# Introduction to Bayesian Optimization

**Javier González**
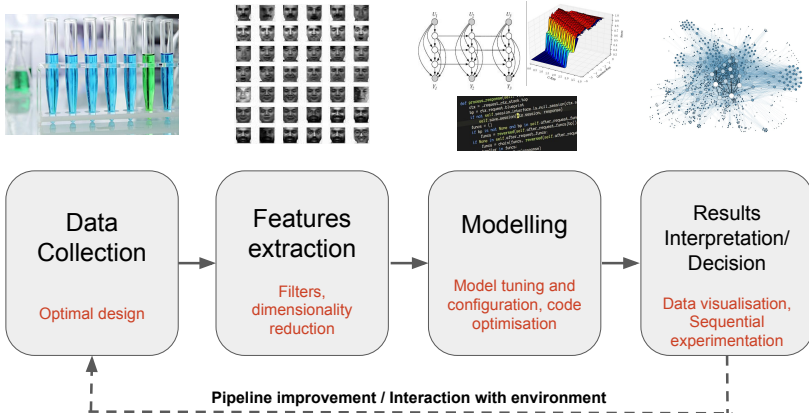
Masterclass, 7-February, 2107 @Lancaster University

Section I: Introduction to Bayesian Optimization

- ▶ What is BayesOpt and why it works?
- ▶ Relevant things to know.

# Data Science pipeline/Autonomous System
Challenges and needs for automation



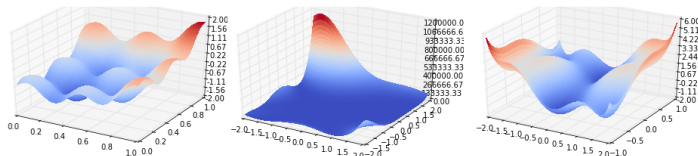| Data Collection | Features extraction | Modelling | Results Interpretation/ Decision |
|---|---|---|---|
| Optimal design | Filters, dimensionality reduction | Model tuning and configuration, code optimisation | Data visualisation, Sequential experimentation |

**Pipeline improvement / Interaction with environment**

Emulator - Simulator - Physical system

# Global optimization

Consider a 'well behaved' function $f : \mathcal{X} \to \mathbb{R}$ where $\mathcal{X} \subseteq \mathbb{R}^D$ is a bounded domain.
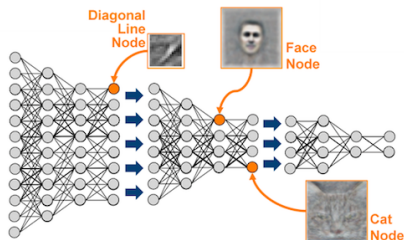
$$x_M = \arg\min_{x \in \mathcal{X}} f(x).$$



- $f$ is explicitly unknown and multimodal.
- Evaluations of $f$ may be perturbed.
- Evaluations of $f$ are expensive.
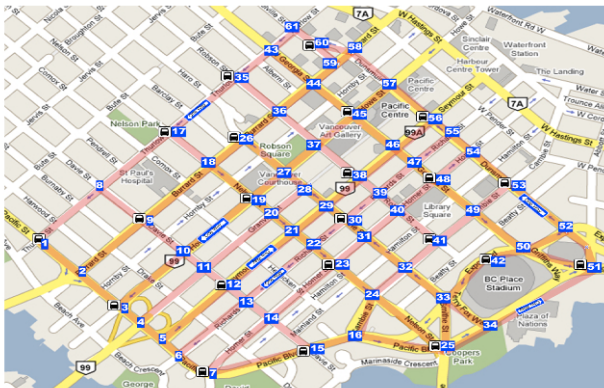
# Expensive functions, who doesn't have one?

**Parameter tuning in ML algorithms.**



- Number of layers/units per layer
- Weight penalties
- Learning rates, etc.

Figure source: http://theanalyticsstore.com/deep-learning

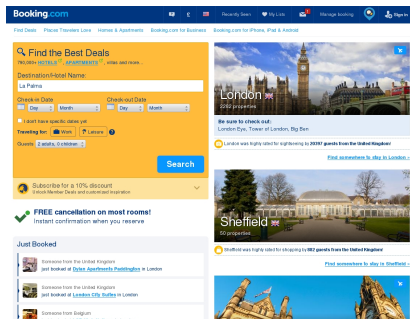# Expensive functions, who doesn't have one?

**Active Path Finding in Middle Level**



Optimise the location of a sequence of waypoints in a map to navigate from a location to a destination.

# Expensive functions, who doesn't have one?
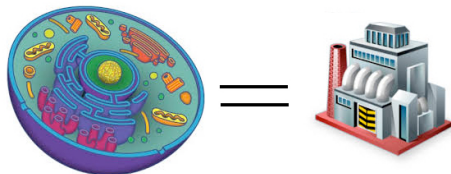
**Tuning websites with A/B testing**



Optimize the web design to maximize sign-ups, downloads, purchases, etc.

**Design of experiments: gene optimization**



- ▶ Use mammalian cells to make protein products.
- ▶ Control the ability of the cell-factory to use synthetic DNA.

Optimize genes (ATTGGTUGA...) to best enable the
cell-factory to operate most efficiently.

# Expensive functions, who doesn't have one?

Many other problems:

- Robotics, control, reinforcement learning.
- Scheduling, planning
- compilers, hardware, software?
- Intractable likelihoods.

# What to do?

**Option 1: Use previous knowledge**

To select the parameters at hand. Perhaps not very scientific but still in use...

