

# Causal Bayesian Optimization

COMP0081 Applied ML

12/03/2024

Virginia Aglietti (Research Scientist)

# Contents

- Bayesian Optimization
- Causal Bayesian Optimization
- Extensions to constrained settings and functional interventions

# Black-box optimization

$$x^* = \operatorname{argmin}_{x \in \mathcal{X}} f(x)$$

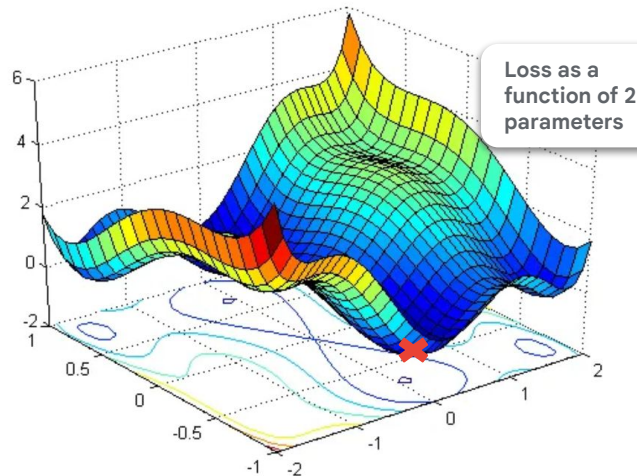
$f : \mathcal{X} \rightarrow \mathbb{R}^P$   
 $\mathcal{X} \subseteq \mathbb{R}^D$

**Goal:** find the global optima  $x^*$  in the smallest number of queries

**Applications:** hyper-parameters optimization, LLMs data mixture, robotics, molecules design, drug design, identification of optimal policies in causal systems etc

## Setting:

- $f$  is explicitly unknown and multimodal.
- Gradients are not available.
- We can query the function but evaluations of  $f$  are expensive.
- Evaluations of  $f$  may be perturbed by noise.



# Bayesian Optimization

- **Surrogate model:** model our belief about the function which gets updated as we sequentially observe function evaluations
  - Gaussian process/BNN/transformer

$$f(x) \sim \mathcal{GP}(m(x), K(x, x'))$$

- **Acquisition function (AF):** determines the sequential acquisition of points thus balancing exploration and exploitation
  - Heuristic, ad-hoc choice, problem specific
  - Generally uses the mean and variance of the prediction to determine next function evaluation

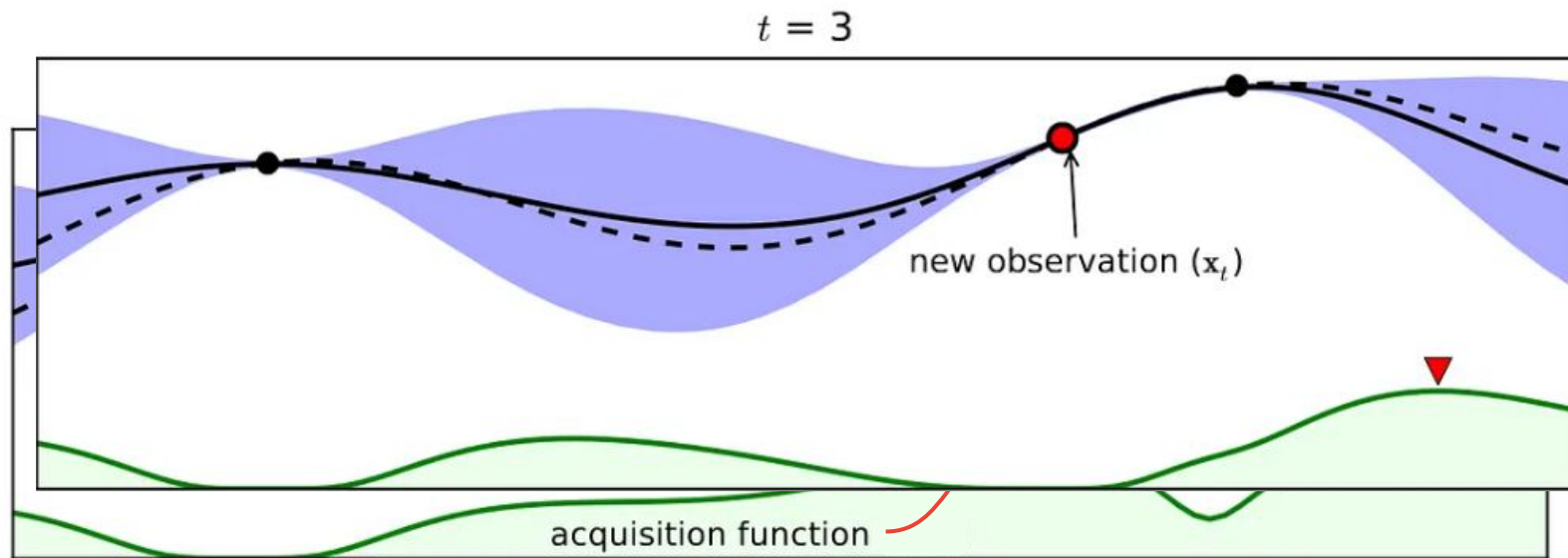
$$x^* = \underset{x \in \mathcal{X}}{\text{argmin}} f(x)$$



