# Data Science Workshop

**Andy Garrett**, EVP Global Scientific Operations, ICON

**Jim Weatherall**, Head Advanced Analytics Centre, AZ

RSS Data Science Section
EFSPI Statistical Leaders' Meeting
4th July 2017

# Agenda

- **Introduction – 15 mins**

LUNCH

- Survey analysis – 15 mins

- Case studies – 10 mins

- Group work: four themes – 30 mins
  - The Internet of Things
  - Big Data: EHRs
  - Decision science
  - Automation and artificial intelligence

- Report back – 25 mins

- Discussion – 25 mins

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS

# 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, pepperdata, Core OS, 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, 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, Deep Reasoning, ORBITAL INSIGHT

**Analytics Platforms:** Microsoft, guavus, Datameer, inter.ana

**Data Science Platforms / context relevant:** CONTINUUM ANALYTICS, DataRobot, Alpine, MODE, plotly, ADATAO, data.iku, 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, kyvos, looker, atscale, ARCADIA, SISENSE

**Statistical Computing:** SAS, SPSS, MATLAB

**Log Analytics:** splunk, sumologic, kibana, CLOUD PHYSICS, loggly

**Social Analytics:** NETBASE, tracx, DATASIFT, 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, Cortana, sentient, VIV, nervana, semantic machines, cortical.io, nara, MindMeld, idibon, yseop, clarifai

**Horizontal AI:** IBM Watson, vicarious, 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, 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, blue yonder, livefyre, Lattice, kahuna, SAILTHRU, persado, infer, 6sense, 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, LiveIntent, dstillery, Data.XU, Appier, TAPAD

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

**Vertical AI Applications:** X., facebook, Clara, KASIST.O, lumiata

**Publisher Tools:** outbrain, mixpanel, Chartbeat, yieldbot, Yieldmo

**Govt/ Regulation:** Socrata, OPENGOV, FN 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, knowre

**Life Sciences:** 23andMe, Counsyl, PATHWAY GENOMICS, Recombine, deep genomics, KYRUUS, FLATIRON, HealthTap, zymergen, 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, RIGHT SIGNATURE, 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, Apache Kafka, Flink, CDAP

**Query / Data Flow:** SLAMDATA, HIVE, APACHE DRILL, Google Cloud Dataflow

**Data Access:** accumulo, cassandra, HBASE, mongoDB, kafka, SciDB, nifi, CouchDB, riak, OPENTSDB

**Coordination:** talend, Apache Zookeeper, APEX, Apache Ambari

**Real-Time:** STORM, Spark, APEX, Flink, TACHYON, druid

**Stat Tools:** R, Scala, NumPy, SciPy

**Machine Learning:** mlib, Apache SINGA, MLlib, mahout, Aerosolve, Caffe, torch, CNTK, TensorFlow, leaflet, FeatureFu, WEKA, VELES, Jupyter, DL4J, 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

# 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

4

# DSS Committee Members

Fran Bennett – Mastodon C

Simon Briscoe (Council representative)

David van Dyk – Imperial / ASA DS Chapter

Andrew Garrett (Chair) - ICON
Martin Goodson – Evolution AI

Mark Girolami – Turing Institute / Imperial

Ioanna Manolopoulou - UCL

Giles Pavey – ex Dunnhumby/Tesco

Harry Powell – Barclays

Richard Pugh (Meetings Secretary) – Mango Solutions

Matthew Upson (Secretary) – Cabinet Office

Leone Wardman  - ONS

James Weatherall (Vice Chair) - AZ

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

# DSS Launch event

*The Industrialisation and Professionalisation of DS* (19th June)

- 12 Questions presented, with three formal responses

- An example topic

- President's response

- Q&A

YouTube: https://m.youtube.com/watch?v=5aH3vVvtOfc

# DSS Social Media

RSS website: landing page

Twitter: @RSS_DSS

GitHub: https://github.com/rssdatascience

LinkedIn:

https://www.linkedin.com/company-beta/111500048/

Slack: https://rssdatascience.slack.com

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Agenda

- Introduction – 15 mins

LUNCH

- Survey analysis – 15 mins
- Case studies – 10 mins
- Group work: four themes  - 30 mins
  - The Internet of Things
  - Big Data: EHRs
  - Decision science
  - Automation and artificial intelligence
- Report back – 25 mins
- Discussion – 25 mins

Please sign up!