**Option 2: Grid search?**

If $f$ is L-Lipschitz continuous and we are in a noise-free domain to guarantee that we propose some $x_{M,n}$ such that

$$f(x_M) - f(x_{M,n}) \leq \epsilon$$

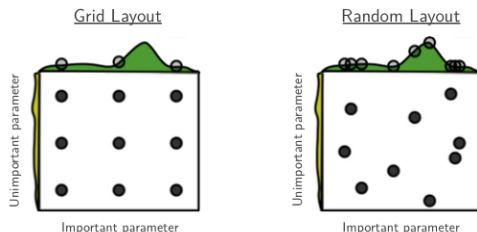we need to evaluate $f$ on a D-dimensional unit hypercube:

$$(L/\epsilon)^D evaluations!$$

**Example**: $(10/0.01)^5 = 10e14...$
... but function evaluations are very expensive!

# What to do?

**Option 3: Random search?**

We can sample the space uniformly [Bergstra and Bengio 2012]



Better than grid search in various senses but still expensive to guarantee good coverage.
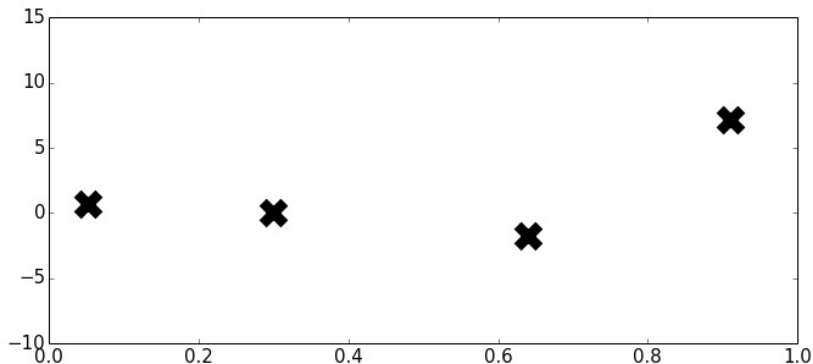
**Key question:**

Can we do better?

# Problem (the audience is encouraged to participate!)

- Find the optimum of some function $f$ in the interval [0,1].

- $f$ is L-Lipchitz continuous and differentiable.
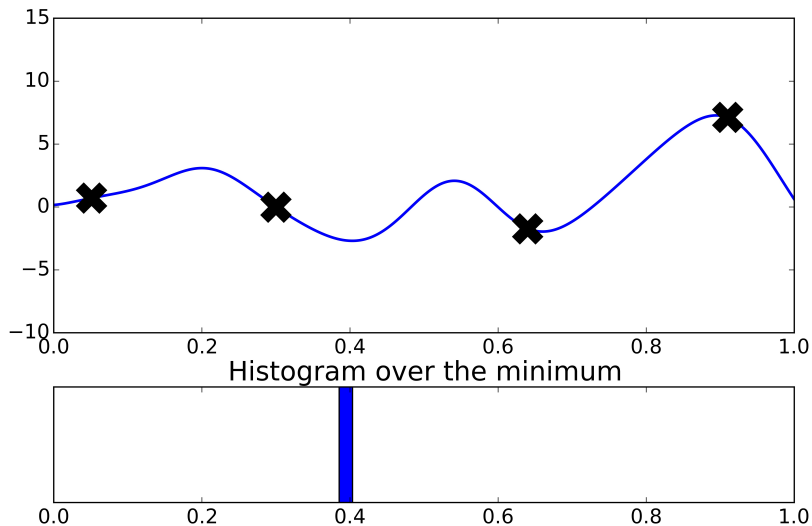
- Evaluations of $f$ are exact and we have 4 of them!

**Where is the minimum of f?**
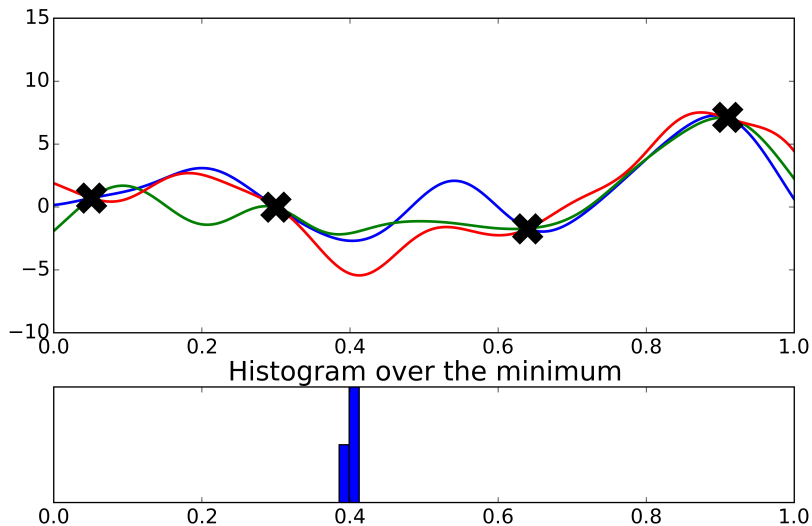**Where should the take the next evaluation?**

# Intuitive solution

One curve

# Intuitive solution

Three curves

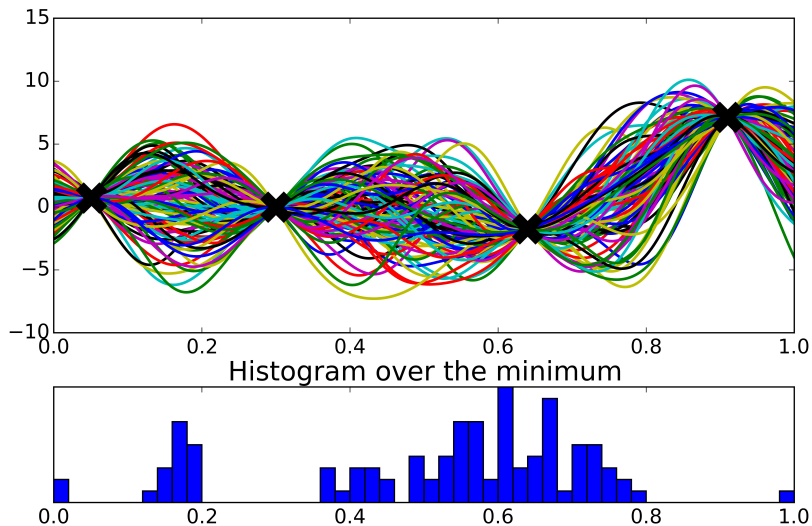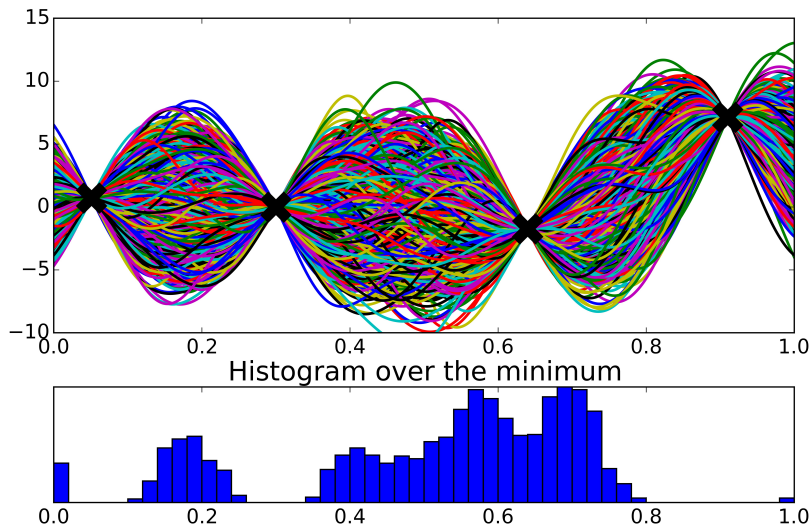# Intuitive solution