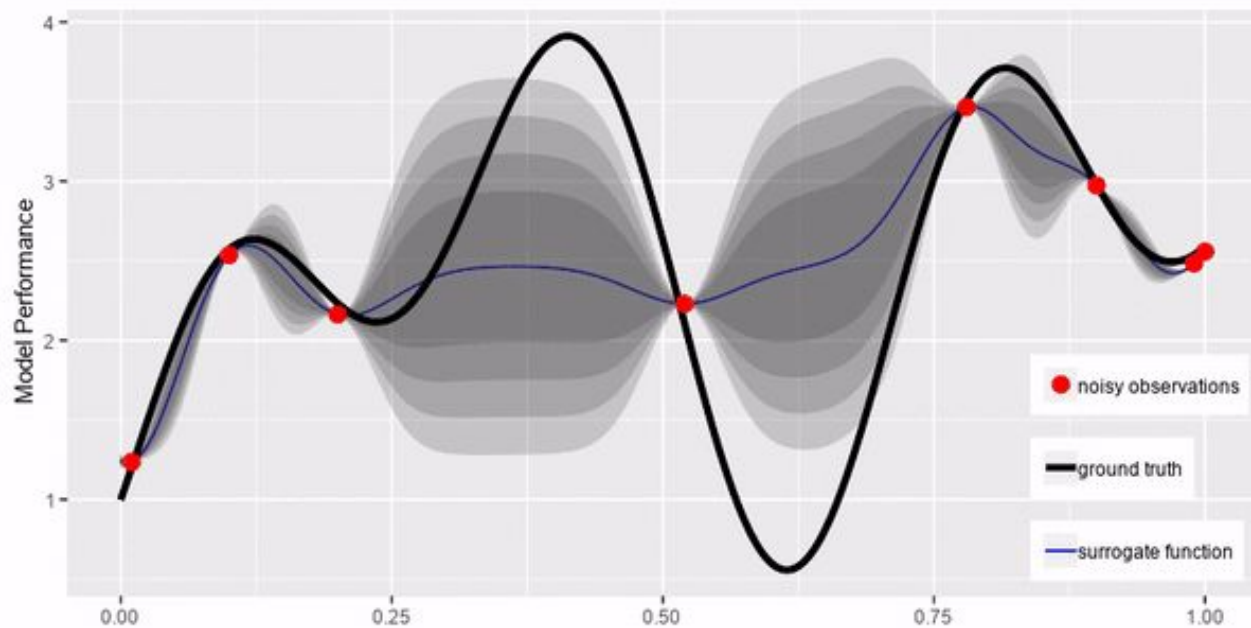
At every trial  $t$ , solve:

$$x_t = \underset{x \in \mathcal{X}}{\text{argmax}} \alpha_{t,\theta}(x)$$

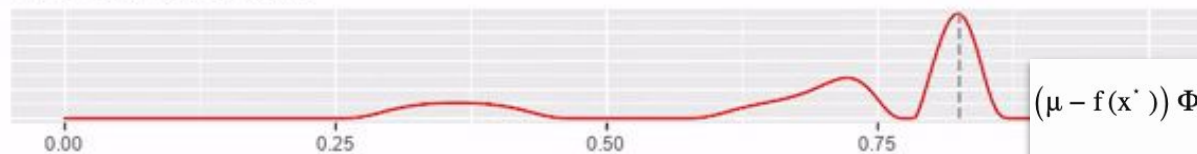
# Bayesian Optimization



# Bayesian Optimization

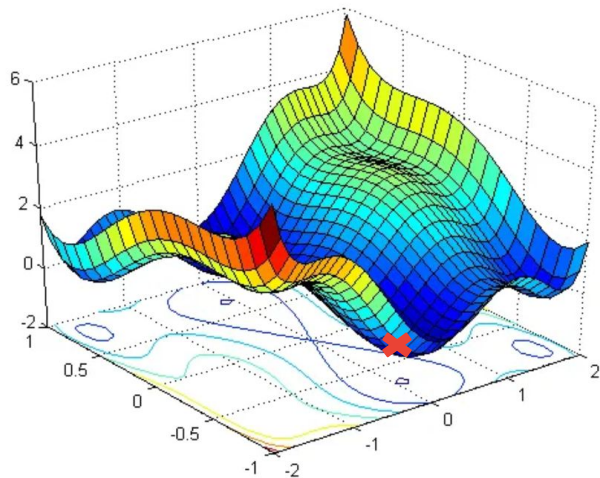


Expected improvement

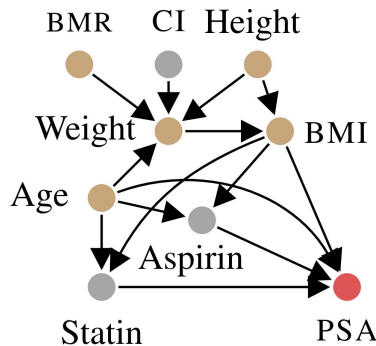


$$(\mu - f(x^*)) \Phi\left(\frac{\mu - f(x^*)}{\sigma}\right) + \sigma \phi\left(\frac{\mu - f(x^*)}{\sigma}\right)$$

# Causal black-box optimization



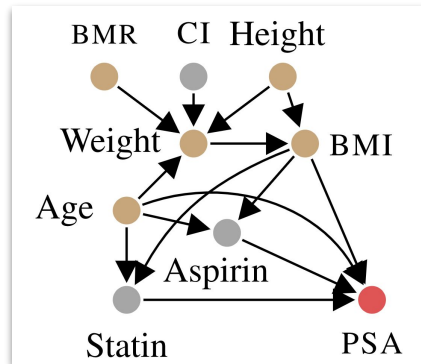
- Loss of a ML model as a function of 2 hyperparameters
- The average causal effect on a target variable given the intervention on 2 different values at different continuous levels