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS
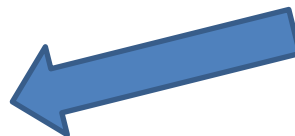
# Agenda

- Introduction – 15 mins

LUNCH

- **Survey analysis – 15 mins**
- Case studies – 10 mins
- Group work: four themes  - 30 mins
    - The Internet of Things
    - Big Data: EHRs
    - Decision science
    - Automation and artificial intelligence
- Report back – 25 mins
- Discussion – 25 mins

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Personal definitions of data science

# There are a wide range of perspectives

Gaining Knowledge and Insights from Data

Data Science is an interdisciplinary field of expertise about processes and systems to extract knowledge or insights from data in various forms, either structured or unstructured, in order to address various kinds of technical, scientific and business needs

Data-driven science based on maths, computer science and domain knowledge

Combination of computational and statistical expertise to access and analyse data

Data visualisation, modelling, simulation and AI technologies are applied in Data science

A multidisciplinary field, merging math/stat skills with computer science and

Evidence that we have failed as a statistical discipline

Database setup/programming, CRF design, data management

A blend of statistics, IT and mathematics for big data

Visualisation skills - usually focussed on a specific domain

# Data science is recognised in most organisations



## What capabilities in data science does your organisation have at present?

Answered: 27    Skipped: 0

- None
- Hobbyists - some who do...
- Recognised practitioner...
- Critical mass - a substant...
- Business as usual - a fu...

(0% 10% 20% 30% 40% 50% 60% 70% 80% 90% 100%)

# Broad range of contributions from data scientists



In which areas are data scientists in your organisation currently working? (please tick all that apply)

Answered: 26    Skipped: 1

Categories: Internet of Things (IoT)..., Applications of machine..., Data visualisatio..., Developing scientific..., Extracting insights fro..., Social media and online..., Data privacy & data ethics, Unstructured data sources, Interoperability - linking..., Modelling & simulation, Bayesian methods, Advanced statistics, Algorithm development, Decision Science, Other (please specify)

ROYAL STATISTICAL SOCIETY

DATA | EVIDENCE | DECISI

# Most believe a more mature data science capability is needed



In your view, where does the data science capability in your organisation need to be in 2 years time?

Answered: 27    Skipped: 0

# Future look: Insights, IoT, visualisation & decision science

# Opportunities for data science throughout development



In which phases of drug development do you see the greatest opportunities for data science? (please select all that apply)

Answered: 26   Skipped: 1

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

16

# Where is the gap?

# There are a wide range of perspectives

we definitely lack people able to assemble or transform the diverse datasets; we also need more associates knowledgeable or experts in Machine learning type of methods

Develop experienced DS teams gathering expertise in technology/mathematics/computer sciences while being open minded and being able to embark and lead DS projects with other scientists (biologists --> clinicians) or internal partners

Organisational boundaries

Complexity of the big data topic and variety of potential applications makes it challenging to focus and join forces between computationally oriented and statistically oriented staff

This is a multidimensional activity needing staff with different skills. Challenge is to have the right balance in the team

Statisticians with an interest in non-traditional data sources people with an interest in non-traditional data sources who understand anything about statistics, uncertainty, randomness

Limited resources/competencies in the critical areas like AI, wearable/sensor technologies

Unstructured data

Strong programming skills

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

18

# Agenda

- Introduction – 15 mins

LUNCH

- Survey analysis – 15 mins

- **Case studies – 10 mins**

- Group work: four themes  - 30 mins
    - The Internet of Things
    - Big Data: EHRs
    - Decision science
    - Automation and artificial intelligence

