



# Data Science, the Key for Competing in Data-driven World

**Setia Pramana**

Politeknik Statistika STIS



# Outline

- Data Driven World
- Big Data Sources
- Data Science
- What we should learn?
- Summary

# The World is Changing

## Big Data

**90%**

of the data created in the last two years alone.



## Mobile

**1 billion (plus)**  
(plus) smart devices shipped in 2013 alone.



## Social

**81%**

of customers depend on social sites for purchasing advice.



## Cloud



**62%**

of total workloads will be in the cloud by 2016.

## Internet of Things



**50 billion**  
devices connected to the internet by 2020.

## API Economy

Global m-commerce sales were

**85 billion**

in 2013 and forecast to rise to \$120 billion by 2015 and an estimated \$1 trillion by 2017





# We can do everything Online!

YESBOSS

traveloka

Qraved

www.BerryKitchen.com  
Fresh • Tasty • Delivered

waze  
OUTSKRIFTING TRAFFIC, TOGETHER

U B E R

Spotify

OPINI.id

HijUp.com

uangteman

Mobil123.com  
PORTAL OTOMOTIF NO. 1

GOJEK

livingsocial

Quipper  
SCHOOL

hukum  
online.com

airbnb

AsmaraKu  
Love • Health • Beauty

bride  
story

GOAL

klikDOKTER®  
MENUJU INDONESIA SEHAT

LewatMania.com

foodpanda

HappyFresh  
JOY DELIVERED!

