

# Causal Bayesian optimization

*'Don't do everything, just do the right thing'*

---

Javier González

July 8, 2020

Microsoft Research Cambridge



Virginia Aglietti



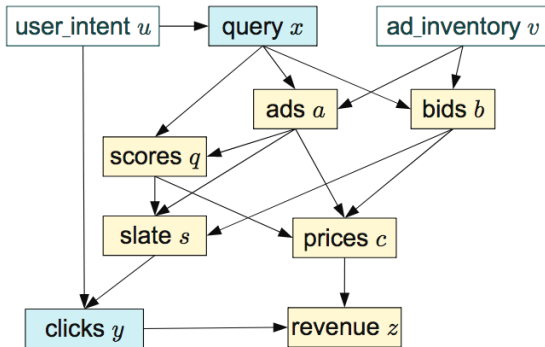
Xiaoyu Lu



Andrei Paleyes

# Systems/processes decompose in sets of interconnected nodes

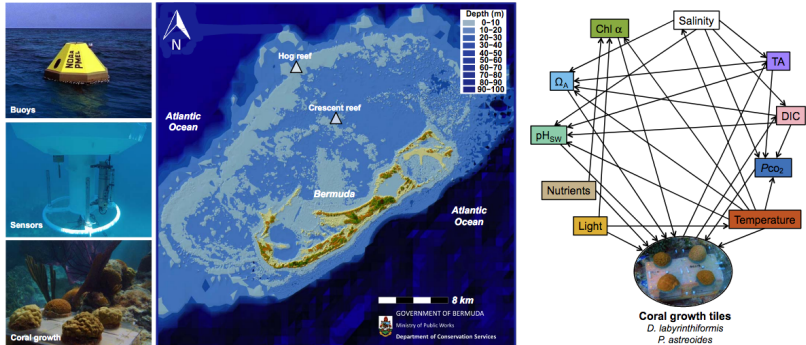
## Example in advertising:



(Bottou et al, 2013) The goal is to design advertising campaigns to *maximize* revenue.

# Systems/processes decompose in sets of interconnected nodes

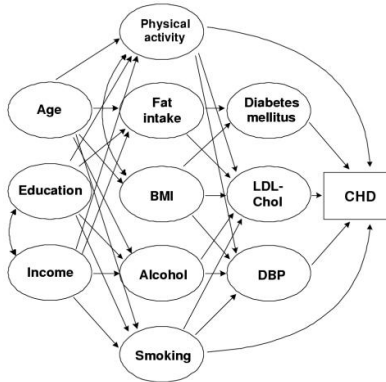
## Example in ecology:



(Courtney et al, 2017) The goal is to apply policies to *improve/maximise* coral calcification.

# Systems/processes decompose in sets of interconnected nodes

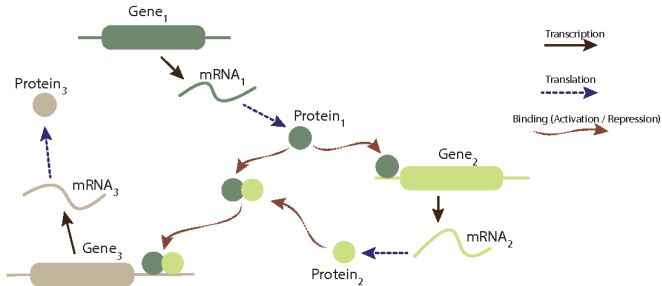
## Example in healthcare:



(Murray et al, 2003) The goal is to define treatments to *minimize* the risk of coronary heart disease (CHD)

# Systems/processes decompose in sets of interconnected nodes

## Example in biology:



(González, 2015; Maksimov, 2015) The goal is to target is to *maximize* the synthetic production of a protein of pharmacological interest.

# Common elements in these problems

- A causal graph.
- Observational data from all (non hidden) nodes.
- Ability of running experiments (in reality or in simulation).
- Cost of experiments depends on the number and type of nodes in which we intervene.

## Common goal

Find the *system/process configuration* that optimises the *target node*.

- *System/process configuration* → actionable variables.
- *Target node* → revenue, coral calcification, risk of disease, etc.