- Report back – 25 mins

- Discussion – 25 mins

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Patient Flows in EHR data
## OncologyFlo



Result: Example showing treatment pathways of Lung cancer patients treated with erlotinib after diagnosis

Total Number by Gender

Female 6006 | 5060 Male

Median Age at First Diagnosis

Female 64 | 65 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 |

ROYAL STATISTICAL SOCIETY
DATA | EVIDENCE | DECISIONS

# Unsupervised machine learning – Insights into healthcare



**Data:** EHRs from Greater Manchester area

# "Seven Ages of Man" healthcare clustering



| | Infant & schoolboy<br>Age 0-17 | Lover<br>Age 18-29 | Solider<br>Age 30-39 | Justice<br>Age 40-59 | Old age<br>Age 60-79 | Incapacity<br>Age > 80 |
|---|---|---|---|---|---|---|

**PCA figures:**
**(1PC vs. 2PC)**

**PCA analysis:**

- Rashes (e.g nappy rashes)
- Acne
- Eczema

- Injuries (e.g sports)
- Pains (knee pain, ankle pain, etc…)
- skin and subcutaneous tissue disease

- Circulatory system disease (hypertension, atrial fibrillation)
- Respiratory system disease (chest infection, throat infection)
- Diabetes

- Falls
- Pains (back pain, pain in limp)
- Urinary system disease (Urinary tract infection)

# What is robotic process automation (RPA)

- **Software that automates repetitive, rules-based tasks to free up your best people to be your best people**

| Robotics (RPA) | Cognitive Automation | | Artificial Intelligence |
|---|---|---|---|
| "Mimics Human Actions" | "Mimics/Augments Quantitative Human Judgment" | "Augments Human Intelligence" | "Mimics Human Intelligence" |
| • Used for rules based processes<br>• Enables:<br>  • Faster processing time<br>  • Higher volumes<br>  • Reduced errors | • Used for judgement based processes<br>• Machine learning capability<br>• Interprets human behavior | • Used for predictive decisioning<br>• Dynamically self-adaptable and managing | Turing Test Definition – "A test for intelligence in a computer, requiring that a human being should be unable to distinguish the machine from another human being by using the replies to questions put to both" |

# Safety data collection via Robotic Process Automation

Metro Map to Data Science — Author: Swami Chandrasekaran

**9. Data Munging** — Denoising, Feature Extraction, Binning Sparse Values, Unbiased Estimators, Handling Missing Values, Data Scrubbing, Normalization, Dimensionality & Numerosity Reduction

**8. Data Ingestion** — Summary of Data Formats, Data Discovery, Data Sources & Acquisition, Data Integration, Data Fusion, Transformation & Enrichment, Data Survey

Sampling, Using ETL, How much Data?, Stratified Sampling, Google OpenRefine, Principal Component Analysis

**5. Text Mining / NLP** — Corpus, Named Entity Recognition, Text Analysis, UIMA, Term Document Matrix, Term Frequency & Weight, Support Vector Machines, Association Rules, Market Based Analysis, Feature Extraction, Using Mahout, Using Weka, Using NLTK

Clustering — Hierarchical Clustering, K-means Clustering, Neural Networks, Sentiment Analysis, Collaborative Filtering, Tagging, Classify Text, Vocabulary Mapping

Regression — Perceptron, Linear Regression, Ranking, Logistic Regression

Classification — K-Nearest Neighbor, Naive Bayes Classifiers, Boosting, Decision Trees, Classification Rate, Trees & Classification, Bias & Variance, Overfitting, Lift, Prediction, Classifier, Training & Test Data, Concepts, Inputs & Attributes, Unsupervised Learning, Supervised Learning, Categorical Var, Numerical Var, What is ML?