## Ten curves



Histogram over the minimum

# Intuitive solution

## Hundred curves

Histogram over the minimum

# General idea: surrogate modelling

1. Use a surrogate model of $f$ to carry out the optimization.

2. Define an utility function to collect new data points satisfying some optimality criterion: *optimization* as *decision*.

3. Study *decision* problems as *inference* using the surrogate model: use a probabilistic model able to calibrate both, epistemic and aleatoric uncertainty.

*Uncertainty Quantification*

# Utility functions

The utility should represent our design goal:.

1. Active Learning and experimental design: Maximize the differential entropy of the posterior distribution $p(f|X, y)$ (D-optimality in experimental design).

2. Minimize the loss in a sequence $x_1, \ldots, x_n$

$$r_N = \sum_{n=1}^{N} f(x_n) - N f(x_M)$$

(1) does to a lot exploration whereas (2) encourages exploitation about the minimum of $f$.

Methodology to perform global optimisation of multimodal black-box functions.

1. Choose some *prior measure* over the space of possible objectives $f$.

2. Combine prior and the likelihood to get a *posterior measure* over the objective given some observations.

3. Use the posterior to decide where to take the next evaluation according to some *acquisition/loss function*.

4. Augment the data.

Iterate between 2 and 4 until the evaluation budget is over.
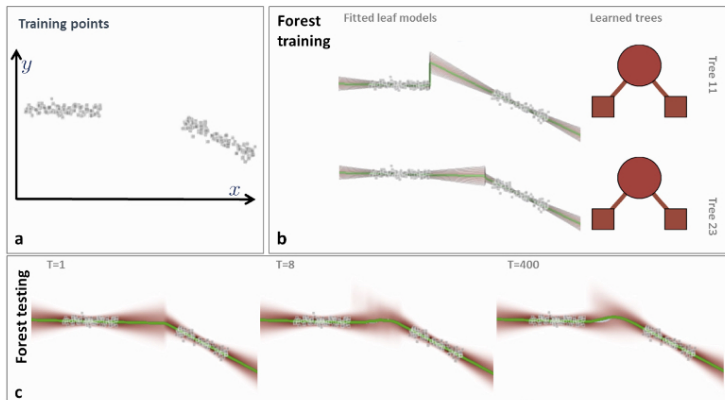
# Surrogate model: Gaussian process

Infinite-dimensional probability density, such that each linear finite-dimensional restriction is multivariate Gaussian.

- Model $f(x) \sim \mathcal{GP}(\mu(x), k(x, x'))$ is determined by the *mean function* $m(x)$ and *covariance function* $k(x, x'; \theta)$.
- Posterior mean $\mu(x; \theta, \mathcal{D})$ and variance $\sigma(x; \theta, \mathcal{D})$ can be *computed explicitly* given a dataset $\mathcal{D}$.

# Other models are also possible: Random Forrest

[Criminisi et al, 2011]

## Student-$t$ Processes as Alternatives to Gaussian Processes

**Amar Shah**
University of Cambridge

**Andrew Gordon Wilson**
University of Cambridge

**Zoubin Ghahramani**
University of Cambridge

### Abstract

We investigate the Student-$t$ process as an alternative to the Gaussian process as a non-parametric prior over functions. We derive closed form expressions for the marginal likelihood and predictive distribution of a Student-$t$ process, by integrating away an

simple exact learning and inference procedures, and impressive empirical performances [Rasmussen, 1996], Gaussian processes as kernel machines have steadily grown in popularity over the last decade.

At the heart of every Gaussian process (GP) is a parametrized covariance kernel, which determines the properties of likely functions under a GP. Typically simple parametric kernels, such as the Gaus-

# Exploration vs. exploitation



Bayesian optimization explains human active search

[Borji and Itti, 2013]
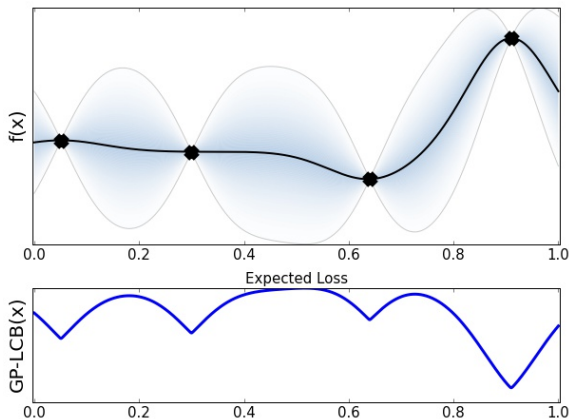
# Exploration vs. exploitation



Picture source: http://peakdistrictcycleways.co.uk

# GP Upper (lower) Confidence Band

[Srinivas et al., 2010]

Direct balance between exploration and exploitation:

$$\alpha_{LCB}(\mathbf{x}; \theta, \mathcal{D}) = -\mu(\mathbf{x}; \theta, \mathcal{D}) + \beta_t \sigma(\mathbf{x}; \theta, \mathcal{D})$$

# GP Upper (lower) Confidence Band

[Srinivas et al., 2010]

- ▶ In noiseless cases, it is a lower bound of the function to minimize.
- ▶ This allows to computer a bound on how close we are to the minimum.
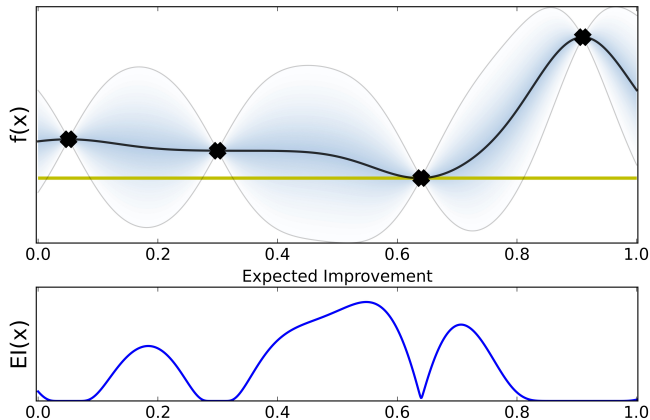- ▶ Optimal choices available for the 'regularization parameter'.

