

Buzzwords surrounding Data Science



Eiko Yoneki

eiko.yoneki@cl.cam.ac.uk

<http://www.cl.cam.ac.uk/~ey204>

*Systems Research Group
University of Cambridge Dept. Computer Science and Technology
Computer Laboratory*

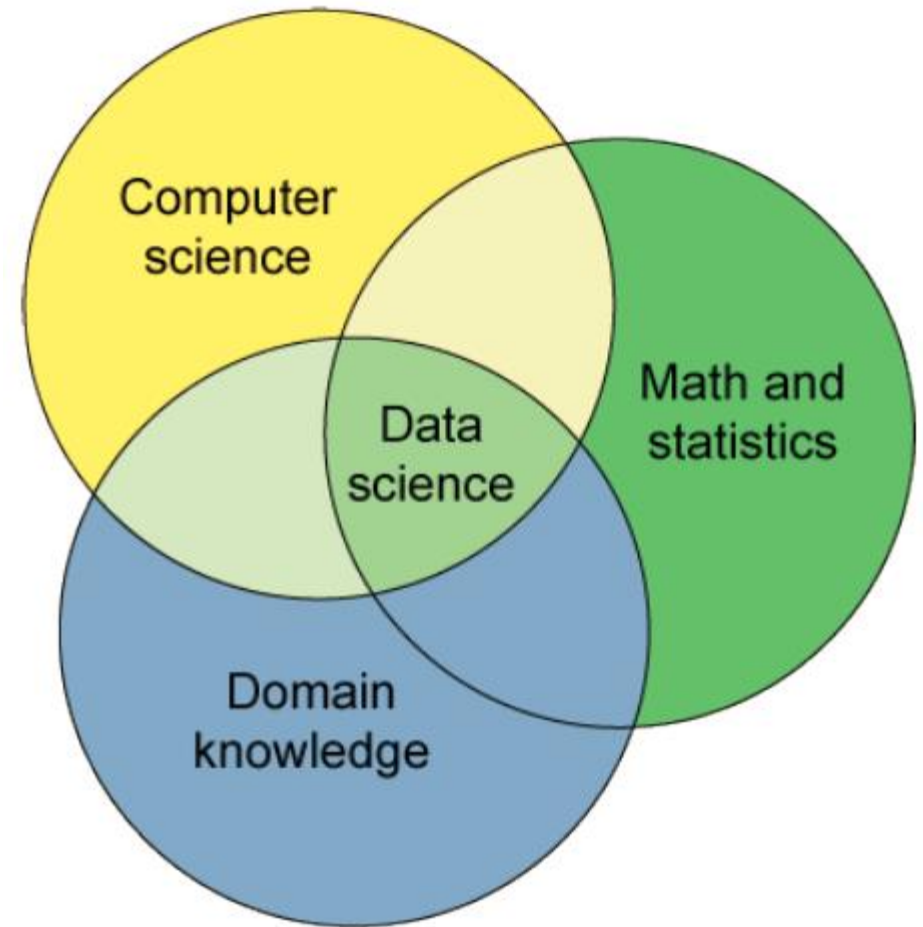


Outline

- Data Science Community
- Landscape of Research in Alan Turing Institute
 - Probabilistic Programming
 - Optimisation
 - Ethics
 - Frameworks
- Becoming Data Scientist?
- Fill the Gap between Research and Practice

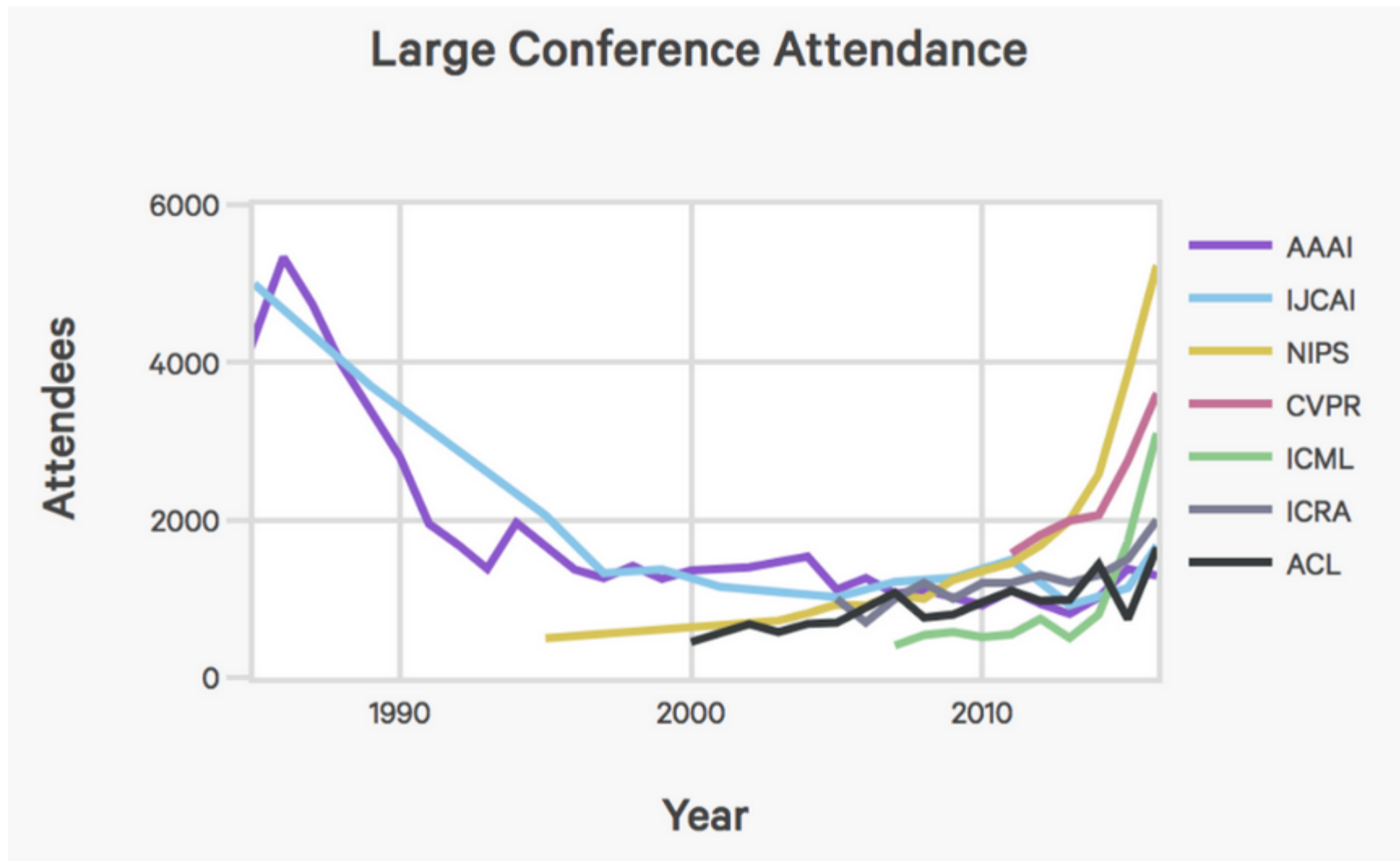
Data Science: Any new intellectual content?

- What does it mean to Computer Science?
- 1970's: EE + Math → Computer Science
- 2010's: CS + Stats + ?? → Data Science
- Is something fundamental emerging here?
- Data Science is a very broad discipline
- Data Science PhD?
 - PhD normally with a narrow field with depth...



based on Drew Conway, NYU

Scale of Community Size in ML/AI





NIPS: 8000 Attendees in 2017

- **Randomness of Paper acceptance?**
- In 2016, 2,406 papers were submitted and 568 were accepted for a 24% acceptance rate. In 2017, 679 papers out of 3,240 submitted were accepted for a 21% acceptance rate.
- In 2014, Corinna Cortes and Neil Lawrence ran the NIPS experiment where 1/10th of papers submitted to NIPS went through the NIPS review process twice, and then the accept/reject decision was compared.
- The 26% disagreement: Given 22% acceptance rate, The immediate implication is that between 1/2 and 2/3 of papers accepted at NIPS would have been rejected if reviewed a second time.
- <http://blog.mrtz.org/2014/12/15/the-nips-experiment.html>
- In particular, about 57% of the papers accepted by the first committee were rejected by the second one and vice versa. In other words, most papers at NIPS would be rejected if one reran the conference review process (with a 95% confidence interval of 40-75%).



SysML Conference spawn in 2018-2019

- SysML is a conference targeting research at the intersection of systems and machine learning
- Aims to elicit new connections amongst these fields, including identifying best practices and design principles for learning systems, as well as developing novel learning methods and theory tailored to practical machine learning workflows