Euclidean Distance, Least Fit, Causation, Pearson Coeff, Correlation, Covariance, Regression, Kernel Density Estimate, MLE, Estimation, Confid Int (CI), Chi-Test, p-Value, Hypothesis Testing, Monte Carlo Method, Central Limit Theorem

**6. Visualization** — Data Exploration in R (Hist, Boxplot etc), Uni, Bi & Multivariate Viz, ggplot2, Histogram & Pie (Uni), Tree & Tree Map, Scatter Plot (Bi), Line Charts (Bi), Spatial Charts, Survey Plot, Timeline, Decision Tree, Tableau, IBM ManyEyes, InfoVis, D3.js

**4. Machine Learning**

**1. Fundamentals** — Matrices & Linear Algebra Fundamentals, Hash Functions, Binary Tree, O(n), Relational Algebra, DB Basics, Inner, Outer, Cross, Theta Join, CAP Theorem, Tabular Data, Data Frames & Series, Sharding, OLAP, Multidimensional Data Model, ETL, Reporting Vs BI Vs Analytics, JSON & XML, NoSQL, Regex, Vendor Landscape, Env Setup, Entropy

Prob Den Fn (PDF), ANOVA, Skewness, Continuos Distributions (Normal, Poisson, Gaussian), Cumul Dist Fn (CDF), Random Variables, Bayes Theorem, Probability Theory, Percentiles & Outliers, Histograms, Exploratory Data Analysis, Descriptive Statistics (mean, median, range, SD, Var), Install Pkgs, Factor Analysis, Pick a Dataset (UCI Repo)

Data Frames, Reading CSV Data, Reading Raw Data, Subsetting Data, Manipulate Data Frames, Functions, Lists, Factors, Arrays, Matrices, Vectors, Variables, Expressions, R Basics, R Setup R Studio

Rapid Miner, IBM SPSS, Python Basics, Working in Excel

**10. Toolbox** — MS Excel w/ Analysis ToolPak, Java, Python, R, R-Studio, Rattle, Weka, Knime, RapidMiner, Hadoop Dist of Choice, Spark, Storm, Flume, Scibe, Chukwa, Nutch, Talend, Scraperwiki, Webscraper, Flume, Sqoop, tm, RWeka, NLTK, RHIPE, D3.js, ggplot2, Shiny, IBM Languageware, Cassandra, MongoDB

**7. Big Data** — Job & Task Tracker, MR Programming, Sqoop, Loading Data in HDFS, Flume, Scribe, Fol Unstruct Data, SQL with Pig, DWH with Hive, Scribe, Chukwa For Weblog, Using Mahout, Zookeeper Avro, Name & Data Nodes, Setup Hadoop (IBM / Cloudera / HortonWorks), Data Replication Principles, HDFS, Hadoop Components, Map Reduce Fundamentals, Storm: Hadoop Realtime, Rhadoop, RHIPE, rmr, Cassandra, MongoDB, Neo4j

