

# 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

4<sup>th</sup> 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

# Big Data Landscape 2016

## Infrastructure

**Hadoop On-Premise**  
 cloudera Hortonworks  
 MAPR Pivotal  
 IBM InfoSphere  
 splince bluedata  
 jethro

**Hadoop in the Cloud**  
 amazon Microsoft Azure  
 Google Cloud Platform  
 IBM InfoSphere  
 CAZENA MESOSPHERE  
 oitiscle  
 Duoble xplenty

**Spark**  
 databricks  
 GridGain  
 TACHYON  
 N E X U S

**Cluster Services**  
 amazon  
 docker  
 HPC SYSTEMS  
 Core OS  
 peppercorn  
 SlackIQ

## Analytics

**Analyst Platforms**  
 Palantir  
 AYASDI  
 Quid enigma  
 Digital Reasoning  
 ORBITALINSIGHT

**Analytics Platforms**  
 Microsoft  
 guavus  
 Datameer  
 Interana

**Data Science Platforms**  
 context relevant  
 CONTINUUM DataRobot  
 Alpine ADAPAT  
 MOBE DORY  
 Adatalki Ionian  
 DOMINO sense  
 what ALGORITHMIA

**Visualization**  
 +tableau  
 Google Cloud Platform  
 Roambi  
 GOMDATA  
 Qlik  
 CHARTIO

**BI Platforms**  
 Power BI  
 amazon  
 Domo  
 Wave Analytics  
 GoodData  
 platforma  
 looker  
 atscale

**Statistical Computing**  
 sas  
 SPSS  
 MATLAB

**Log Analytics**  
 Splunk  
 LogLogic  
 Loggly

**Social Analytics**  
 NETBASE  
 DATASIFT  
 track  
 bitly  
 synthesis  
 bottlenose  
 simple reach

## Applications

**Sales & Marketing**  
 RADIUS Gainsight  
 bloomreach Zeta  
 livefyre blueyonder  
 kahuna Lattice  
 SAILTHRU  
 persado infer  
 AVISO sense  
 ACTIONIQ  
 QUANTIFIND JENGA GIG

**Customer Service**  
 MEDALLIA  
 ATTENITY CLARABRIDGE  
 STELLAService  
 NGDATA  
 DigitalGenius  
 Preact  
 Wiseio  
 appurri  
 fusemachines

**Human Capital**  
 gold  
 Connectifier  
 textic  
 entelo  
 hiIQ

**Legal**  
 RAVEL  
 JUDICATA  
 Everlaw  
 Brevia  
 PRESENTION

**NoSQL Databases**  
 amazon  
 dynamoDB  
 Google Cloud Platform  
 ORACLE  
 Microsoft Azure  
 MarkLogic  
 mongoDB  
 DATASTAX  
 KEROPIKE  
 Couchbase  
 SequoiaDB  
 redislabs  
 Influxdata

**NewSQL Databases**  
 SAP  
 Clustrix Pivotal  
 paradigm4  
 memsql  
 nuODB  
 MariaDB  
 VOLTDB  
 citusdata  
 deepdb  
 Trafletion  
 Cockroach LABS

**Graph Databases**  
 neo4j  
 Infogreph

**MPP Databases**  
 TERADATA  
 VERTICA  
 NETEZZA  
 OrientDB  
 Infogreph

**Cloud EDW**  
 amazon  
 Microsoft Azure  
 Pivotal  
 snowflake  
 WATERSHOLE  
 Infoworks

**Data Transformation**  
 alteryx  
 TRIFACTA  
 tamr  
 StreamSets  
 Alation

**Data Integration**  
 Informatica  
 MuleSoft  
 snapLogic  
 Bedrock Data

**Real-Time**  
 amazon  
 METAMARKETS  
 confluent  
 DATATOURNMENT  
 dataArtisans

**Machine Learning**  
 Azure  
 H2O  
 SKYTRIE  
 rapidminer  
 DEEPQUEST  
 PredictionIO  
 glowfish

**Speech & NLP**  
 NarrativeScience  
 api.ai  
 NUANCE  
 semantic machines  
 cortico.io  
 YippeeSense  
 Mindfield  
 IDIBON  
 yseop

**Horizontal AI**  
 IBM Watson  
 Cortana  
 sentient  
 viv  
 Numenta  
 MetaMind  
 clarifai

**Publisher Tools**  
 Outbrain  
 mixpanel  
 Chartbeat  
 yieldbot  
 Yieldmo

**Govt/Regulation**  
 Socrata  
 OPENGOV  
 EN FiscalNote  
 enigma  
 PREDPOL  
 OpenDataSoft

**Finance**  
 affirm  
 OnDeck  
 LendingClub  
 Kreditech  
 finance LendUp  
 Kabbage  
 tldemark  
 INSIKT  
 ZUORA  
 Dataminr  
 Lenddo  
 KENSHC  
 AIDYA  
 ISENTIUM  
 Quantopian  
 sentient

