

# Challenges in practical data science

Philip Jonathan

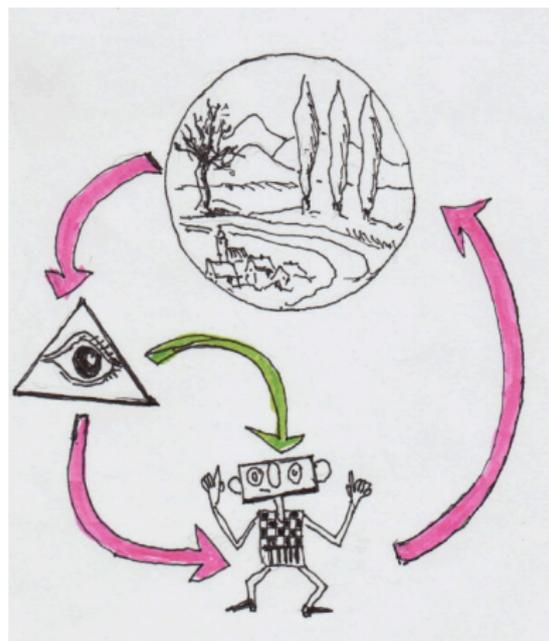
Lancaster University, Department of Mathematics & Statistics, UK.  
Shell Research Ltd., London, UK.

Seminar, Data Science Institute  
(slides at [www.lancs.ac.uk/~jonathan](http://www.lancs.ac.uk/~jonathan))



# Outline

- The digital transformation
- Representative applications
- Assessing empirical models
- Assurance in data science
  
- Acknowledgement
  - Colleagues in Shell
  - Academic partners



# Big picture

# Delivering a 'digital transformation'

- Digital roadmap
  - All parts of organisation involved (e.g. upstream, downstream, retail, finance, HR)
  - Convergence to common way of working
- Agile system development
  - Self-organising, cross-disciplinary teams
  - Common platform, governance, replication
  - Fail fast
- Delivery model
  - New technologies and R&D
  - Proof of concept studies, minimum viable products
  - Accelerators
  - Replication

# What's changed?

## Basics

- Scale, speed, connectedness
- Numeric, text, image, sound
- Parallelism, cores, clusters, cloud
- Freeware (PYTHON)
- Data engineering systems

## Cross-discipline

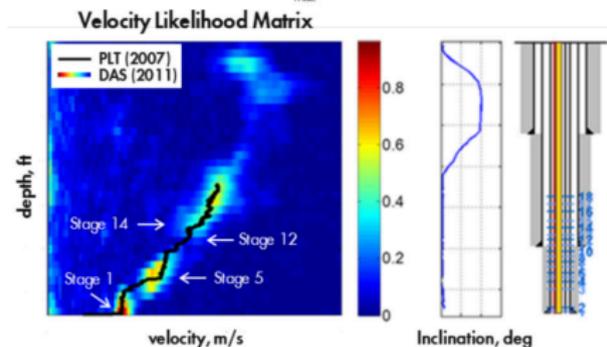
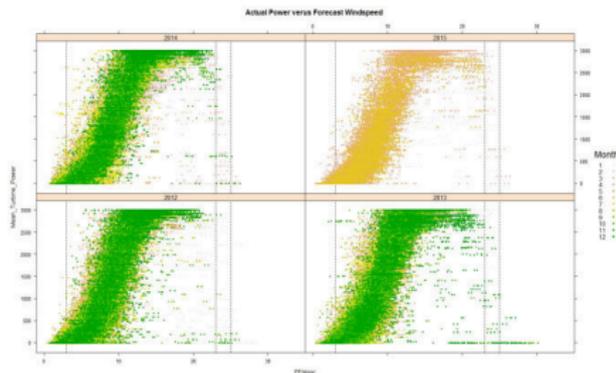
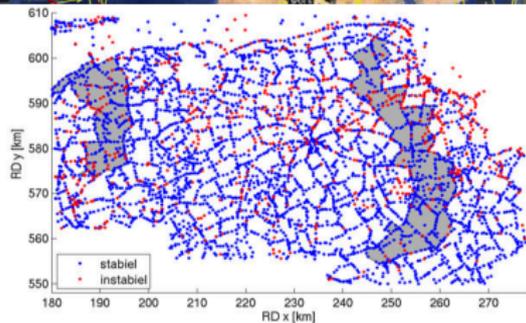
- Software engineering
- Computer science
- Data engineering
- Statistics & data science
- Communication