## Take home messages:

---

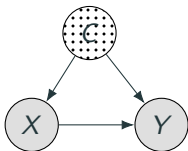
1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.



# Crash course in causal models and do-calculus (1 of 4)

**Causal model:** Directed acyclic graph  $\mathcal{G}$  + four-tuple  $\langle \mathbf{U}, \mathbf{V}, F, P(\mathbf{U}) \rangle$

- $\mathbf{U}$ : independent *exogenous* background variables.
- $P(\mathbf{U})$  distribution of  $\mathbf{U}$ .
- $\mathbf{V}$ : *endogenous* variables (non-manipulative  $\mathbf{C}$ , treatment  $\mathbf{X}$ ).
- $F = \{f_1, \dots, f_{|\mathbf{V}|}\}$ : functions  $v_i = f_i(pa_i, u_i)$ ,  $pa_i$  are the parents of  $V_i$ .



$$C = f_c(U_c), U_c \sim \mathcal{N}(0, \sigma_c^2)$$

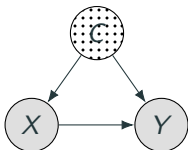
$$X = f_x(C, U_x), U_x \sim \mathcal{N}(0, \sigma_x^2)$$

$$Y = f_y(X, C, U_y), U_y \sim \mathcal{N}(0, \sigma_y^2)$$

## Crash course in causal models and do-calculus (2 of 4)

**Intervention:** Setting a manipulative variable  $X$  to a value  $x$ ,  $do(X = x)$ .

*Observed universe*



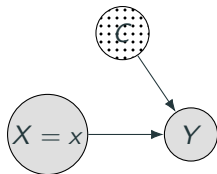
$$C = f_c(U_c)$$

$$X = f_x(C, U_x)$$

$$Y = f_y(X, C, U_y)$$

$$P(X, C, Y)$$

*Post-interventional universe*



$$C = f_c(U_c)$$

$$X = x$$

$$Y = f_y(x, C, U_y)$$

$$P^{do(X=x)}(C, Y)$$

$$P(Y|do(X = x)) := P^{do(X=x)}(Y|X = x)$$

# Crash course in causal models and do-calculus (3 of 4)

**Key question:** *How to do inference in the post-interventional universe with data from the observed universe.*

Observing vs. doing:

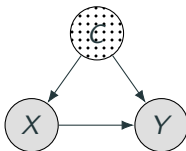
- $P(Y|\text{do}(X = x), C)$  requires change the way the universe works.
- $P(Y|X = x, C)$  only requires 'observing' the universe.

**do-calculus:** algebra to emulate the post-intervention universe in terms of conditionals in the observed universe (experiments emulation).

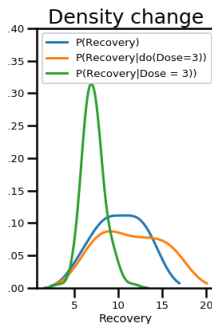
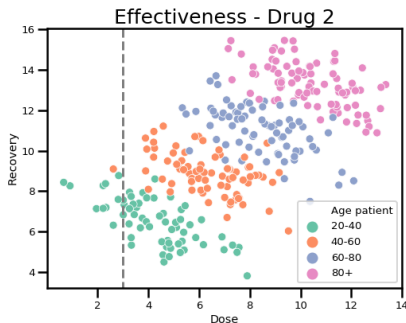
$$P(\mathbf{V}) = \prod_{i=1}^{|\mathbf{V}|} p(V_i | pa_i) \quad (\text{Markov condition})$$

# Crash course in causal models and do-calculus (4 of 4)

Do-calculus: back-door adjustment:



$$p(Y|do(X = x)) = \int P(Y|C, X)P(C)dC$$



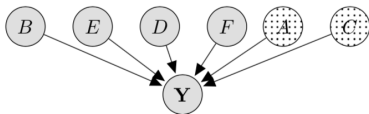
## Take home messages:

---

1. Many real systems decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.

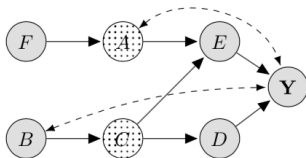
# Gobal optimization vs. Causal optimization