Steering Committee

Jennifer Chayes

Bill Dally

Jeff Dean

Michael I. Jordan

Yann LeCun

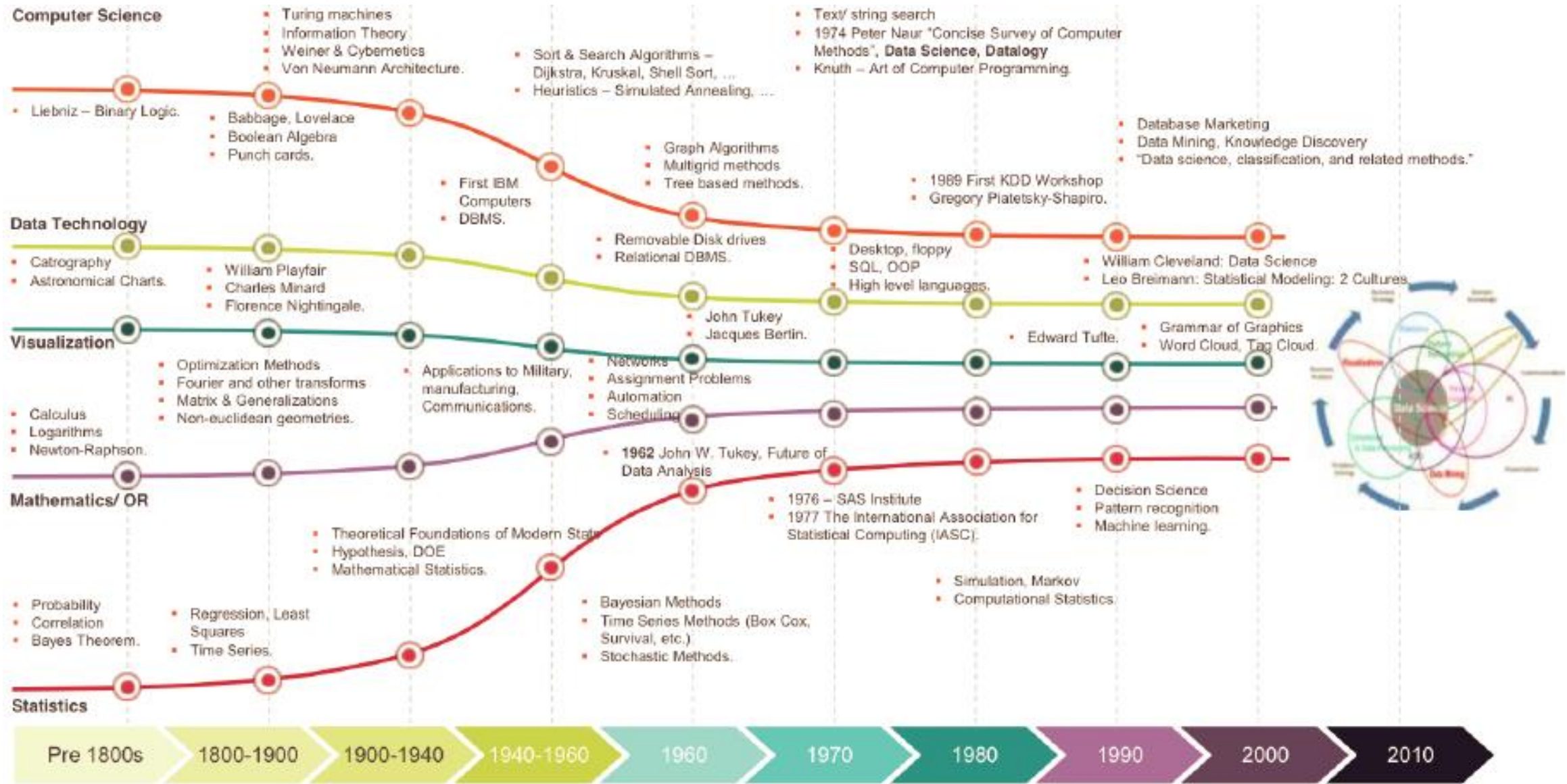
Fei-Fei Li

Alex Smola

Dawn Song

Eric Xing

Trajectory

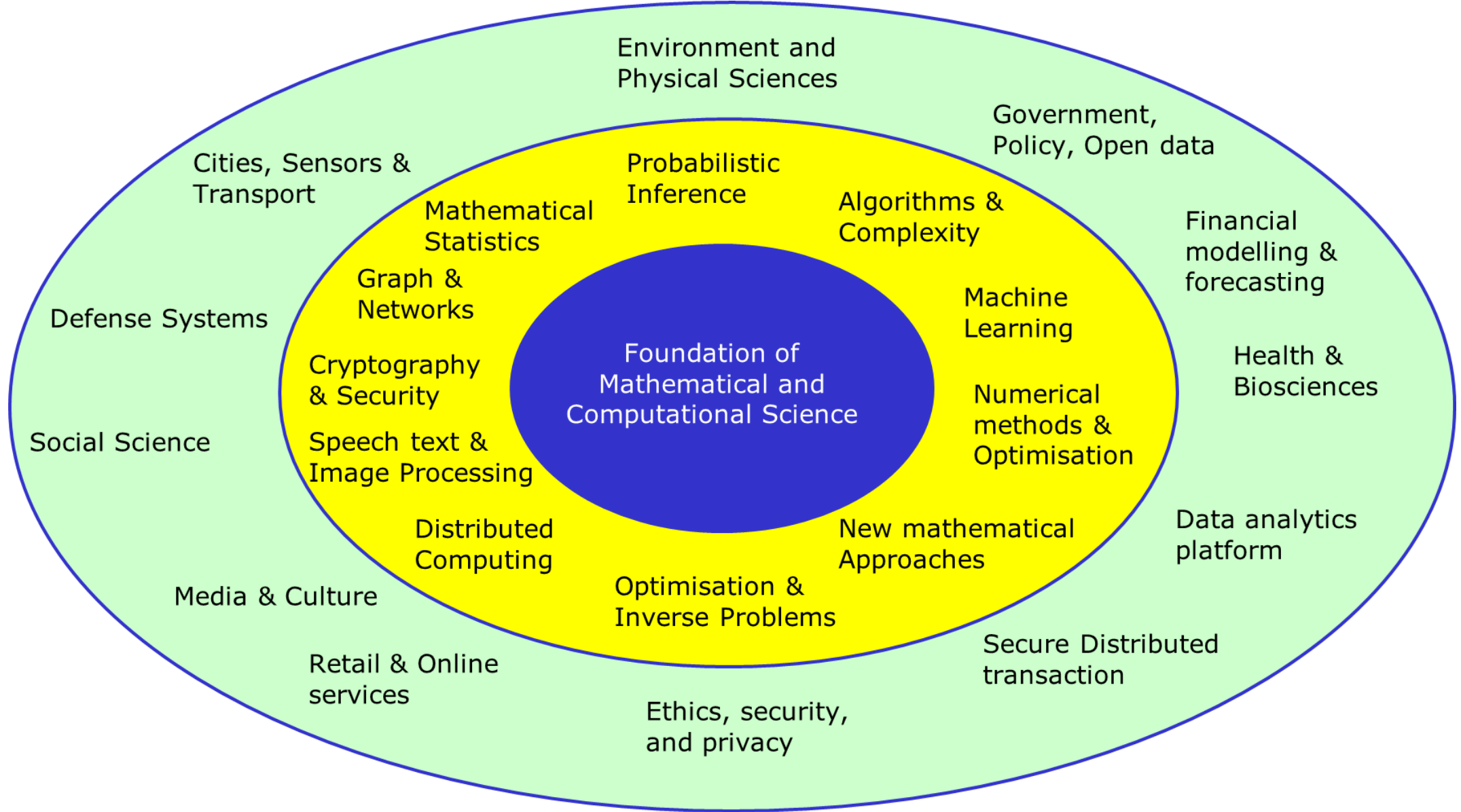




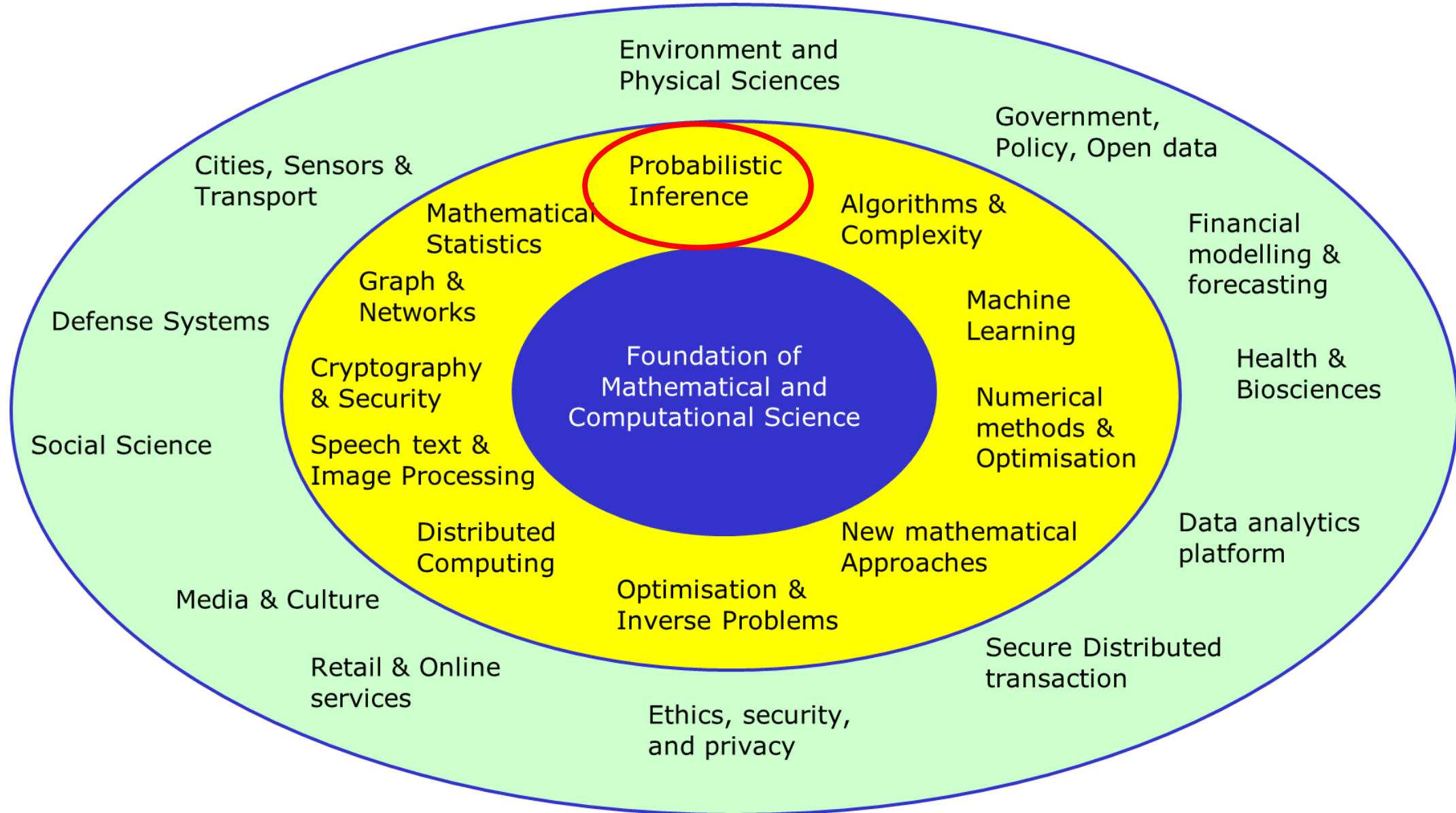
Alan Turing Institute (ATI)

- Established in 2015 in London as a National Institute for Data Science
- >£20M Capital Investment from Government
- Originally 5 Universities formed core body (UCL, Warwick, Edinburgh, Oxford and Cambridge) and now expanded to 13 universities
- Goal: Data Science and after 2018 Artificial Intelligence
- Translating output into practice

Broad Landscape of Research in ATI in 2015



Broad Landscape of Research in ATI in 2015



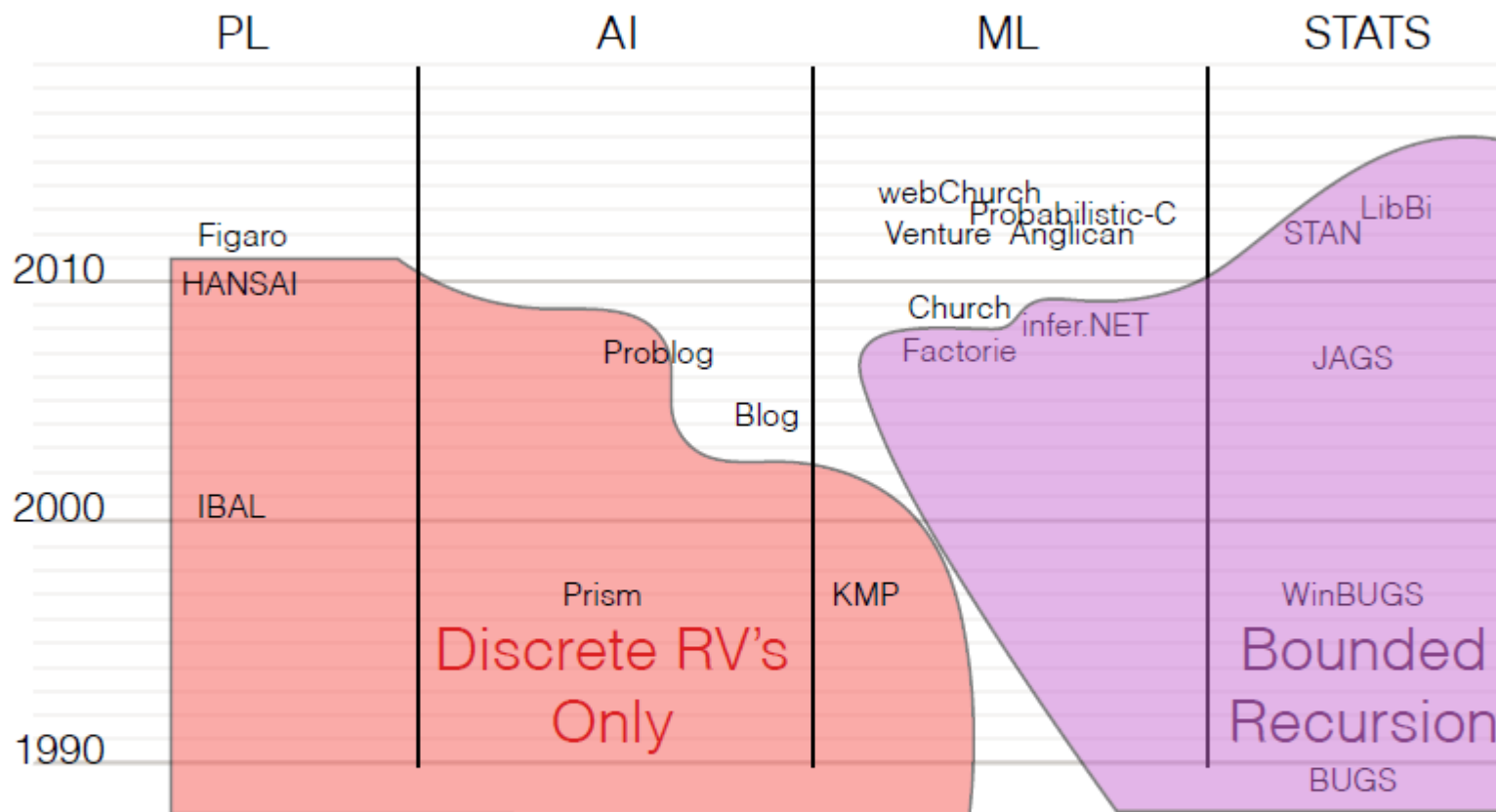


Probabilistic Model

- Probabilistic models incorporate random variables and probability distributions into the model
 - Deterministic model gives a single possible outcome
 - Probabilistic model gives a probability distribution
- Used for various probabilistic logic inference (e.g. MCMC-based inference, Bayesian inference...)



Probabilistic Programming



Edward based on Python

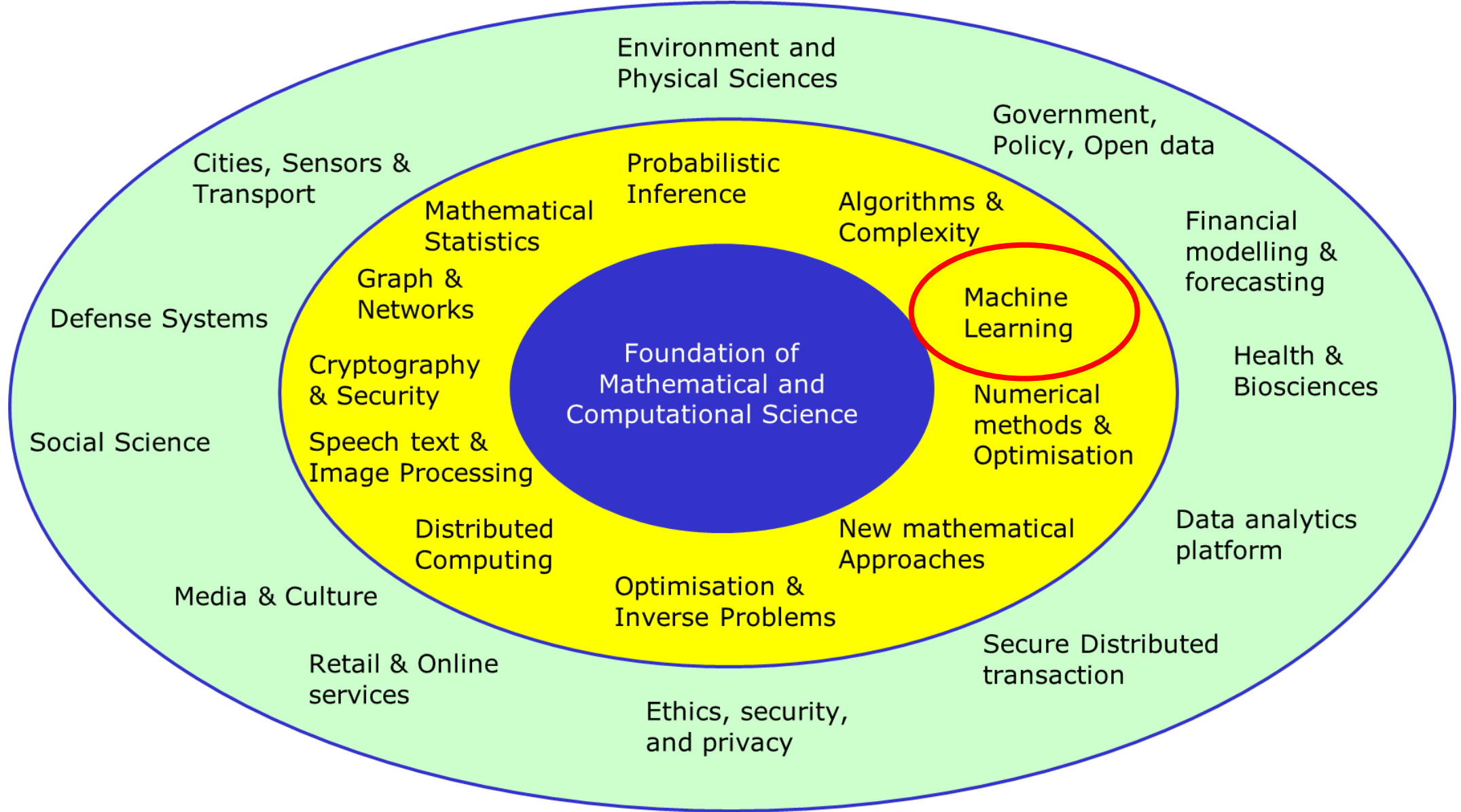
Probabilistic C++

Simula

Prolog

Improbable - Java version

Broad Landscape of Research in ATI in 2015



4 Great Pictures Illustrating Machine Learning Concepts

- Neural Networks: The Backpropagation algorithm
- Cheat Sheet on Probability
- 24 Neural Network Adjustments
- Matrix Multiplication in Neural Networks

See <https://www.datasciencecentral.com>.

Neural Networks: The Backpropagation algorithm

Backpropagation Algorithm

Before Weight Adjustment

Parameters

For $w_2 = 5 \wedge x = 2 \wedge w_1 = 3.5$

Where $MAE_1 = w_1x - y \wedge MAE_2 = w_2x - y$

For $w_2x = 10 \wedge w_1x = 7 \wedge y = 4$

$$f(x) = \frac{MAE_2^2}{2}$$

$$g(x) = \frac{MAE_1^2}{2}$$

Backpropagation of Error = $g'(x)$

Chain Rule

$$\frac{\partial g}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x} = g'(x) = g'(f(x)) \cdot f'(x)$$

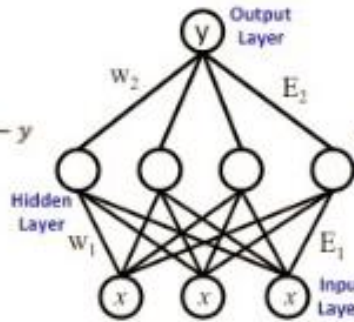
$$\frac{\partial f}{\partial x} = f'(x) = 2 \cdot \frac{E_2}{2} = E_2 = w_2x = 10$$

$$\frac{\partial g}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x} = g'(x) = g'(f(w_1x - y)) \cdot f'(x)$$

$$g'(x) = g'(f(3)) \cdot 10$$

$$g'(x) = \frac{1}{2} \cdot 3^2 \cdot 10 = \frac{90}{2} = 45$$

Derivative of error



After Weight Adjustment

Weight Adjustment

Adjust w_2 from 5 to 4 \wedge w_1 from 3.5 to 2.5

Goal is to decrease derivative of error $g'(x) \rightarrow 0$

For $w_2x = 8 \wedge w_1x = 5 \wedge y = 4$

$$f(x) = \frac{E_2^2}{2}$$

$$\frac{\partial f}{\partial x} = f'(x) = 2 \cdot \frac{E_2}{2} = E_2 = w_2x = 8$$

$$g(x) = \frac{E_1^2}{2}$$

Backpropagation of Error = $g'(x)$

$$\frac{\partial g}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x} = g'(x) = g'(f(x)) \cdot f'(x)$$

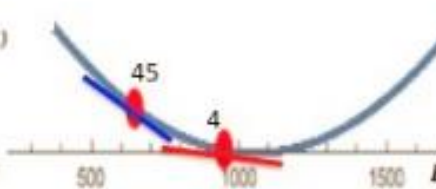
$$g'(x) = g'(f(w_1x - y)) \cdot f'(x)$$

$$g'(x) = g'(f(1)) \cdot 8$$

Derivative of error

$$g'(x) = \frac{1}{2} \cdot 1^2 \cdot 8 = \frac{8}{2} = 4$$

After Backprop



Rubens
Zimbres

Cheat Sheet on Probability

Rubens Zimbres

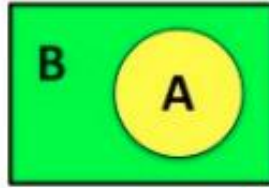
Probability

Marginal Probability

long hair

$$\sum Prob = 1 \quad P(A) = \frac{P(A)}{\sum P(A, B)}$$

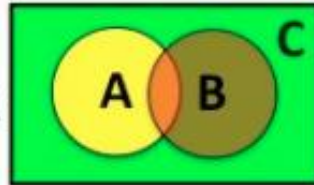
$$0 < Prob < 1 \quad P(\bar{A}) = 1 - A$$



Conditional Probability (Bayes)

long hair, given that is woman

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{P(A \cap B)}{P(B)}$$



Independent events

coins

$$P(A \cap B) = P(A) \cdot P(B)$$



Dependent events

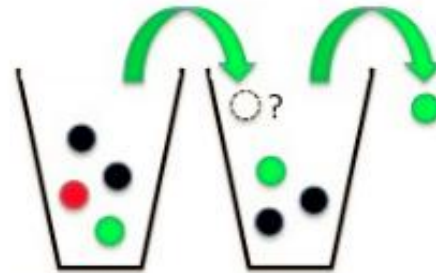
cards

$$P(A \cap B) = P(A) \cdot P(B|A)$$

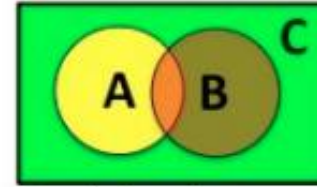
Total Probability

jar

$$P(\text{2nd Green}) = P(\text{Green|1st Black}) + P(\text{Green|1st Green}) + P(\text{Green|1st Red})$$