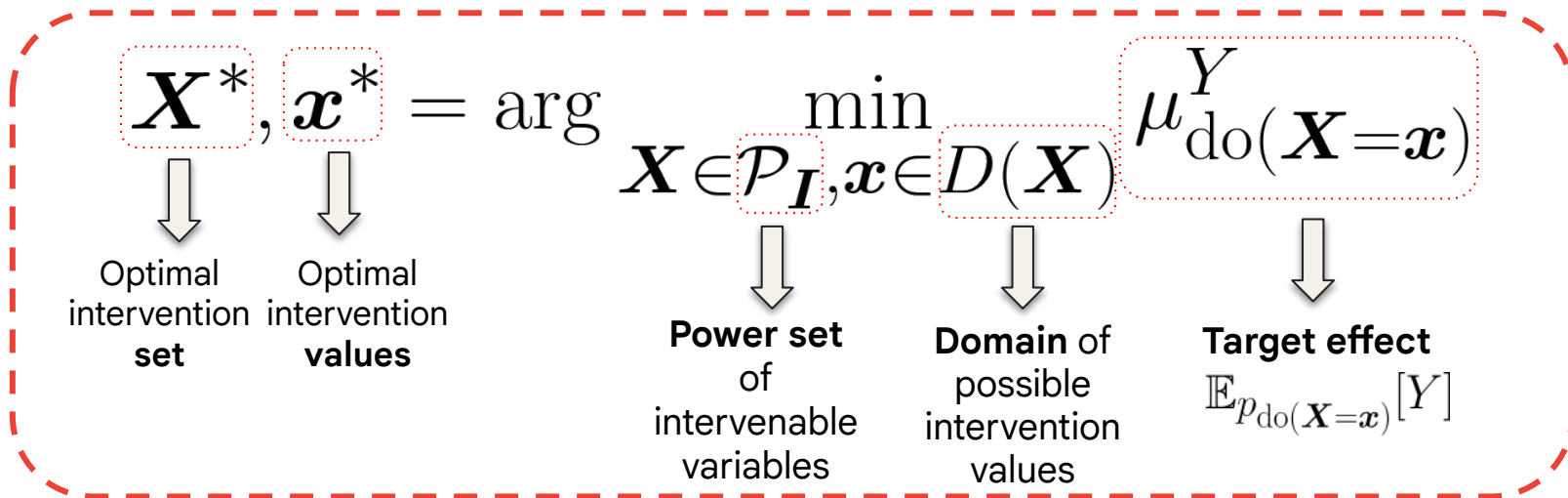
- Target variable  $Y = \text{PSA}$
- Intervenable variables  $I = \{\text{CI}, \text{Statin}, \text{Aspirin}\}$
- $\mathcal{P}_I = \{\emptyset, \{\text{CI}\}, \{\text{Statin}\}, \{\text{Aspirin}\}, \{\text{CI}, \text{Statin}\}, \{\text{CI}, \text{Aspirin}\}, \{\text{Statin}, \text{Aspirin}\}, \{\text{CI}, \text{Statin}, \text{Aspirin}\}\}$

# Causal black-box optimization

- A causal graph  $\mathcal{G}$  with nodes  $V$
- Target variable  $Y \in V$
- Intervenable variables  $I \subseteq V \setminus Y$
- Interventional domain  $D(\mathbf{X}) = \times_{X \in \mathbf{X}} D(X)$



**E.g.**  $\mathbf{X}^* = \{\text{CI}, \text{Statin}, \text{Aspirin}\}$   
 $\mathbf{x}^* = (\text{CI}=1, \text{Statin}=1, \text{Aspirin}=0)$





# Non-causal vs Causal Bayesian Optimization

## Non-causal

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in D(\mathbf{I})} \mu_{\text{do}(\mathbf{I}=\mathbf{x})}^Y$$

- Target function is explicitly unknown and multimodal
- Evaluations are perturbed by noise
- Evaluations are expensive



**Bayesian Optimization**

## Causal

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y$$

...

**+ Causal Graph**



**Causal Bayesian Optimization**

# Causal Bayesian Optimization

1

**Restrict the search space**  
by exploiting  
redundancies

2

**Model the target effects using GPs**  
and exploiting  
observational and  
interventional data

3

**Define an acquisition function**  
that allows to explore  
the interventions  
space



[Causal Bayesian Optimization.](#)

[V. Aglietti, X. Lu, A. Paleyes, & J. González](#)

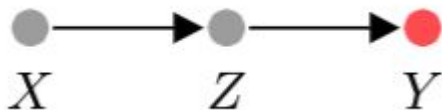
# Causal Bayesian Optimization: search space

1

**Restrict the search space**

by exploiting redundancies in  $\mathcal{G}$

$$X^*, x^* = \arg \min_{\substack{X \in \mathcal{P}_I, x \in D(X)}} \mu_{\text{do}(X=x)}^Y$$



$$\mu_{\text{do}(X=x, Z=z)}^Y = \mu_{\text{do}(Z=z)}^Y$$

**Assumption:** preference for sets of smaller cardinality (e.g. due the intervention cost)

$$\mathcal{P}_I = \{\emptyset, \{X\}, \{Z\}, \{X, Z\}\} \xrightarrow{\text{Minimal intervention sets [1]}} \mathbb{M}_{Y, \mathcal{G}} = \{\emptyset, \{X\}, \{Z\}\}$$