*Global optimization*



$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in D(\mathbf{X})} \mathbb{E}_{P(\mathbf{Y}|\text{do}(\mathbf{X}=\mathbf{x}), \mathbf{C})}[\mathbf{Y}]$$

*Causal optimization*



$$\mathbf{X}_s^*, \mathbf{x}_s^* = \arg \min_{\substack{\mathbf{X}_s \in \mathcal{P}(\mathbf{X}) \\ \mathbf{x}_s \in D(\mathbf{X}_s)}} \mathbb{E}_{P(\mathbf{Y}|\text{do}(\mathbf{X}_s=\mathbf{x}_s), \mathbf{C})}[\mathbf{Y}]$$

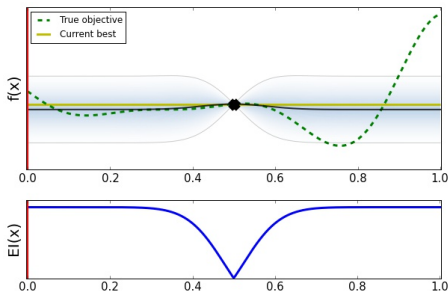
## Global optimization

- $f$ , the objective function, is explicitly unknown and multimodal.
- Evaluations of  $f$  may be perturbed by noise.
- Evaluations of  $f$  are expensive.

Standard method for this scenario → *Bayesian optimization*

# Bayesian optimization

- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$

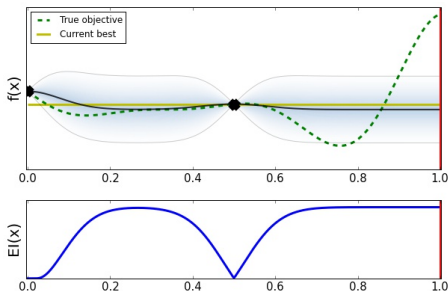


Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$



# Bayesian optimization

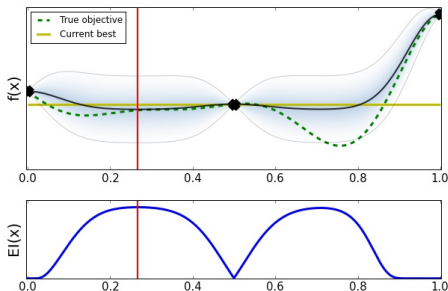
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

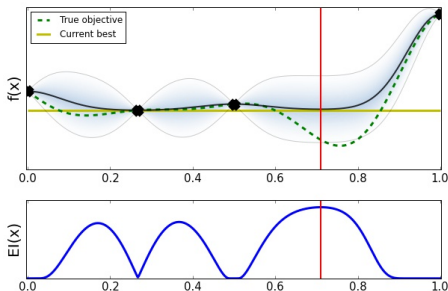
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

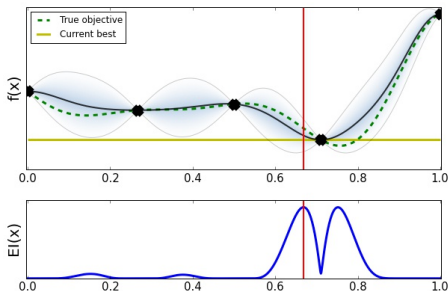
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

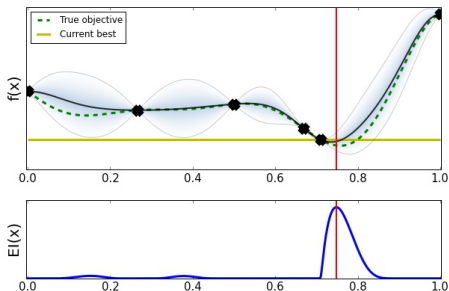
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

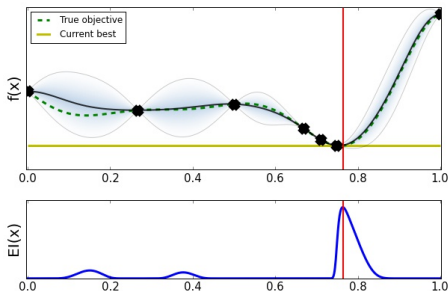
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

