



Institución
Universitaria
Reacreditada en Alta Calidad

Optimización

Black-box Optimization

Optimización Bayesiana

Docente: Cristian Guarnizo Lemus

Somos Innovación Tecnológica con *Sentido Humano*



Alcaldía de Medellín



Institución
Universitaria
Reacreditada en Alta Calidad

Contenido

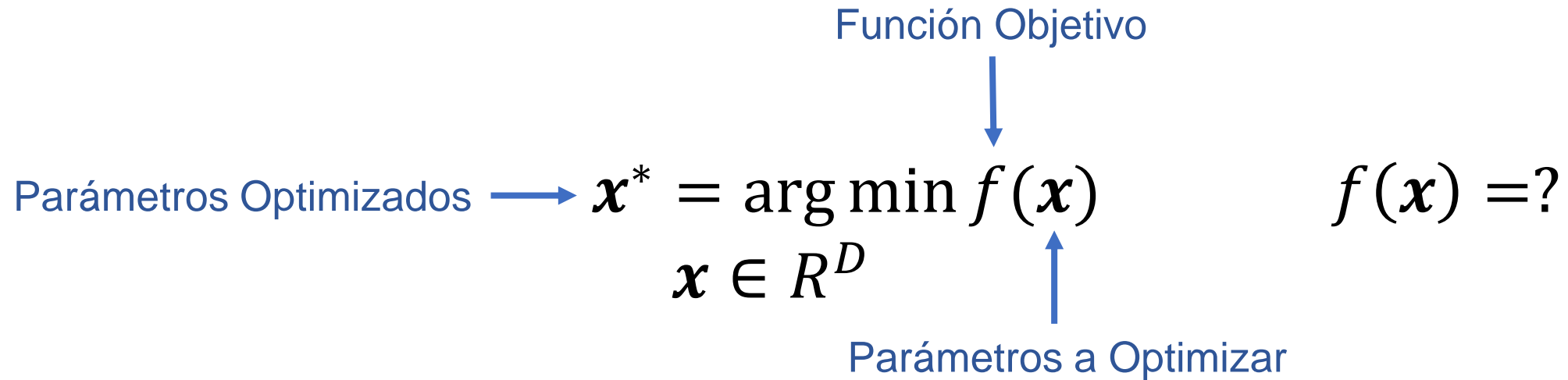
1. **Introducción.**
2. Procesos Gaussianos.
3. Técnicas para NLPs.



Objetivo de la presentación

- Entender los conceptos básicos de la optimización de caja negra (black-box).
- Entender el funcionamiento básico de la optimización Bayesiana.
- Cuando es útil emplear este tipo de optimización.

Black-box Optimization



Tipos de funciones objetivos

Un solo mínimo
(p.e. funciones convexas)

De primer orden
(podemos calcular gradientes)

Sin ruido
(una evaluación repetida
entrega el mismo resultado)

Evaluación barata
(se permiten una cantidad
de evaluaciones infinitas)

Fácil de resolver
(p.e. con gradientes)

Múltiples mínimos
(optimización global)

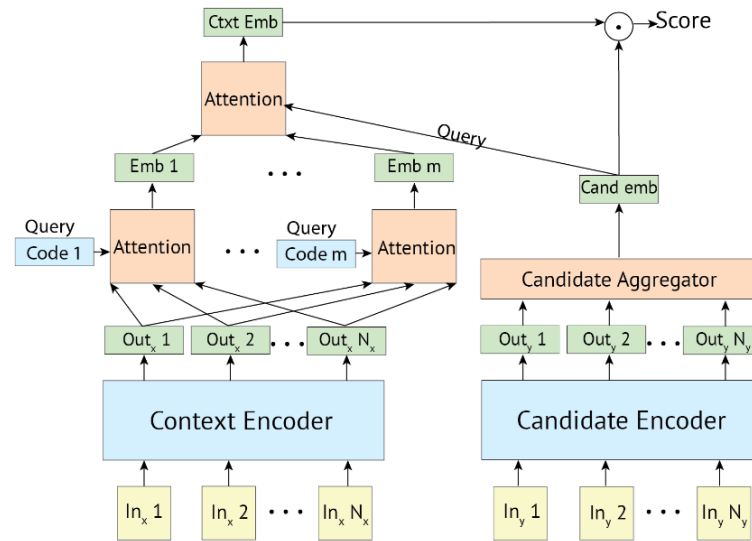
Orden cero
(sin gradientes)