telunjuk.com  
Sahabat Kita Saat Berjaya

SETIPE.COM

car(e) sharing  
nebengers  
.com

pakdok

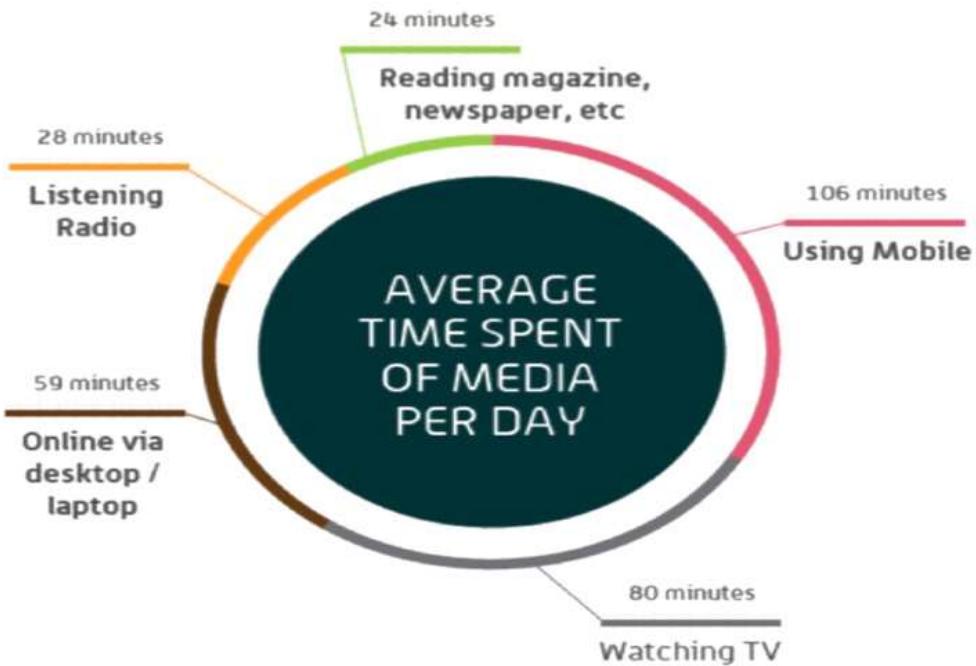
PROSEHAT  
Aplikasi Kesehatan Indonesia



JobStreet.com



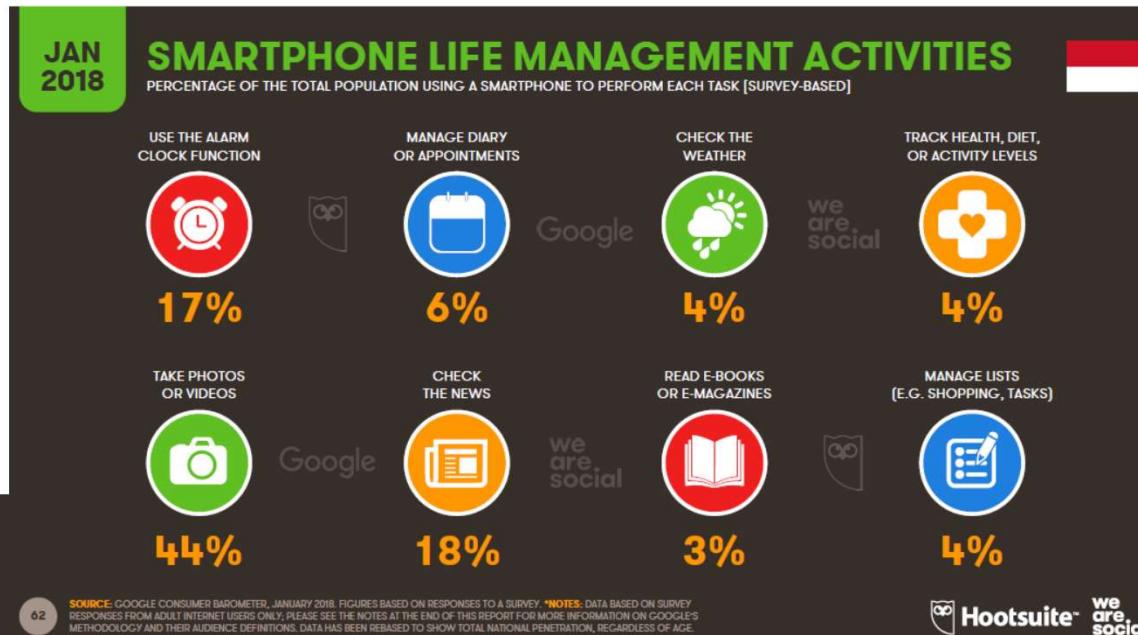
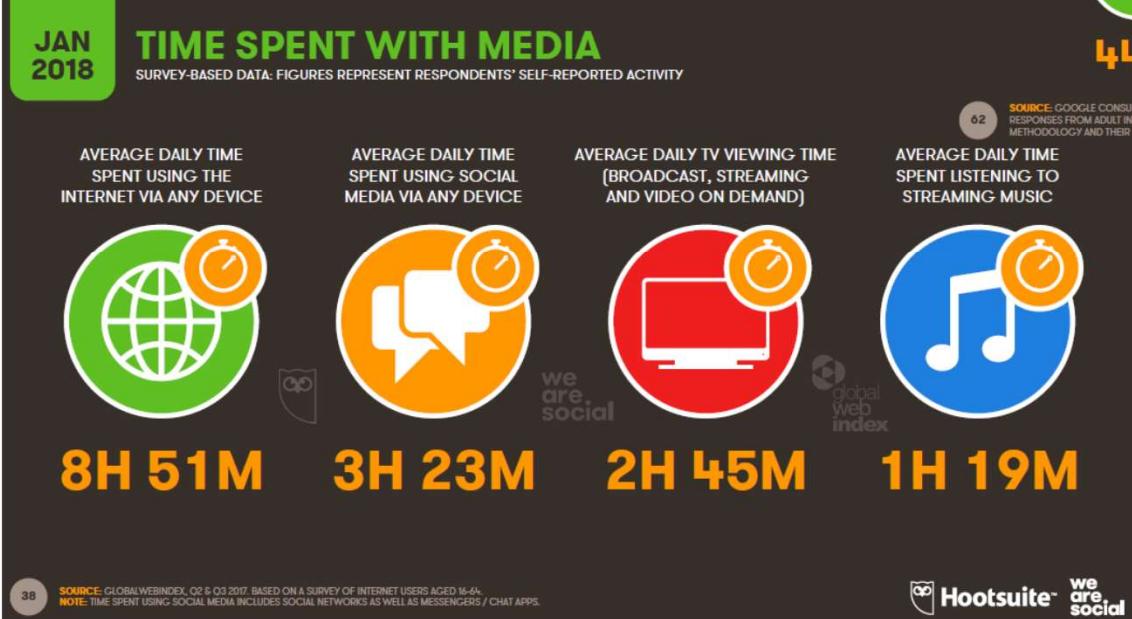
## What we “produce”?



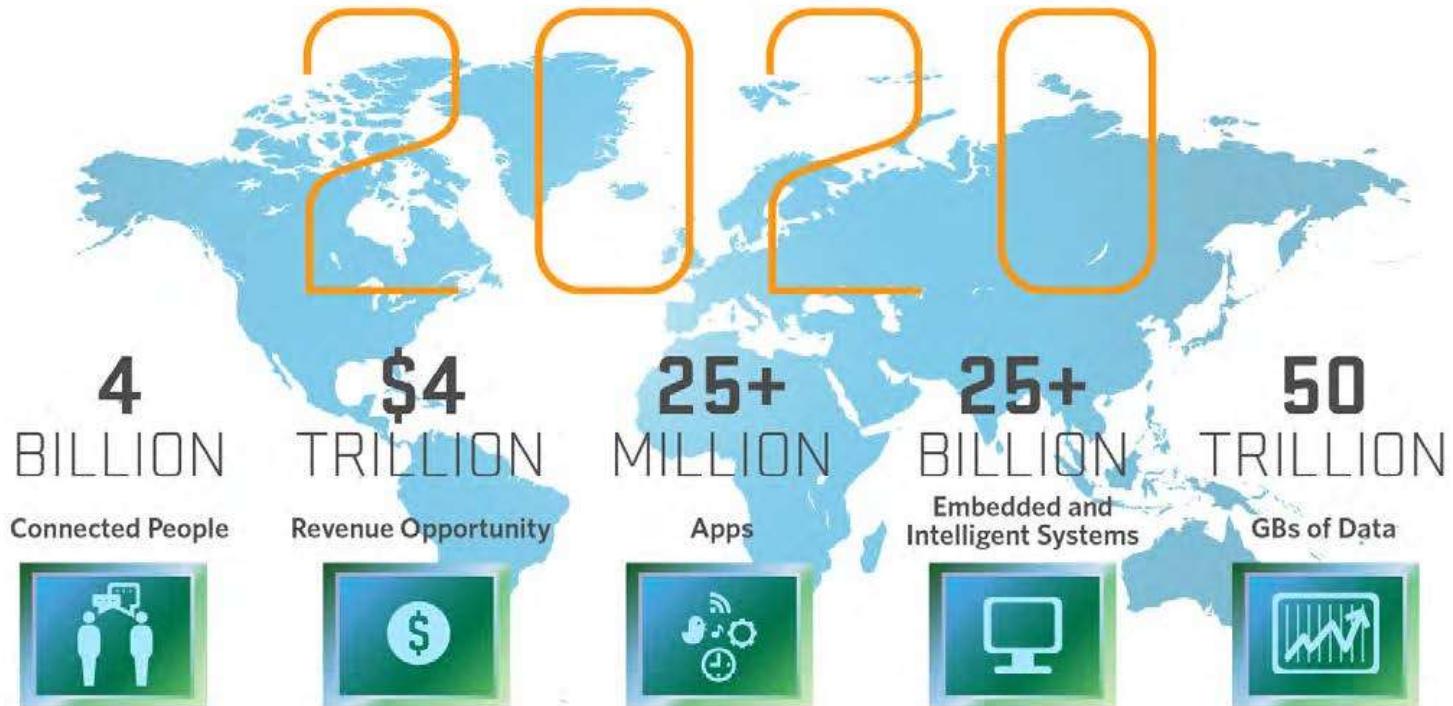
## 2018 This Is What Happens In An Internet Minute