**Education/Learning**  
 KNEWTON  
 Clever  
 Oeclara  
 PANORAMA  
 knowTR

**Life Sciences**  
 23andMe  
 Counsyl  
 Recombine  
 KYRUS FLATIRON  
 zymogen HealthTop  
 METABIOTA ZEPHYR  
 ovia  
 Gingerio  
 transcriptic  
 Glow  
 enlitic  
 AiCure  
 Atomix

**Industries**  
 POWER eHarmony  
 RetailNext  
 duoetto  
 STITCH FIX  
 WorkFusion  
 TACHYUS  
 SwiftKey  
 Seeq  
 FarmLogs  
 HowGood  
 select  
 statmuse  
 BEXEVER

## Cross-Infrastructure/Analytics

amazon Google Microsoft IBM SAP SAS VMware talend TIBCO TERADATA ORACLE NetApp

## Open Source

**Framework**  
 HADOOP HDFS  
 YARN Spark  
 MESOS TEZ  
 Flink CDAP

**Query / Data Flow**  
 SLAMDATA  
 DFRILL  
 couchDB  
 rick

**Data Access**  
 HERSE  
 mongoDB  
 kaffka  
 SCIDB  
 OPENSTACK  
 nifi

**Coordination**  
 Apache Zookeeper  
 Apache Ambari

**Real-Time**  
 STORM Spark  
 APEX Flink  
 TACHYON

**Stat Tools**  
 R  
 Scala  
 SciPy

**Machine Learning**  
 mlib  
 Apache SINGA  
 MADlib  
 CNTK  
 TensorFlow  
 DL4J  
 WEKA  
 DIMSUM

**Search**  
 elasticsearch  
 Solr  
 Lucene

**Security**  
 Apache Ranger  
 Visualization  
 Kogniton

## Data Sources & APIs

**Health**  
 JAWBONE GARMIN  
 practicefusion fitbit  
 Withings VALIDIC netatmo  
 kinsa Human API

**IOT**  
 UPTAKE  
 ThingWorx  
 helium samsara

**Financial & Economic Data**  
 Bloomberg DOW JONES  
 YDLEE PREMISE  
 quandl xignite  
 mattermark  
 estimizez  
 FLAID

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

**Location/People/Entities**  
 GARMIN foursquare  
 InsideView  
 STREETLINE  
 CARTODB  
 factual PlaceIQ  
 plancemeter  
 BASIS  
 Sense

**Other**  
 qualtrics  
 panjiva  
 DATA.GOV

**Incubators & Schools**  
 DataCamp  
 INSIGHT  
 METIS  
 DataElite  
 The Data Incubator

# 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



# 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 Upton (Secretary) – Cabinet Office

Leone Wardman - ONS

James Weatherall (Vice Chair) - AZ

# DSS Launch event

---

*The Industrialisation and Professionalisation of DS (19<sup>th</sup> 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>