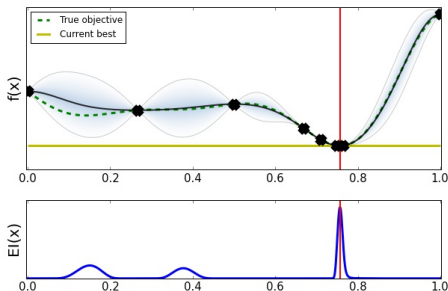
- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Bayesian optimization

- **Goal:** Collect data  $x_1, \dots, x_n$  to find the optimum as fast as possible.
- **Model:** Gaussian process  $f(x) \sim \mathcal{GP}(\mu(x), k_\theta(x, x'))$ .
- **Acquisition:**  $\alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{best} - y) p(y|\mathbf{x}; \theta, \mathcal{D}) dy$



Each point  $x_{n+1}$  is collected as  $x_{n+1} = \arg \max \alpha_{EI}(\mathbf{x}; \theta, \mathcal{D}_n)$

# Assumptions

## Causal optimization

- $f$  is explicitly unknown and multimodal.
- Evaluations of  $f$  may be perturbed by noise.
- Evaluations of  $f$  are expensive.

+

- Causal graph
- Cost associate to experiment with each variable.

New method for this scenario → ***Causal Bayesian optimization***



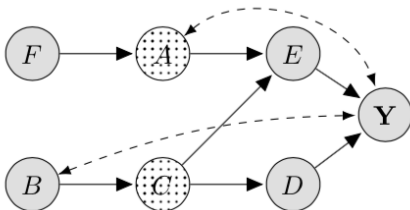
## Take home messages:

---

1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.
4. Standard Bayesian Optimization ignores causal assumptions.
5. Causal Optimization requires a new approach.

# Goal for the new setup

- $\mathbf{X} = \{B, E, D, F\}$ : treatment variables
- $\mathcal{P}(\mathbf{X})$ , all possible combinations of interventions.
- $\mathbf{X}_s, \mathbf{x}_s$ , intervention set and its value.
- $\mathbf{X}_s^*, \mathbf{x}_s^*$ , optimal intervention set and value.



**Goal:** Run interventions  $(\mathbf{X}_{s_1}, \mathbf{x}_{s_1}), \dots, (\mathbf{X}_{s_n}, \mathbf{x}_{s_n})$  to find the optimum as fast as possible.

# Do we need to find $\mathbf{X}_s^*$ in the $2^{|\mathbf{X}|}$ sets in $\mathcal{P}(\mathbf{X})$ ? NO!

## Minimal Intervention Set (MIS, $\mathbb{M}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ )

Given  $\langle \mathcal{G}, \mathbf{Y}, \mathbf{X}, \mathbf{C} \rangle$ , a set  $\mathbf{X}_s \in \mathcal{P}(\mathbf{X})$  is said to be a MIS if there is no  $\mathbf{X}'_s \subset \mathbf{X}_s$  such that  $\mathbb{E}[Y | \text{do}(\mathbf{X}_s = \mathbf{x}_s), \mathbf{C}] = \mathbb{E}[Y | \text{do}(\mathbf{X}'_s = \mathbf{x}'_s), \mathbf{C}]$ .

## Possibly-Optimal Minimal Intervention set (POMIS, $\mathbb{P}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ )

Let  $\mathbf{X}_s \in \mathbb{M}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ .  $\mathbf{X}_s$  is a POMIS if there exists a sem conforming to  $\mathcal{G}$  such that  $\mathbb{E}[Y | \text{do}(\mathbf{X}_s = \mathbf{x}_s^*), \mathbf{C}] > \forall \mathbf{w} \in \mathbb{M}_{\mathcal{G}, \mathbf{Y} \setminus \mathbf{x}_s}^{\mathbf{C}} \mathbb{E}[Y | \text{do}(\mathbf{W} = \mathbf{w}^*), \mathbf{C}]$  where  $\mathbf{x}^*$  and  $\mathbf{w}^*$  denote the optimal intervention values.