Estocástico
(evaluaciones repetidas
entregan diferentes resultados)

Evaluación costosa
(limitado a decenas o cientos
de evaluaciones)

Difícil de Optimizar

Tipos de funciones objetivas





Optimización de Hiperparámetros

- Los modelos de aprendizaje automático están creciendo más y más complejos (Deep learning).
- Modelos de aprendizaje profundo modernos tienen docenas de parámetros que requieren ser optimizados (tasas de aprendizaje, número de capas, tamaño del lote). (p.e. Optuna)
- Para alcanzar resultados del estado del arte, encontrar buenos hiperparámetros es fundamental.
- Inclusive para expertos encontrar buenos hiperparámetros puede ser difícil y consumir mucho tiempo.

Métodos de optimización tradicionales

- Sintonización manual (requiere del experto).
- Grid search (no escala bien a espacio de parámetros grandes).
- Random Search (mejor que Grid, pero aun así requiere muchas evaluaciones).
- Gradiente descendente (solo factible para funciones de primer orden).
- Búsqueda Heurística (requiere de miles de evaluaciones)

Para docenas de parámetros, con correlaciones complejas, y evaluaciones costosas estos métodos se vuelven imprácticos.

Intuición detrás de la Optimización Bayesiana

Muchos optimizadores capturan solamente la información local de la función objetivo

$$\mathbf{x}^{(t+1)} = g(\mathbf{x}^{(t)}, f(\mathbf{x}^{(t)}))$$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha g(\nabla f(\mathbf{x}^{(t)}))$$

Cómo podemos usar toda la información (evaluaciones) recolectadas hasta ahora para tomar decisiones más informadas, por ende mejorando eficiencia de los datos?

$$\mathbf{x}^{(t+1)} = g(D)$$

$$D = \{\mathbf{x}_i, f(\mathbf{x}_i)\}, i = 1, \dots, n$$

Como hacer esto en la practica?

Podemos construir un modelo sustituto
(Surrogate)