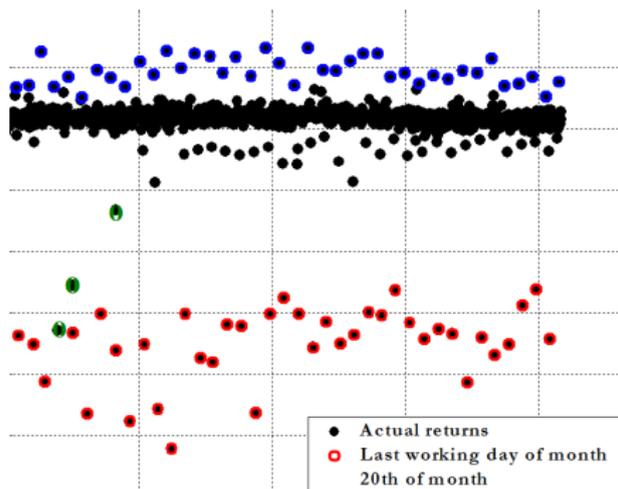
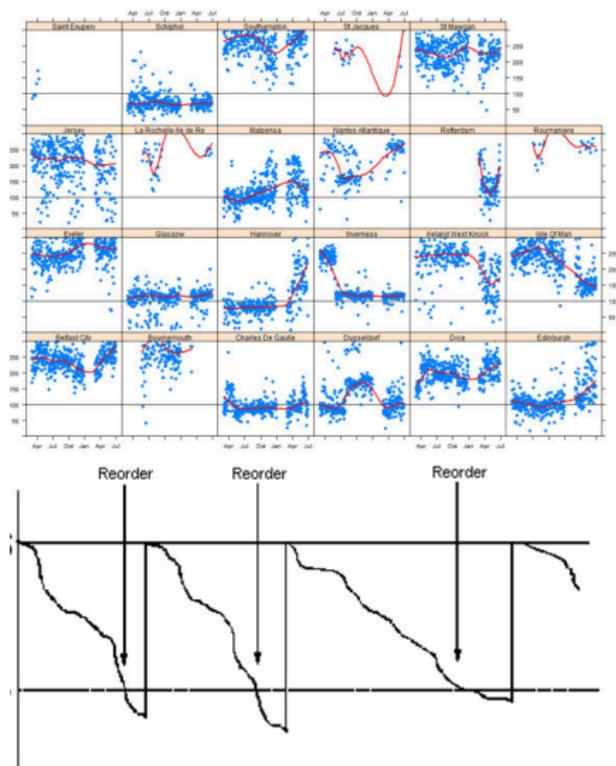
[Microsoft]

# Typical applications

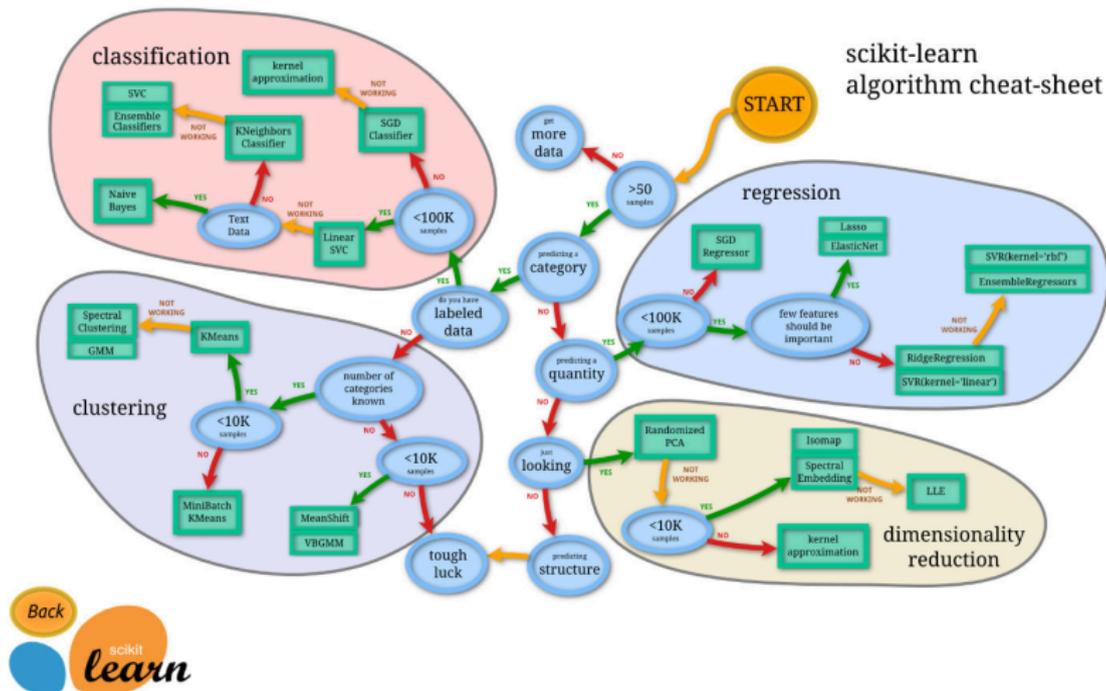
# Physical environment



# Business, finance



# Standard toolkit



[sklearn]