## Causal Bayesian Optimization: surrogate models

2

**Model the target effects using GPs**  
and exploiting  
observational and  
interventional data

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y$$

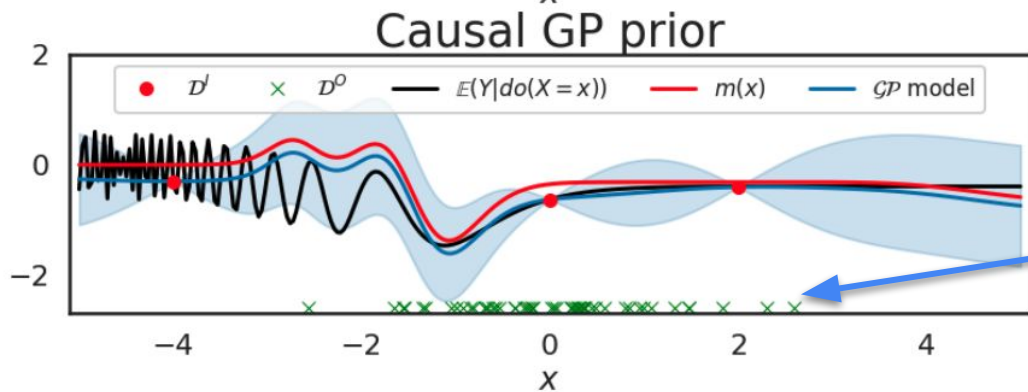
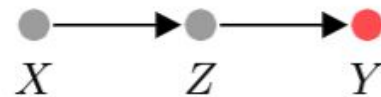
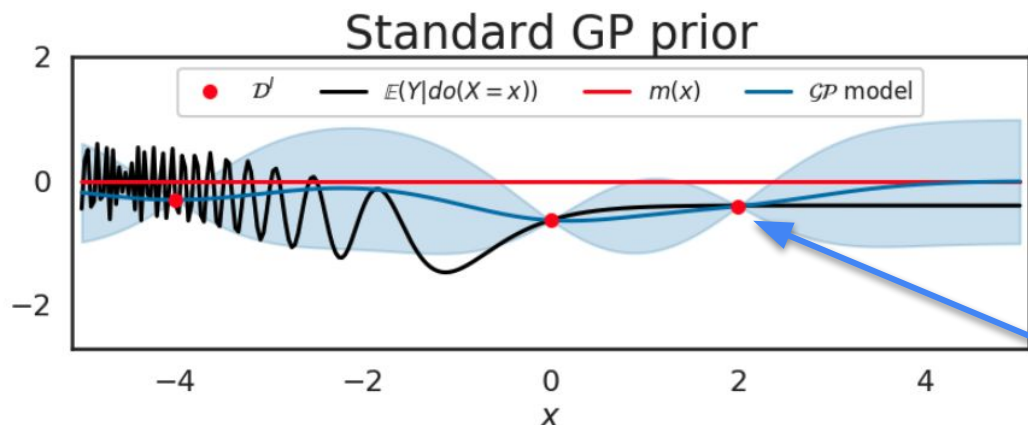
$$g_{\mathbf{X}}^Y \sim \mathcal{GP}(m_{\mathbf{X}}^Y(\mathbf{x}'), S_{\mathbf{X}}^Y(\mathbf{x}, \mathbf{x}'))$$

$$m_{\mathbf{X}}^Y(\mathbf{x}') = \hat{\mu}_{\text{do}(\mathbf{X}=\mathbf{x})}^Y \rightarrow$$

Estimated using observational data  
and the do-calculus

$$S_{\mathbf{X}}^Y(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right) + \hat{\sigma}_{\text{do}(\mathbf{X}=\mathbf{x})}^Y \times \hat{\sigma}_{\text{do}(\mathbf{X}=\mathbf{x}')}^Y$$

# Causal Bayesian Optimization: surrogate models

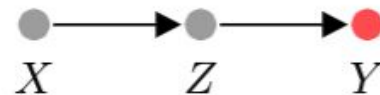


# Causal Bayesian Optimization: acquisition function

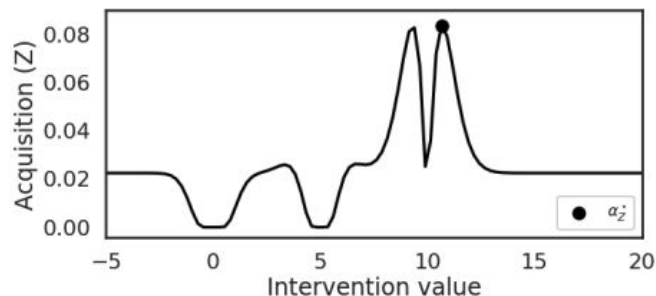
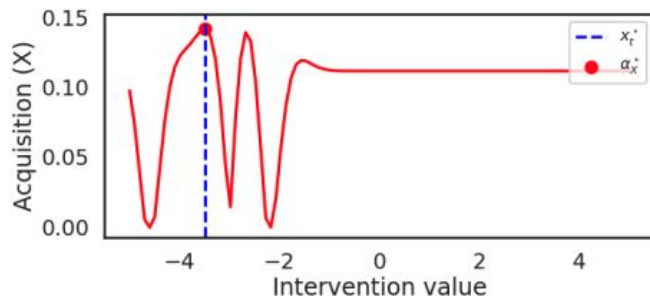
# 3

**Define an acquisition function** that allows to explore the interventions space

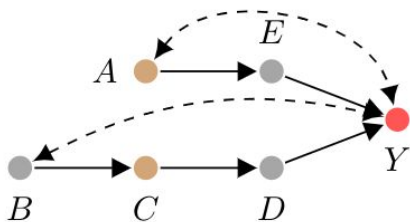
In CBO we optimize the **expected improvement per unit of cost** for every set in  $\mathbb{M}_{Y,\mathcal{G}}$  and select the intervention set and intervention values giving the highest expected improvement.