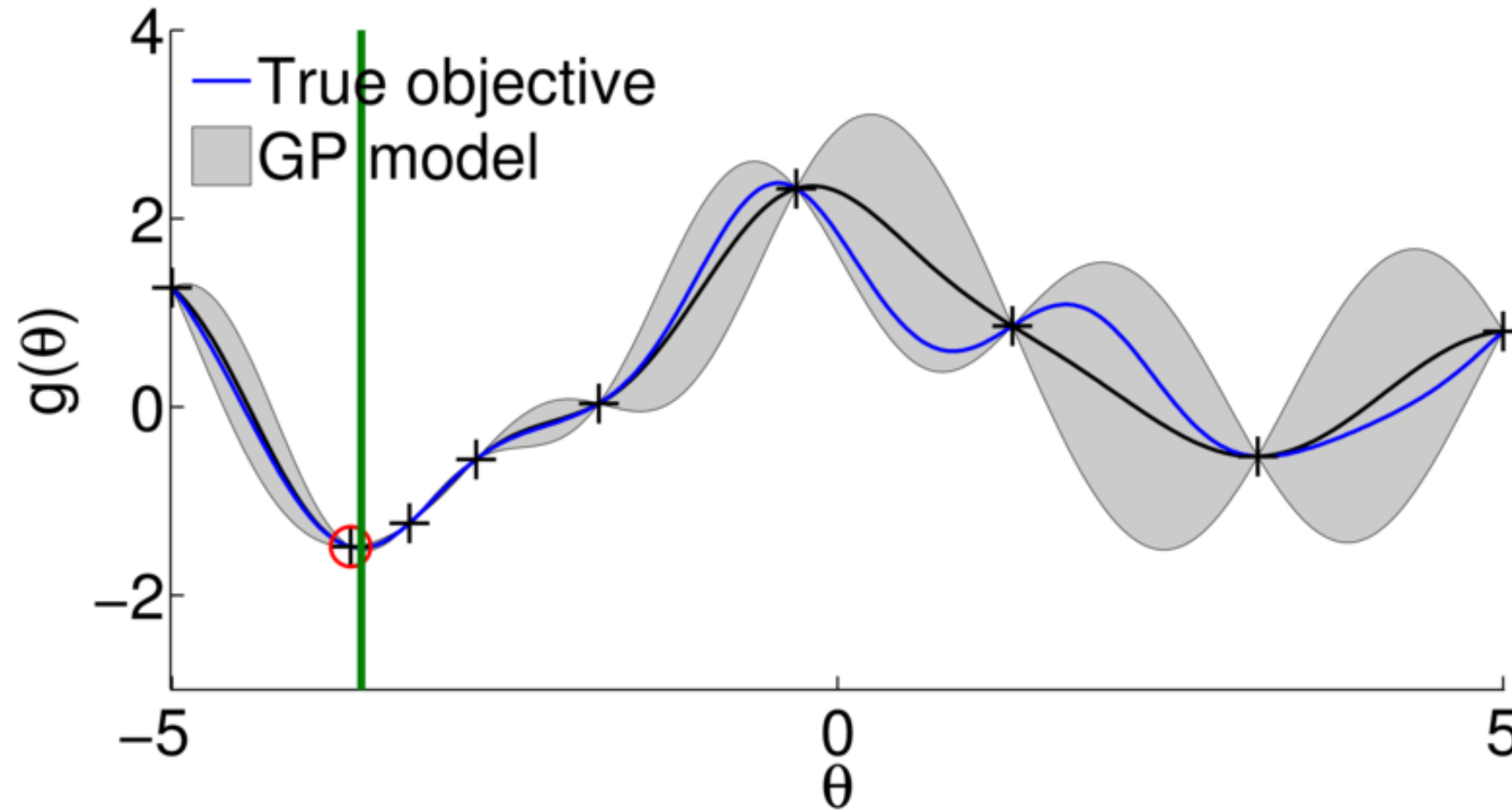
$$\tilde{f}(x_i) \Big|_D \sim f(x)$$

$$\mathbf{x}^* = \arg \min \tilde{f}(\mathbf{x})$$

Optimización Bayesiana

- Aprender la superficie de respuesta $\tilde{f}(x)$
- Basado en la superficie, seleccione los siguientes parámetros a evaluar $x^{(t+1)}$
- Evaluar $x^{(t+1)}$ de la función objetivo
- Repetir hasta que se cumpla un criterio de parada

Optimización Bayesiana





Contenido

1. **Introducción.**
2. Procesos Gaussianos – Superficie de respuesta.
3. Técnicas para NLPs.

Superficie de respuesta

Modelo sustituto (surrogate) que se requiere para aproximar la función objetivo a partir de los datos disponibles

$$D = \{\mathbf{x}_i, f(\mathbf{x}_i)\}, i = 1, \dots, n \quad \tilde{f}(\mathbf{x}_i) \Big|_D \sim f(\mathbf{x})$$

Existen una gran cantidad de técnicas en la literatura:

- Funciones polinomiales.
- Random forests.
- Bayesian neural networks.
- Deep neural networks.
- Gaussian Processes.
- etc