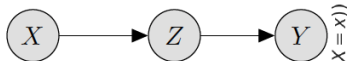
- BO,  $\mathbb{B}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ : all treatment variables
- MIS,  $\mathbb{M}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ : set of variables 'worth' intervening.
- POMIS,  $\mathbb{P}_{\mathcal{G}, \mathbf{Y}}^{\mathbf{C}}$ : set of variables in MIS that always improve  $Y$ .

## Take home messages:

---

1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.
4. Standard Bayesian Optimization ignores causal assumptions.
5. Causal Optimization requires a new approach.
6. No need to explore  $\mathcal{P}(\mathbf{X})$  to solve Causal Optimization problems.

# Toy example



$$X = \epsilon_X$$

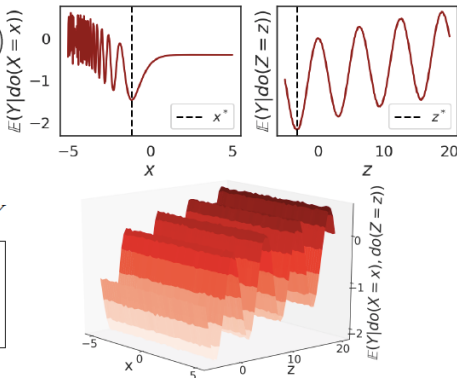
$$Z = \exp(-X) + \epsilon_Z$$

$$Y = \cos(Z) - \exp\left(-\frac{Z}{20}\right) + \epsilon_Y$$

$$\mathbb{M}_{\mathcal{G}, Y} = \{\emptyset, \{X\}, \{Z\}\}$$

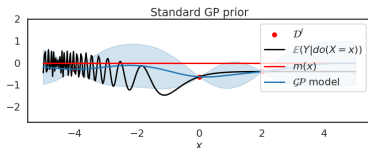
$$\mathbb{P}_{\mathcal{G}, Y} = \{\{Z\}\}$$

$$\mathbb{B}_{\mathcal{G}, Y} = \{\{X, Z\}\}$$



# Modelling $\mathbb{E}_{P(Y|\text{do}(\mathbf{X}_s=\mathbf{x}_s))}[\mathbf{Y}]$ for each $\mathbf{X}_s$

## Standard Gaussian process



$$f(\mathbf{x}_s) \sim \mathcal{GP}(m(\mathbf{x}_s), k(\mathbf{x}_s, \mathbf{x}'_s))$$

$$m(\mathbf{x}_s) = \mathbf{0}$$

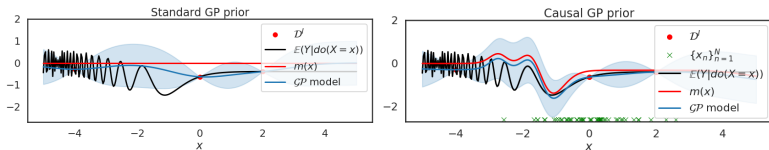
$$k(\mathbf{x}_s, \mathbf{x}'_s) = k_{RBF}(\mathbf{x}_s, \mathbf{x}'_s)$$

where

- $k_{RBF}(\mathbf{x}_s, \mathbf{x}'_s) := \exp\left(-\frac{\|\mathbf{x}_s - \mathbf{x}'_s\|^2}{2l^2}\right)$

# Modelling $\mathbb{E}_{P(Y|\text{do}(\mathbf{X}_s=\mathbf{x}_s))}[Y]$ for each $\mathbf{X}_s$

## Causal Gaussian process



$$f(\mathbf{x}_s) \sim \mathcal{GP}(m(\mathbf{x}_s), k(\mathbf{x}_s, \mathbf{x}'_s))$$

$$m(\mathbf{x}_s) = \hat{\mathbb{E}}[Y|\text{do}(\mathbf{X}_s = \mathbf{x}_s)]$$

$$k(\mathbf{x}_s, \mathbf{x}'_s) = k_{RBF}(\mathbf{x}_s, \mathbf{x}'_s) + \sigma(\mathbf{x}_s)\sigma(\mathbf{x}'_s)$$

where

- $k_{RBF}(\mathbf{x}_s, \mathbf{x}'_s) := \exp(-\frac{\|\mathbf{x}_s - \mathbf{x}'_s\|^2}{2l^2})$
- $\sigma(\mathbf{x}_s) = \sqrt{\hat{\mathbb{V}}(Y|\text{do}(\mathbf{X}_s = \mathbf{x}_s))}$  with  $\hat{\mathbb{V}}$  is the variance of the causal estimated from observational data.