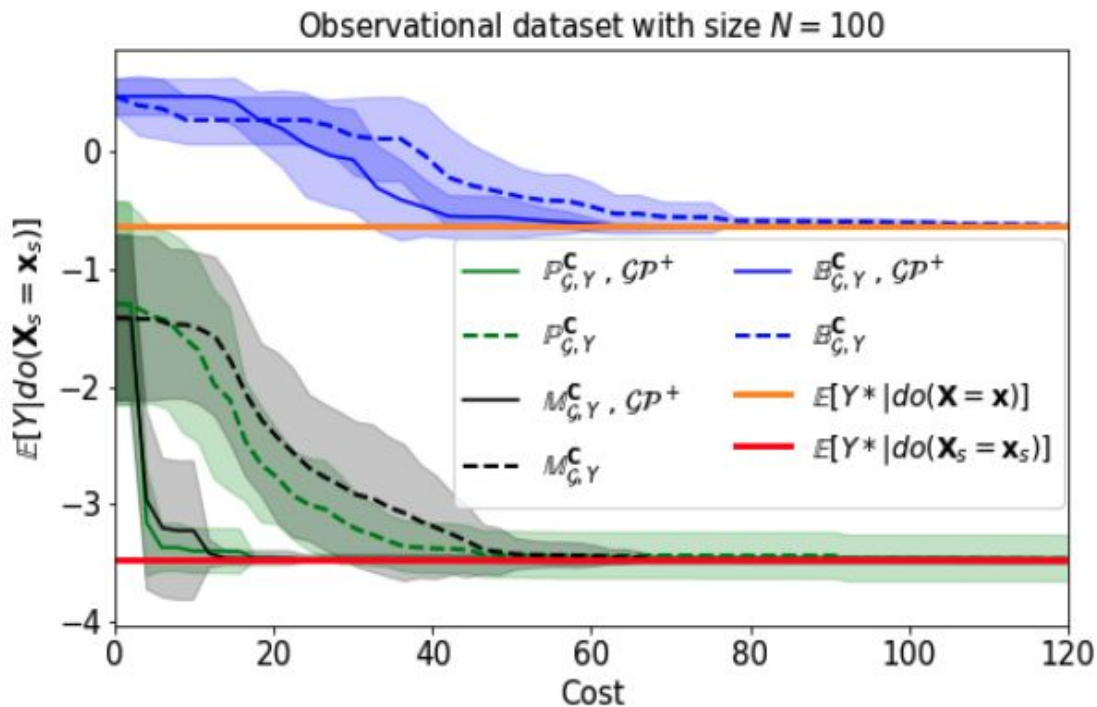
$$\text{EI}_{\mathbf{X}}(\mathbf{x}) = \mathbb{E}_{p(g_{\mathbf{X}}^Y | \mathcal{D}^I)} [\max(g_{\mathbf{X}}^Y(\mathbf{x}) - y^*, 0)] \setminus \text{Co}(\mathbf{X}, \mathbf{x})$$



# Causal Bayesian Optimization: experimental results



- BO is slower and identifies a suboptimal intervention
- CBO achieves the best result when using the Causal GP model



# Causal Bayesian Optimization: limitations

- The number of models GPs we require is determined by the number of sets to explore which is potentially large.
- We don't transfer interventional information across GPs e.g. we don't account for the fact that intervening on e.g.  $X$  might give us some information about an intervention on  $X$  and  $Z$ .
- We do not account for time and dynamic changes in the causal effects.
- We assume the causal graph to be known.
- We do not account for the existence of constrained variables.
- We only consider hard interventions.
- ...



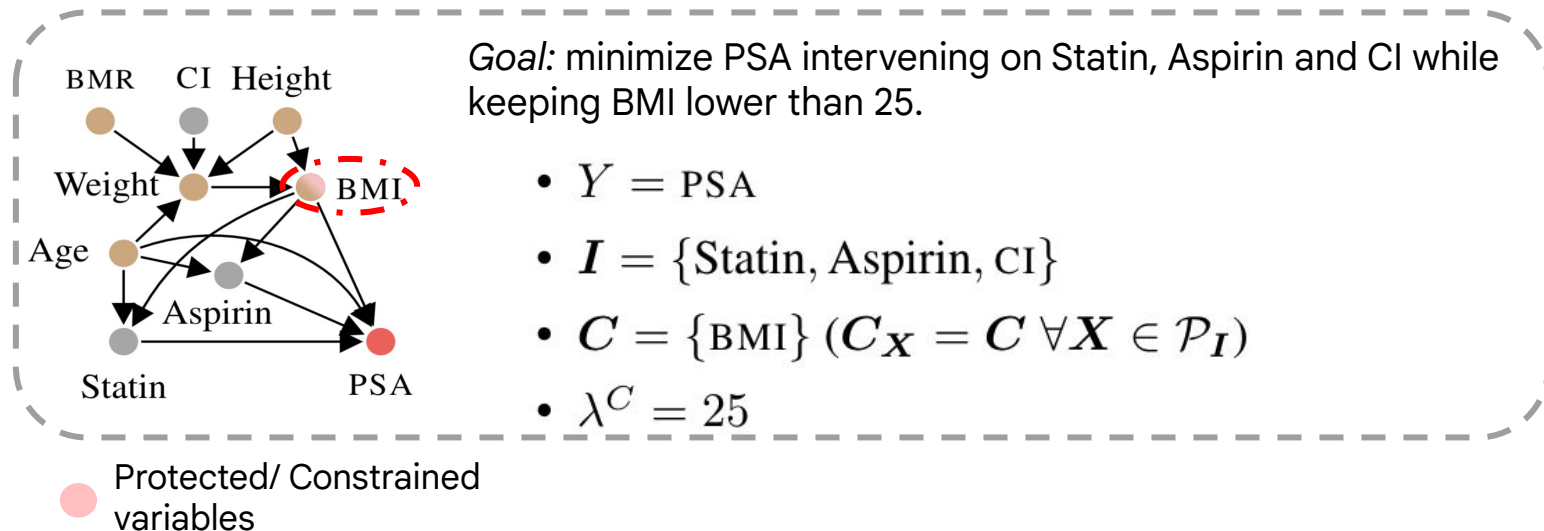
## Causal Bayesian Optimization: limitations

- The number of models GPs we require is determined by the number of sets to explore which is potentially large.
- We don't transfer interventional information across GPs e.g. we don't account for the fact that intervening on e.g.  $X$  might give us some information about an intervention on  $X$  and  $Z$ .
- We do not account for time and dynamic changes in the causal effects.
- We assume the causal graph to be known.
- **We do not account for the existence of constrained variables.**
- **We only consider hard interventions.**
- ...