El más utilizado en la actualidad

Procesos Gaussianos

- Modelo de regresión flexible (no paramétrico)
- Distribución sobre funciones $\tilde{f} \sim GP(m_f, k_f)$
- Modelo probabilístico

$$y = \tilde{f}(X) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

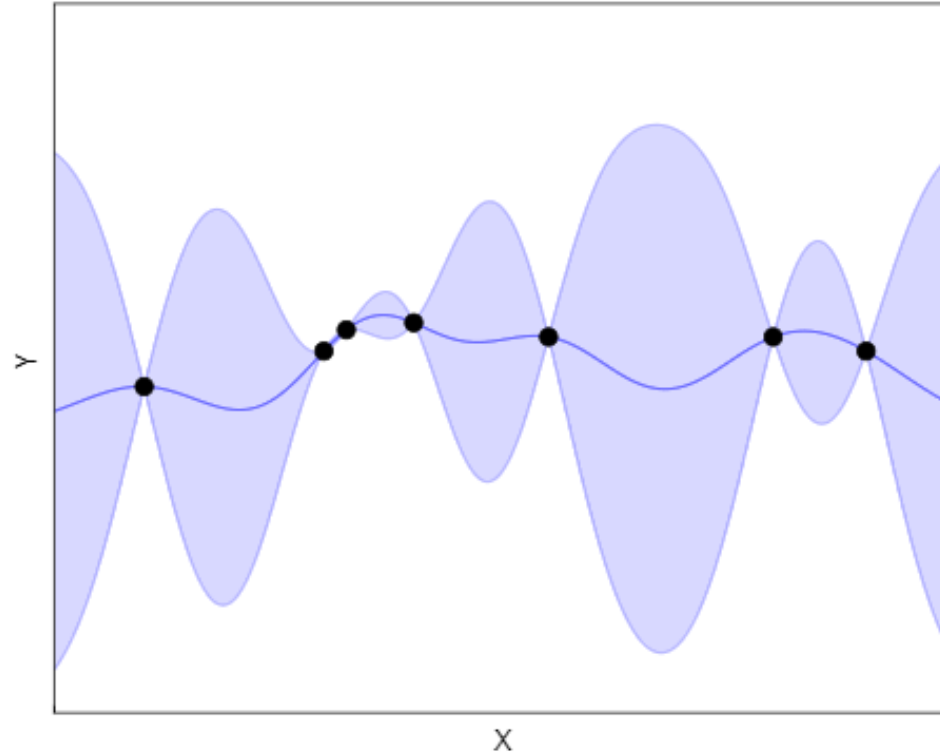
- El posterior de la distribución predictiva para una entrada arbitraria

$$p(\tilde{f}(x^*) | D, x^*) \sim N(\mu, \sigma^2)$$

$$\mu|_{x^*} = k(X, x^*)^T k(X, X)^{-1} y$$

$$\sigma^2|_{x^*} = k(x^*, x^*) - k(X, x^*)^T k(X, X)^{-1} k(X, x^*)$$

Procesos Gaussianos



Procesos Gaussianos

$$k(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{(x_i - x_j)^2}{2l^2}\right) + \delta_{ij} \sigma_n^2$$

Parámetros del GP
(hiperparametros)

Por que Procesos Gaussianos?

Pro:

- Matemática bien definida.
- Incertidumbres estimadas.
- Posibilidad de incluir información previa.
- Fácil forzar la suavidad.
- Buenas capacidades de modelamiento en regímenes de baja cantidad de datos.

Cons:

- Difícil escalar a entradas con alto dimensionalidad.
- Computacionalmente costoso.
- La calidad del modelo depende del uso apropiado del kernel.

Función de adquisición

Como seleccionamos los siguientes parámetros a evaluar $x^{(t+1)}$?

$$x^* = \arg \min \alpha(\tilde{f}(x), D)$$

Intuición: Una buena función de adquisición $\alpha()$ necesita tener un balance entre exploración y explotación.