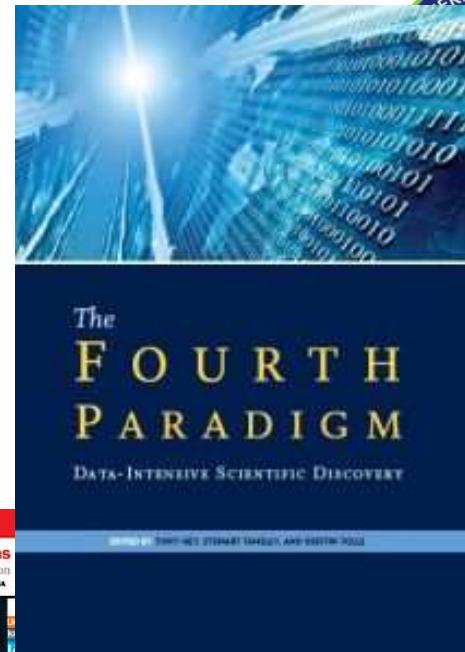
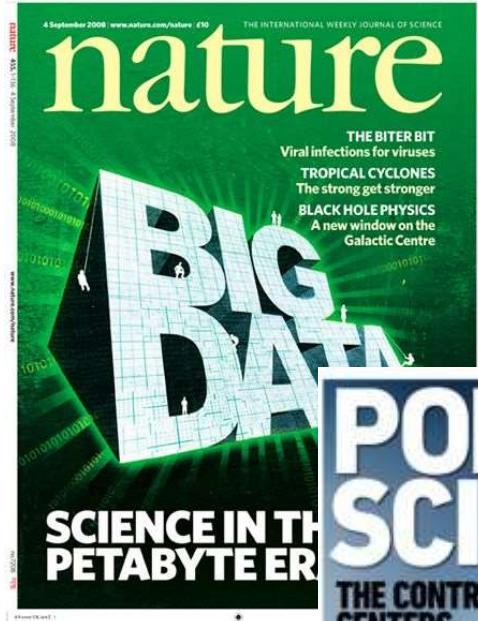
# How about Indonesia?