# Constrained Causal Global Optimization

- A causal graph  $\mathcal{G}$  with nodes  $V$
- Target variable  $Y \in V$
- Intervenable variables  $I \subseteq V \setminus Y$
- Interventional domain  $D(\mathbf{X}) = \times_{X \in \mathbf{X}} D(X)$

(A set of protected variables  
 $C \subseteq V \setminus Y$ )



# Constrained Causal Global Optimization

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y,$$

$$\text{s.t. } \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^{C_X} \geq \lambda^{C_X}$$



Set of **constraint effects** with  
 $C_X := C \setminus (C \cap \mathbf{X})$



**Threshold values**

# Constrained vs Unconstrained vs non-causal Bayesian Optimization

## Non-causal

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in D(\mathbf{I})} \mu_{\text{do}(\mathbf{I}=\mathbf{x})}^Y$$

- Target function is explicitly unknown and multimodal
- Evaluations are perturbed by noise
- Evaluations are expensive

**Bayesian  
Optimization**

## Causal unconstrained

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y$$

...

+ Causal Graph

**Causal Bayesian  
Optimization**

## Causal constrained

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y, \text{ s.t. } \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^{C_X} \geq \lambda^{C_X}$$

...

+ Causal Graph  
+ Unknown constraints

**Constrained Causal Bayesian  
Optimization**

# Constrained Causal Bayesian Optimization (cCBO)

1

## **Restrict the search space**

by exploiting redundancies for target *and* constraint effects

2

**Model the target effects using GPs** and exploiting observational and interventional data and **capturing the correlation** between target and constraint effects.

3

**Define an acquisition function** that allows to explore the interventions space accounts for both the target and constraint effects.

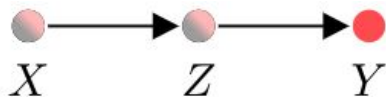


[Constrained Causal Bayesian Optimization.](#)  
[V. Aglietti, A. Malek, S. Chiappa.](#)

# Constrained Causal Bayesian Optimization: surrogate models

2

**Model the target effects using GPs** and exploiting observational and interventional data and **capturing the correlation** between target and constraint effects.



$$X = \alpha U_X$$

$$Z = \gamma X + U_Z$$

$$Y = \beta Z + U_Y$$

$$U_X, U_Z, U_Y \sim \mathcal{N}(0, 1)$$

When intervening on  $X$ :

- Target effect  $\mu_{\text{do}(X=x)}^Y$
- Constraint effect  $\mu_{\text{do}(X=x)}^Z$



$$\mu_{\text{do}(X=x)}^Y = \beta \mu_{\text{do}(X=x)}^Z$$

## Constrained Causal Bayesian Optimization: surrogate models

$$\mathbf{X}^*, \mathbf{x}^* = \arg \min_{\mathbf{X} \in \mathcal{P}_I, \mathbf{x} \in D(\mathbf{X})} \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^Y, \quad \text{s.t.} \quad \mu_{\text{do}(\mathbf{X}=\mathbf{x})}^{C_{\mathbf{X}}} \geq \lambda^{C_{\mathbf{X}}}$$

We model each effect  $\mu_{\text{do}(\mathbf{X})}^{V_k} \forall V_k \in C_{\mathbf{X}} \cup Y$  with a GP:

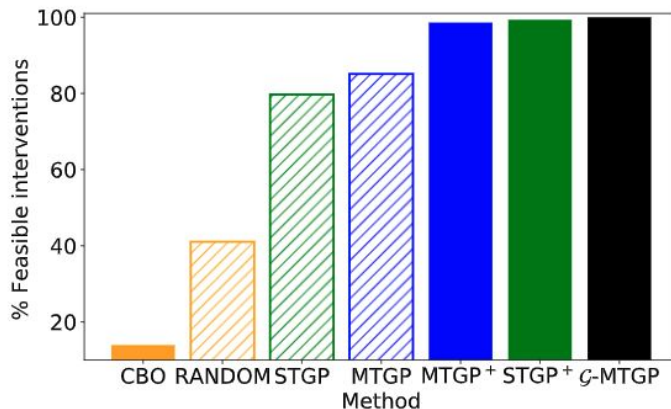
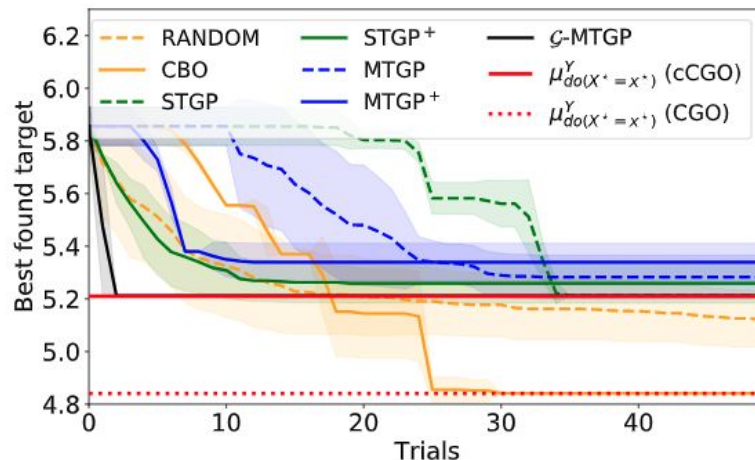
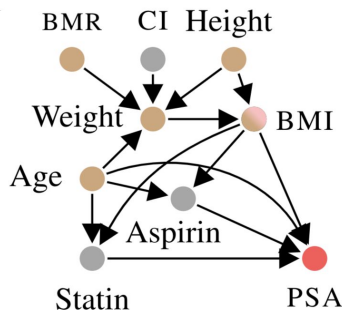
$$g_{\mathbf{X}}^{V_k}(\mathbf{x}) \sim \mathcal{GP}(\underbrace{m_{\mathbf{X}}^{V_k}(\mathbf{x})}_{\text{mean}}, \underbrace{S_{\mathbf{X}}^{V_k}(\mathbf{x}, \mathbf{x}')}_{\text{covariance}})$$