# 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!



# 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

# 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

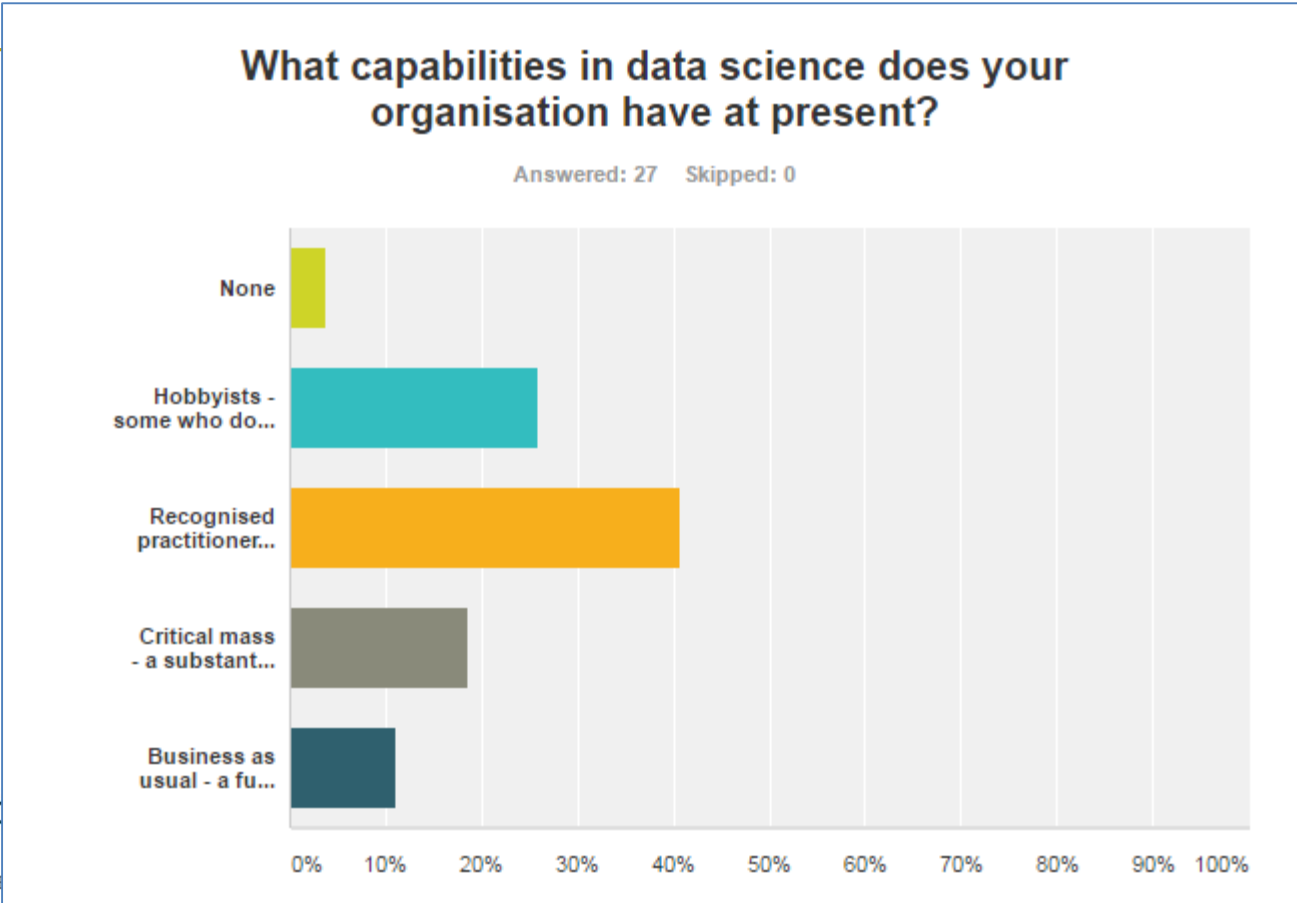
Evidence that we have failed as a statistical discipline

Database setup/programming, CRF design, data management

A blend of statistics, IT and mathematics for big data

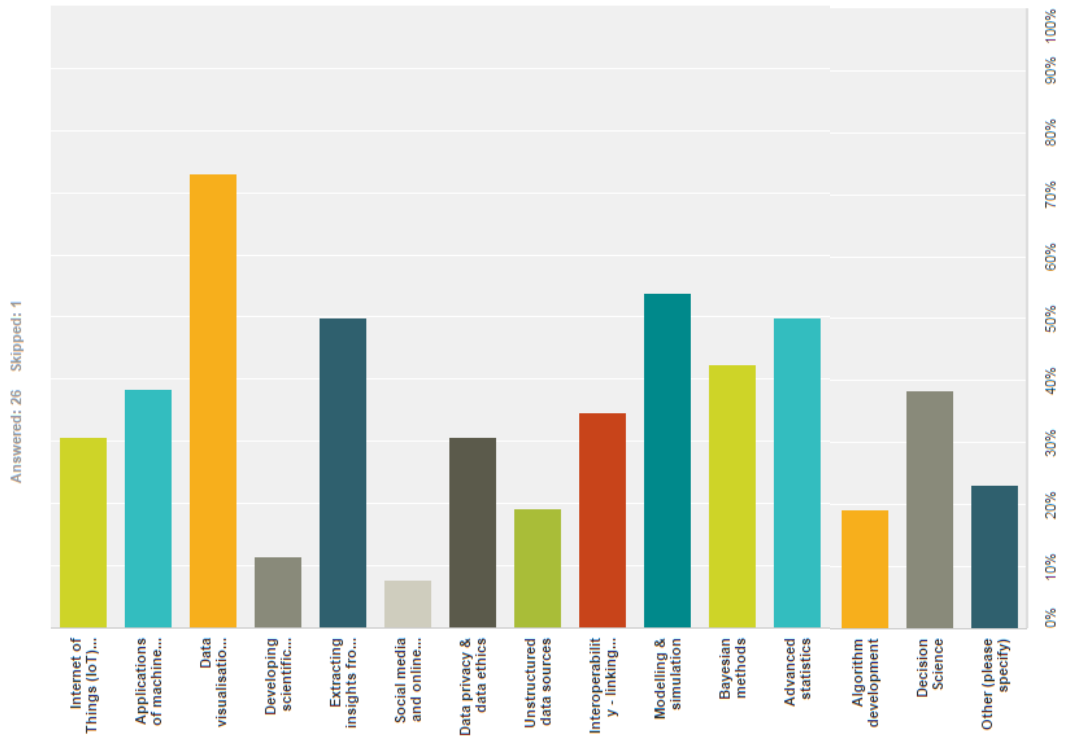
Visualisation skills - usually focussed on a specific domain