# Typical process application

- Liquefied natural gas (LNG)
  - Source gas → Liquid for transport → Gas for use
- World-wide facilities
  - Brunei, Oman, Nigeria, Australia, Qatar, Russia, Trinidad&Tobago, Egypt
- Nigerian plant
  - Six processing units
  - Total processing capacity: 22 million tonnes of LNG p.a.
  - Accounted for approximately 7% of global LNG supply in 2017

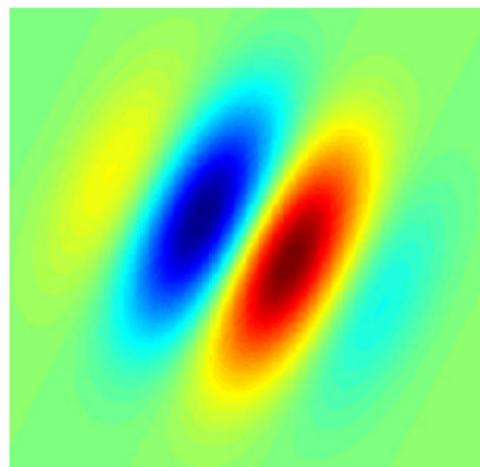
(Information from [www.shell.com](http://www.shell.com))

# Typical process application

- System: input  $\rightarrow$  process  $\rightarrow$  output
  - Flows, Temperatures, Pressures, Compositions, Valve settings, Level settings
  - Recycles, automatic control, constraints
- Data
  - Multivariate time-series,  $\log_{10}(p) \in (2, 4)$ ,  $\mathbf{X}(t), \mathbf{Y}(t), t \in \mathcal{T}$
  - Some of  $\mathbf{X}$  manipulable ( $\mathbf{X}'$ ), others not ( $\mathbf{X}''$ )
  - $\mathcal{T} \approx 5$  years, sampling  $\approx 1$  Hz,  $\log_{10}(n) \in (3, 8)$
- Goal: optimisation
  - $f_1(\mathbf{X}'_{\mathcal{T}}, \mathbf{X}''_{\mathcal{T}}) = 0$ ,  $\mathbf{Y}_{\mathcal{T}} = f_2(\mathbf{X}'_{\mathcal{T}}, \mathbf{X}''_{\mathcal{T}})$ ,  $\mathbf{X}''_{\mathcal{T}} = f_3(\mathbf{X}'_{\mathcal{T}})$
- Issues:
  - Data quality (redundancy, reconciliation)
  - Complex time-series dependence structure (models are 'static')
  - Existing models at play (APC, engineering:  $f = f_{\text{ENG}} + f_{\text{Emp}}$ )
  - In-situ performance
- Most popular models (linear & regularised regression)

# Computer vision, NLP

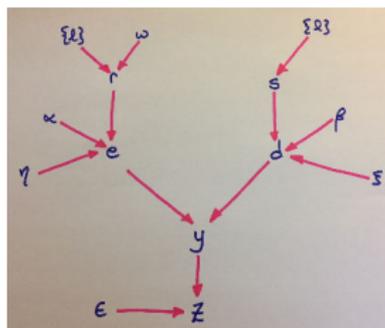
- Raw input (images, text, sound)  
 $X_0$
- Multiple filters  $X = F(X_0, \theta)$
- Inference on processed  $X$
- Optimal choice of filter (kernels) for prediction
  
- Sparse + low rank
- Higher criticism (Donoho)



Gabor filter [wiki]

(More illustrations at [www.shell.ai](http://www.shell.ai))