**Theorem 1** *Let* $\delta \in (0,1)$ *and* $\beta_t = 2\log(|D|t^2\pi^2/6\delta)$. *Running* GP-UCB *with* $\beta_t$ *for a sample* $f$ *of a GP with mean function zero and covariance function* $k(\boldsymbol{x}, \boldsymbol{x}')$, *we obtain a regret bound of* $\mathcal{O}^*(\sqrt{T\gamma_T\log|D|})$ *with high probability. Precisely, with* $C_1 = 8/\log(1+\sigma^{-2})$ *we have*

$$\Pr\left\{R_T \le \sqrt{C_1 T\beta_T\gamma_T} \quad \forall T \ge 1\right\} \ge 1-\delta.$$

# Expected Improvement

$$\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$$

# Expected Improvement

[Jones et al., 1998]

- ► Perhaps the most used acquisition.
- ► Explicit for available for Gaussian posteriors.
- ► It is too greedy in some problems. It is possible to make more explorative adding a 'explorative' parameter

$$\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \sigma(\mathbf{x}; \theta, \mathcal{D})(\gamma(x)\Phi(\gamma(x))) + \mathcal{N}(\gamma(x); 0, 1).$$

where

$$\gamma(x) = \frac{f(x_{best}) - \mu(\mathbf{x}; \theta, \mathcal{D}) + \psi}{\sigma(\mathbf{x}; \theta, \mathcal{D})}.$$

# Maximum Probability of Improvement

$$\gamma(\mathbf{x}) = \sigma(\mathbf{x}; \theta, \mathcal{D})^{-1}(\mu(\mathbf{x}; \theta, \mathcal{D}) - y_{best})$$

$$\alpha_{MPI}(\mathbf{x}; \theta, \mathcal{D}) = p(f(\mathbf{x}) < y_{best}) = \Phi(\gamma(\mathbf{x}))$$

# Maximum Probability of Improvement

- First used acquisition: very intuitive.
- Less used in practice.
- Explicit for available for Gaussian posteriors.
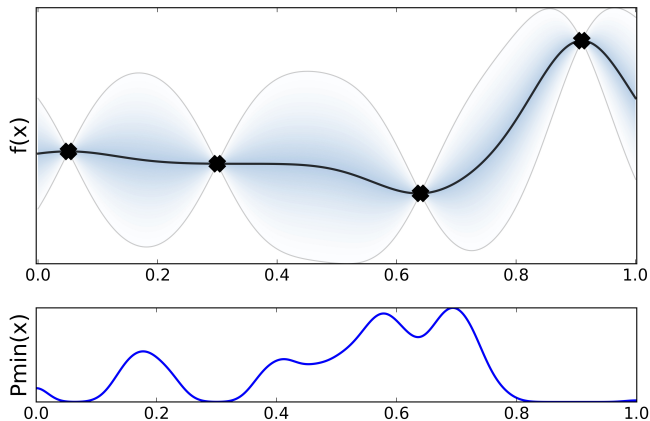
$$\alpha_{MPI}(\mathbf{x}; \theta, \mathcal{D}) = \Phi(\gamma(x))).$$

where

$$\gamma(x) = \frac{f(x_{best}) - \mu(\mathbf{x}; \theta, \mathcal{D}) + \psi}{\sigma(\mathbf{x}; \theta, \mathcal{D})}.$$

# Information-theoretic approaches

[Hennig and Schuler, 2013; Hernández-Lobato et al., 2014]

$$\alpha_{ES}(\mathbf{x}; \theta, \mathcal{D}) = H[p(x_{min}|\mathcal{D})] - \mathbb{E}_{p(y|\mathcal{D},\mathbf{x})}[H[p(x_{min}|\mathcal{D} \cup \{\mathbf{x}, y\})]]]$$

# Information-theoretic approaches

Uses the distribution of the minimum

$$p_{min}(x) \equiv p[x = \arg \min f(x)] = \int_{f:I \to \Re} p(f) \prod_{\substack{\tilde{x} \in I \\ \tilde{x} \neq x}} \theta[f(\tilde{x}) - f(x)] df$$

where $\theta$ is the Heaviside's step function. No closed form!

Use Thomson sampling to approximate the distribution.
Generate many sample paths from the GP, optimize them to
take samples from $p_{min}(x)$.

# Thomson sampling
## Probability matching

$$\alpha_{THOMSON}(\mathbf{x}; \theta, \mathcal{D}) = g(\mathbf{x})$$

$g(\mathbf{x})$ is sampled form $\mathcal{GP}(\mu(x), k(x, x'))$

# Thompson sampling
Probability matching [Rahimi and B. Recht, 2007]

- It is easy to generate posterior samples of a GP at a finite set of locations.
- More difficult is to generate 'continuous' samples.

Possible using the Bochner's lemma: existence of the Fourier dual of $k$, $s(\omega)$ which is equal to the spectral density of $k$

$$k(x, x') = \nu \mathbb{E}_\omega \left[ e^{-i\omega^T(x-x')} \right] = 2\nu \mathbb{E}_{\omega,b} \left[ \cos(\omega x^T + b) \cos(\omega x^T + b) \right]$$

With sampling and this lemma (taking $p(w) = s(\omega)/\nu$ and $b \sim \mathcal{U}[0, 2\pi]$) we can construct a feature based approximation for sample paths of the GP.

$$k(x, x') \approx \frac{\nu}{m} \sum_{i=1}^{m} e^{-i\omega^{(i)T}x} e^{-i\omega^{(i)T}x'}$$

# The choice of utility matters

The choice of the utility may change a lot the result of the optimisation.