# Data science is recognised in most organisations



# Broad range of contributions from data scientists

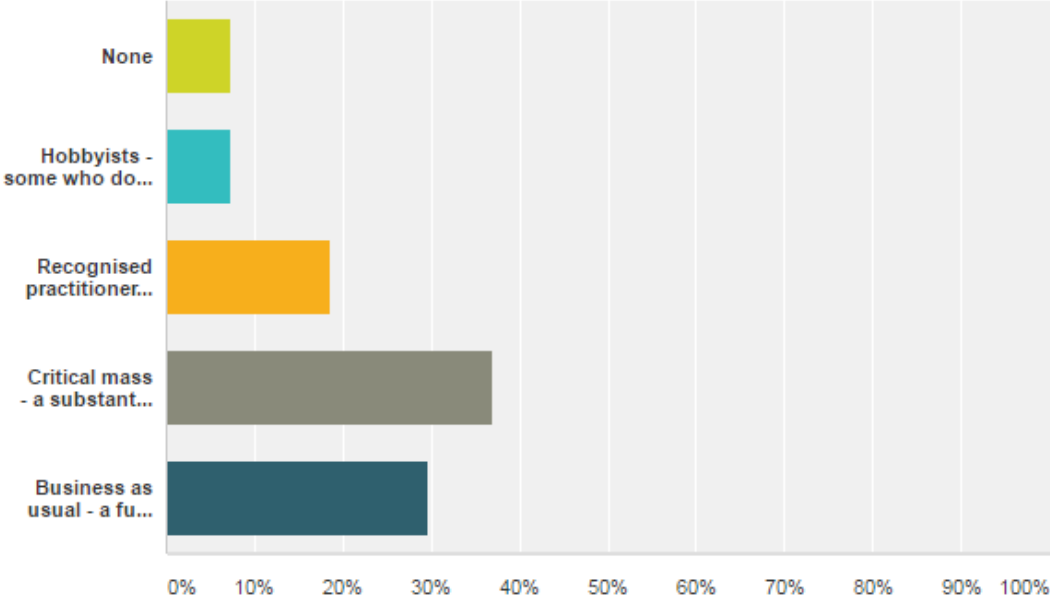
In which areas are data scientists in your organisation currently working? (please tick all that apply)