# Uncertainty quantification, emulation



A simple system model

- Flexible framework, Bayes linear
- Optimal design
- Probabilistic ODEs, Bayesian optimisation

$$\text{Obs} : z(x) = y(x) + \epsilon$$

$$\text{Sys} : y(x) = e(x) + d(x)$$

$$\text{Emul} : e(x) = \alpha'g(x) + r(x, \omega) + \eta$$

$$\text{Disc} : d(x) = \beta'h(x) + s(x) + \xi$$

$e$  : emulator or 'process' model

$d$  : discrepancy model

$g, h$  : bases for covariate space

$r, s$  : Gaussian process residuals

Priors : all Gaussian

Data : emulator  $E$ , measured  $Z$

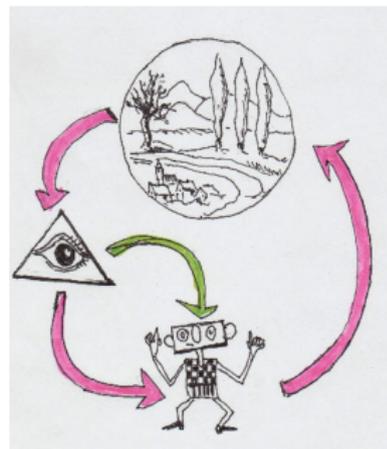
Estimation :  $f(\alpha, \beta, \{l_r\}, \{l_s\}, \omega | E, Z)$

Prediction :  $f(y(x) | E, Z)$

# Robotics, automatic decisions

Markov decision = 'Reinforcement learning'

- 'Agent' observes state  $s_t$
- Takes action  $a_t$
- Leads to state  $s_{t+1}$  with  $\mathbb{P} p(s_{t+1}|s_t, a_t)$
- Agent receives reward  $r_t = r(s_{t+1}, a_t, s_t)$
- Policy  $\pi(a|s)$  governs behaviour
- History  $h = [s_1, a_1, s_2, a_2, \dots, s_T, a_T]$
- Discounted return
 
$$R(h) = \sum_t \gamma^t r_t, \gamma \in [0, 1)$$
- $p_\pi(h) = [\prod_t p(s_{t+1}|s_t, a_t) \pi(a_t|s_t)] p(s_1)$
- Optimal  $\pi^* = \operatorname{argmax}_\pi \mathbb{E}_{p_\pi(h)} R(h)$



- Boston Dynamics (dogs opening door, robot somersault; YouTube)

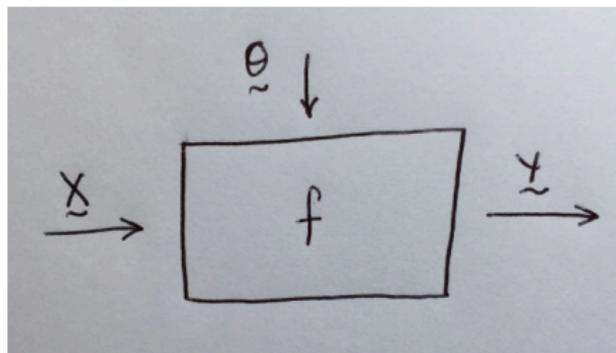
# Assessing empirical models

# Fundamentals

- Think stand-alone empirical modelling (not robotics)
- Data not from a designed experiment
  - Interpolation (failure to fit)
  - Extrapolation (failure to generalise)
  - Curse of dimensionality: often can't tell one from the other
- Data preparation and preprocessing
- Naive presumptions
  - IID, ignore dependence structure, time-series!
- Use cases
  - Non-expert users
  - Automatic, embedded predictions

# Assessing predictive algorithms

- $Y = f(X, \theta)$
- Black-box
  - Manipulate inputs,  $X$
  - Manipulate some tuning parameters,  $\theta$
- Predictive performance
- Prediction uncertainty
- Exceptionality



# In-vitro performance

- Competitions
  - e.g. Makridakis (M3, M4) forecasting, Kaggle, NN3, imagenet
  - e.g. Fernandez-Delgado (~ 200 classifiers), OpenML (~ 100 classifiers)
  - e.g. MLaut regression and classification (Kiraly), AutoML (Freiburg)
- Findings for time-series (M3) and general prediction (Kiraly): ensemble methods best, GPs/SVMs close, ML (inc. 'deep learning') poor
- Lazy focus on 'predictive  $R^2$ ' not prediction uncertainty

Table 10. Features of various Artificial Intelligence (AI) applications.