# The choice of utility in practice
[Hoffman, Shahriari and de Freitas, 2013]



The best utility depends on the problem and the level of
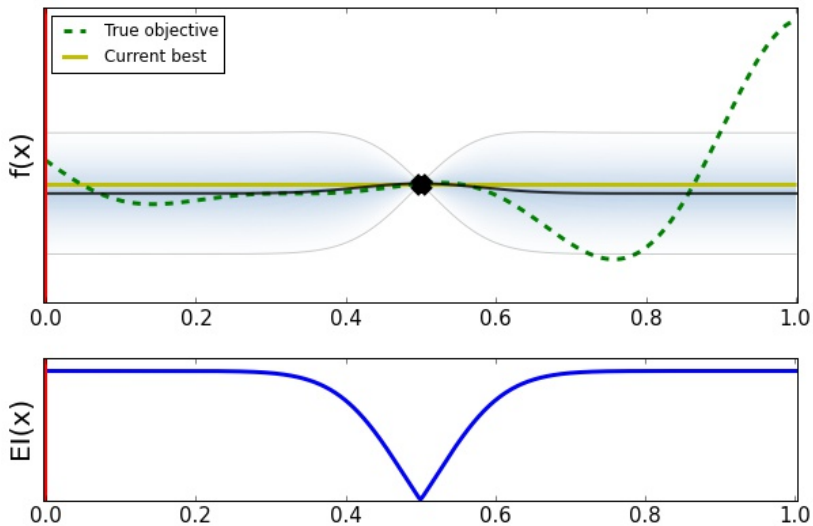exploration/exploitation required.
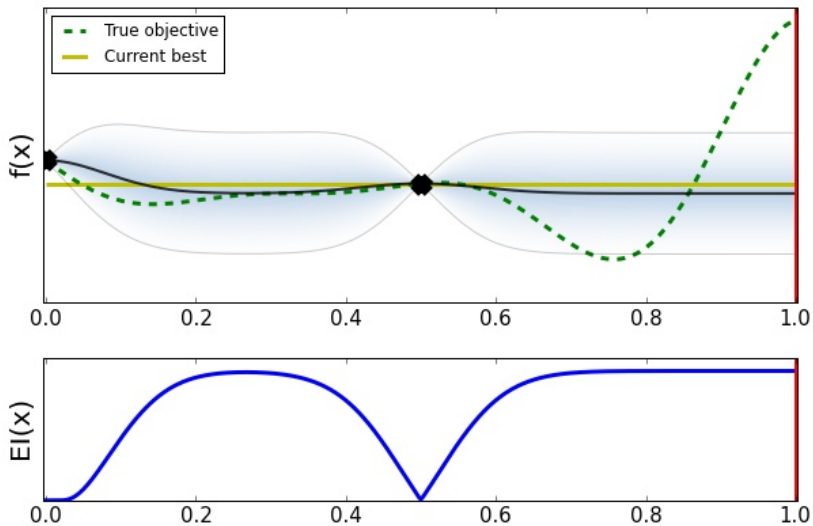
# Illustration of BO