# 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?

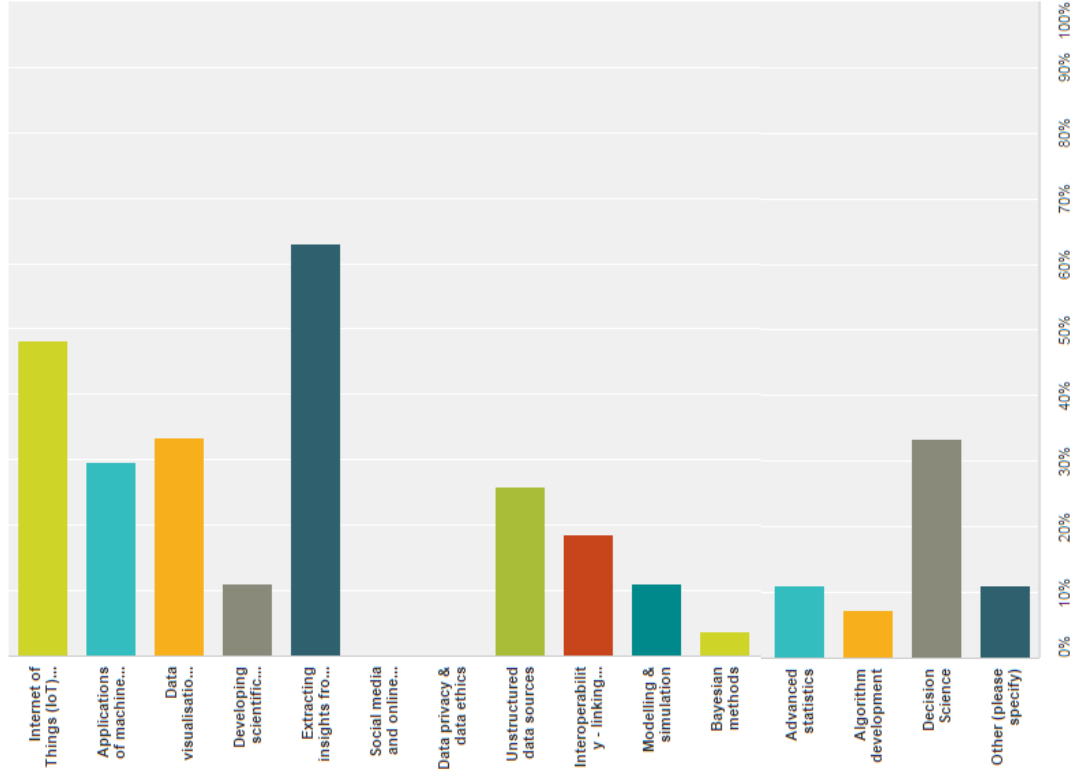
Answered: 27 Skipped: 0



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

What are the top 3 key future opportunities you see for data science in pharma?  
(please select your top 3)

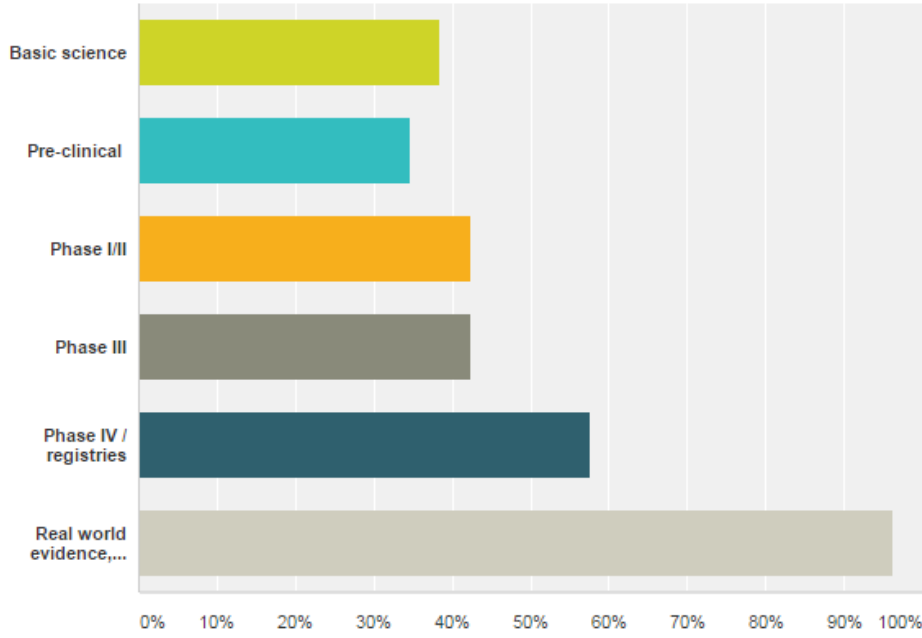
Answered: 27 Skipped: 0



# 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







# There are a wide range of perspectives

we definitely lack people able to assemble or transform the diverse datasets; we also need more 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