Type of Application	Rules are known and do not change	The environment is known and stable	Predictions can influence the future	Extent of Uncertainty (or amount of noise)	Examples
Games	Yes	Yes	No	None	Chess, GO
Image and speech recognition	Yes	Yes	No	Minimal (can be minimized by big data)	Face Recognition, Siri, Cortana, Google AI
Predictions based on the Law of large numbers	Yes	Yes	Minimally	Measurable (Normally distributed)	Forecasting the sales of beer, coffee, soft drinks, weather etc.
Autonomous Functions	Yes	Yes	No	Can be assessed and minimized	Self-Driving Vehicles
Strategy, Competition, Investments	No	No	Yes, often to a great extent	Cannot be measured (fat tails)	Decisions, Anticipations, Forecasts
Combinations of the above	It may be the ultimate challenge moving towards GAI (General AI) but also increasing the level of complexity and sophistication of algorithms				Eventually it can cover everything

# In-vivo performance

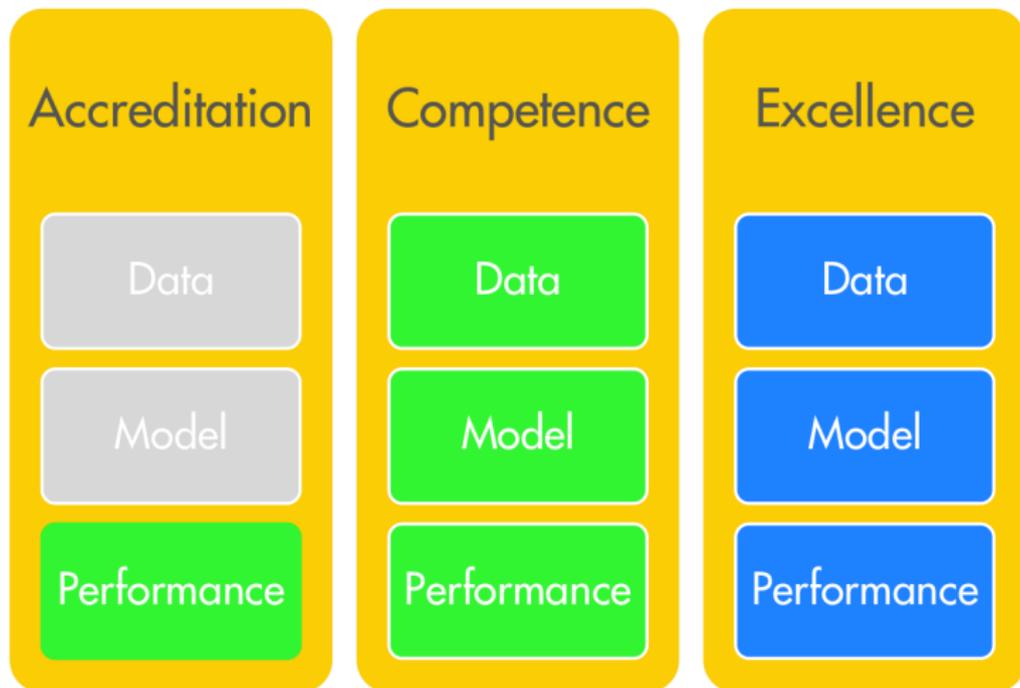
- Interpretability
  - Functional form for  $\hat{f}$
  - Variable importance (easy for regression, trees; difficult for others)
- Parsimony
  - Variable selection not regularisation
- 'Human in the loop'
  - 'What legal scholars should learn about machine learning' (Lehr, Ohm)
  - Algorithms making decisions with no human supervision
  - Business- and safety-critical
  - Reinforcement learning of 'human experience' (Faisal)
- High level architectures enable but constrain
  - e.g. Azure DataBricks (database)
  - e.g. TensorFlow (images)
- Solution maintenance
  - Version control
  - Documentation

# Performance, uncertainty, exceptionality

- Predictive performance
  - Cross-validation,  $\hat{\theta} = \operatorname{argmin}_{\theta} L(\theta)$
  - Preserve dependence, choice of partition, bias, nested
  - Problem : does not provide  $\operatorname{var}(\hat{\theta}) \Rightarrow$  use BS also
- Prediction uncertainty
  - Bootstrapping (BS),  $p_{\text{BS}}(\hat{\theta})$
  - Preservation of dependence structure
  - Problem : not inherently out-of-bag  $\Rightarrow$  use CV also
- Exceptionality
  - Randomised permutation (RP) testing
  - $E = \mathbb{P}_{p_{\text{RP}}} [L(\hat{\theta}) \leq L(\hat{\theta}_{\text{RP}})]$
- Model-agnostic assessment ('wrapper')
  
- c.f. Bayes: posterior, evidence, posterior predictive, estimate or minimise generalisation loss
- Huge model-specific applied literature
- Conformal inference (CMU)