$$y = \alpha(\mu(x^*), \sigma(x^*))$$

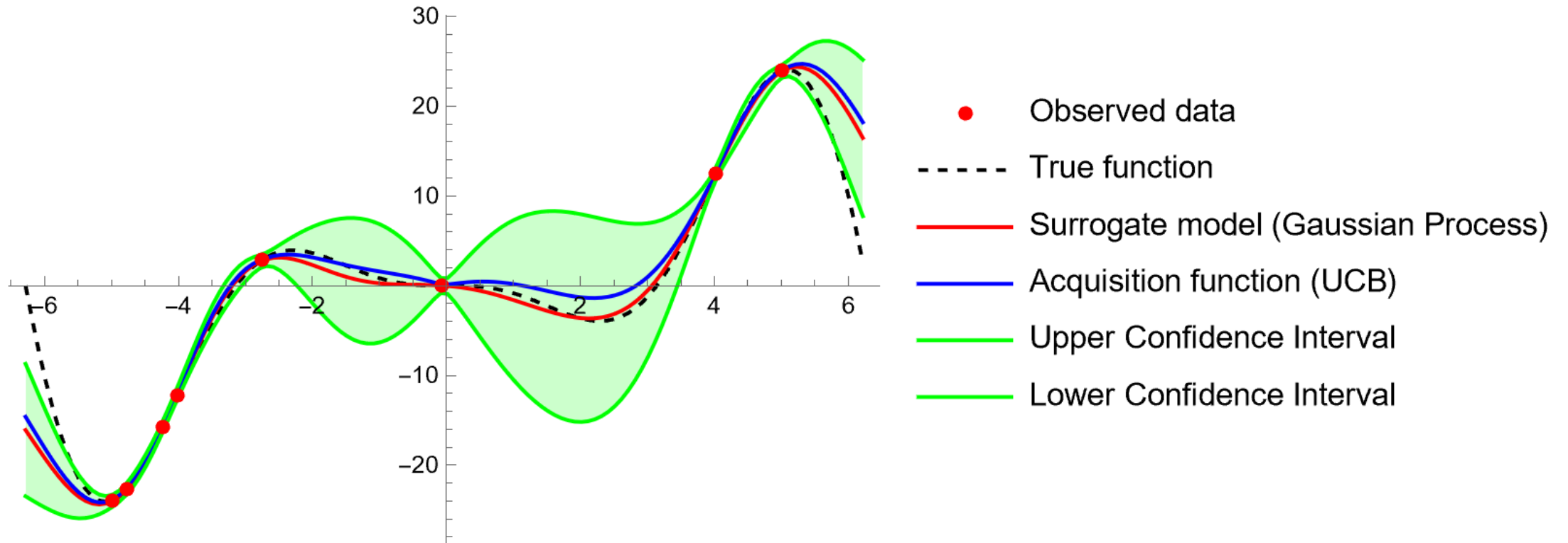
Muchas funciones de adquisición en la literatura:

- Probabilidad de mejoramiento.
- Mejoramiento esperado.
- Limite de confianza superior.
- Búsqueda de entropía.
- Búsqueda de entropía predictiva

Limite de confianza superior

$$UCB(x; \lambda) = \mu(x) - \lambda \sigma(x)$$

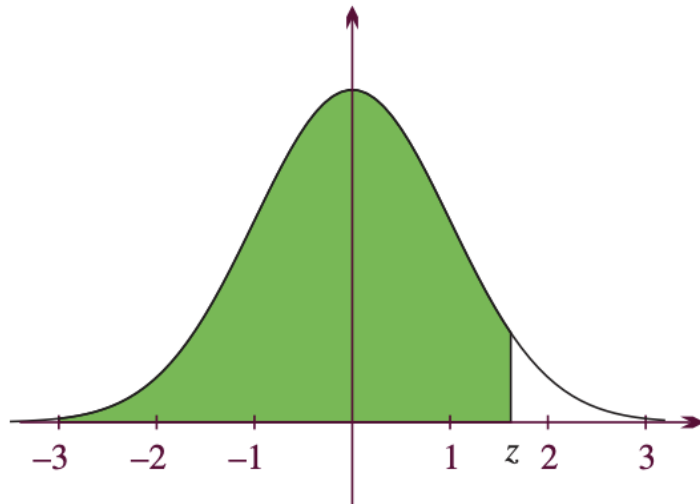
$UCB_{\lambda}(x) = \mu(x) + \lambda \sigma(x)$ for $\lambda=0.2$