hyperparameters construction that accounts for the correlation induced by the structure of the graph (**multi-task** GP models).

# Constrained Causal Bayesian Optimization: Experimental results

- High level of noise in the observational data leads to less accurate prior formulation thus penalizing STGP<sup>+</sup>, MTGP<sup>+</sup> and  $\mathcal{G}$ -MTGP.
- Capturing the correlation is very important in this setting and leads G-MTGP to outperform all other methods according to both metrics.





# Why only Hard Interventions?

## Soft/Contextual/Functional Intervention

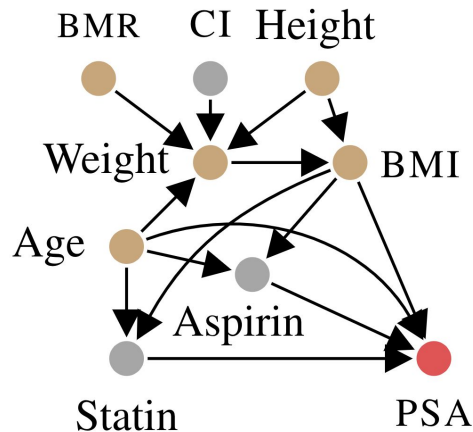
Often the decision maker has the ability to perform a **conditional/contextual** replacement of the existing causal mechanism, i.e. replace  $p(\mathbf{X} \mid pa_{\mathcal{G}}(\mathbf{X}))$  with another conditional distribution  $\pi_{\mathbf{X}} \mid \mathbf{C}_X$

↑  
New parents called  
**contexts**

E.g.

When finding an optimal value for Statin, we would likely want to take Age and BMI levels into account, as those hold information about the outcome node

Replace  $p(\text{Statin} \mid \text{Age}, \text{BMI})$  with  $\pi_{\text{Statin} \mid \text{Age}, \text{BMI}}$



- Targeted, more personalized treatment
- Subgroup optimality
- Lower treatment cost
- Hard interventions are special cases of soft, so no loss by considering soft

# Causal Bayesian Optimization with functional interventions

## Mixed Policy Scope (MPS) $\mathcal{S}$

Collection of tuples  $\langle X, \mathbf{C}_X \rangle$  where

- ❖  $X$  is an intervenable node  $X \in \mathbf{I}$
- ❖  $\mathbf{C}_X$  is associated set of **contexts** for intervention  $\pi_X | \mathbf{C}_X$
- ❖  $\langle X, \mathbf{C}_X \rangle$  does not introduce cycles in the graph

$$\boxed{\mathcal{S}^*}, \boxed{\pi_{\mathcal{S}^*}^*} = \arg \min_{\mathcal{S} \in \Sigma, \pi_{\mathcal{S}} \in \Pi_{\mathcal{S}}} \mu_{\pi_{\mathcal{S}}}^Y$$

Optimal **mixed policy scope (MPS)**

Optimal **MPS realization** (collection of functions)

Set of **all MPSs**

Set of **all possible MPS realizations** (interventions)



[Functional Causal Bayesian Optimization.](#)

[L. Gultchin, V. Aglietti, A. Bellot, I. Ktena, S. Chiappa](#)

# Causal Bayesian Optimization with functional interventions

$$\boxed{\mathcal{S}^*}, \boxed{\pi_{\mathcal{S}^*}} = \arg \min_{\mathcal{S} \in \Sigma, \pi_{\mathcal{S}} \in \Pi_{\mathcal{S}}} \mu_{\pi_{\mathcal{S}}}^Y$$

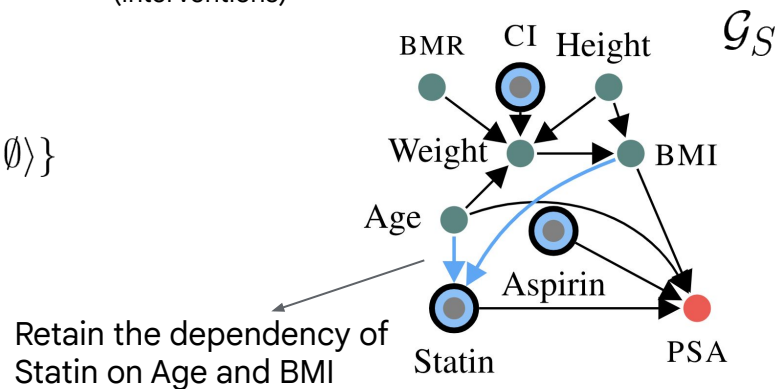
Optimal **mixed policy scope (MPS)**      Optimal **MPS realization** (collection of functions)