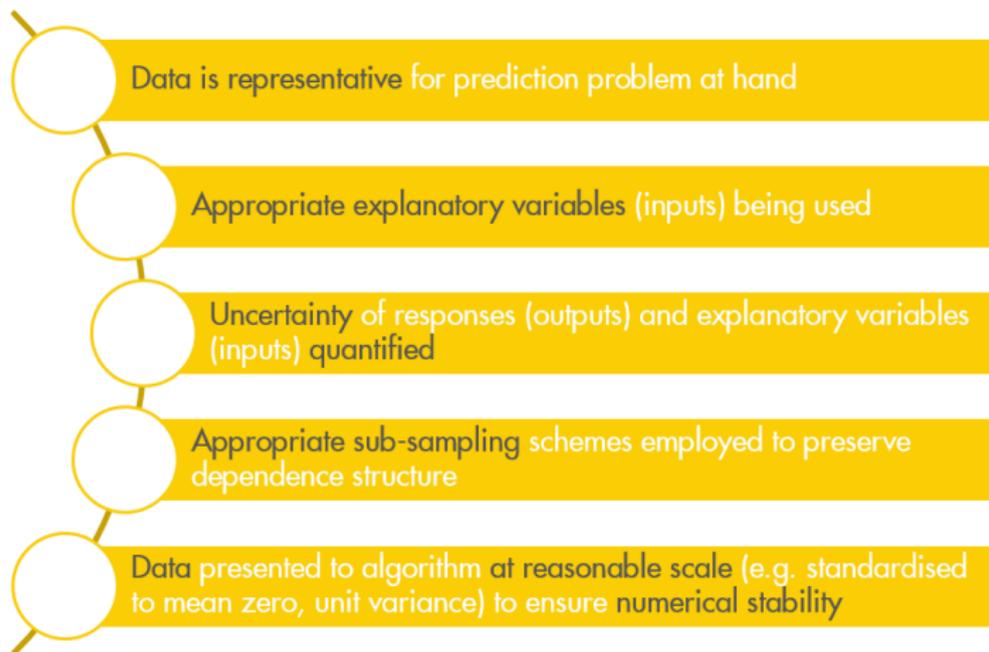
# Assurance in data science

# Fundamentals



Accreditation, Competence, Excellence

# Data

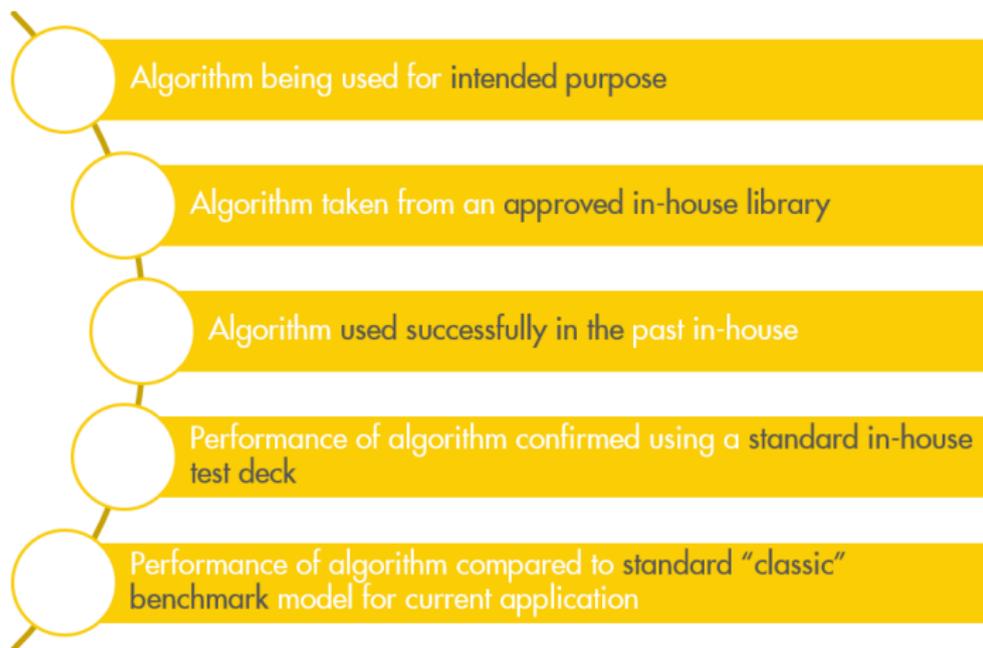


## Data

### Assessing the data

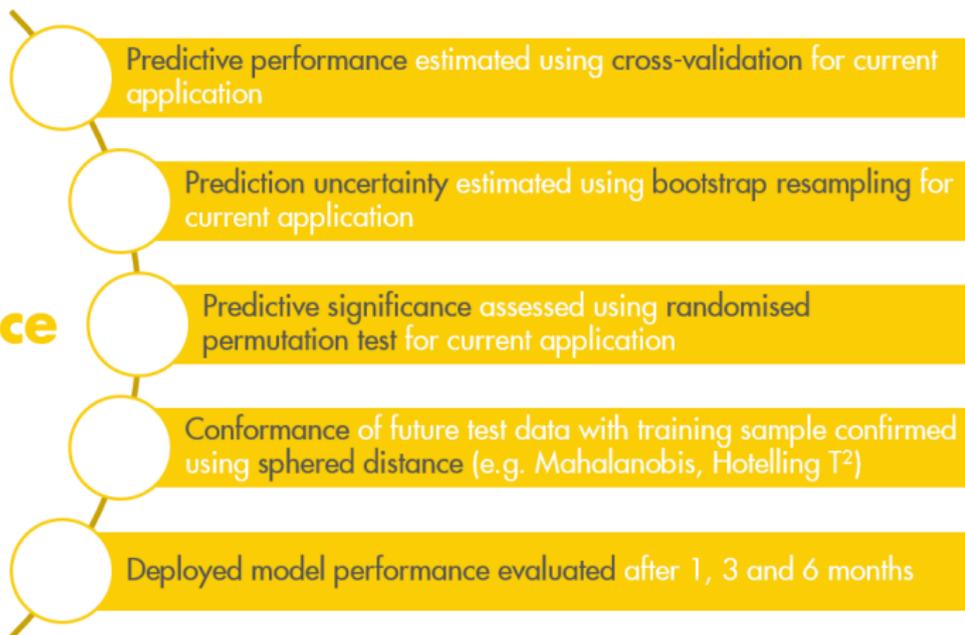
# Algorithm

## Algorithm



Assessing the algorithm

# Performance



Assessing performance

# Long-term

- Good governance
- Standardisation
  - Approved algorithms
  - Test decks
  - Reporting framework
- Continuous improvement
  - Monitoring of applications in development
  - Monitoring of deployments
- Learning
  - Data scientists
  - Expert reviewers
- Impact from data science

# Reasons for concern, optimism

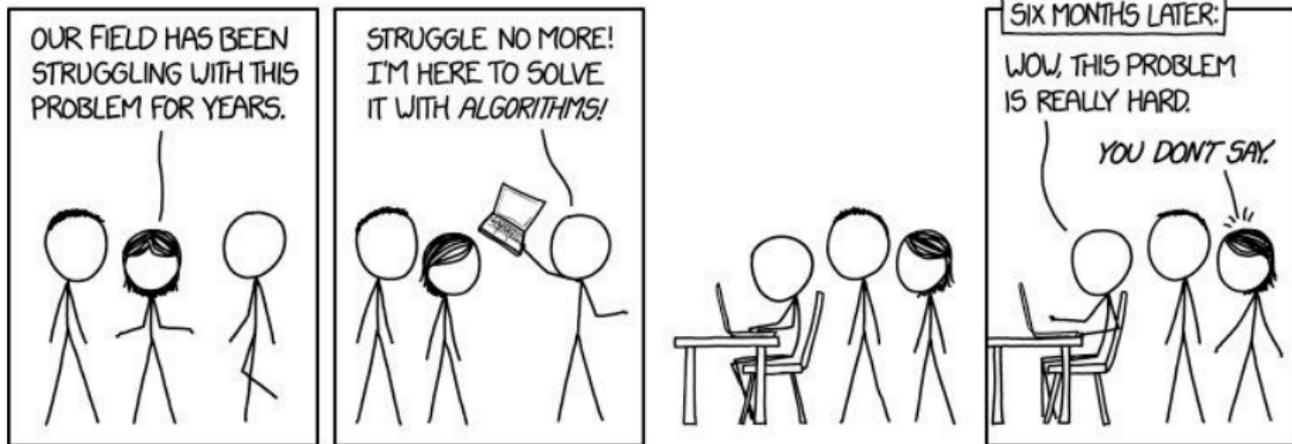
# Concern

- Data-driven models everywhere
  - Don't collect data until you know why you're collecting it
- Bias and uncertainty
- 'Lazy' application, just because it's possible
  - Current criticism of 5G
- Intrusive application
  - Spam emails, phone-calls, inadvertent discrimination
  - Snooping by third parties, phishing, cameras, voice monitoring
  - Terrorism, autonomous weapons
- Erroneous application
  - Boeing 737Max8
  - Autonomous vehicles
- Uncontrolled application
  - 'Black Mirror', AI / AGI takes over (Berners-Lee on www: 'downward plunge to dysfunctional future')
  - Nascent ethical, legal, governmental frameworks
  - Moral choices, dilemmas; intelligence vs. wisdom