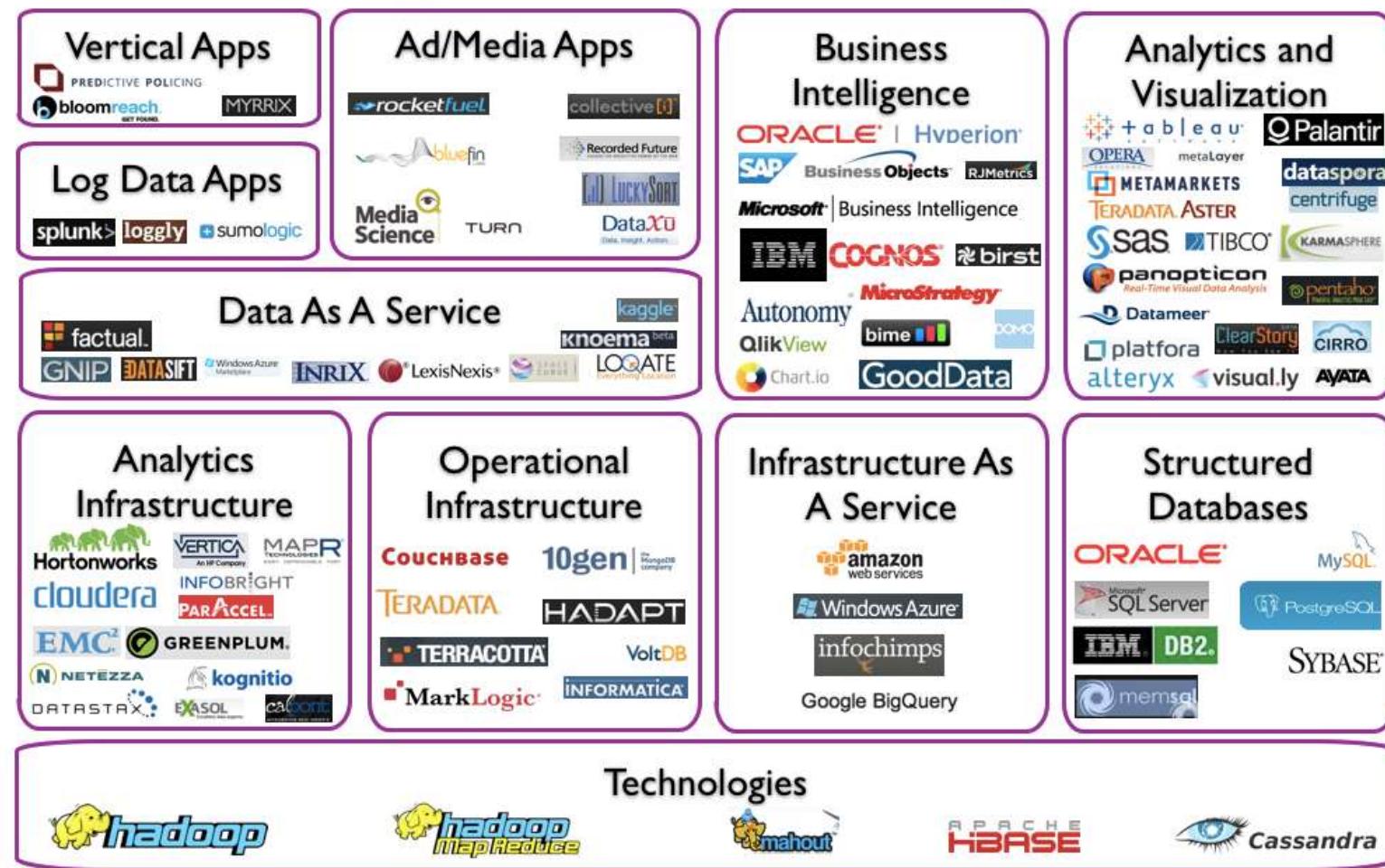
# Data Explosion



Source: Mario Morales, IDC



# Big Data Landscape



## Big Data Landscape 2016 (Version 3.0)





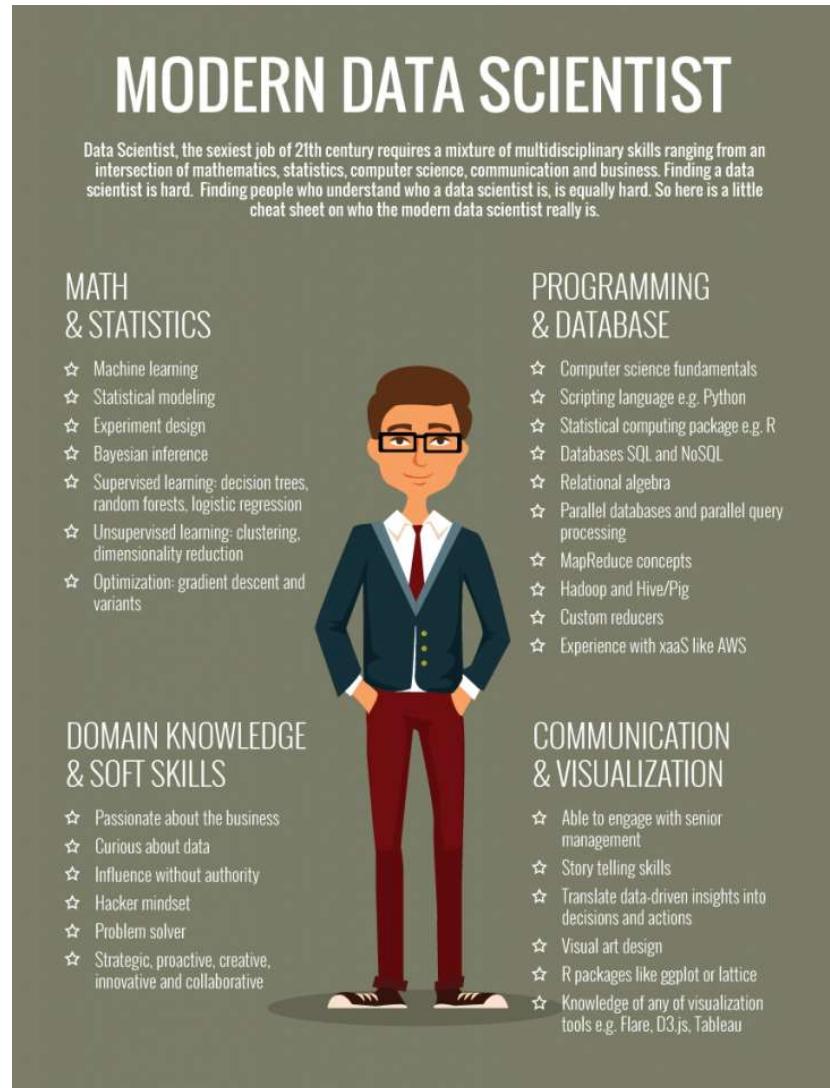
# Analytics Approaches

- Descriptive: What happened or what is happening now?
- Diagnostic: Why did it happen or Why is it happening now?
- Predictive: What will happen next? What will happen under various conditions?
- Prescriptive: What are the options to create the most optimal/high value result/outcome?

# Expertise Needed

**“The ability to take data — to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it — that’s going to be a hugely important skill in the next decades.”**

Hal Varian, chief economist at Google and UC Berkeley professor of information sciences, business, and economics



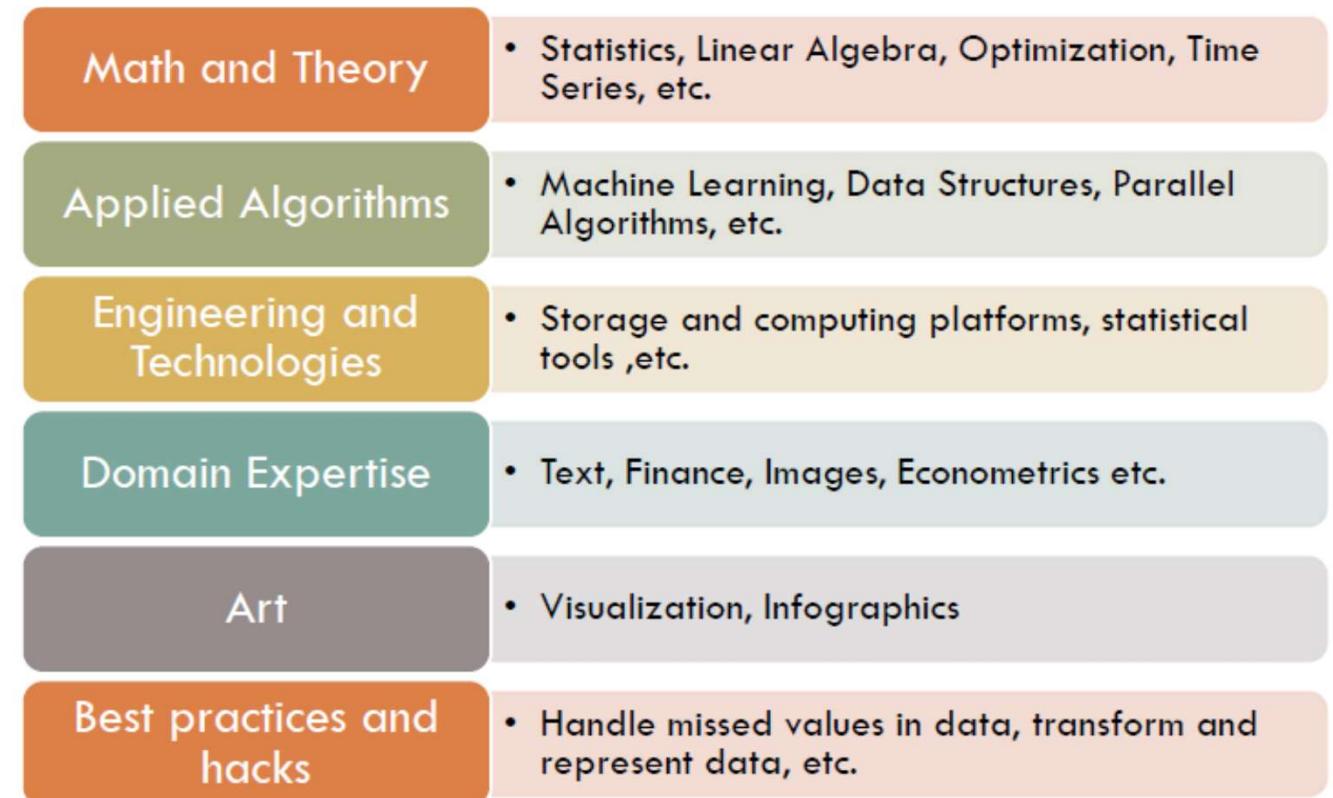
RESOURCES:

