6. **(*) How should DS fit into the structure of an organisation**

7. **How should DS business practices change to make a success of DS**

8. **(*) What do executives and managers need to know about DS**

9. **How can an organisation build a coherent DS capability from a collection of DS projects**

10. **What career paths are available to a Data Scientist?**

11. **How can Data Scientists measure the value they create?**

12. **(*) What is a Data Scientist's responsibility to wider society?**
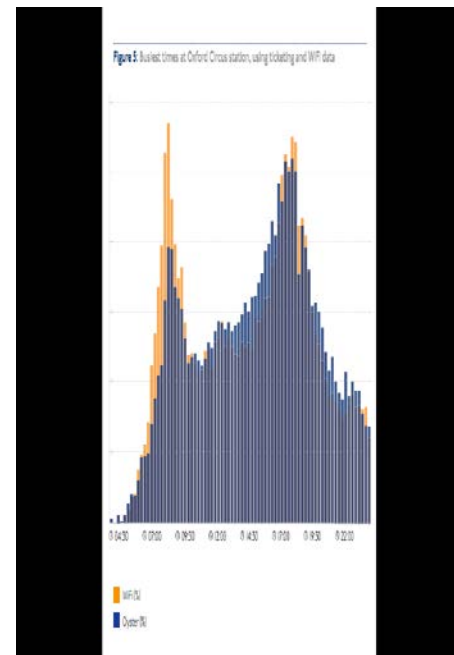
ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS

# DS is great, DS is rubbish (this week)

Neural network to generate Pub names

- 1053 NE England Pub names fed in, network creates its own rules on how to make names, through iteration

  - *Mingside Arms, Castle Stan, Burn Horse Hotel…..*

TfL WiFi tracking trial (Oyster cards don't give route)

- King's Cross → Waterloo (32% via Oxford Circus) – investigate over-crowding in stations



ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Data Science – Our Challenge

- Engage and be open – DS <u>is</u> a different mind set to traditional "industrial" statistics

- Call out bad practices and encourage good practices – there are some really neat things being done

- Be enablers, not gatekeepers – we should be pleased with the raised profile, and we don't have all the answers

- Bring rigour, not rigidity – statistical under-pinning is key

- You don't have to call yourself a data scientist to do data science – but it might help

**ROYAL STATISTICAL SOCIETY**

DATA | EVIDENCE | DECISIONS