## Take home messages:

---

1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.
4. Standard Bayesian Optimization ignores causal assumptions.
5. Causal Optimization requires a new approach.
6. No need to explore  $\mathcal{P}(\mathbf{X})$  to solve Causal Optimization problems.
7. Causal GPs to merge observational and interventional data.



# Causal Expected Improvement (CEI)

Expected improvement in the intervention set:

- $y_s = \mathbb{E}[Y | \text{do}(\mathbf{X}_s = \mathbf{x}_s), \mathbf{C}]$
- $y^* = \max_{\mathbf{X}_s \in \text{es}, \mathbf{x} \in D(\mathbf{X}_s)} \mathbb{E}[Y | \mathbf{X}_s = \mathbf{x}_s, \mathbf{C}],$

$$EI^s(\mathbf{x}) = \mathbb{E}_{p(y_s)}[\max(y_s - y^*, 0)] / Co(\mathbf{x}).$$

- $\alpha_1, \dots, \alpha_{|\text{es}|}$ : solutions of optimizing  $EI^s(\mathbf{x})$  for each set in **es** and

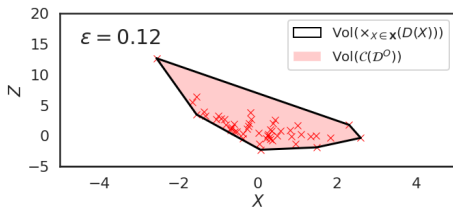
## New intervention set and value

$$\alpha^* := \max\{\alpha_1, \dots, \alpha_{|\text{es}|}\}$$

$$s^* = \underset{s \in \{1, \dots, |\text{es}|\}}{\operatorname{argmax}} \alpha_s.$$

# Intervention-observation trade off

$\epsilon$ -greedy criteria to balance interventions and observations



$$\epsilon = \frac{\text{Vol}(C(D^O))}{\text{Vol}(X_{X \in \mathbf{x}}(D(X)))} \times \frac{N}{N_{\max}},$$

- $\text{Vol}(C(D^O))$ : volume of the convex hull for observational data
- $\text{Vol}(X_{X \in \mathbf{x}}(D(X)))$ : volume of the interventional domain.

# Causal Bayesian Optimization (CBO)

---

**Algorithm 1:** Causal Bayesian Optimization-CBO

---

**Data:**  $\mathcal{D}^O$ ,  $\mathcal{D}^I$ ,  $\mathcal{G}$ ,  $\mathbf{ES}$ , number of steps  $T$

**Result:**  $\mathbf{X}_s^*$ ,  $\mathbf{x}_s^*$ ,  $\hat{\mathbb{E}}[\mathbf{Y}^* | \text{do}(\mathbf{X}_s^* = \mathbf{x}_s^*), \mathbf{C}]$

**Initialise:** Set  $\mathcal{D}_0^I = \mathcal{D}^I$  and  $\mathcal{D}_0^O = \mathcal{D}^O$

**for**  $t=1, \dots, T$  **do**

    Compute  $\epsilon$  and sample  $u \sim \mathcal{U}(0, 1)$

**if**  $\epsilon > u$  **then**

        (Observe)

1. Observe new observations  $(\mathbf{x}_t, c_t, \mathbf{y}_t)$ .
2. Augment  $\mathcal{D}^O = \mathcal{D}^O \cup \{(\mathbf{x}_t, c_t, \mathbf{y}_t)\}$ .
3. Update prior of the causal GP (Eq. (2)).

**end**

**else**

        (Intervene)

1. Compute  $EI^s(\mathbf{x})/Co(\mathbf{x})$  for each element  $s \in \mathbf{ES}$  (Eq. (5)).
2. Obtain the optimal interventional set-value pair  $(s^*, \alpha^*)$ .
3. Intervene on the system.
4. Update posterior of the interventional GP.