<https://insidebigdata.com/2017/08/05/benefits-data-scientist-career/>  
[https://www.glassdoor.com/Salaries/us-data-scientist-salary-SRCH\\_JL,0,2\\_IN1\\_KO3,17.htm](https://www.glassdoor.com/Salaries/us-data-scientist-salary-SRCH_JL,0,2_IN1_KO3,17.htm)  
<https://blog.udacity.com/2014/11/data-science-job-skills.html>



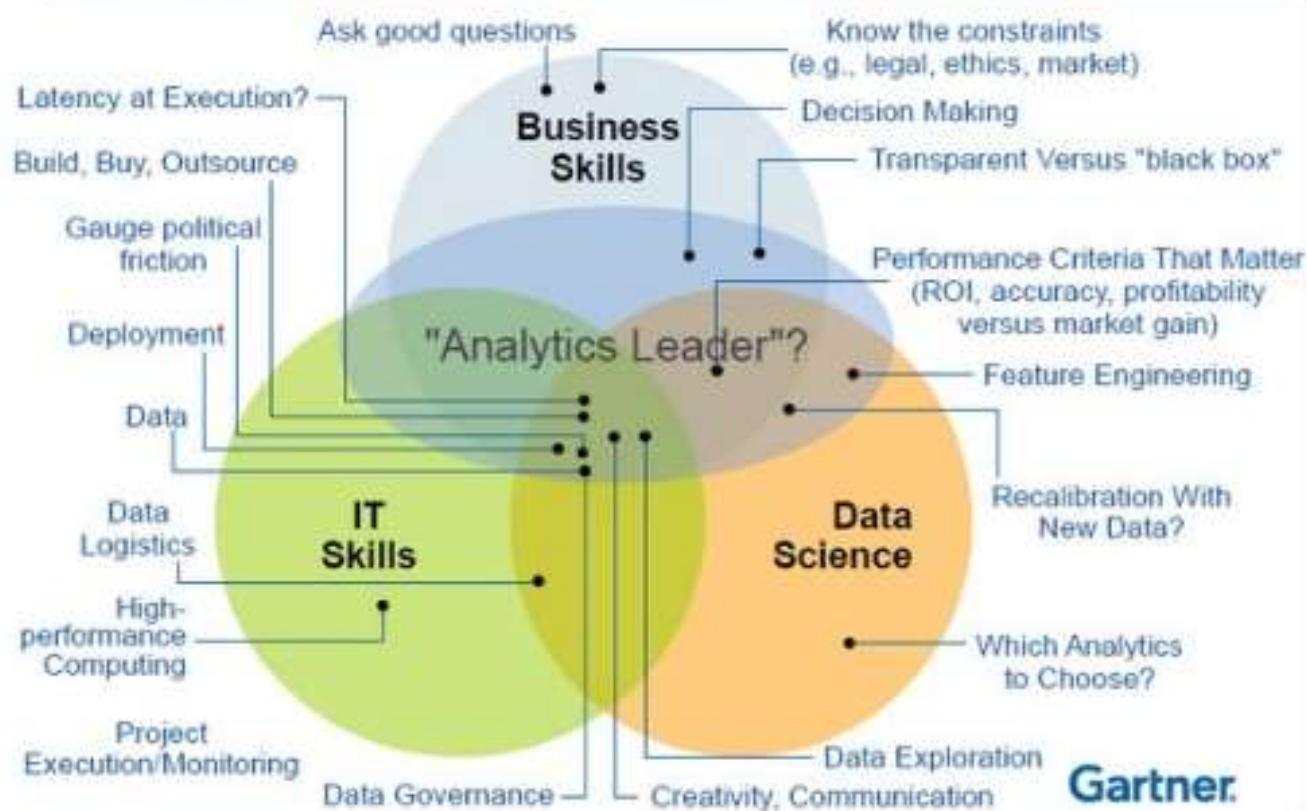
# Data Science

- A Mashed Up Discipline



# Data Science

## Driving the Success of Data Science Solutions: Skills, Roles and Responsibilities ...



# Huge Demands

**28%**

Demand Increase by  
2020

**4,524**

Number of Job  
Openings

**\$120,931**

Average Base Salary

**#1**

Best Job in America  
2016, 2017, 2018

Sources: [Glassdoor](#) and [Forbes](#)

# Data Scientist Roles

- **Math and Statistics**
- Programming and Database
- Domain Knowledge
- Communication and Visualization