Strong programming skills

Unstructured data



# Agenda

---

- Introduction – 15 mins

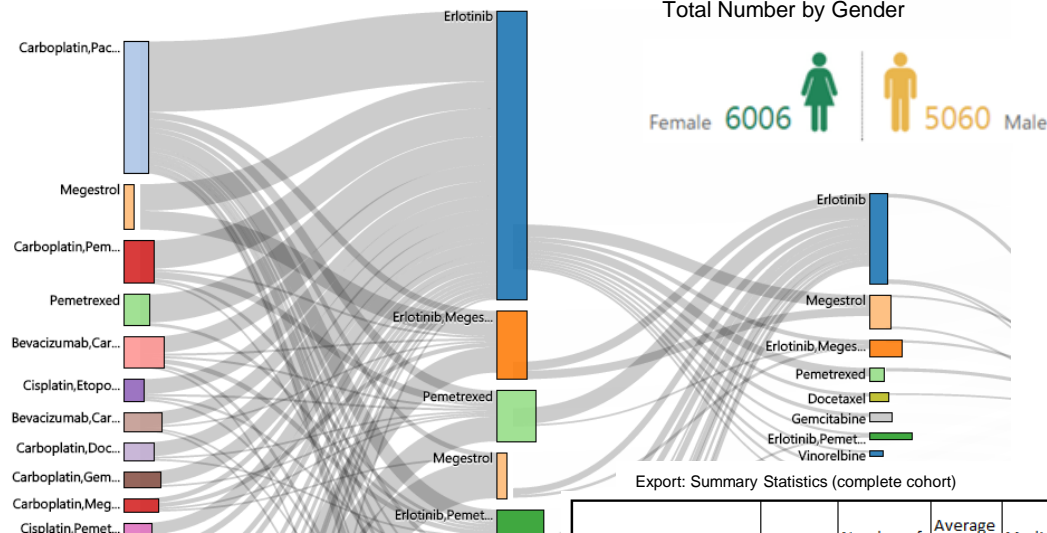
## LUNCH

- 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

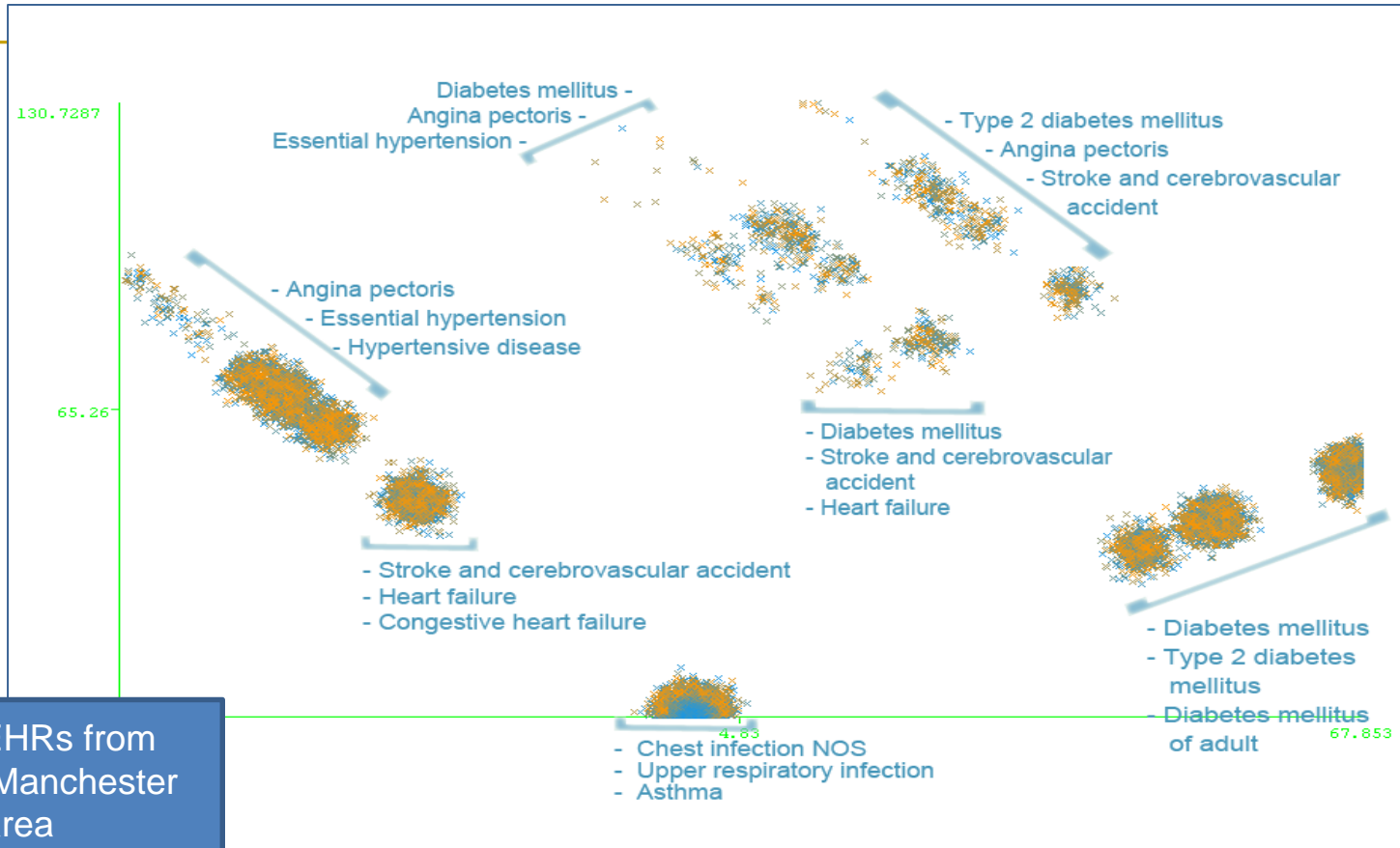
# Patient Flows in EHR data

## OncologyFlo

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



# Unsupervised machine learning – Insights into healthcare



# “Seven Ages of Man” healthcare clustering

Infant & schoolboy  
Age 0-17

Lover  
Age 18-29

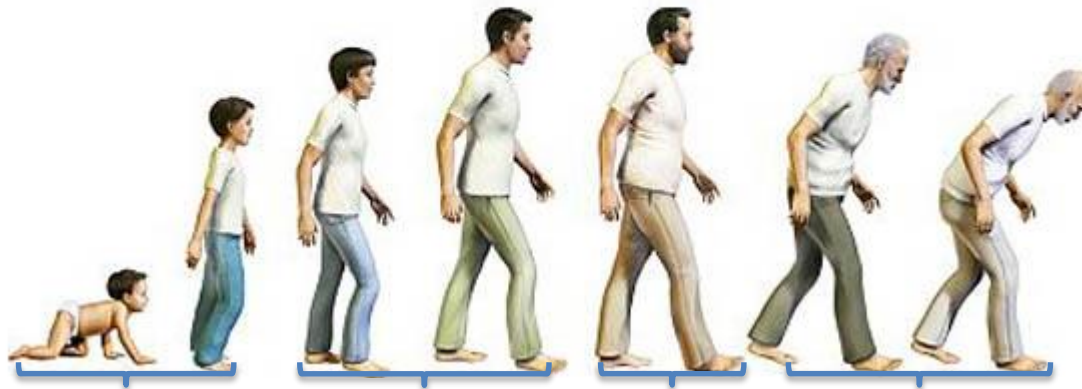