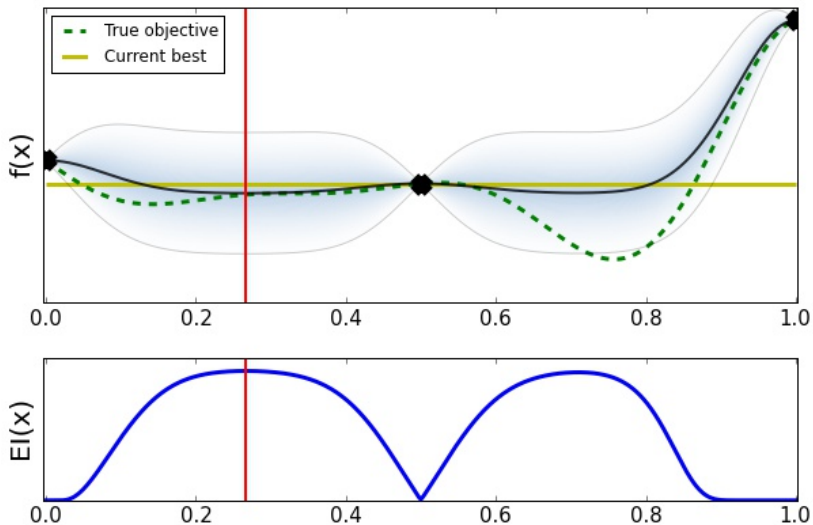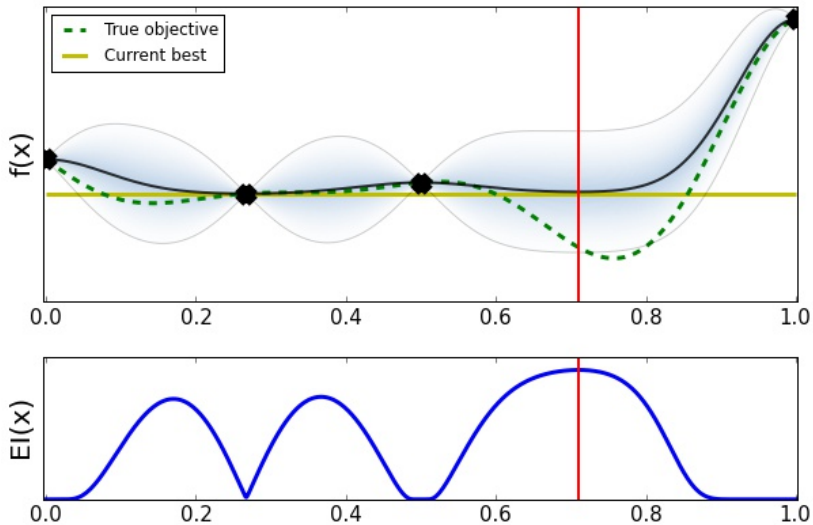# Illustration of BO
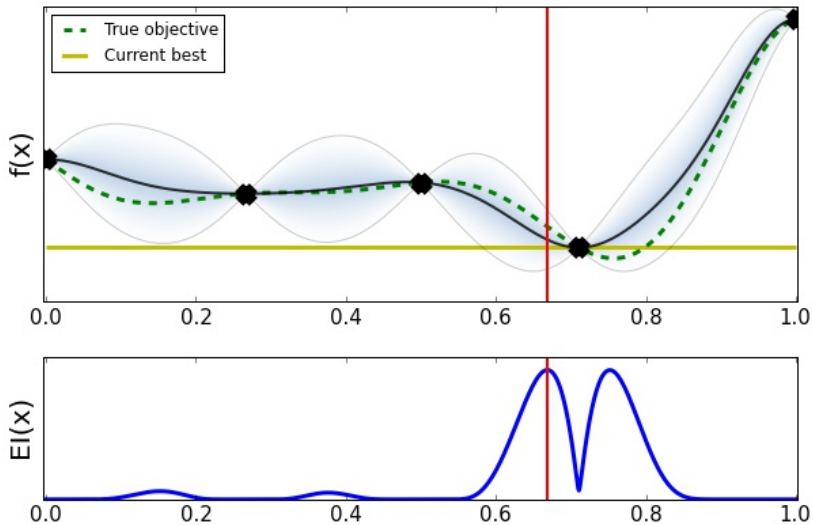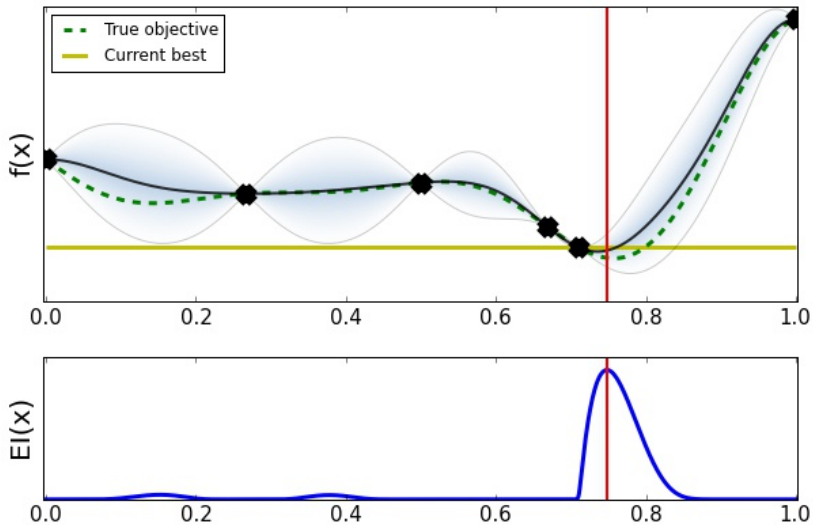
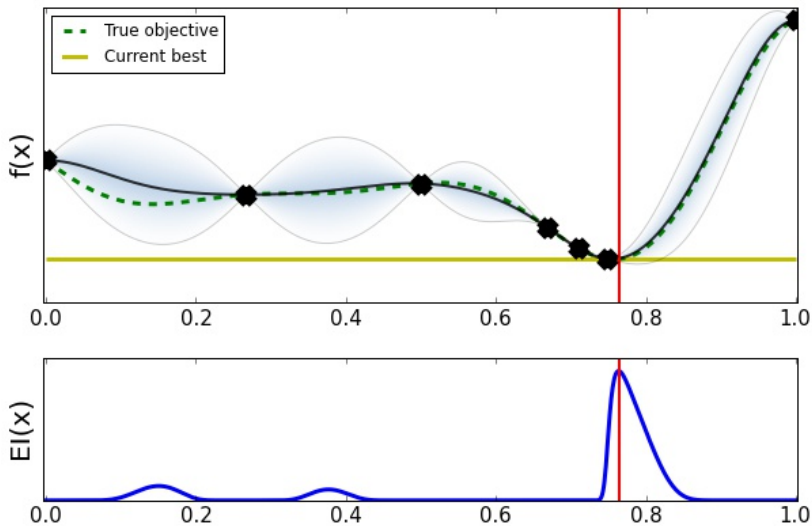# Illustration of BO

# Illustration of BO

# Illustration of BO

# Illustration of BO

# Illustration of BO

BO is an strategy to transform the problem

$$x_M = \arg\min_{x \in \mathcal{X}} f(x)$$
$$\textit{solvable}!$$

into a series of problems:

$$x_{n+1} = \arg\max_{x \in \mathcal{X}} \alpha(x; \mathcal{D}_n, \mathcal{M}_n)$$
$$\textit{solvable}!$$

where now:

- $\alpha(x)$ is inexpensive to evaluate.
- The gradients of $\alpha(x)$ are typically available.
- Still need to find $x_{n+1}$.

# BO vs other methods

[Osborne et al, 2009]

Bayesian optimization works better in practice!

| | EGO | RBF | DIRECT | GPGO 1-Step Non-Periodic | Periodic | GPGO 2-Step Non-Periodic |
|------|-------|-------|--------|--------|--------|--------|
| Br | 0.943 | 0.960 | 0.958 | **0.980** | — | — |
| C6 | 0.962 | 0.962 | 0.940 | 0.890 | — | **0.967** |
| G–P | 0.783 | 0.815 | **0.989** | 0.804 | — | **0.989** |
| H3 | 0.970 | 0.867 | 0.868 | **0.980** | — | — |
| H6 | 0.837 | 0.701 | 0.689 | **0.999** | — | — |
| Sh5 | 0.218 | 0.092 | 0.090 | **0.485** | — | — |
| Sh7 | 0.159 | 0.102 | 0.099 | **0.650** | — | — |
| Sh10 | 0.135 | 0.100 | 0.100 | **0.591** | — | — |
| GK2 | 0.571 | 0.567 | 0.538 | **0.643** | — | — |
| GK3 | 0.519 | 0.207 | 0.368 | **0.532** | — | — |
| Shu | **0.492** | 0.383 | 0.396 | 0.437 | 0.348 | 0.348 |
| G2 | 0.979 | **1.000** | 0.981 | **1.000** | **1.000** | — |
| G5 | **1.000** | 0.998 | 0.908 | 0.925 | 0.957 | — |
| A2 | 0.347 | 0.703 | 0.675 | 0.606 | 0.612 | **0.781** |
| A5 | 0.192 | **0.381** | 0.295 | 0.089 | 0.161 | — |
| R | 0.652 | 0.647 | 0.776 | 0.675 | **0.933** | — |
| mean | 0.610 | 0.593 | 0.604 | **0.705** | — | — |

# Recap

- Bayesian optimization is a way of encoding our beliefs about a property of a function (the minimum)

- Two key elements: the model and the acquisition function.