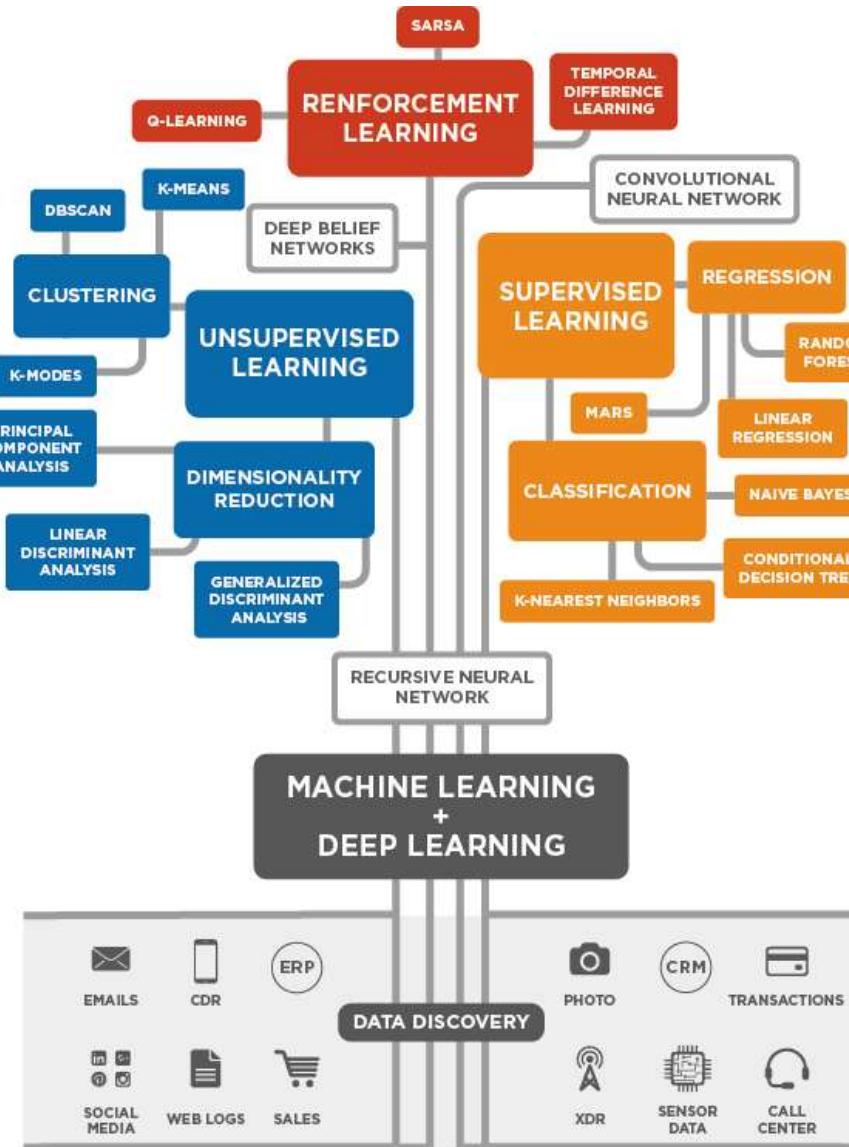
$$L(G|D) = P(G)P(D|G) = \prod_{f_i \in \{fragments\}} P(f_i|G)$$

Likelihood for the genotype      Prior for the genotype      Likelihood of the data given the genotype      Inference from reads and bases to sequenced DNA fragment to chromosomes

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma^2), \quad i = 0, \dots, K, \quad j = 1, 2, \dots, n_i,$$

$$\begin{aligned} Y_{ij} &\sim N(\mu_i, \sigma_1^{-2}) && \text{likelihood,} \\ \mu_i &\sim N(\eta_{\mu_i}, \sigma_{\mu_i}^{-2}) I(\mu_{i-1}, \mu_{i+1}) && \text{prior,} \end{aligned}$$

# Numerous (New) Algorithms





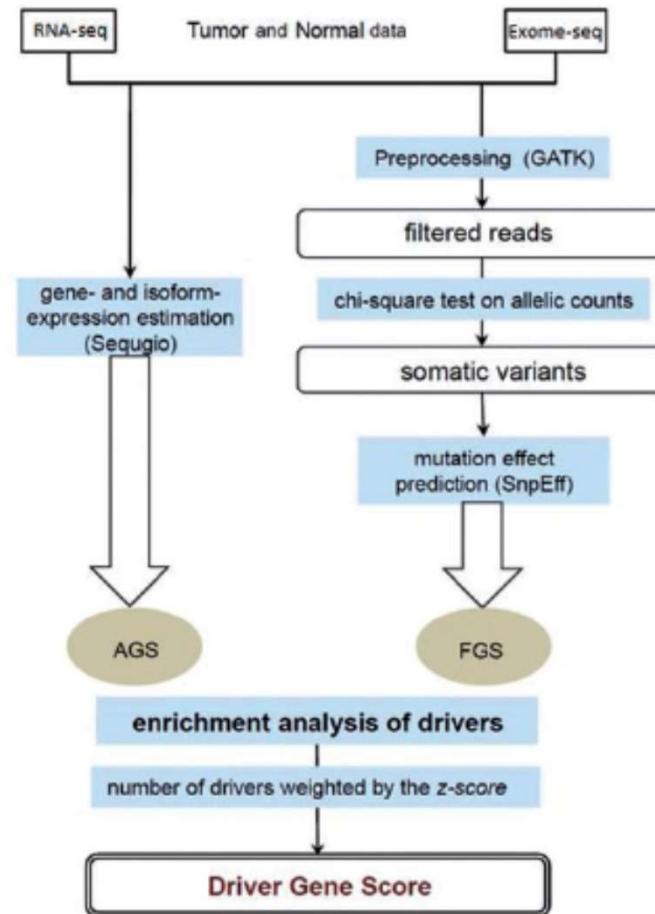
# Numerous (New) Algorithms

### Classification algorithms considered in the benchmarking study.

	BM selection	Classification algorithm	Acronym					
Individual classifier	n.a.	Bayesian Network	B-Net	n.a.	Simple average ensemble	AvgS		
		CART	CART				AvgW	
		Extreme learning machine	ELM				Stack	
		Kernalized ELM	ELM-K		Complementary measure	CompM		
		k-nearest neighbor	kNN				EPVRL	
		J4.8	J4.8				GASEN	
		Linear discriminant analysis <sup>a</sup>	LDA				HCES	
		Linear support vector machine	SVM-L				HCES-B	
		Logistic regression <sup>a</sup>	LR				MPOE	
		Multilayer perceptron artificial neural network	ANN				Top-T	
Classification models from individual classifiers	n.a.	Naive Bayes	NB	Heterogeneous ensembles	Static direct	Hill-climbing ensemble selection		
		Quadratic discriminant analysis <sup>a</sup>	QDA				CuCE	
		Radial basis function neural network	RbfNN				k-Mean	
		Regularized logistic regression	LR-R				KaPru	
		SVM with radial basis kernel function	SVM- Rbf				MDM	
		Voted perceptron	VP		Static indirect	Margin distance minimization	UWA	
<b>Classification models from individual classifiers</b>		<b>16</b>						
Homogenous ensembles	n.a.	Alternating decision tree	ADT	Dynamic	Probabilistic model for classifier competence	PMCC		
		Bagged decision trees	Bag				kNORA	
		Bagged MLP	BagNN		k-nearest oracle			
		Boosted decision trees	Boost					
		Logistic model tree	LMT					
		Random forest	RF					
		Rotation forest	RotFor					
		Stochastic gradient boosting	SGB					
					<b>Classification models from heterogeneous ensembles</b>		<b>17</b>	