Joint Probability



long hair and woman

$$P(A \cap B) = P(A) \cdot P(B)$$

long hair or woman

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

not long hair and not woman

$$P(\bar{A} \cap \bar{B}) = 1 - P(A) \cdot P(B)$$

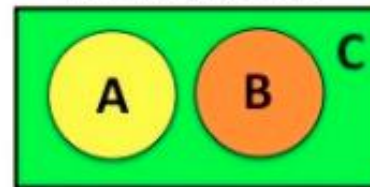
neither long hair nor woman

$$P(\overline{A \cup B}) = 1 - (P(A) + P(B) - P(A \cap B))$$

Disjoint Probability

Mutually Exclusive

weather and coins



$$P(A \cap B) = \{ \}$$

$$P(A \cup B) = P(A) + P(B)$$

24 Neural Network Adjustments

ARCHITECTURE

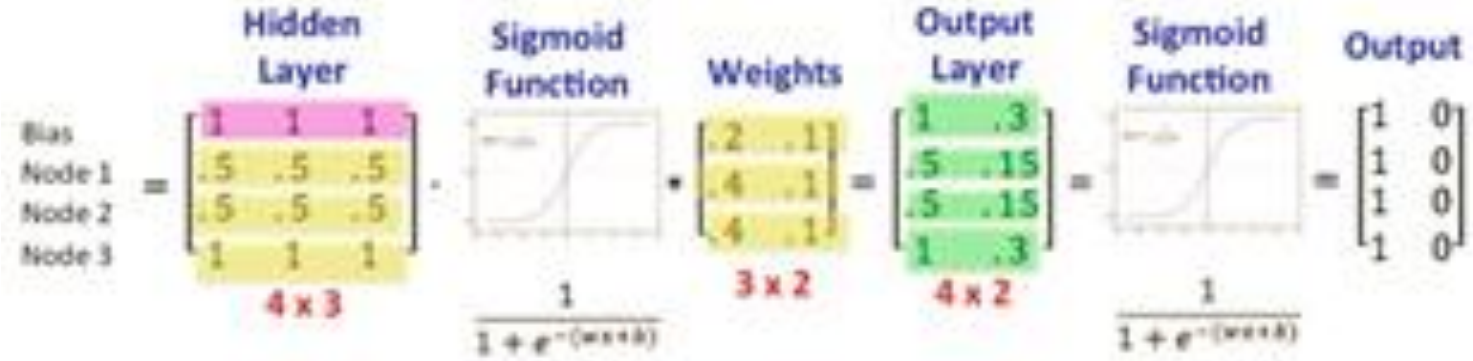
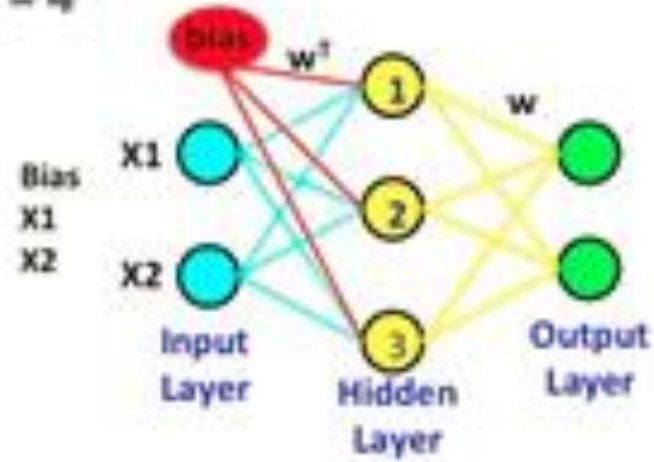
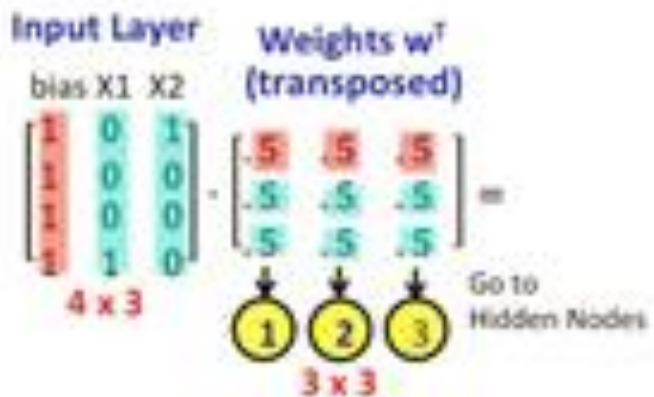
- Variables type
- Variable scaling
- Cost function
- Neural Network type:
 - RBM,FFN,CNN,RNN...
- Number of layers
- Number of hidden Layers
- Number of nodes
- Type of layers:
 - LSTM, Dense, Highway
 - Convolutional, Pooling...
- Type of weight initialization
- Type of activation function
 - Linear, sigmoid, relu...
- Dropout rate (or not)
- Threshold

HYPERPARAMETER TUNING

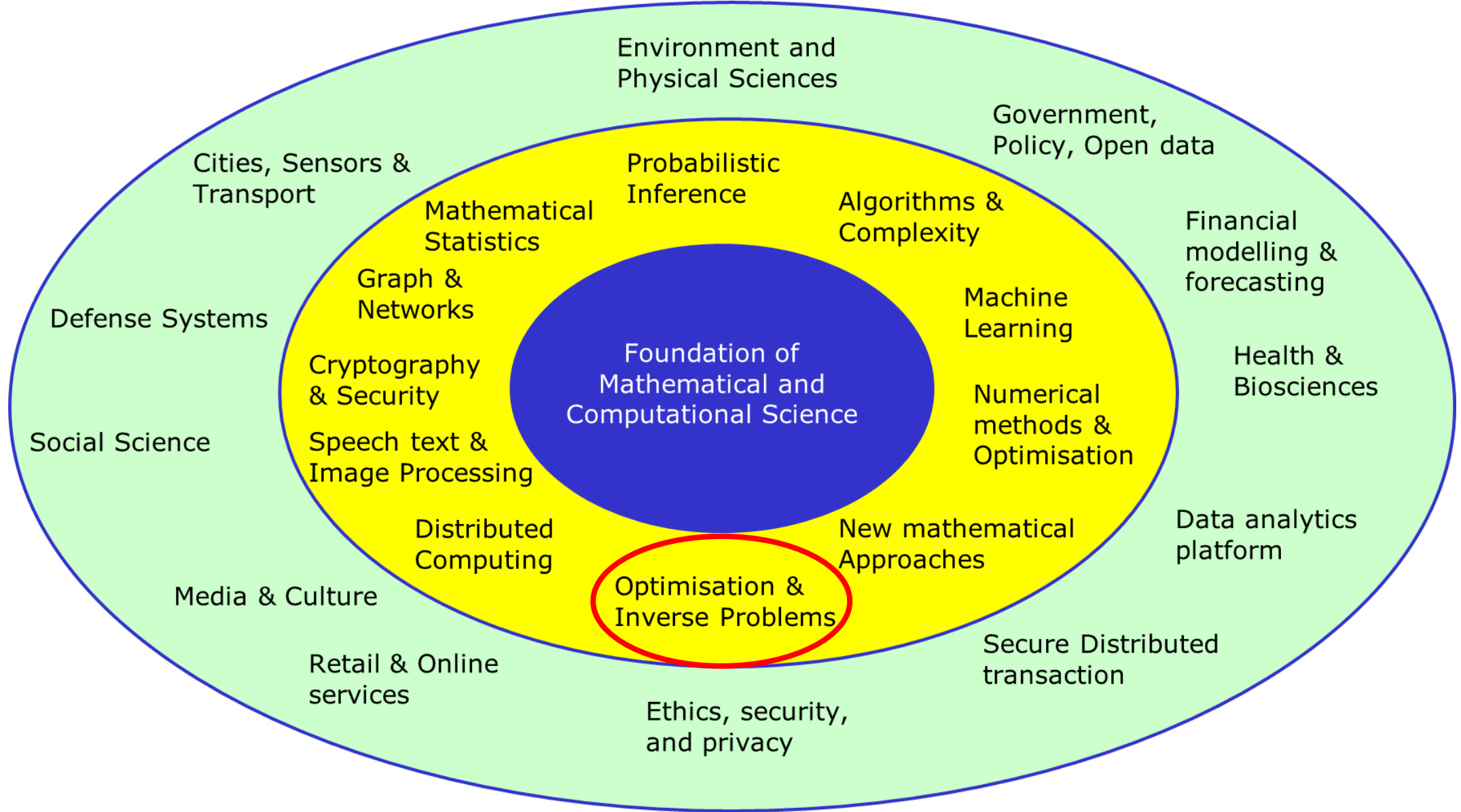
- Type of optimizer
- Learning rate (fixed or not)
- Regularization rate (or not)
- Regularization type: L1, L2, ElasticNet
- Type of search for local minima:
 - Gradient descent, simulated
 - annealing, evolutionary...
- Batch size
- Nesterov momentum (or not)
- Decay rate (or not)
- Momentum (fixed or not)
- Type of fitness measurement:
 - MSE, accuracy, MAE, cross-entropy,
 - precision, recall
- Epochs
- Stop criteria

Matrix Multiplication in Neural Networks

Color Guided Matrix Multiplication for a Binary Classification Task with N = 4



Broad Landscape of Research in ATI in 2015



Tuning Computer Systems

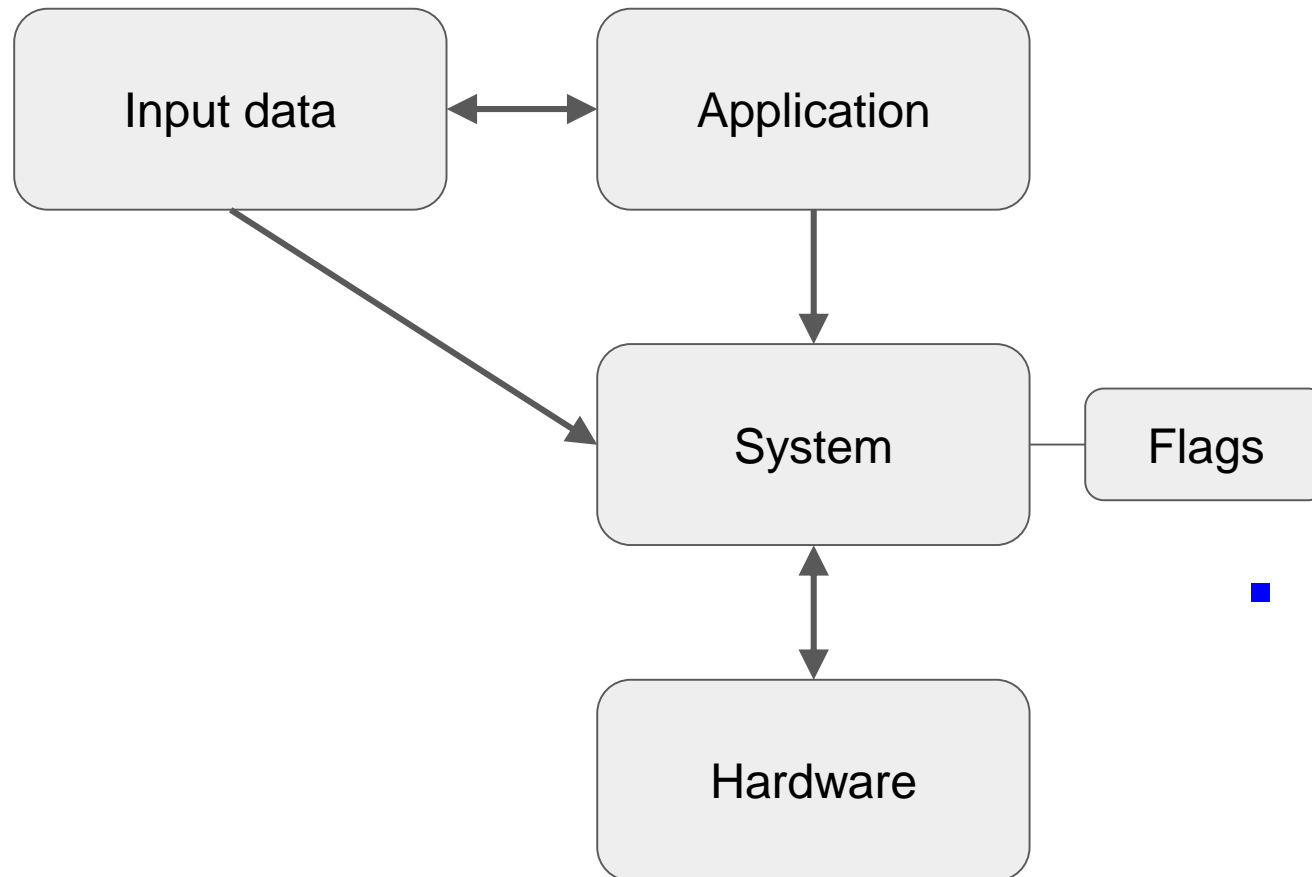
- Complex configuration parameter space and increasing number of parameters
- Configurations need tuning to optimise resource utilisation, minimise cost
- Containerisation means software is frequently restarted and deployed in new contexts
- Hand-crafted solutions impractical, often left static or configured through extensive offline analysis
- Auto-tuning: automatically determine parameters, e.g. resources for Hadoop job

Optimisation of Complex Data Processing in Computer Systems

- Auto-tuning to deal with complex parameter space using machine-learning
 - Structured Bayesian Optimisation, Reinforcement Learning
 - Build a solid auto-tuning platform in a complex and large parameter space
- e.g. Tuning Cluster task scheduling, ML framework, JVM garbage collector, NN model, Compiler, DB indexing...

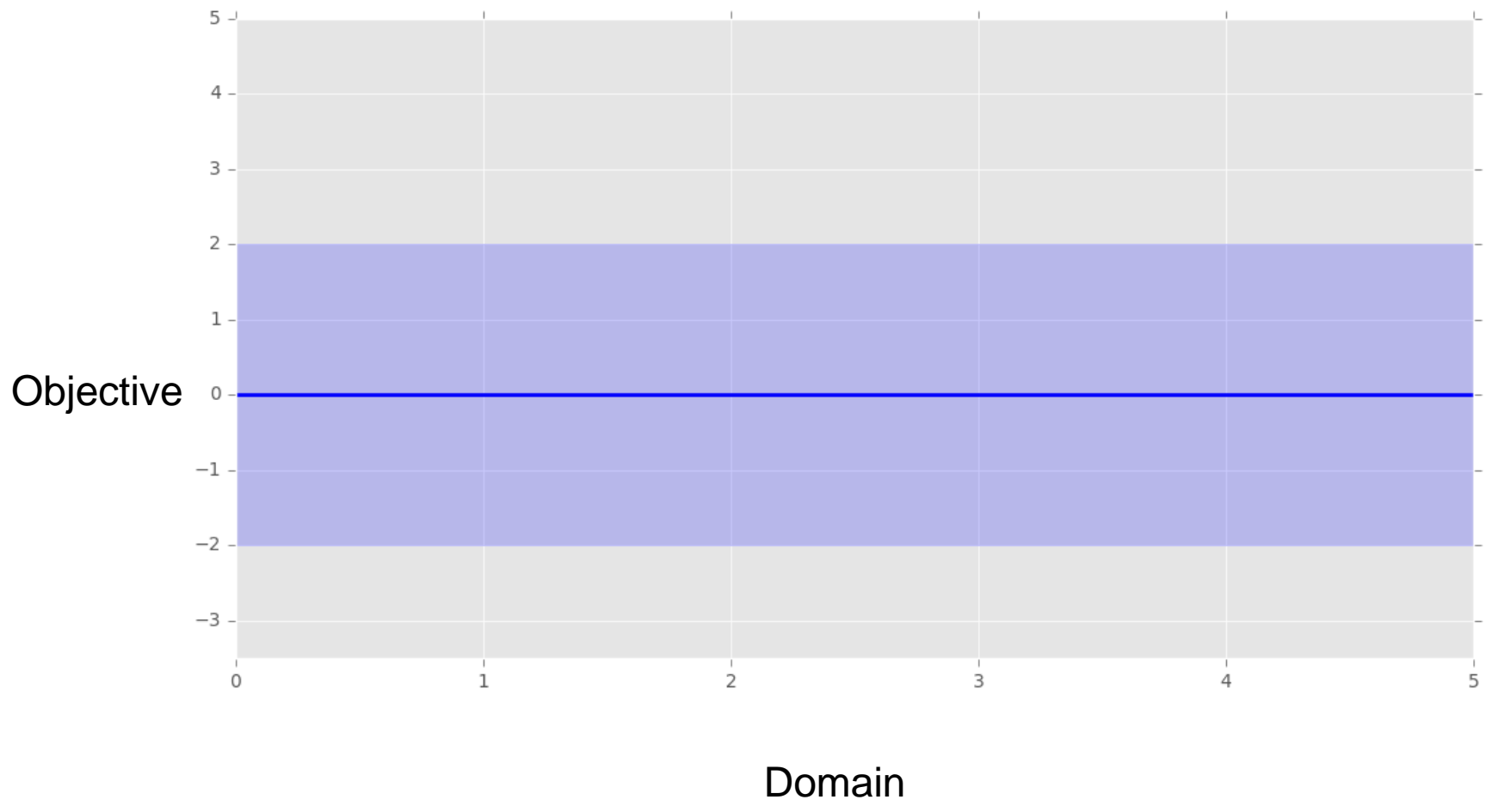


Auto-tuning systems

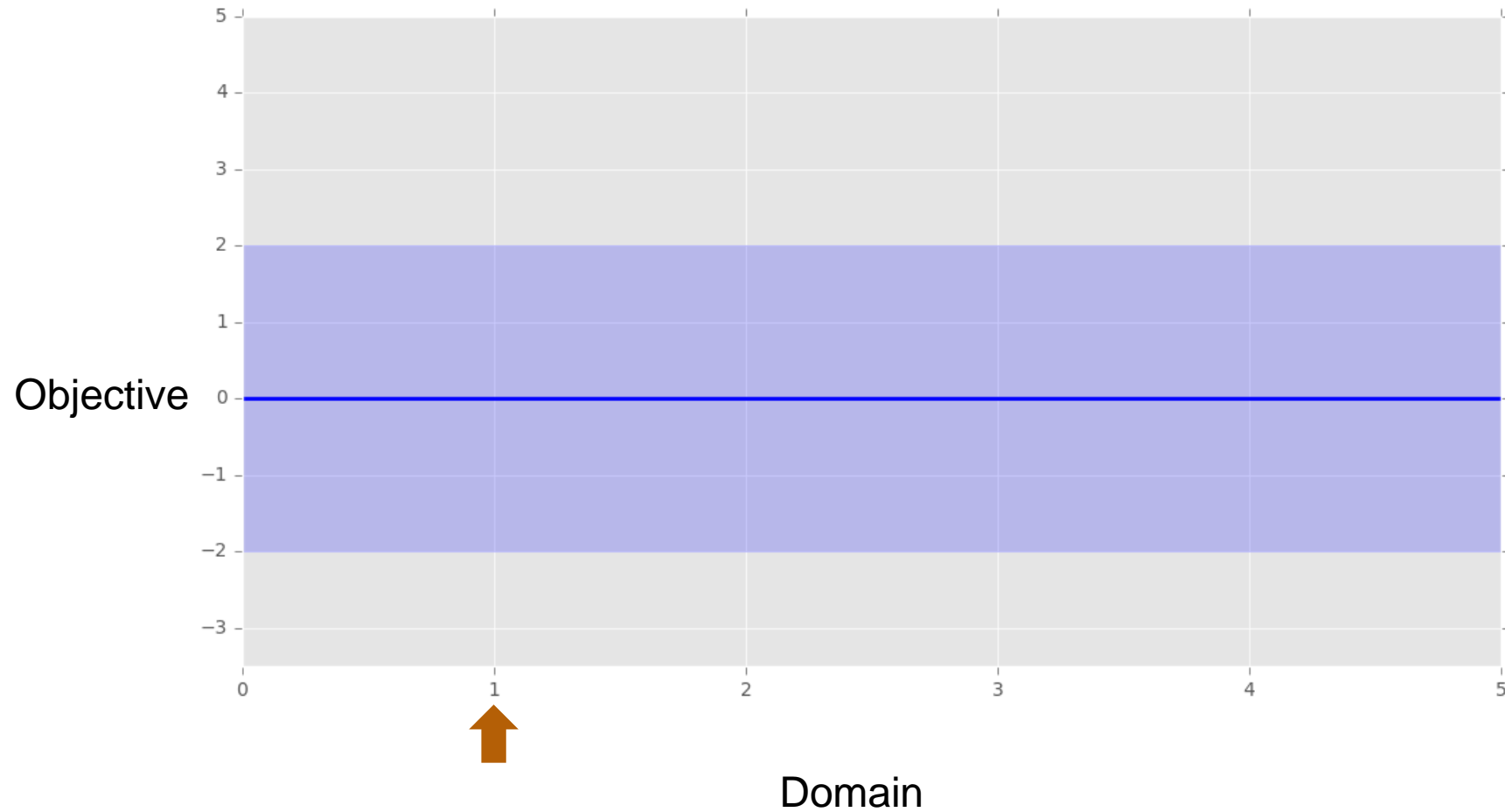


- Properties:
 - Many dimensions (30+)
 - Expensive objective function
 - Understanding of the underlying behaviour
- Bayesian Optimisation is popular tool