- Many choices in both cases, especially in terms of the acquisition function used.

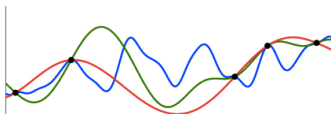- The key is to find a good balance between exploration and exploitation.

- ▶ What to do with the hyper-parameters of the model?

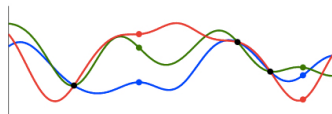- ▶ How to select points to initialize the model?

- ▶ How to optimize the acquisition function?

# BO independent of the parameters of the GP.

[Snoek et al. 2012]

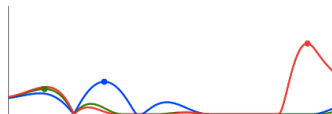Integrate out across parameter values or location outputs.



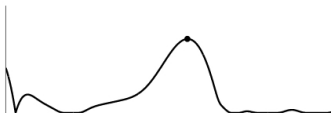(a) Posterior samples under varying hyperparameters
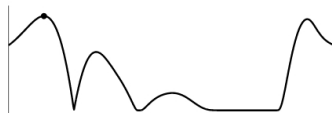
(a) Posterior samples after three data

(b) Expected improvement under varying hyperparameters

(b) Expected improvement under three fantasies

(c) Integrated expected improvement

(c) Expected improvement across fantasies

# How to initialise the model?

- One point in the centre of the domain.
- Uniformly selected random locations.
- Latin design.
- Halton sequences.
- Determinantal point processes.

The idea is always to start at some locations trying to minimise the initial model uncertainty.

# Latin design

$n \times n$ array filled with $n$ different symbols, each occurring
exactly once in each row and exactly once in each column.

| A | B | F | C | E | D |
|---|---|---|---|---|---|
| B | C | A | D | F | E |
| C | D | B | E | A | F |
| D | E | C | F | B | A |
| E | F | D | A | C | B |
| F | A | E | B | D | C |

# pyDOE

Python framework for standard experimental design

# Latin design

Window honors Ronald Fisher. Fisher's student, A. W. F. Edwards, designed this window for Caius College, Cambridge.

# Halton sequences

- ► Used to generate points in $(0, 1) \times (0, 1)$
- ► Sequence that is constructed according to a deterministic method that uses a prime number as its base.
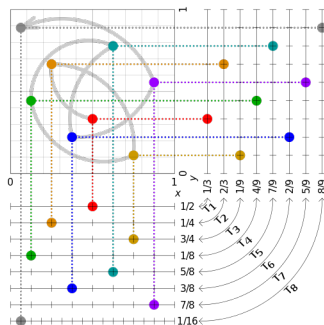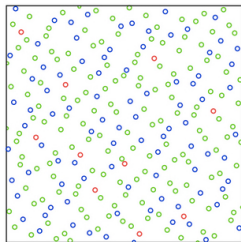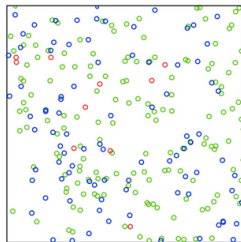


Figure source: Wikipedia

# Halton sequences
[Halton, 1964]

Better coverage than random.



*Halton*          *Random*

Figure source: Wikipedia

We say that $X$ is a 'determinantal point process' on $\Lambda$ with kernel $K$ if it is a simple point process on $\Lambda$ with a joint intensity or 'correlation function' given by

$$\rho_n(x_1, \ldots, x_n) = det(K(x_i, x_j)_{1 \le i,j \le n})$$

- Probability measures over subsets.
- Possible to characterise the samples in terms of quality and diversity.

# Determinantal point processes
Kulesza and Taskar, [2012]



DPP      Independent

Key idea:

$$\mathcal{P}(i,j \in \boldsymbol{Y}) = \left| \begin{array}{cc} K_{ii} & K_{ij} \\ K_{ji} & K_{jj} \end{array} \right|$$
$$= K_{ii}K_{jj} - K_{ij}K_{ji}$$
$$= \mathcal{P}(i \in \boldsymbol{Y})\mathcal{P}(j \in \boldsymbol{Y}) - K_{ij}^2.$$

# Methods to optimise the acquisition function

This may not be easy.

- ▶ Gradient descent methods: Conjugate gradient, BFGS, etc.
- ▶ Lipschitz based heuristics: DIRECT.
- ▶ Evolutionary algorithms: CMA.

Some of these methods can also be used to directly optimize $f$

# Gradient descent

[Avriel,2013], but many others

---

**Algorithm 2**: Gradient Descent

**input** : $f : \mathbb{R}^n \to \mathbb{R}$ a differentiable function
  $\mathbf{x}^{(0)}$ an initial solution

**output**: $\mathbf{x}^\star$, a local minimum of the cost function $f$.

1 **begin**
2    $k \leftarrow 0$ ;
3    **while** STOP-CRIT **and** $(k < k_{max})$ **do**
4      $\mathbf{x}^{(k+1)} \leftarrow \mathbf{x}^{(k)} - \alpha^{(k)} \boldsymbol{\nabla} f(\mathbf{x})$ ;
5      with $\alpha^{(k)} = \arg \min_{\alpha \in \mathbb{R}_+} f(\mathbf{x}^{(k)} - \alpha \boldsymbol{\nabla} f(\mathbf{x}))$ ;
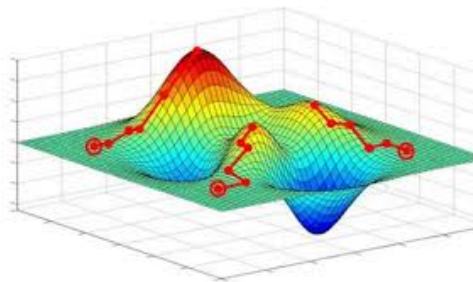6      $k \leftarrow k + 1$ ;
7    **return** $\mathbf{x}^{(k)}$
8 **end**

---

We need to know the gradients. This is the case for most
acquisitions but not for all of them (PES for instance).