# Data Scientist Roles

- Math And Statistics
- **Programming and Database**
- Domain Knowledge
- Communication and Visualization



Vu et. al, 2016

# Technical Skills

R

Python

Apache Hadoop

MapReduce

Apache Spark

NoSQL databases

Cloud computing

D3

Apache Pig

Tableau

iPython notebooks

GitHub

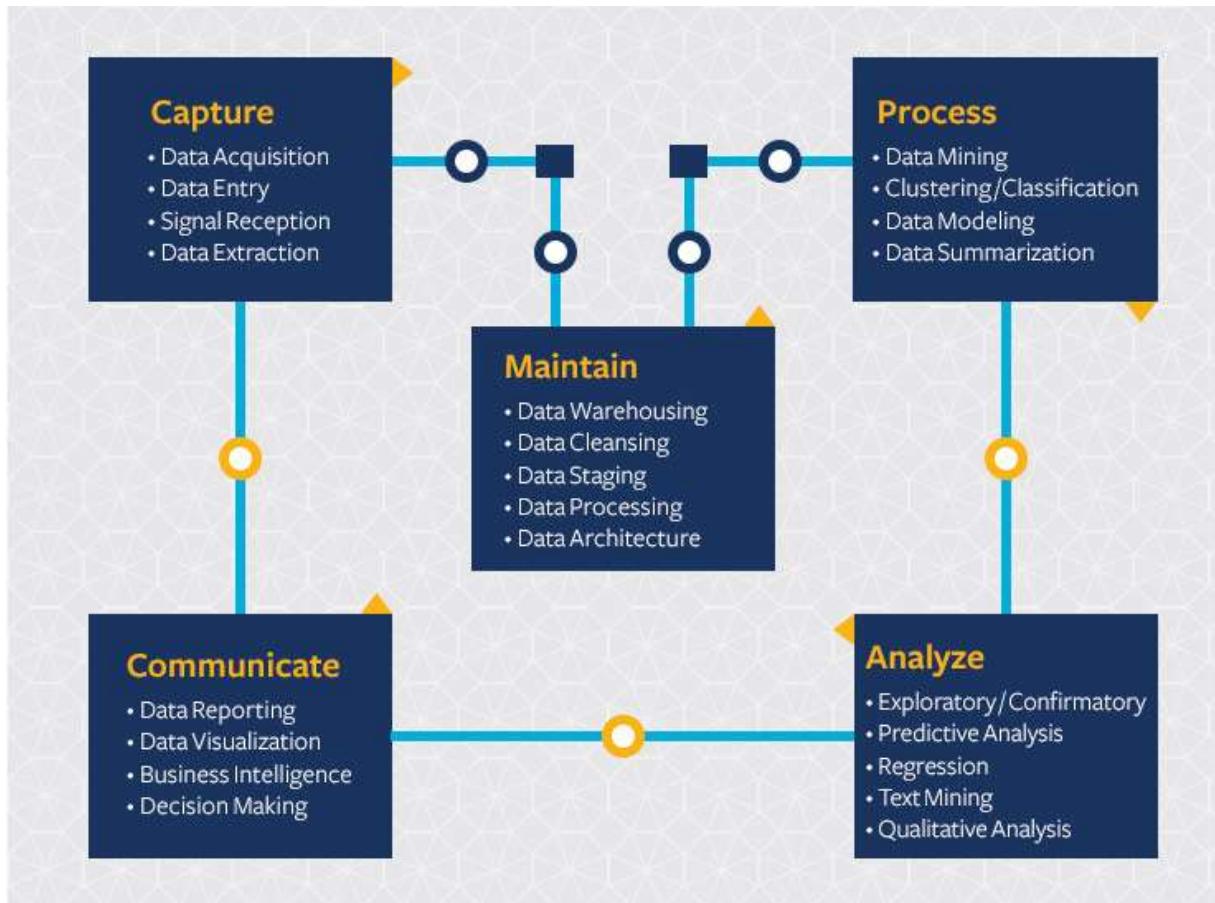
# Data Scientist Roles

- Math And Statistics
- Programming and Database
- **Domain Knowledge**
- Communication and Visualization

# Data Scientist Roles

- Math And Statistics
- Programming and Database
- Domain Knowledge
- **Communication and Visualization**

# Data Science Life Cycle



[datascience.berkeley.edu](http://datascience.berkeley.edu)

# Which one?

- **Data Scientist:** examine which questions need answering and where to find the related data. Have analytical skills as well as the ability to mine, clean, and present data, then synthesize and communicate to key stakeholders.
  - **Skills needed:** Programming skills (SAS, R, Python), statistical and mathematical skills, storytelling and data visualization, Hadoop, SQL, machine learning
- **Data Analyst:** are responsible for translating technical analysis to qualitative action items and effectively communicating their findings to diverse stakeholders.
  - **Skills needed:** Programming skills (SAS, R, Python), statistical and mathematical skills, data wrangling, data visualization
- **Data engineers** manage exponential amounts of rapidly changing data and focus on the development, deployment, management, and optimization of data pipelines and infrastructure to transform and transfer data to data scientists for querying.
  - **Skills needed:** Programming languages (Java, Scala), NoSQL databases (MongoDB, Cassandra DB), frameworks (Apache Hadoop)



# Data Science, What should we learn?

- 70% Statistical Machine Learning (7 weeks)
  - Focus on practical aspects
  - Classes
    - Necessary theoretical background
    - Basic R programming lab
- 20% Big Data Algorithms (2 weeks)
  - Focus on algorithms not on big data technologies
- 10% Data Visualization (1 weeks)
  - Grammar of graphics in R



# Data Science, What should we learn?

- Appropriate skills in statistical computing and computational topics relevant to statistics
- Introduce multivariate data early
- Teach about, and with, interactive graphics
- Work with multiple data sources (e.g, Open data etc.,) and types (text, images)
- Teach with, and about, familiar technologies: smartphones, fitness trainers and tablets offer opportunities to provoke statistical thinking
- Use Internet resources to revitalize teaching