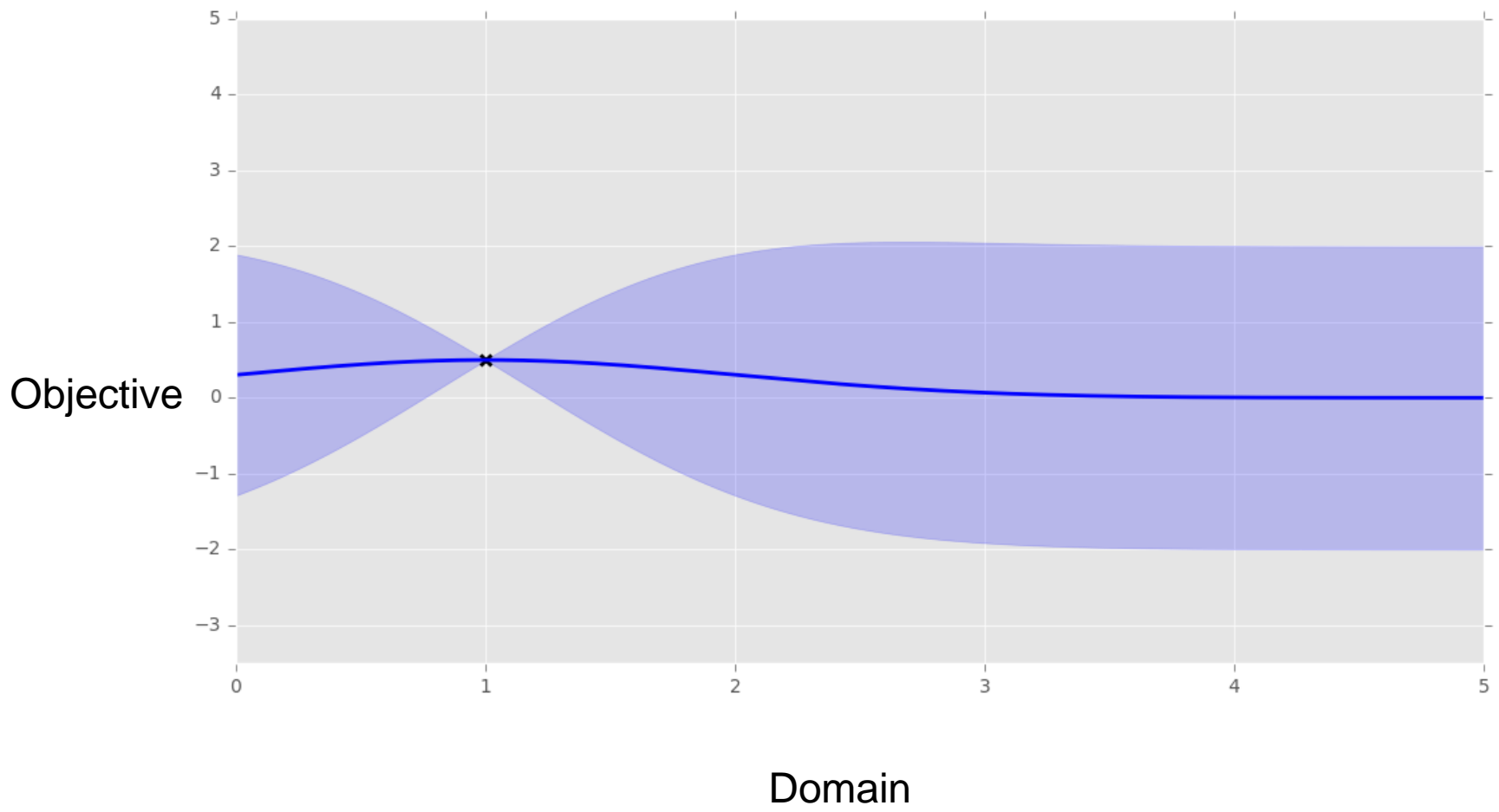
Bayesian optimisation



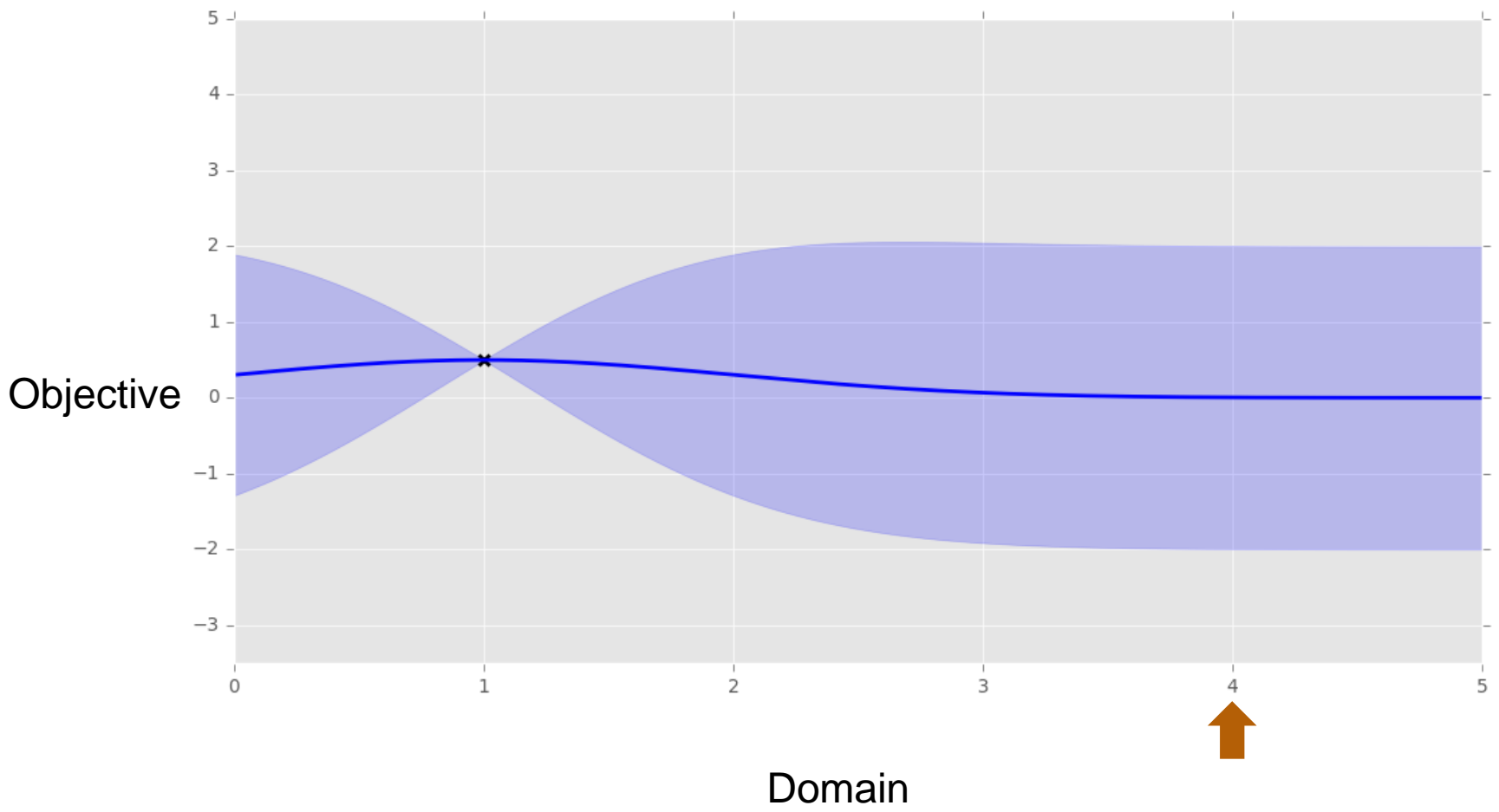
Bayesian optimisation



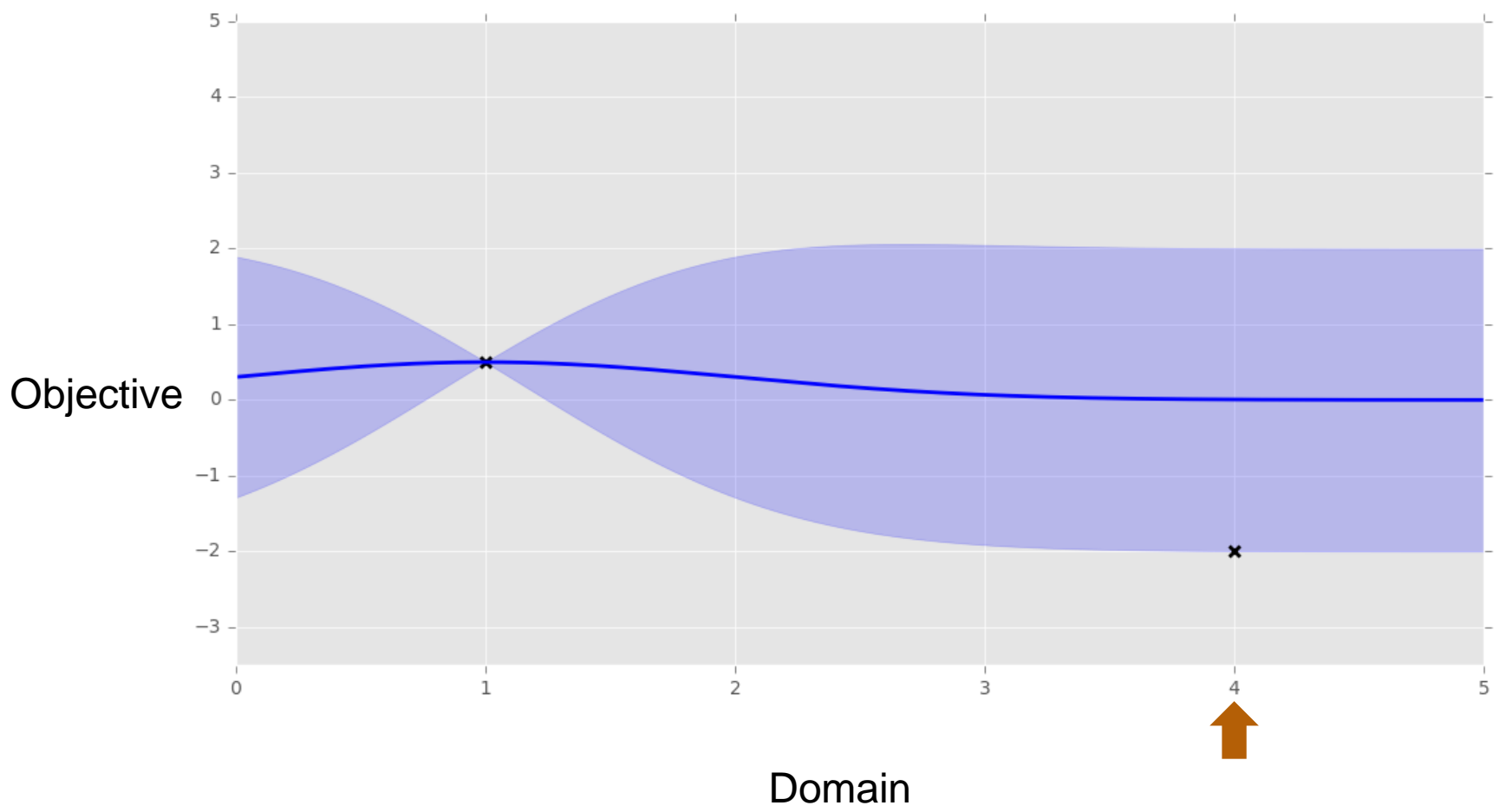
Bayesian optimisation



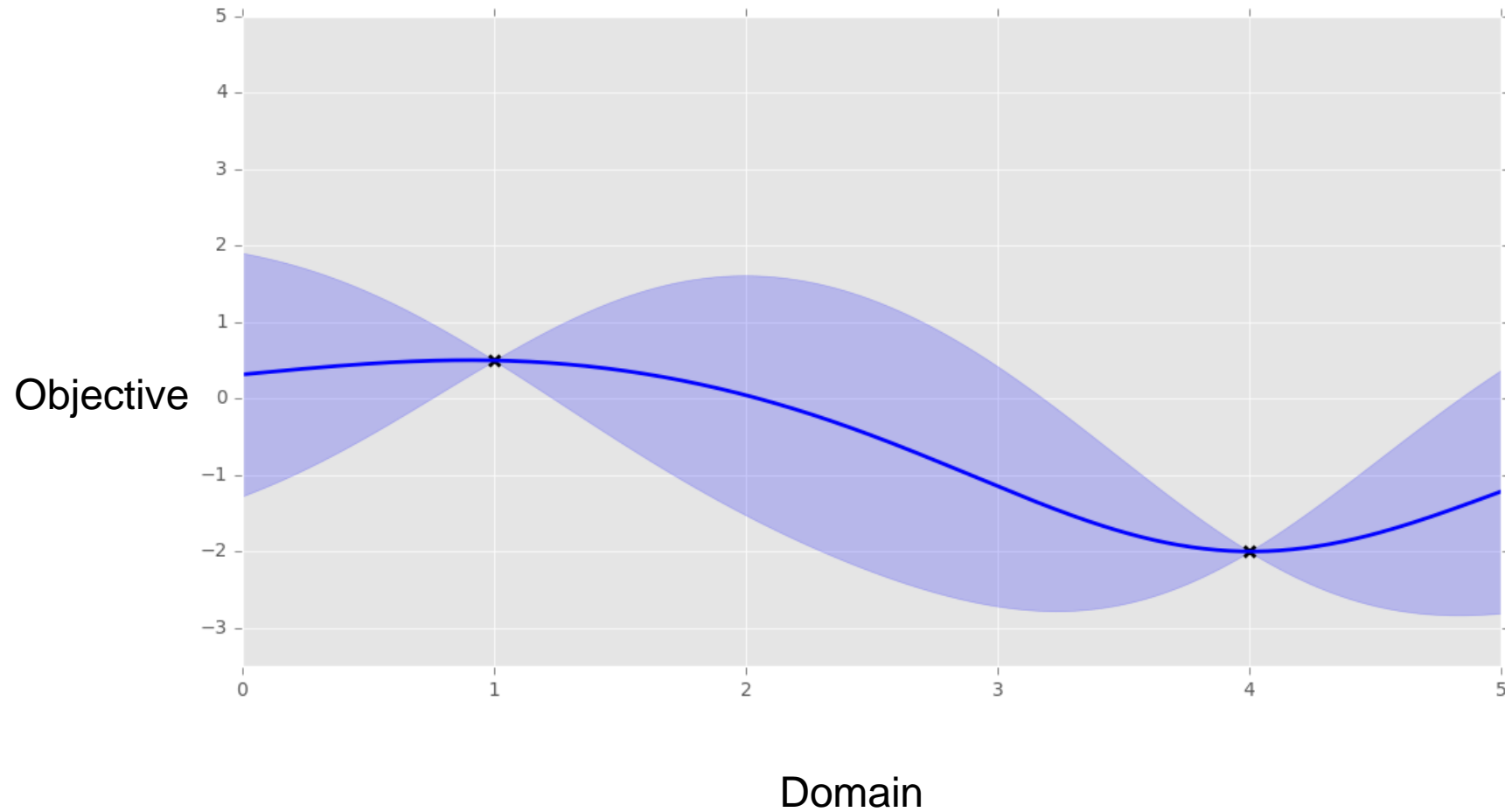
Bayesian optimisation



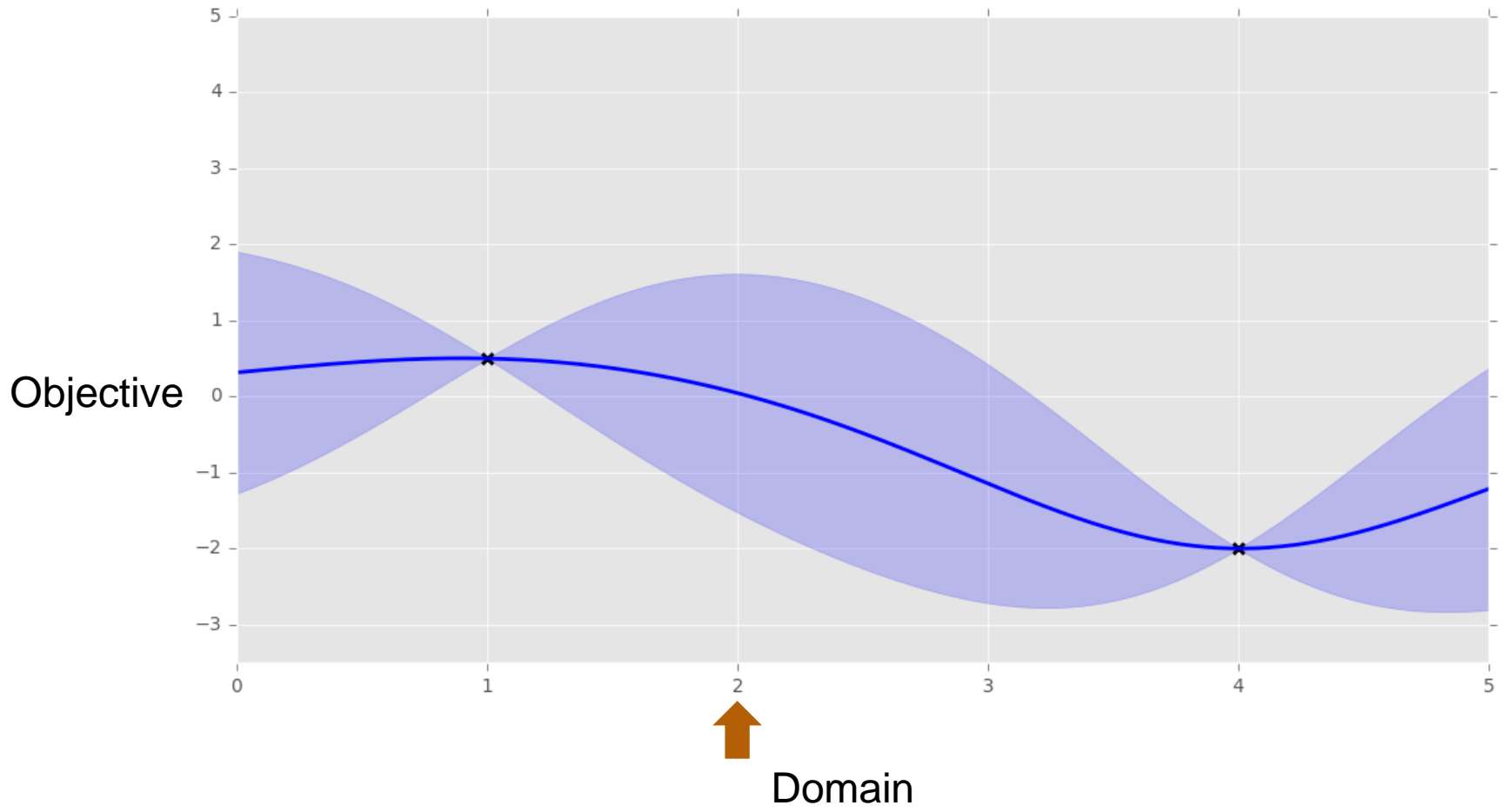
Bayesian optimisation



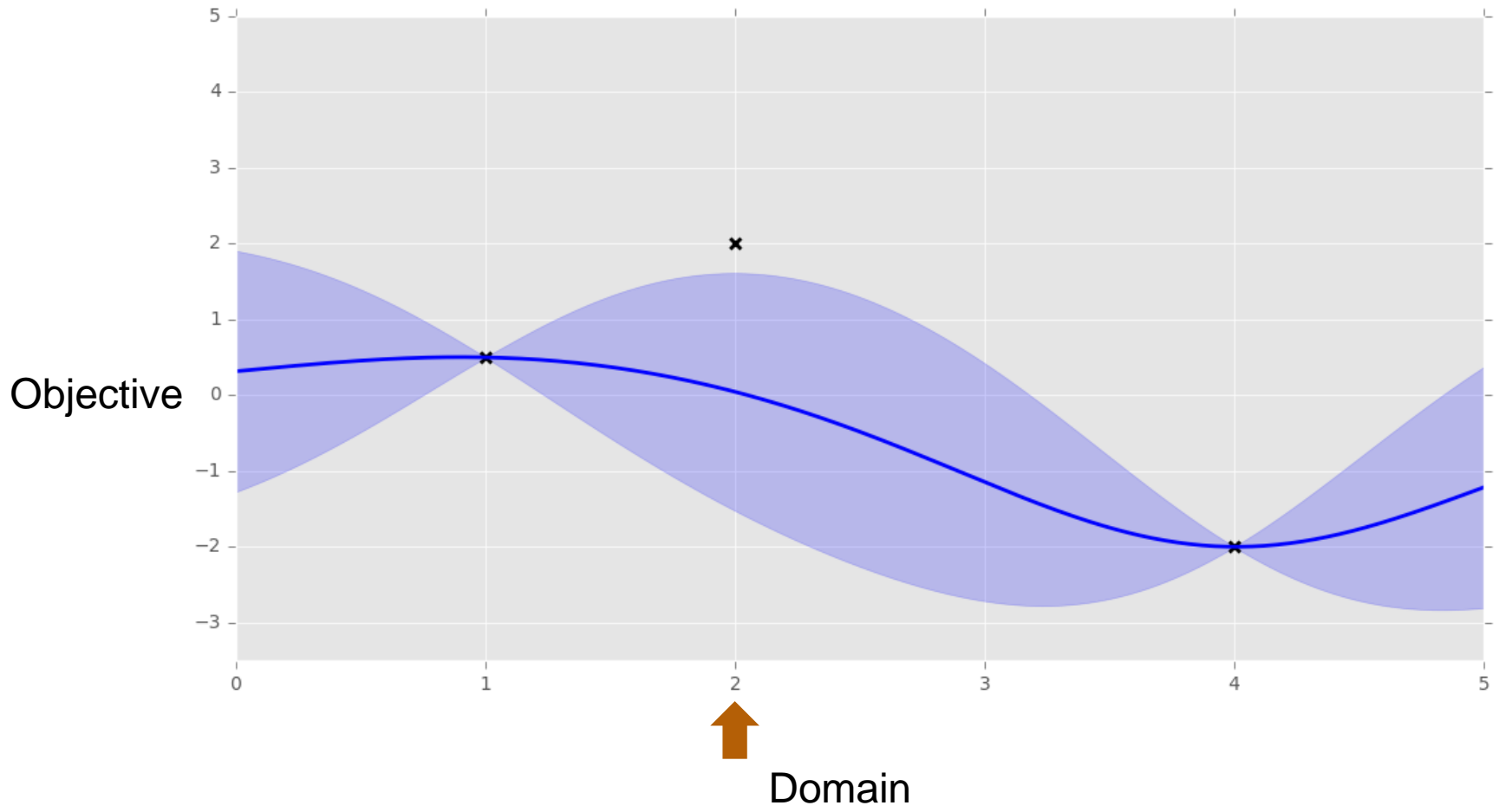
Bayesian optimisation



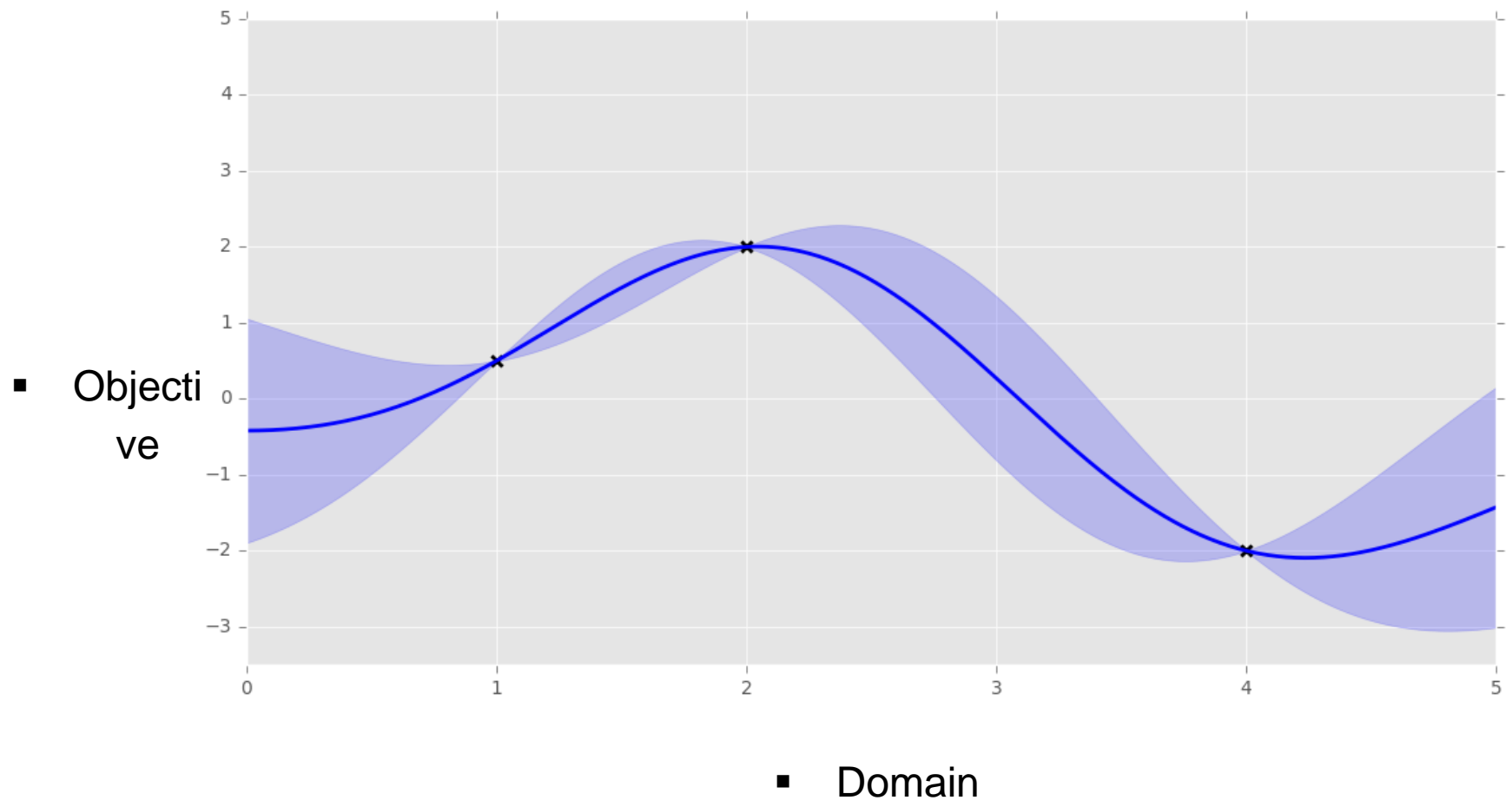
Bayesian optimisation



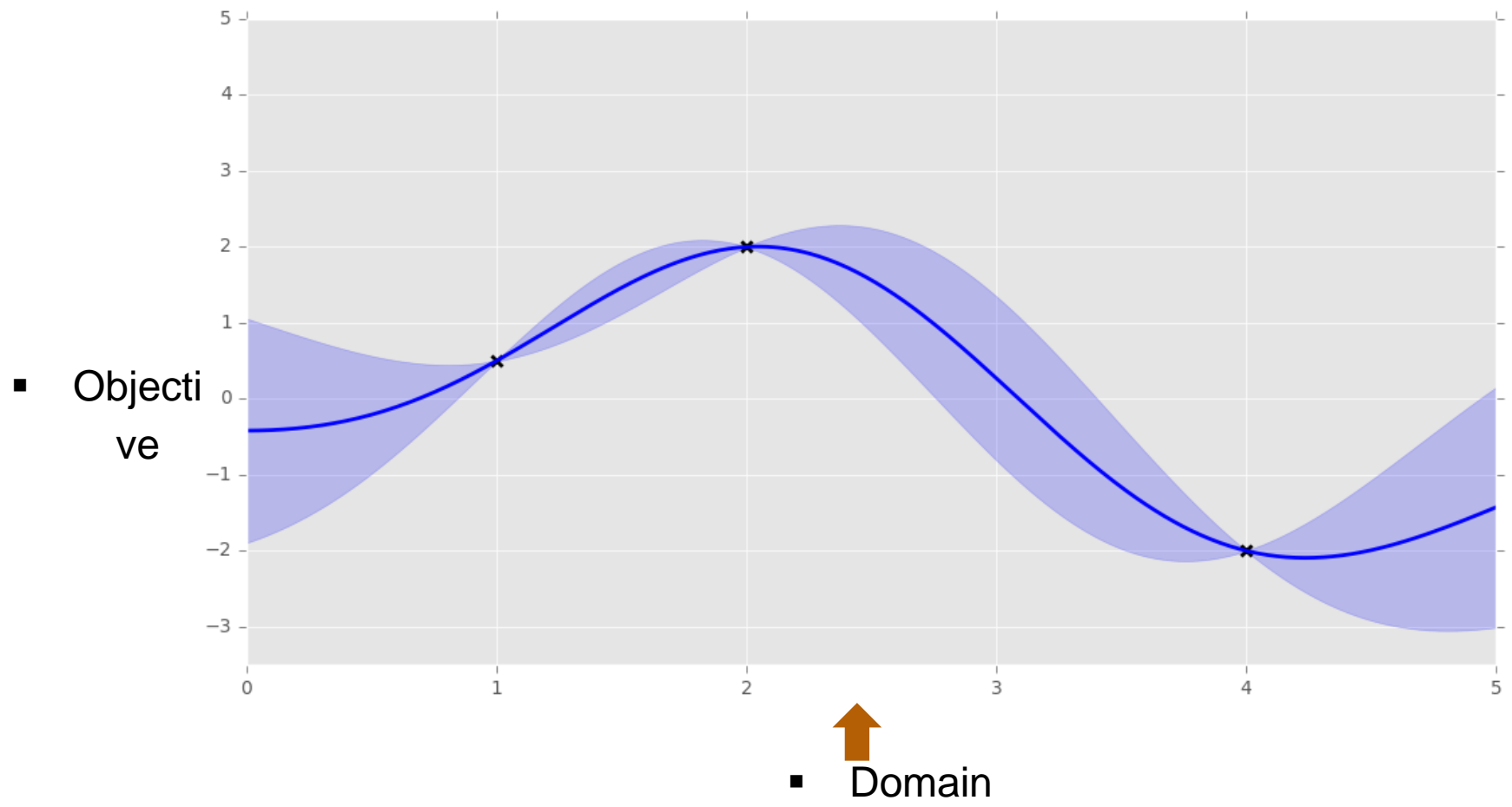
Bayesian optimisation



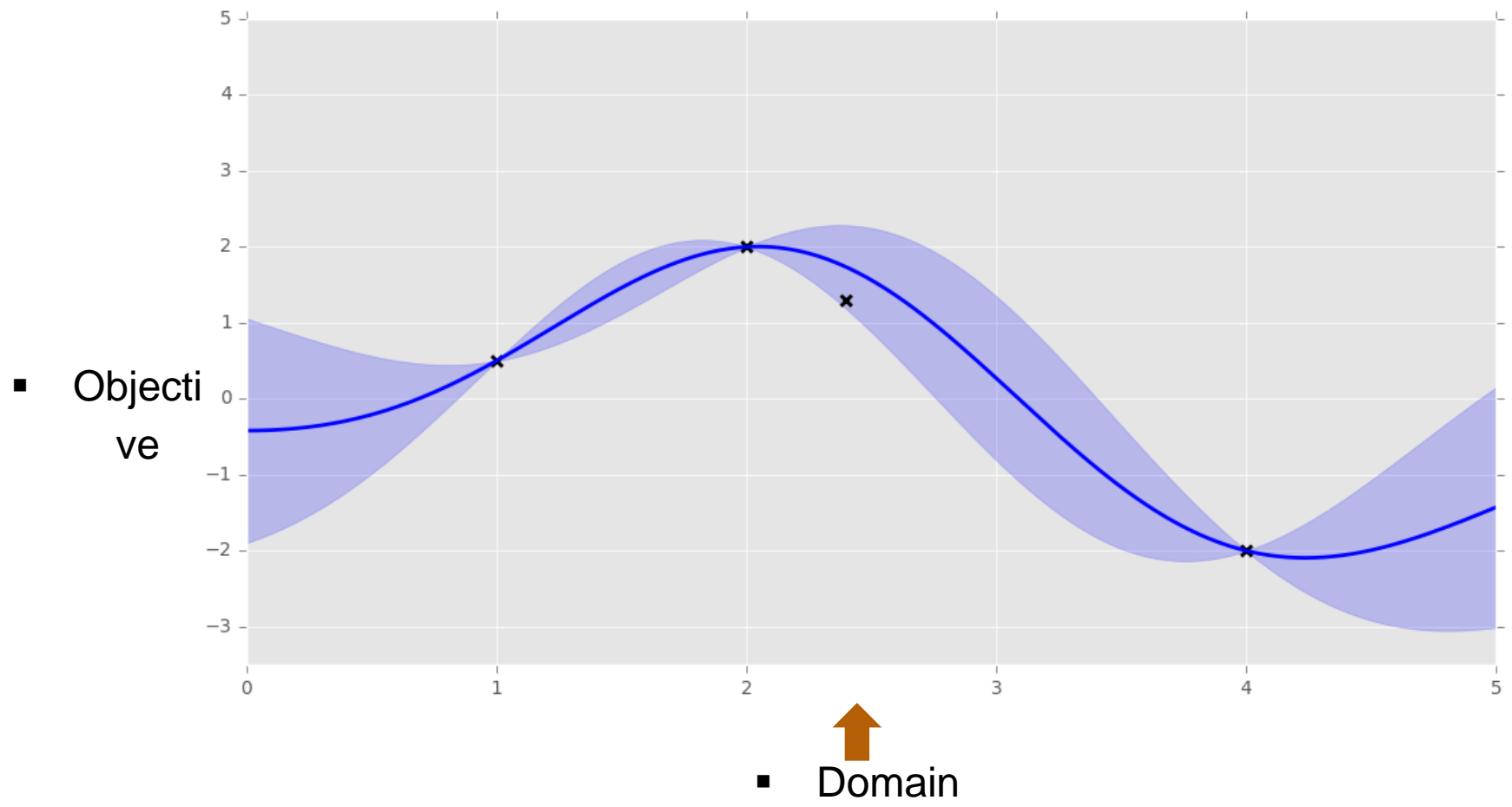
Bayesian optimisation