# Gradient descent



May fall in local minima if the function is multimodal: multiple initializations.
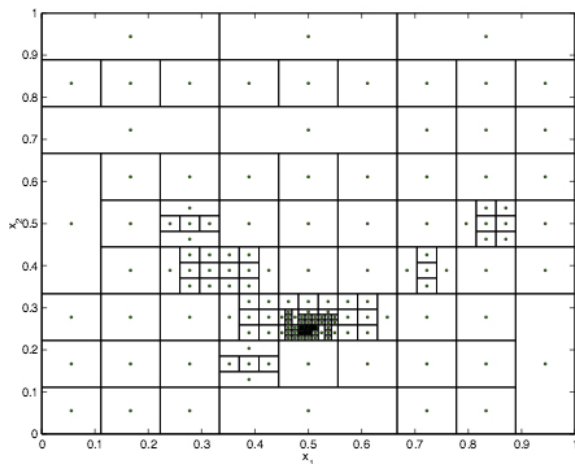
# 'DIviding RECTangles', DIRECT

---

**Algorithm** DIRECT('myfcn',bounds,opts)

1: Normalize the domain to be the unit hyper-cube with center $c_1$
2: Find $f(c_1)$, $f_{min} = f(c_1)$, $i = 0$ , $m = 1$
3: Evaluate $f(c_1 \pm \delta e_i$ , $1 \leq i \leq n$, and divide hyper-cube
4: **while** $i \leq maxits$ and $m \leq maxevals$ **do**
5:     Identify the set $S$ of all pot. optimal rectangles/cubes
6:     **for** all $j \in S$
7:         Identify the longest side(s) of rectangle $j$
8:         Evaluate myfcn at centers of new rectangles, and divide $j$ into smaller rectangles
9:         Update $f_{min}$, $xatmin$, and $m$
10:     **end for**
11:     $i = i + 1$
12: **end while**

---

Minimal hypothesis about the acquisition

# 'DIviding RECTangles', DIRECT
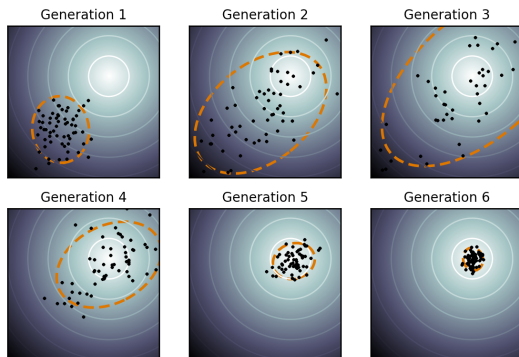
[Perttunen at al. 1993]



Finds good solution in general and doesn't need gradient. Not generalizable to non-squared domains.

# Covariance Matrix Adaptation, CMA

[Hansen and Ostermeier, 2001].

- ▶ Sample for a Gaussian with some mean $\mu$ and covariance matrix $\Sigma$.
- ▶ Select the best points and use them to update $\mu$ and $\Sigma$.
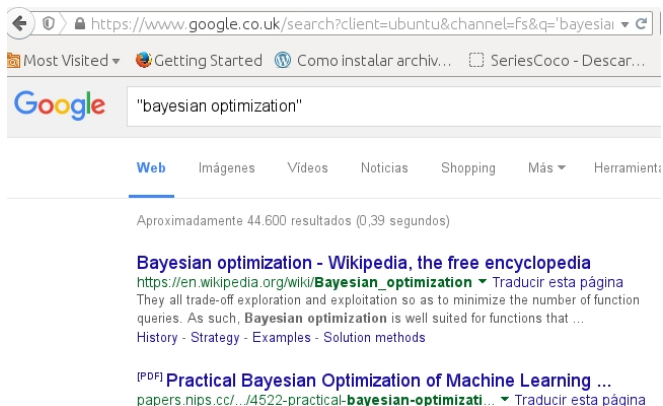- ▶ Sample form the new Gaussian.

# Took a while to start using these ideas in ML

Although in the stats community have been there for a while

- BO depends on its own parameters.

- Lack of software to apply these methods as a black optimization boxes.

- Reduced scalability in dimensions and number of evaluations (this is still a problem).

*Practical Bayesian Optimization of Machine Learning Algorithms.* Snoek, Larochelle and Adams. NIPS 2012 (Spearmint)

# Increasing popular field



- Hot topic in Machine Learning.
- The BO workshop at NIPS is well stablished and it is a mini-conference itself.

*It has become increasingly popular since it allows to configure algorithms without human intervention.*

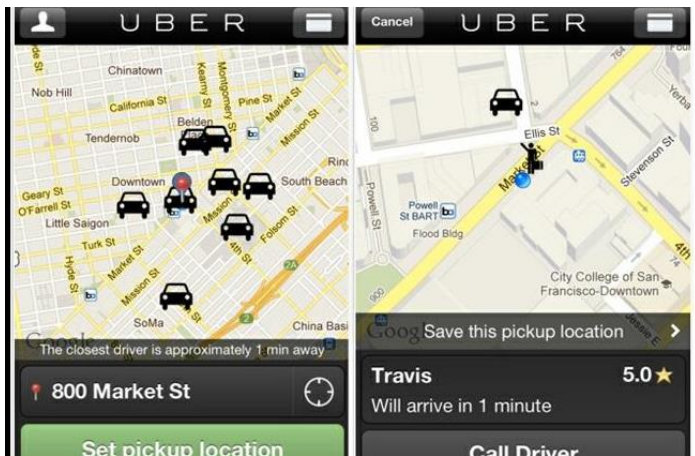BO takes to human out of the loop!

We have joined forces with Twitter!

W + 🐦

We have created a technology to make machine learning better and faster for companies, automatically. Twitter is *the* platform for open communication on the internet and we believe that Whetlab's technology can have a great impact by accelerating Twitter's internal machine learning efforts.

Questions?