Expected Improvement

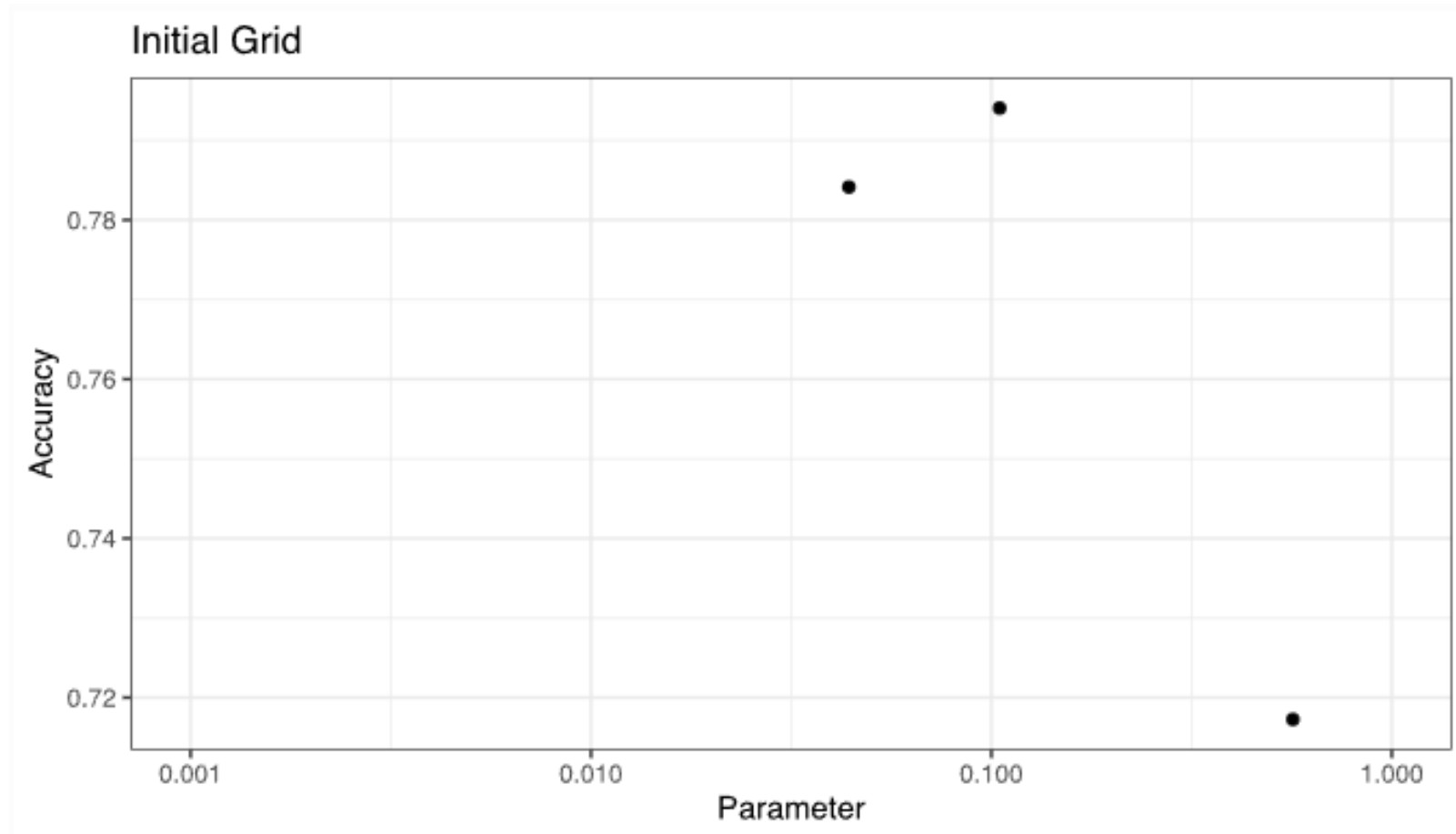
$$EI(x; \xi) = (\mu - \tilde{f}(x) - \xi) \Phi \left(\frac{\mu - \tilde{f}(x) - \xi}{\sigma} \right) + \sigma \phi \left(\frac{\mu - \tilde{f}(x) - \xi}{\sigma} \right)$$