# Optimism

- Data-driven models everywhere
  - Broad range of application complexity
- Bias and uncertainty
  - What statisticians and data scientists know how to do well
- Game-changing
  - Better science
  - Better systems (bigger, faster, more accurate, more integrated, more efficient, less wasteful)
  - Better management of global resources (food, water, energy, people time, personalised health care)
  - Healthier, longer, more enriching lives (Morris, ‘News from nowhere’)

Hofstadter's Law: It always takes longer than you expect, even when you take into account Hofstadter's Law



[indico.cern.ch]

Diolch yn fawr!

# Backup

# Regression

Linear regression:  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$

- $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I})$ , prior:  $\boldsymbol{\beta} \sim N(0, \boldsymbol{\Sigma})$ ,  $\mathbf{X}$  is  $n \times p$
- Negative log posterior:  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^2 / \sigma^2 + \boldsymbol{\beta}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}$
- $\mathbb{E}[\boldsymbol{\beta} | \mathbf{y}] = (\mathbf{X}' \mathbf{X} / \sigma^2 + \boldsymbol{\Sigma}^{-1})^{-1} \mathbf{X}' \mathbf{y} / \sigma^2$
- Prediction:  $\mathbb{E}[\mathbf{X}_* \boldsymbol{\beta} | \mathbf{y}] = \mathbf{X}_* (\mathbf{X}' \mathbf{X} / \sigma^2 + \boldsymbol{\Sigma}^{-1})^{-1} \mathbf{X}' \mathbf{y} / \sigma^2 = \mathbf{S} \mathbf{y}$
- Regression is a **smoother**, invert  $p \times p$  matrix  $\mathbf{X}' \mathbf{X}$
- But:  $(\mathbf{X}' \mathbf{X} / \sigma^2 + \boldsymbol{\Sigma}^{-1}) \boldsymbol{\Sigma} \mathbf{X}' = \mathbf{X}' (\mathbf{X} \boldsymbol{\Sigma} \mathbf{X}' + \sigma^2 \mathbf{I}) / \sigma^2$
- Prediction:  $\mathbb{E}[\mathbf{X}_* \boldsymbol{\beta} | \mathbf{y}] = \mathbf{X}_* \boldsymbol{\Sigma} \mathbf{X}' (\mathbf{X} \boldsymbol{\Sigma} \mathbf{X}' + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$
- Invert  $n \times n$  matrix  $\mathbf{X} \boldsymbol{\Sigma} \mathbf{X}'$ .  $n$  can be big,  $> 10^5$

Kernel regression:  $\mathbf{y} = \Phi(\mathbf{X})' \boldsymbol{\beta} + \boldsymbol{\epsilon}$

- Inner product  $\mathbf{X}_* \boldsymbol{\Sigma} \mathbf{X}'$  generalised to  $\Phi(\mathbf{X}_*)' \boldsymbol{\Sigma} \Phi(\mathbf{X})$ , where  $\Phi(\mathbf{X})$  is a **kernel function**

# Kernel

- **Kernel trick:** Never write kernel functions explicitly. Only need inner-product expression  $k(\mathbf{X}_*, \mathbf{X}) = \Phi(\mathbf{X}_*)' \Sigma \Phi(\mathbf{X})$
- Arbitrary (smooth) relationships can be estimated using kernel regression
- Product form:  $k(\mathbf{X}_*, \mathbf{X}) = \prod_{j=1}^p k_j(\mathbf{X}_{*j}, \mathbf{X}_j)$  for covariate  $j$
- Squared exponential:  $k_j(x_{*hj}, x_{ij}) = \exp(-(x_{*hj} - x_{ij})^2 / (2\ell_j^2))$
- Correlation lengths:  $\{\ell_j\}_{j=1}^p$  must be estimated using cross-validation. Also estimate nugget variances
- Slick leave-one-out cross-validation available (cf 'hat matrix')
- Slick Kronecker product form (fast inverse!) when  $\mathbf{X}$  is of blocked form:  $k(\mathbf{X}, \mathbf{X}) = k(\mathbf{X}_1, \mathbf{X}_1) \otimes k(\mathbf{X}_2, \mathbf{X}_2)$