**end**

**end**

Return the optimal value  $\hat{\mathbb{E}}[\mathbf{Y}^* | \text{do}(\mathbf{X}_s^* = \mathbf{x}_s^*), \mathbf{C}]$   
in  $\mathcal{D}_T^I$  and the corresponding  $\mathbf{X}_s^*$ ,  $\mathbf{x}_s^*$ .

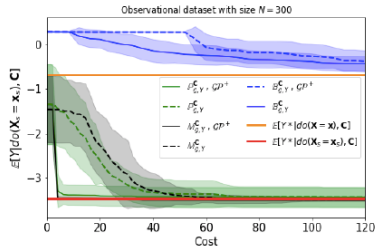
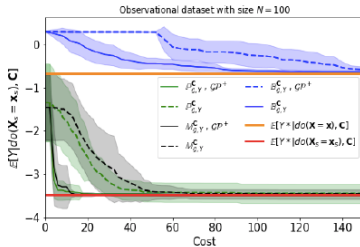
---

## Take home messages:

---

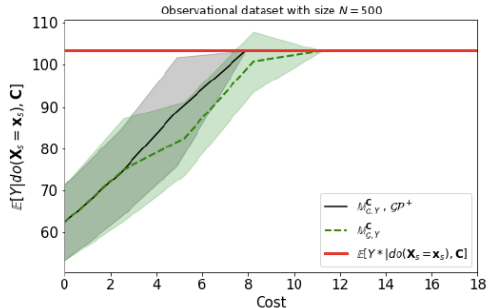
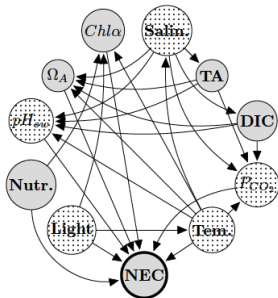
1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.
4. Standard Bayesian Optimization ignores causal assumptions.
5. Causal Optimization requires a new approach.
6. No need to explore  $\mathcal{P}(\mathbf{X})$  to solve Causal Optimization problems.
7. Causal GPs to merge observational and interventional data.
8. CBO optimizes systems/processes with interconnected nodes.

# Toy example - simulation analysis

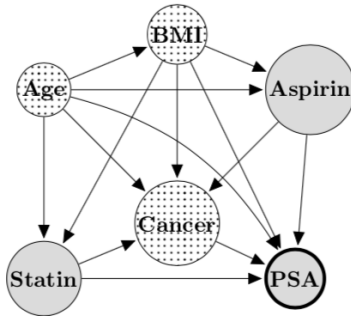


- Equal cost of intervening each variable.
- Results are consistent with what is expected.
- Better results than BO: propagation of effect beyond default domain.

# CBO in Ecology



- Goal: maximising the net coral ecosystem calcification (NEC).
- Five manipulative variables: cardinality of MIS is 25.
- Fast convergence to the optimum.



- Goal is minimize prostate specific antigen (PSA) providing aspirin and/or statin (domain  $[0,1]$ ).
- Optimal found solution is to only provide statin.
- This agrees with the general practice in medicine.

## Take home messages:

---

1. Many systems/processes decompose in interconnected nodes.
2. Optimization requires 'intervening' in the actionable nodes.
3. Do-calculus: 'emulating' experiments with observational data.
4. Standard Bayesian Optimization ignores causal assumptions.
5. Causal Optimization requires a new approach.
6. No need to explore  $\mathcal{P}(\mathbf{X})$  to solve Causal Optimization problems.
7. Causal GPs to merge observational and interventional data.
8. CBO optimizes systems/processes with interconnected nodes.
9. CBO improves BO when causal information is available.



‘Causal decision making’ takes the best of two worlds:

		<i>Obs. data + Causal assumptions</i>	
		No	Yes
<i>Int. data</i>	No	Mechanical models (OR, control, etc.)	Causal inference (PO, do-calculus etc.)
	Yes	Sequential decision making (AL, BayesOpt, etc.)	Causal Decision Making (This work!)

*Causal Bayesian Optimization.* Virginia Aglietti, Xiaoyu Lu, Andrei Paleyes and Javier González. AISTATS 2020.