Set of **all MPSs**      Set of **all possible MPS realizations** (interventions)

$$\mathcal{S} = \{ \langle \text{Statin}, \{ \text{Age}, \text{BMI} \} \rangle, \langle \text{Aspirin}, \{ \text{Age}, \text{BMI} \} \rangle, \langle \text{CI}, \emptyset \rangle \}$$

$$\pi_{\mathcal{S}} = \{ \pi_{\text{Statin} | \text{Age}, \text{BMI}}, \pi_{\text{Aspirin} | \text{Age}, \text{BMI}}, \pi_{\text{CI}} \}$$

$$\pi_{\text{Statin} | \text{Age}, \text{BMI}} = \delta_{\text{Statin}}(\alpha * \text{Age} + \beta * \text{BMI})$$



# Functional Causal Bayesian Optimization: Surrogate models

2

Model the target effects using GPs

$$g_{\mathcal{S}}(\pi) \sim \mathcal{GP}(m_{\mathcal{S}}(\pi), K_{\mathcal{S}}^{\theta}(\pi, \pi'))$$

Surrogate model for the target effect  $\mu_{\mathcal{S}}^Y$  under possible interventions on MPS  $\mathcal{S}$

- Prior mean functional  $m_{\mathcal{S}}(\pi)$ , initialized at 0
- Prior covariance functional  $K_{\mathcal{S}}^{\theta}(\pi, \pi')$ , RBF kernel with hyperparameters  $\theta$
- Functional objective from the space  $\Pi_{\mathcal{S}}$  of all bounded (vector-valued) functions on to the reals  $C_{\mathcal{S}} = \bigcup_{(X, C_X) \in \mathcal{S}} C_X$

$$K_{\mathcal{S}}^{\theta} : \Pi_{\mathcal{S}} \times \Pi_{\mathcal{S}} \rightarrow \mathbb{R}$$

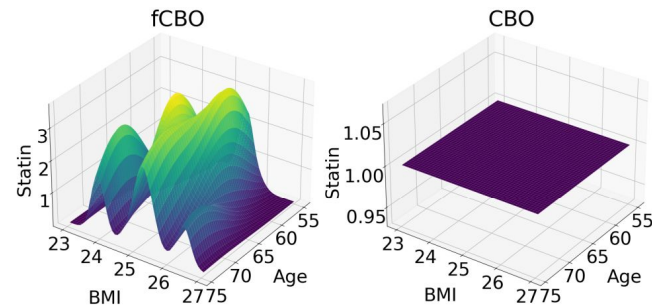
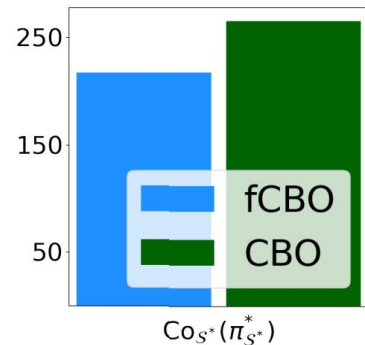
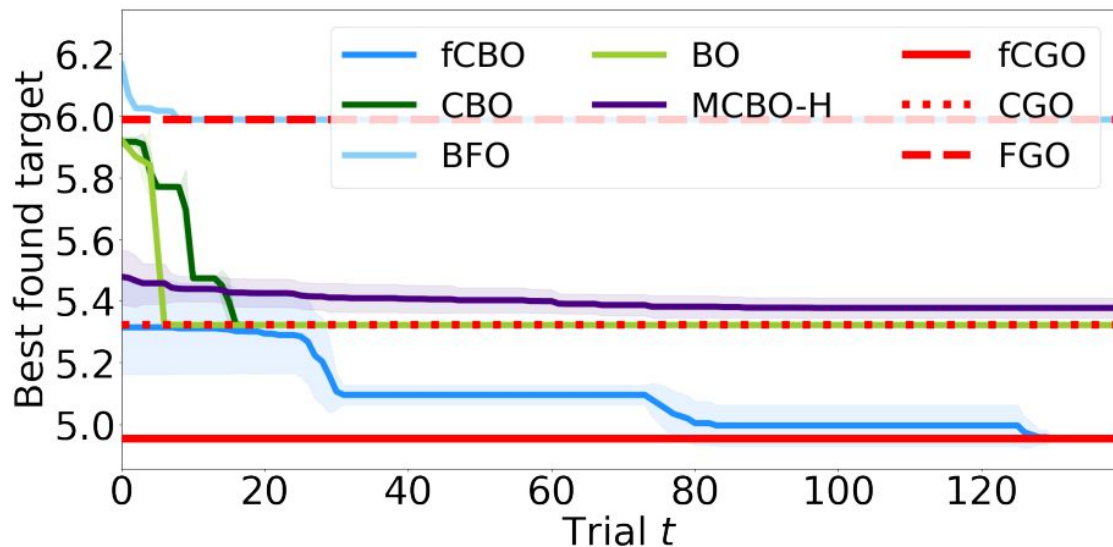
$\Downarrow$

$$\sigma_f^2 \exp\left(-\frac{\|\pi - \pi'\|^2}{2\ell^2}\right) \Rightarrow \langle \pi_{\text{func}} - \pi'_{\text{func}}, \pi_{\text{func}} - \pi'_{\text{func}} \rangle_{H_{\kappa}^{\mathcal{S}}}$$

# Functional Causal Bayesian Optimization: Healthcare experiment

**CBO**  $\mathbf{X}^* = \{\text{CI, Statin, Aspirin}\}$   $\mathbf{x}^* = (1, 1, 0)$

**fcBO**  $\mathcal{S}^* = \{\langle \text{CI}, \emptyset \rangle, \langle \text{Statin}, \{\text{Age, BMI}\} \rangle, \langle \text{Aspirin}, \emptyset \rangle\}$





# Thank you.

Get in touch at [aglietti@google.com](mailto:aglietti@google.com)

Google DeepMind