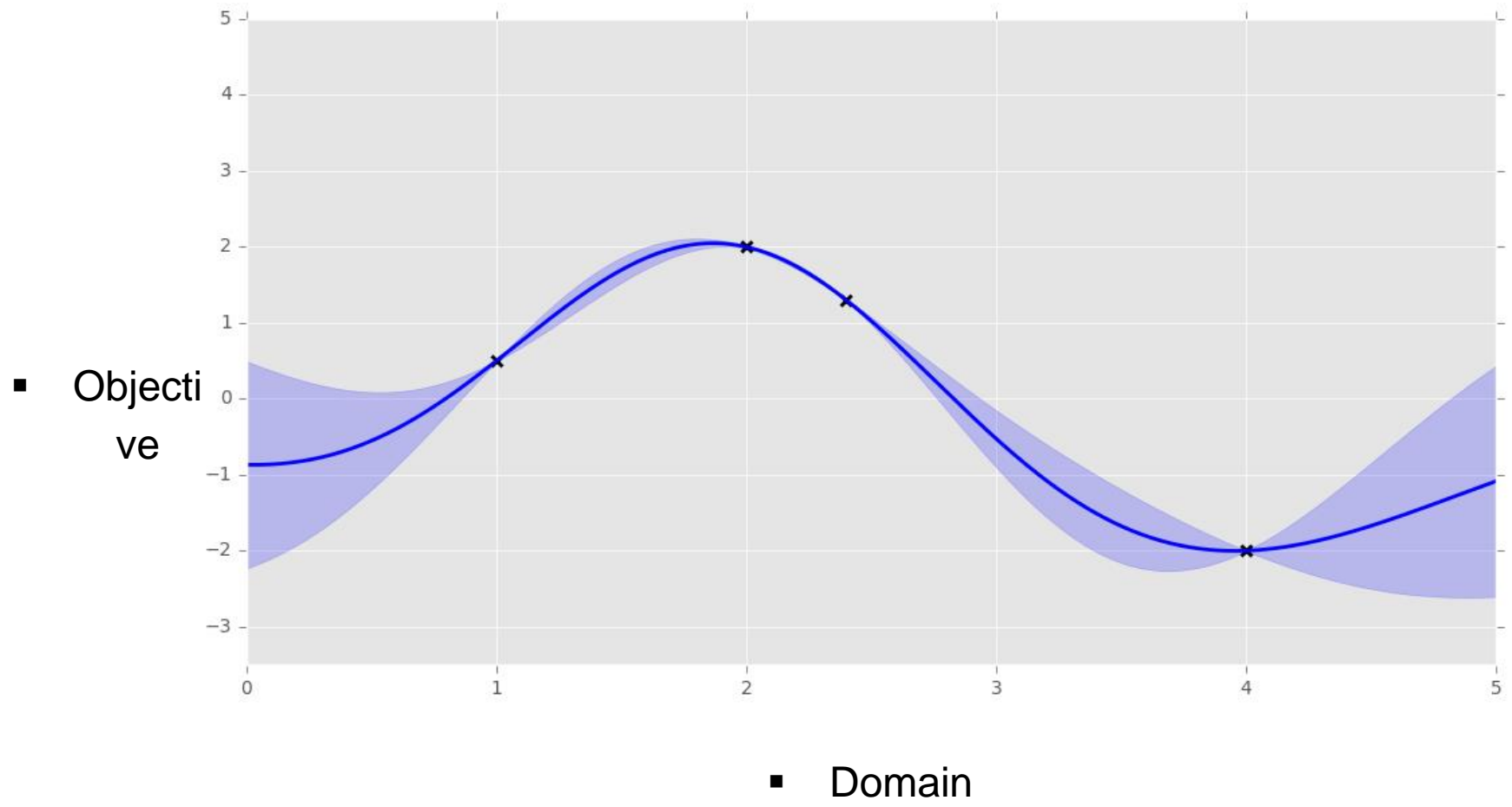
Bayesian optimisation



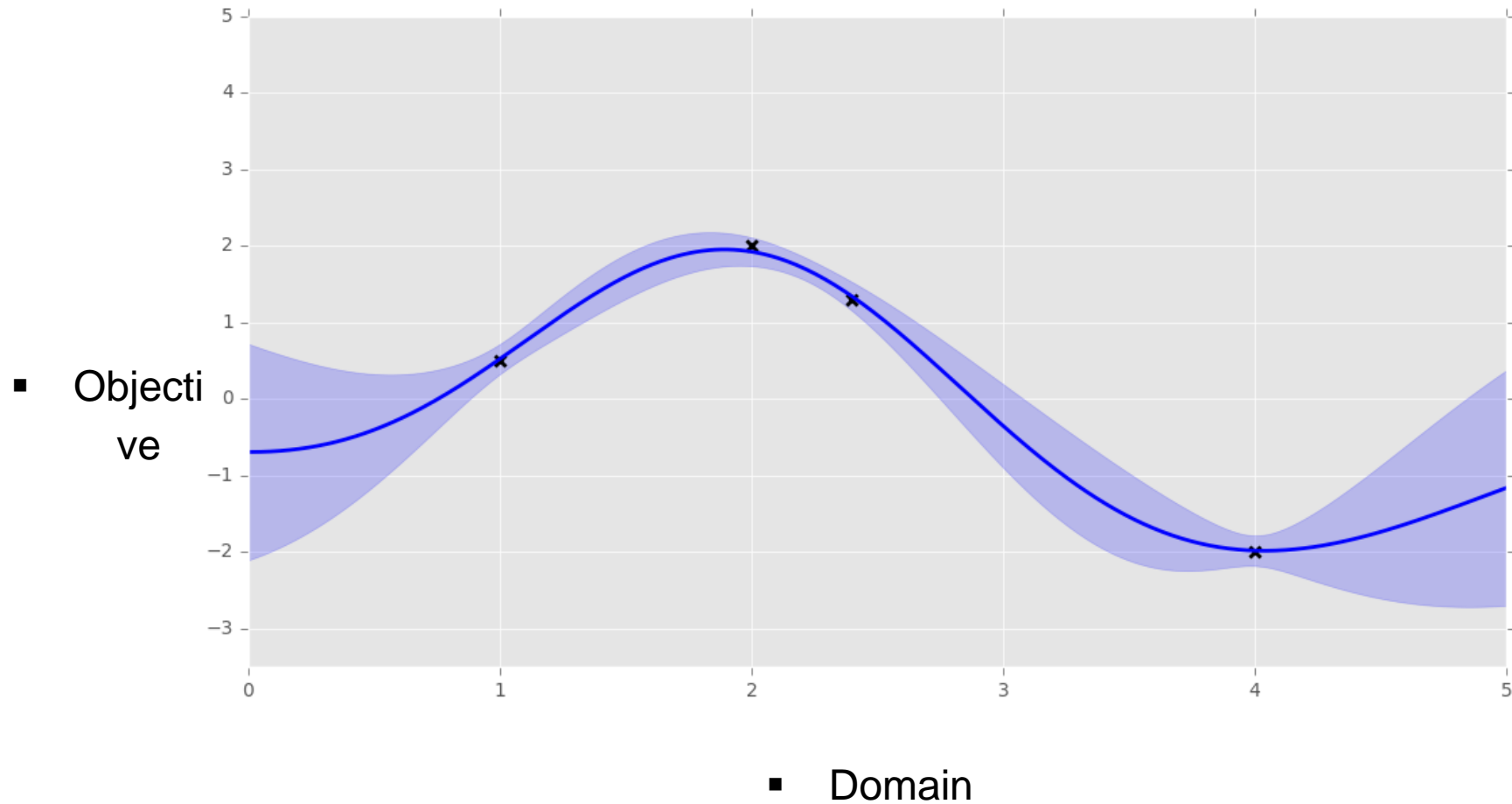
Bayesian optimisation



Bayesian optimisation

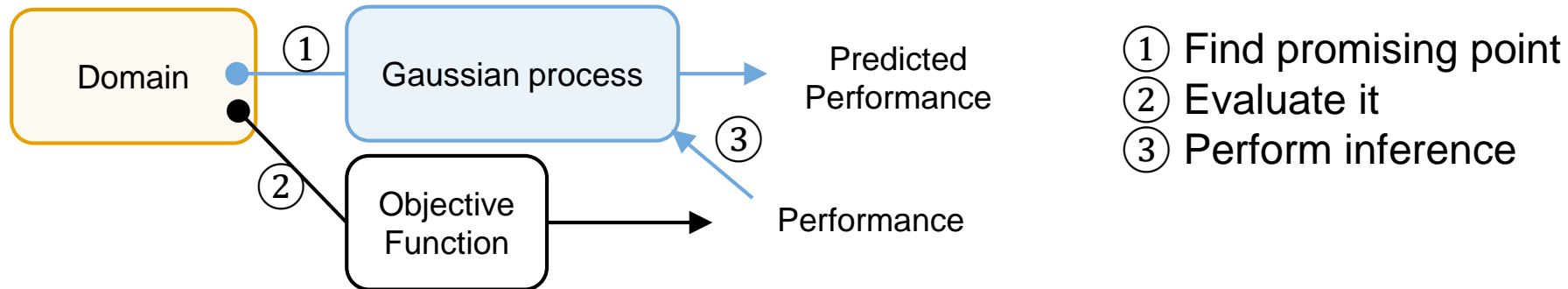


Bayesian optimisation

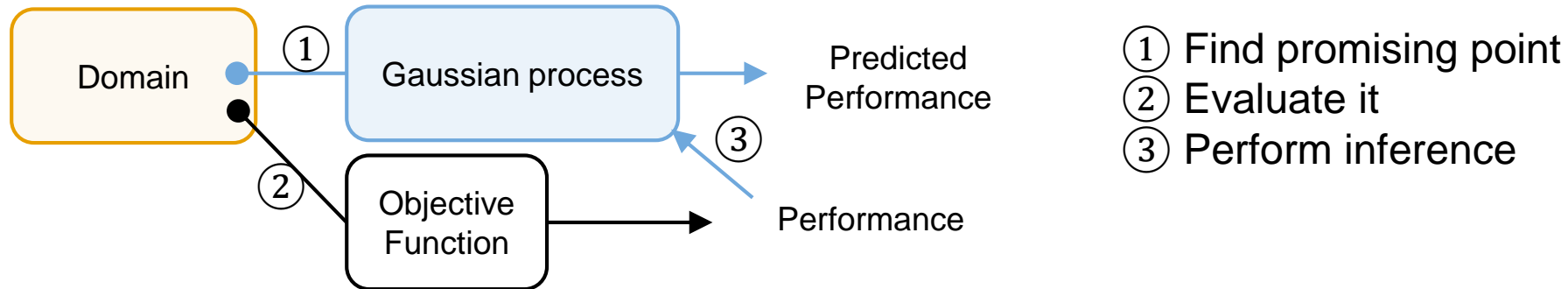


Bayesian optimisation

- Iteratively build a model of the objective function



Bayesian optimisation



Pros:

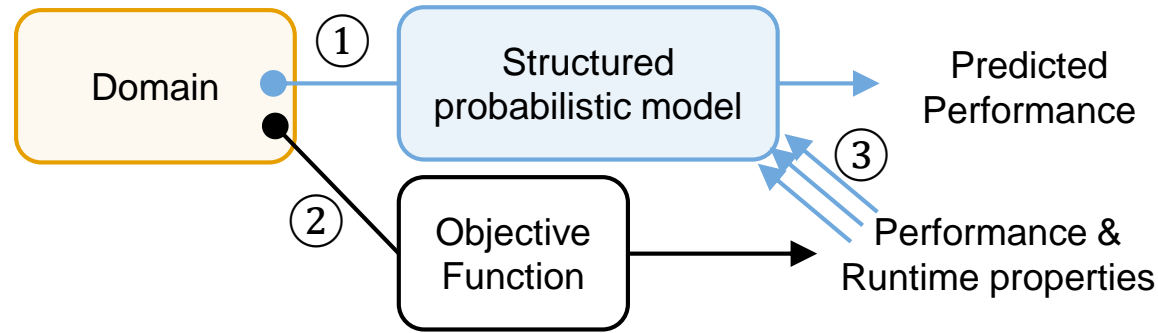
- ✓ Data efficient: converges in few iterations
- ✓ Able to deal with noisy observations

Cons:

- ✗ In many dimensions, model does not converge to the objective function