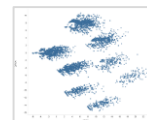
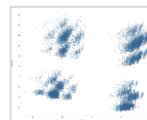
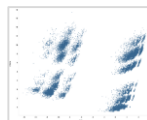
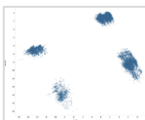
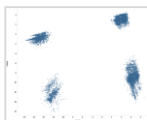
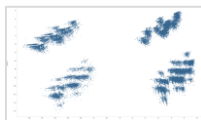
Solider  
Age 30-39

Justice  
Age 40-59

Old age  
Age 60-79

Incapacity  
Age > 80

PCA figures:  
(1PC vs. 2PC)



PCA analysis:

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

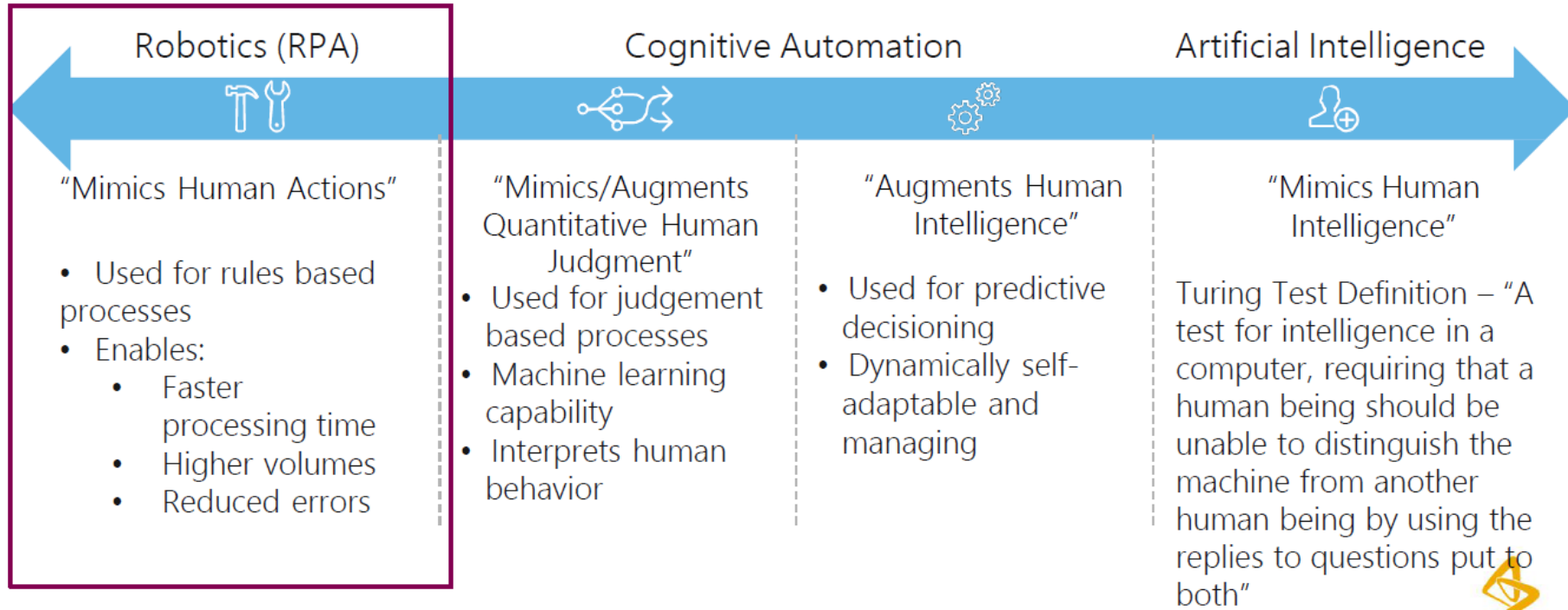
- Injuries (e.g sports)
- Pains (knee pain, ankle pain, etc...)
- skin and subcutaneous tissue disease