# Office for National Statistics

**“Data Science for public good”**



**Data Science**

Applying the tools, methods and practices of the digital and data age to create new understanding and improve decision-making

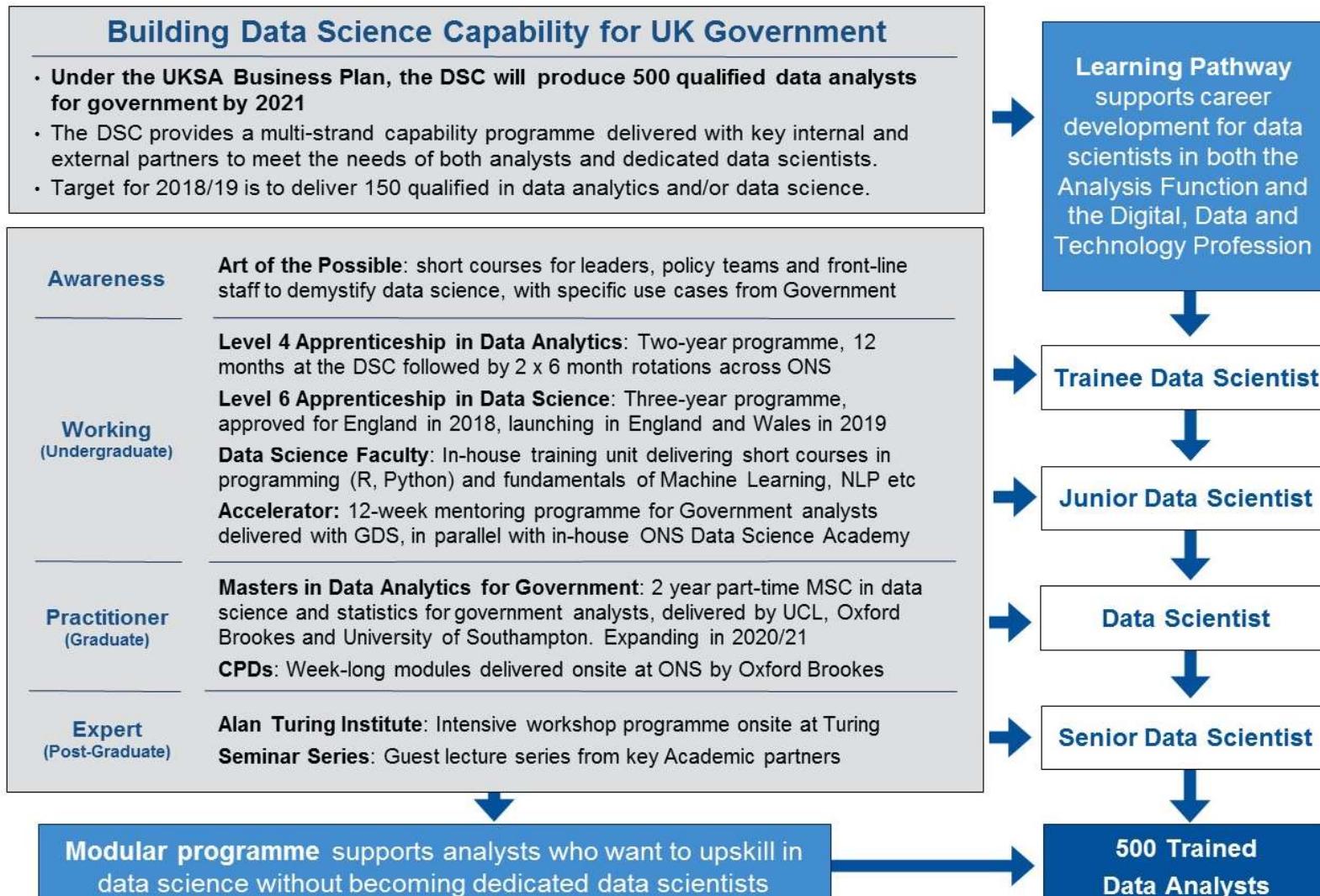
**Purpose**

We apply data science, and build skills, for public good across the UK and internationally

**Mission**

We work at the frontier of data science and AI – building skills and applying tools, methods and practices – to create new understanding which improves decision-making for public good

# Learning Pathway



# Knowledge Exchange

Data Science Faculty				
Academic Manager (G7) Solange Correa-Onel	Senior Data Scientist (G7) Rob Breton	Academic Manager (G7) Ceri Regan	Academic Manager (G7) Alison Adams	Partnerships Manager (SEO) Jane Crowe
	Data Science Lecturer (SEO) Sonia Mazzi	Data Science Trainer (SEO) Julie Owens		
	Data Science Lecturer (SEO) - Oct 1st Paraskevi Pericleousi	Data Science Training Associate (HEO) - TBD		
	Data Science Lecturer (SEO) – Oct 29th Daniel Lewis	Data Science Training Associate (HEO) - TBD		

Learning and Development Programmes	Partnership Programmes
<ul style="list-style-type: none"><li>Externally delivered programmes</li><li>MDataGov: Masters in Data Analytics for Government</li><li>CPDs: Continuous Professional Development Modules</li></ul>	<ul style="list-style-type: none"><li>In-house programme development and delivery</li><li>Data Science and AI Curriculum</li><li>R and Python building blocks</li><li>Art of the Possible</li><li>Learning Pathways</li><li>Accelerator and Data Science Academy mentoring programmes</li><li>Partnerships with ONS Learning Academy and GDS Academy</li></ul>



# Take Home Message

- Data Size has grown Exponentially and introduce a new type of data
- More Sophisticated Algorithms have been developed
- Computational Power and Storage have been improved
- Challenges in:
  - Extracting Value from Data
  - Get the right talent for Data Scientist and Business translator
- Universities (MSc Statistics and others) have to revolutionize their courses
- Needs Statistical Literacy in all Aspects



Thank you!

[setia.pramana@stis.ac.id](mailto:setia.pramana@stis.ac.id)