**2. Statistics**

**3. Programming**

Author: Swami Chandrasekaran

25

# Agenda

- Introduction – 15 mins

LUNCH

- Survey analysis – 15 mins

- Case studies – 10 mins

- **Group work: four themes – 30 mins**
    - **The Internet of Things**
    - **Big Data: EHRs**
    - **Decision science**
    - **Automation and artificial intelligence**
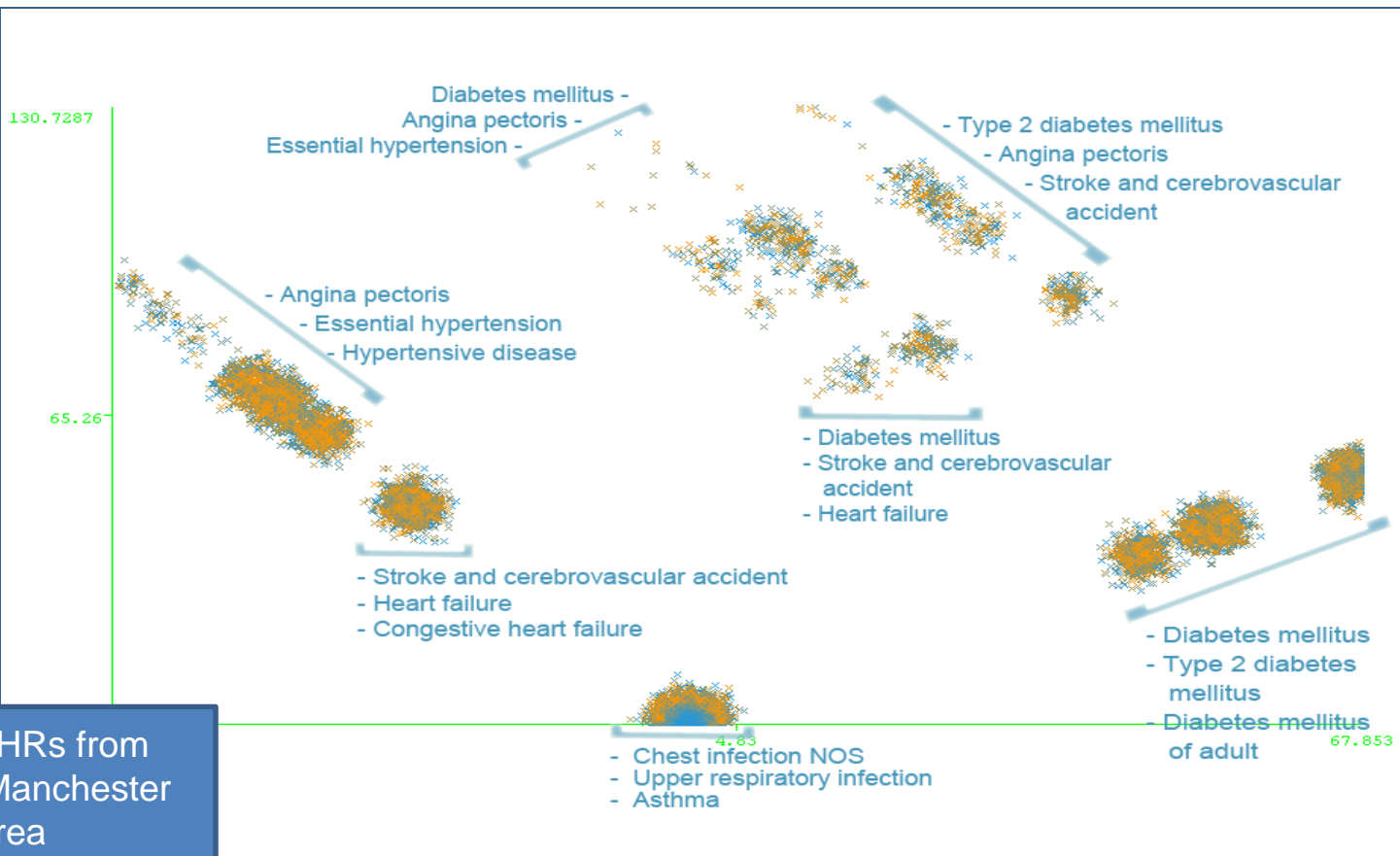
- Report back – 25 mins

- Discussion – 25 mins

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS

# Group work – three key questions

1. Brainstorm: what are the main opportunities and challenges

2. What are the top 3 areas we should address as statistical leaders

3. What immediate action should we take next?

# Agenda

- Introduction – 15 mins

LUNCH

- Survey analysis – 15 mins

- Case studies – 10 mins

- Group work: four themes - 30 mins
  - The Internet of Things
  - Big Data: EHRs
  - Decision science
  - Automation and artificial intelligence

- Report back – 25 mins

- Discussion – 25 mins

ROYAL
STATISTICAL
SOCIETY
DATA | EVIDENCE | DECISIONS

# Agenda

- Introduction – 15 mins

LUNCH

- Survey analysis – 15 mins

- Case studies – 10 mins

- Group work: four themes  - 30 mins

  - The Internet of Things

  - Big Data: EHRs

  - Decision science

  - Automation and artificial intelligence

- Report back – 25 mins

- Discussion – 25 mins

# The End

# Thank you!

ROYAL
STATISTICAL
SOCIETY

DATA | EVIDENCE | DECISIONS