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

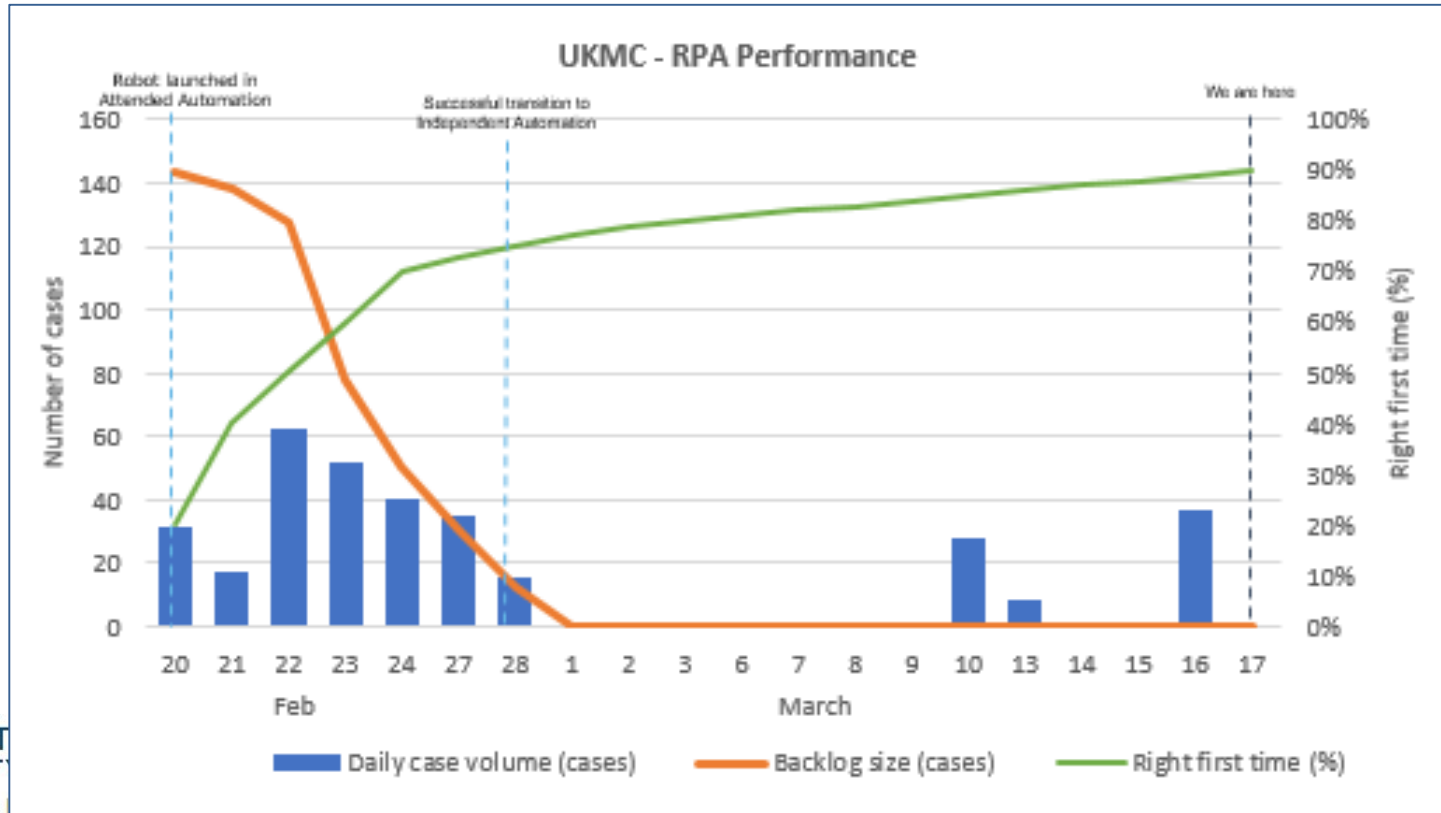
- Falls
- Pains (back pain, pain in limp)
- Urinary system disease (Urinary tract infection)

# What is robotic process automation (RPA)

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



# Safety data collection via Robotic Process Automation







# 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

# 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

# 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!