Solution: Use the known structure of the optimisation problem

Structured Bayesian optimisation



Three desirable properties:

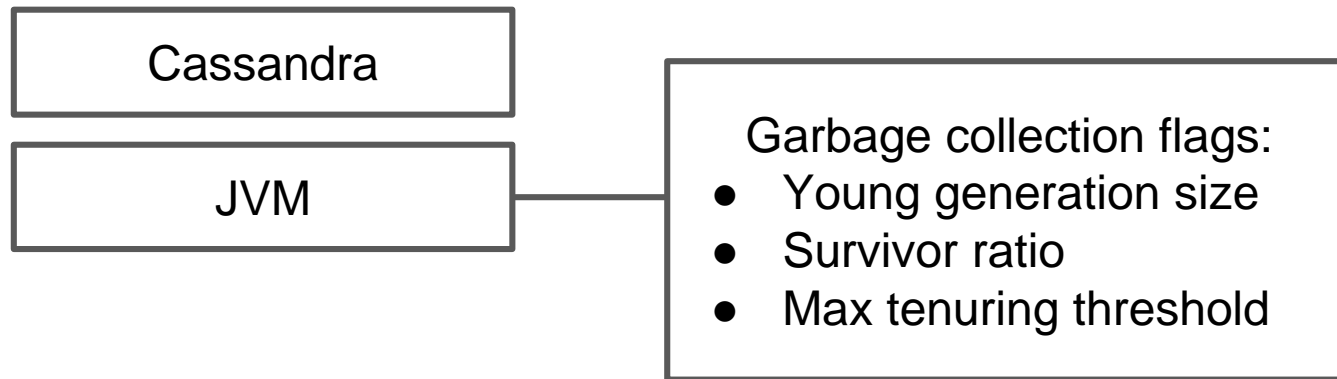
- Able to use many measurements
- Understand the trend of the objective function
- High precision in the region of the optimum

- ✓ Better convergence
- ✓ Use all measurements

- **BOAT**: a framework to build **BespOke Auto-Tuners**
- It includes a probabilistic library to express these models
- V. Dalibard, M. Schaarschmidt, and E. Yoneki: BOAT: Building Auto-Tuners with Structured Bayesian Optimization, WWW 2017. (Morning Paper on May 18, 2017)

Example:

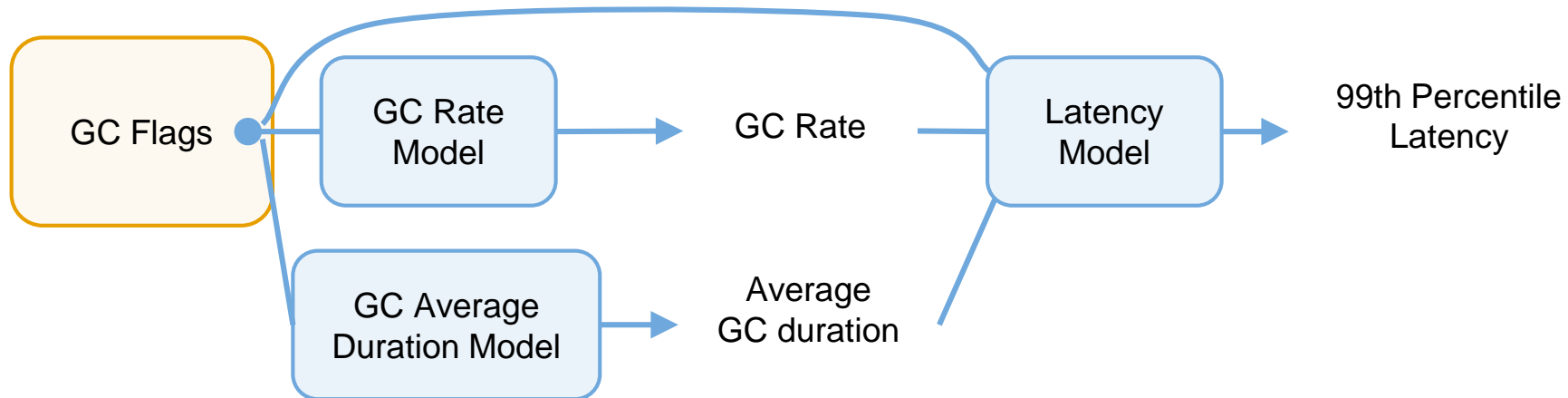
- Cassandra's garbage collection



- Minimise 99th percentile latency of Cassandra

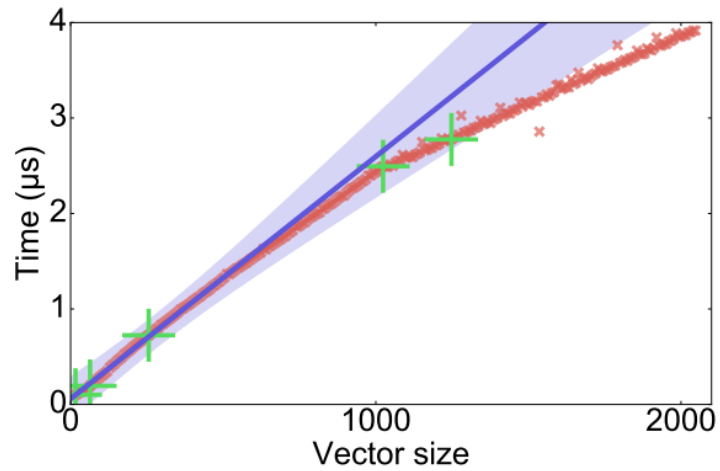
Using many measurements

- Define a directed acyclic graph (DAG) of models

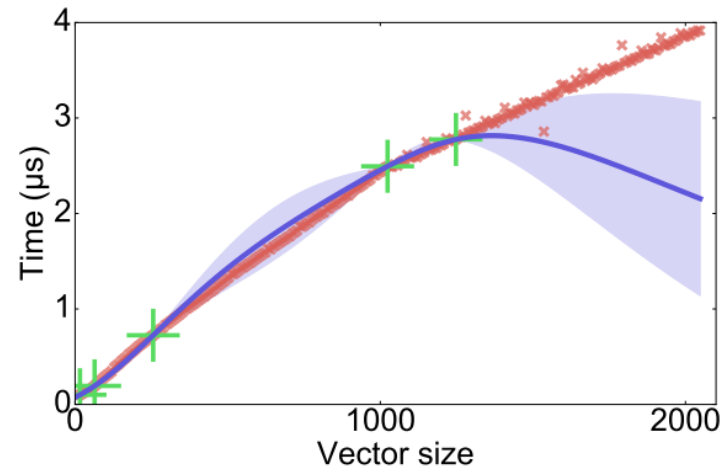


Tune three JVM parameters of a database (Cassandra) to minimise latency

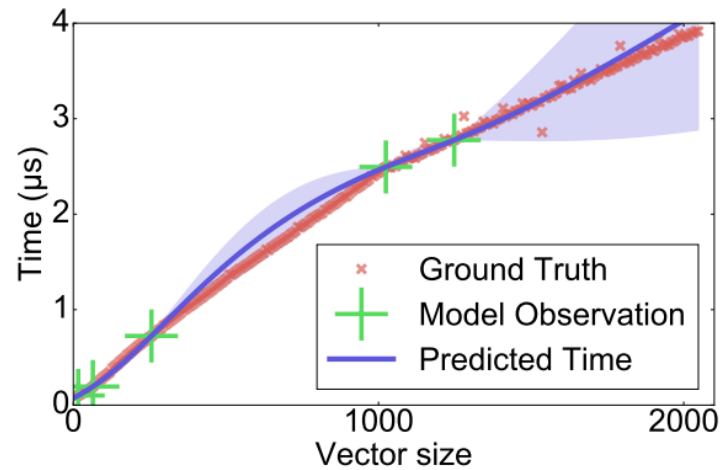
Semi-parametric models



(a) Parametric (Linear regression)



(b) Non-parametric (Gaussian process)



(c) Semi-parametric (Combination)



Computer Systems Optimisation Models

- **Long-term planning:** requires model of how actions affect future states. Only a few system optimisations fall into this category, e.g. network routing optimisation.
- **Short-term dynamic control:** major system components are under dynamic load, such as resource allocation and stream processing, where the future load is not statistically dependent on the current load. Bayesian optimisation is sufficient to optimise distinct workloads. For dynamic workload, Reinforcement Learning would perform better.
- **Combinatorial optimisation:** a set of options must be selected from a large set under potential rules of combination. For this situation, one can either learn online if the task is cheap via random sampling, or via RL and pre-training if the task is expensive, or massively parallel online training given sufficient resources.

Towards Reinforcement Learning

- Given a set of actions with some unknown reward distributions, maximise the cumulative reward by taking the actions sequentially, one action at each time step and obtaining a reward immediately.
- To find the optimal action, one needs to explore all the actions but not too much. At the same time, one needs to exploit the best action found so-far by exploring.
- What makes reinforcement learning different from other machine learning paradigms?
 - There is no supervisor, only a reward signal
 - Feedback is delayed, not instantaneous
 - Time really matters (sequential)
 - Agent's actions affect the subsequent data it receives

AlphaGo defeating the Go World Champion



Practical Issues: Scalable action spaces

Many systems problems are combinatorial in nature:

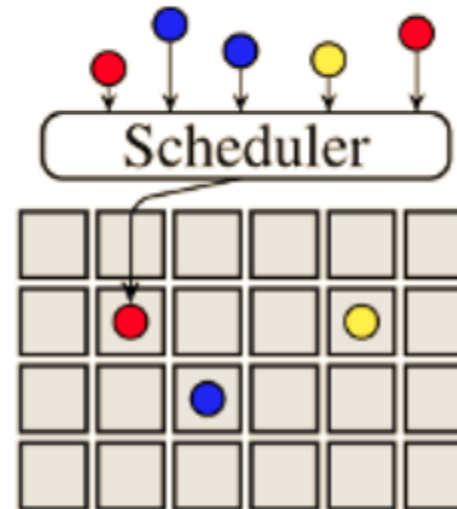
- Mapping tasks to resources
- Selecting a subset of different configuration options
- Large number of options, e.g. thousands of devices in a cluster

‘Standard RL’ cannot deal with this – training data requirements grow exponentially

Problem: Controlling dynamic behaviour

Assume workload dynamic,
e.g. seasonality, load spikes,
shared resources, failures..

- Algorithm: workload \rightarrow behavior **distribution**
- Involves approximations to NP-complete problems, e.g. bin packing, sub-graph isomorphism, ..



Source: firmament.io

Trade-offs in iterative optimisation

Grid search: Discretised sweep

Random search: No risk of 'getting stuck',
potentially many samples required

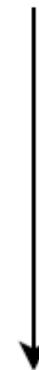
Evolution strategies: Evaluate
permutations against fitness function

Bayes Opt: Sample efficient, requires
continuous function, some configuration

More
requirements

Computation
more expensive

Fewer samples





Trade-offs in dynamic control

Single static configuration/rule: E.g.
FIFO scheduler

Online estimate of distributions: E.g.
join-order in query planning

Workload clustering: Identify distinct
classes, e.g. write-heavy, read-heavy
workloads, per-class decisions

Fully adaptive: Optimal per-task
behavior, unrealistic in practice

Robust, predictable,
low deployment cost

Analytical overhead

Training/deployment
cost factor

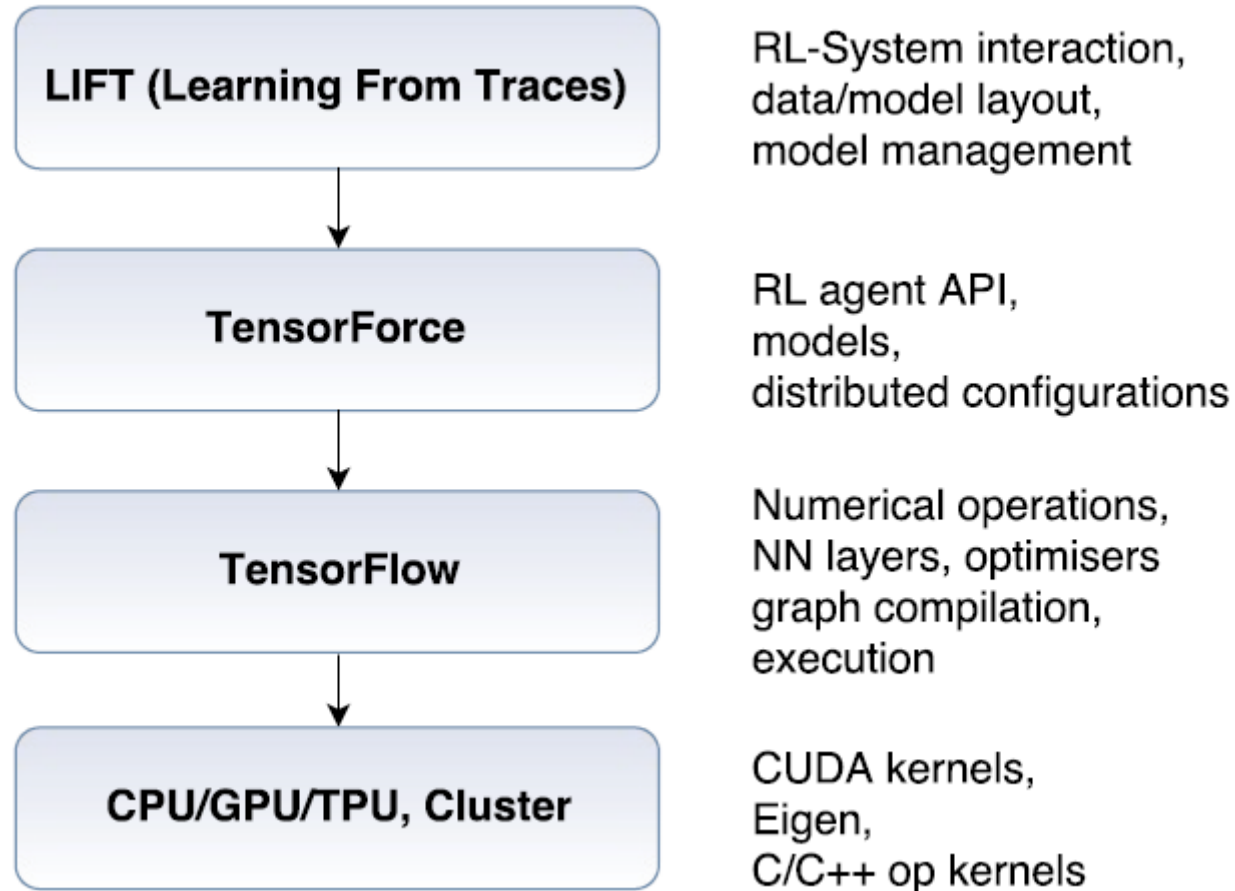




Practical Issues continued...

- Many deep learning tools, **no standard library** for modern RL (~2014-2018)
- **Exploration** in production system not a good idea
 - Unstable, unpredictable
- **Simulations** can oversimplify problem
 - Expensive to build, not justified versus gain
- **Online steps take too long**

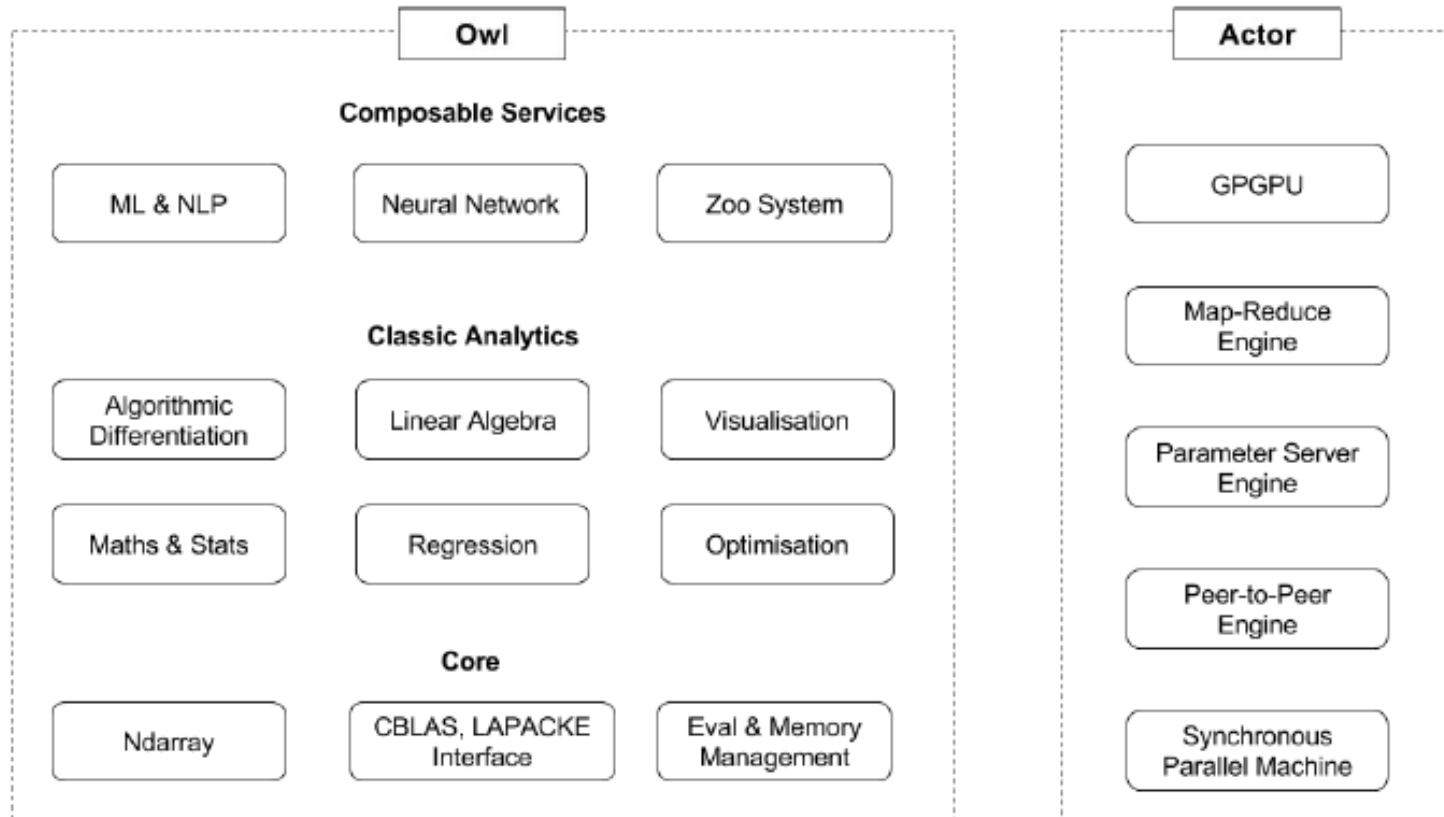
LIFT: A software stack for deep RL



Dynamic Control Flow in Current Frameworks

- There are static computation frameworks WITHOUT dynamic control flow (mxnet, cntk) -> dynamic control flow is in the out of graph host program.
- There are dynamic computation graph frameworks WITH dynamic control flow (PyTorch, DyNet) -> graph is only implicitly defined via imperative operations, cannot do static graph optimisations, typically slower but easier to use.
- There is static computation with dynamic differentiable control flow in the graph -> only TensorFlow offers this among modern deep learning frameworks.

OWL Architecture for OCaml



Owl + Actor = Distributed & Parallel Analytics

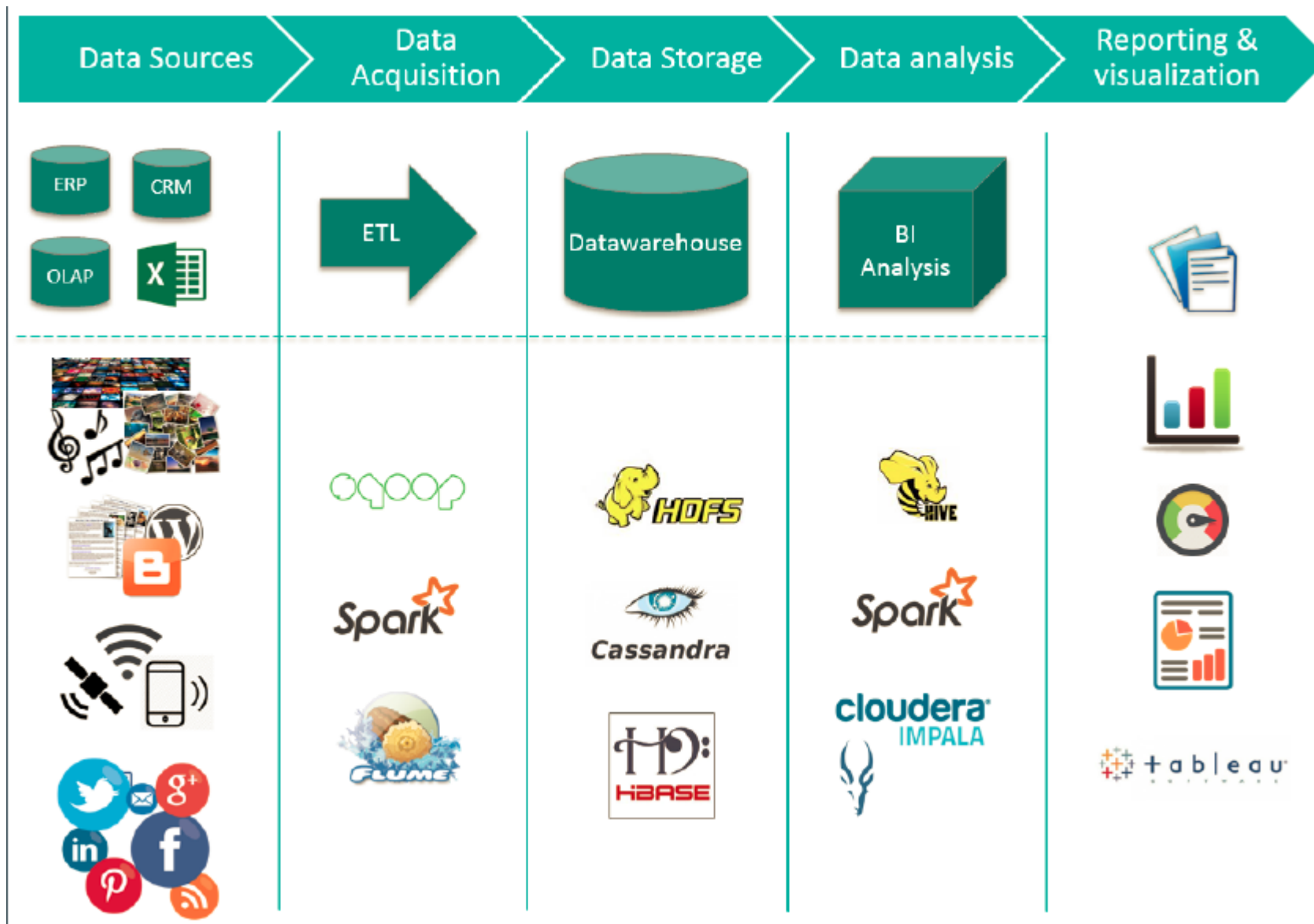
Owl provides numerical backend; whereas Actor implements the mechanisms of distributed and parallel computing. Two parts are connected with functors.

Various system backends allows us to write code once, then run it from cloud to edge devices, even in browsers.

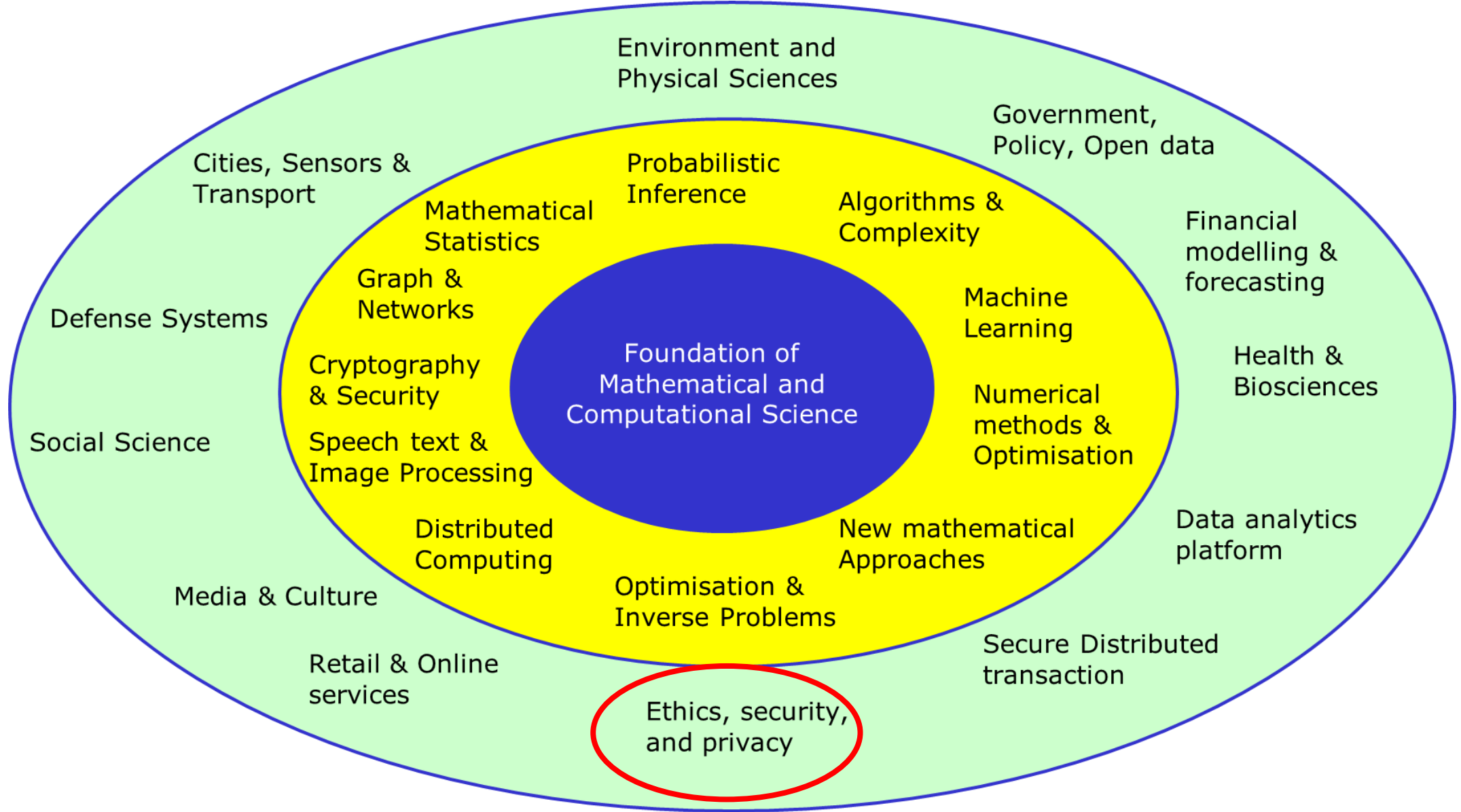
Same code can run in both sequential and parallel mode with Actor engine.



Pipeline of Data Processing...



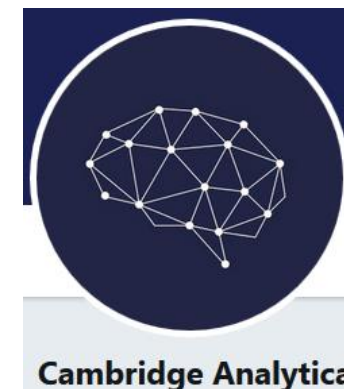
Broad Landscape of Research in ATI in 2015





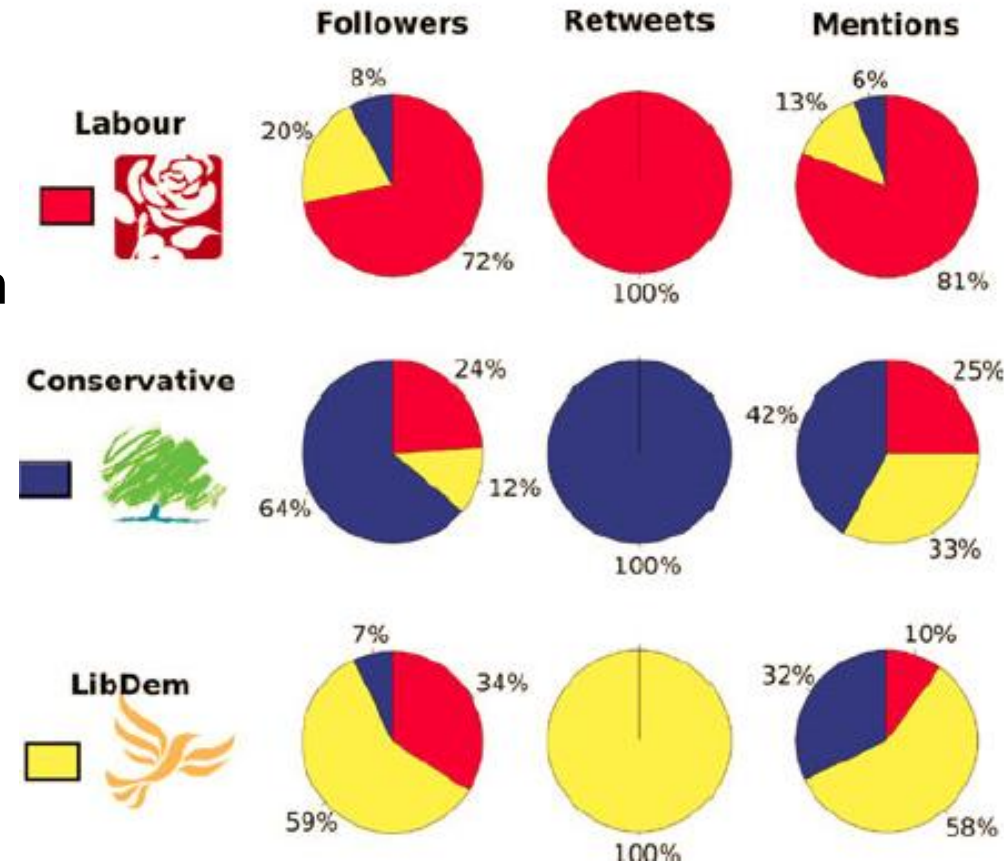
Network Structures

- Eight Friends Are Enough: Social Graph Approximation via Public Listings by Joseph Bonneau et al. 2009.
- Many interesting properties can be accurately approximated. This has disturbing implications for online privacy, since leaking graph information enables transitive privacy loss.
- Cambridge Analytica was caught tampering with elections by exploiting Facebook, but chances are that this is the tip of the iceberg, and that many others.



Social Network Analysis

- What's in Twitter: I Know What Parties are Popular and Who You are Supporting Now!
- 1,150,000 messages from the main stream of Twitter related to the 2010 UK General Election between the 5th and the 12th of May from about 220,000 users in Twitter.
- With party characteristics, classification algorithm based on Bayesian framework to compute political preferences of users is developed. Accuracy of 86 % without any training and the network topology information.



Modern Data Scientist: The sexiest job of 21th century



MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing package e.g. R
- ☆ Databases SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

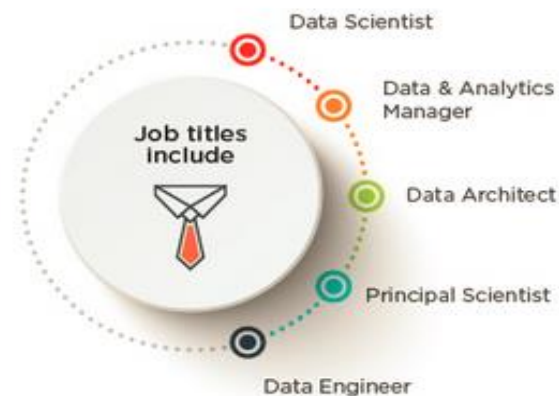
Many Courses offered: e.g. Master Certificate

- Data scientist is the pinnacle rank in an analytics organisation. You will be required to understand the business problem, design the analysis, collect and format the required data, apply algorithms or techniques using the correct tools, and finally make recommendations backed by data.

£ 1799

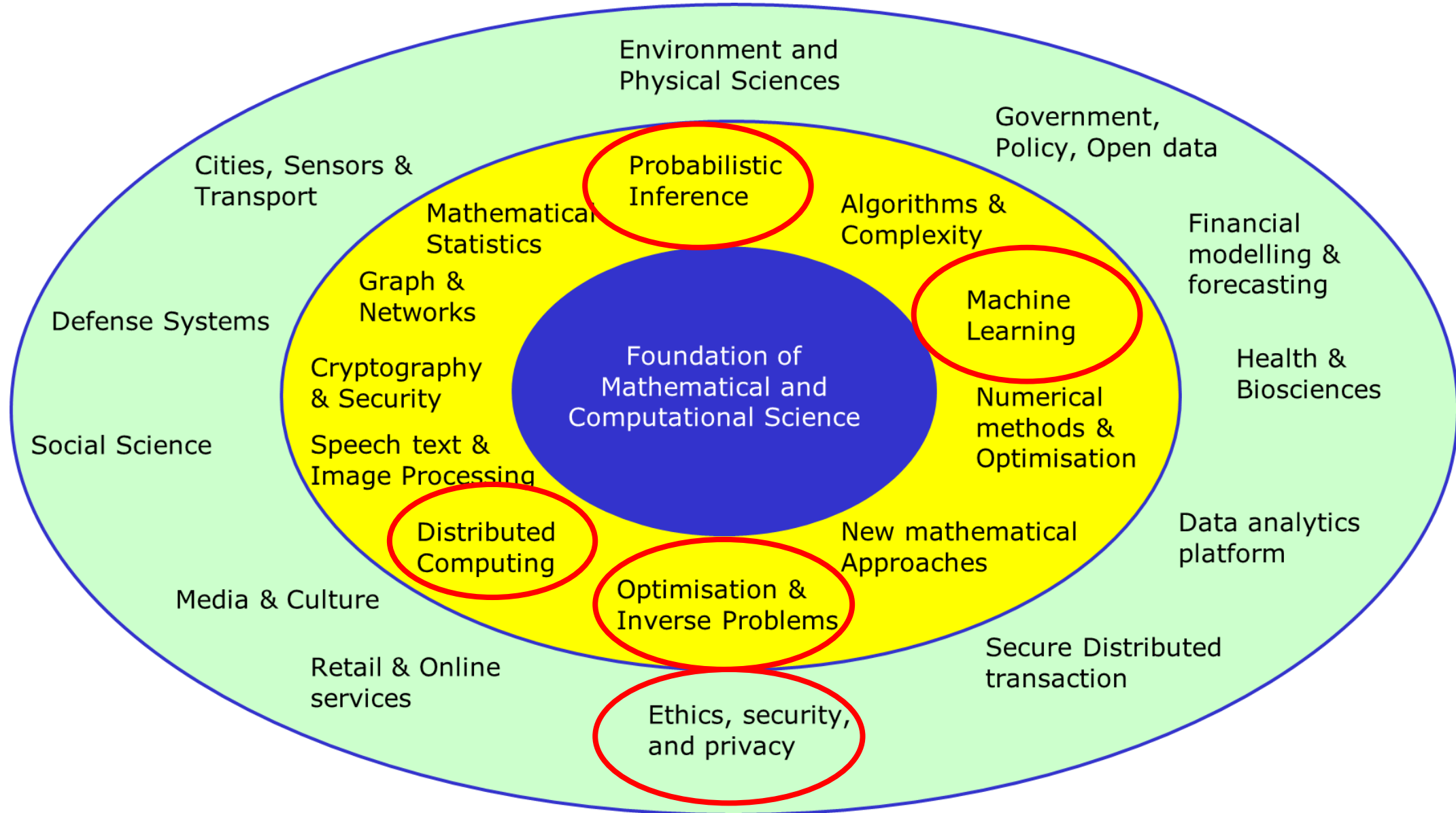
ENROLL NOW

* VAT Included



- Course 1
Data Science with SAS Training
- Course 2
Data Science Certification Training - R Programming
- Course 3
Big Data Hadoop and Spark Developer
- Course 4
Data Science with Python
- Course 5
Business Analytics with Excel
- Course 6
Machine Learning
- Course 7
Deep Learning with TensorFlow
-  Masters Certificate
*You will get individual certificates for each course.

Broad Landscape of Research in ATI in 2015



Modern Theory of Deep Learning

Why does it work so well?

- On the Expressive Power of Deep Neural Networks PMLR 2017: understanding of how and why neural network architectures achieve their empirical successes is still lacking.
- Ali Rahimi's talk at NIPS(NIPS 2017 Test-of-time award presentation)
- Deep Learning works in Practice. But Does it work in Theory? By L. Hoang and R. Guerraoui. (<https://arxiv.org/pdf/1801.10437.pdf>)
- Understanding deep learning requires rethinking generalisation
- Fundamental theory behind the paradoxical effectiveness of deep learning. Still open research problem...

Gap between Research and Practice

Device Placement Optimization with Reinforcement Learning

Azalia Mirhoseini^{*12} Hieu Pham^{*12} Quoc V. Le¹ Benoit Steiner¹ Rasmus Larsen¹ Yuefeng Zhou¹
Naveen Kumar³ Mohammad Norouzi¹ Samy Bengio¹ Jeff Dean¹

20H with 80GPUs!

Research opportunities ahead!

<http://www.cl.cam.ac.uk/~ey204/>

Acknowledgements: Michael Schaarschmidt (RL) and Valentin Dalibard (BOAT).