$$\Phi(z) = \int_{-\infty}^z \phi(t) dt.$$

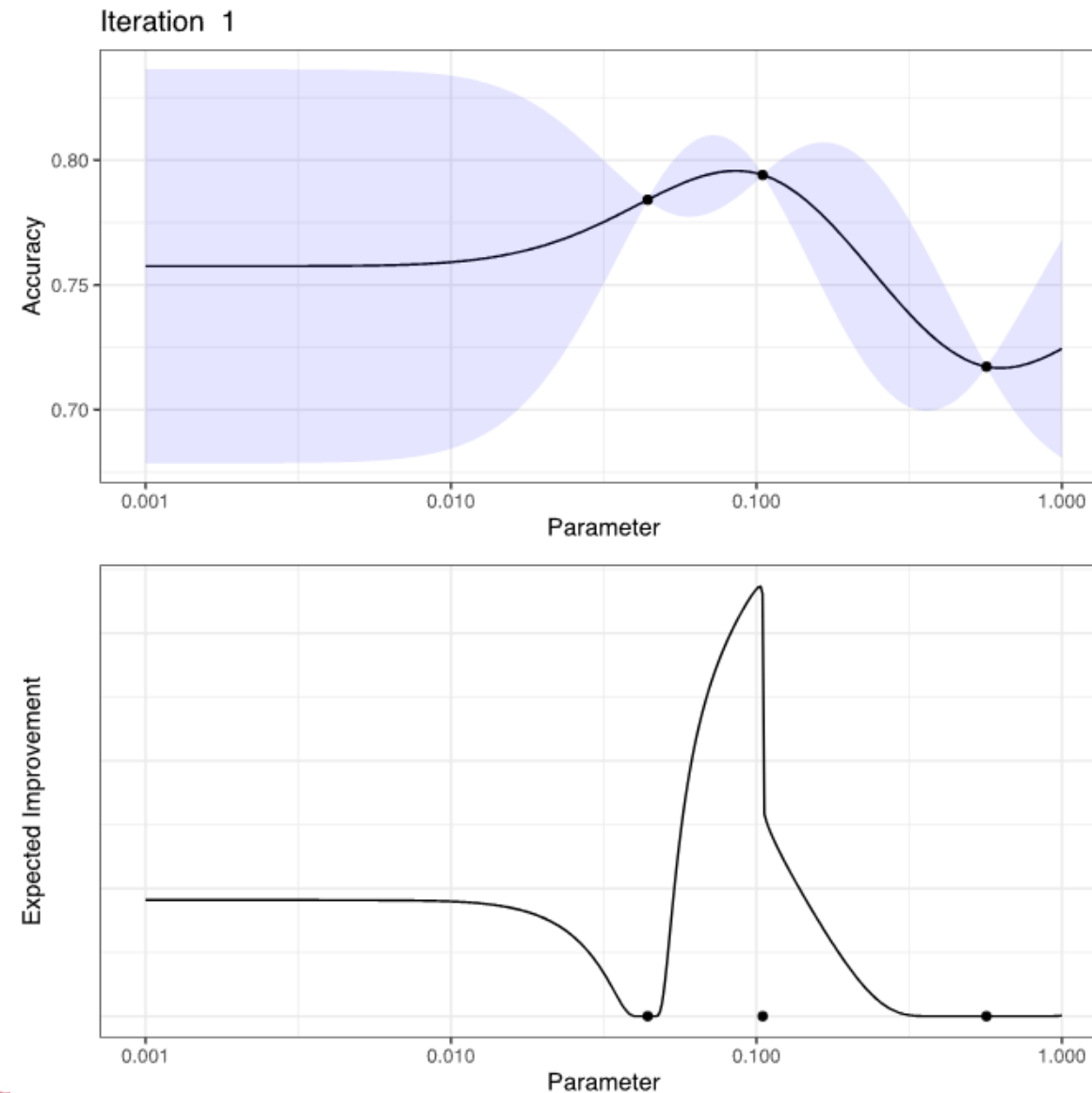


ξ igual a 0 realiza explotación.
 ξ grande realiza exploración.

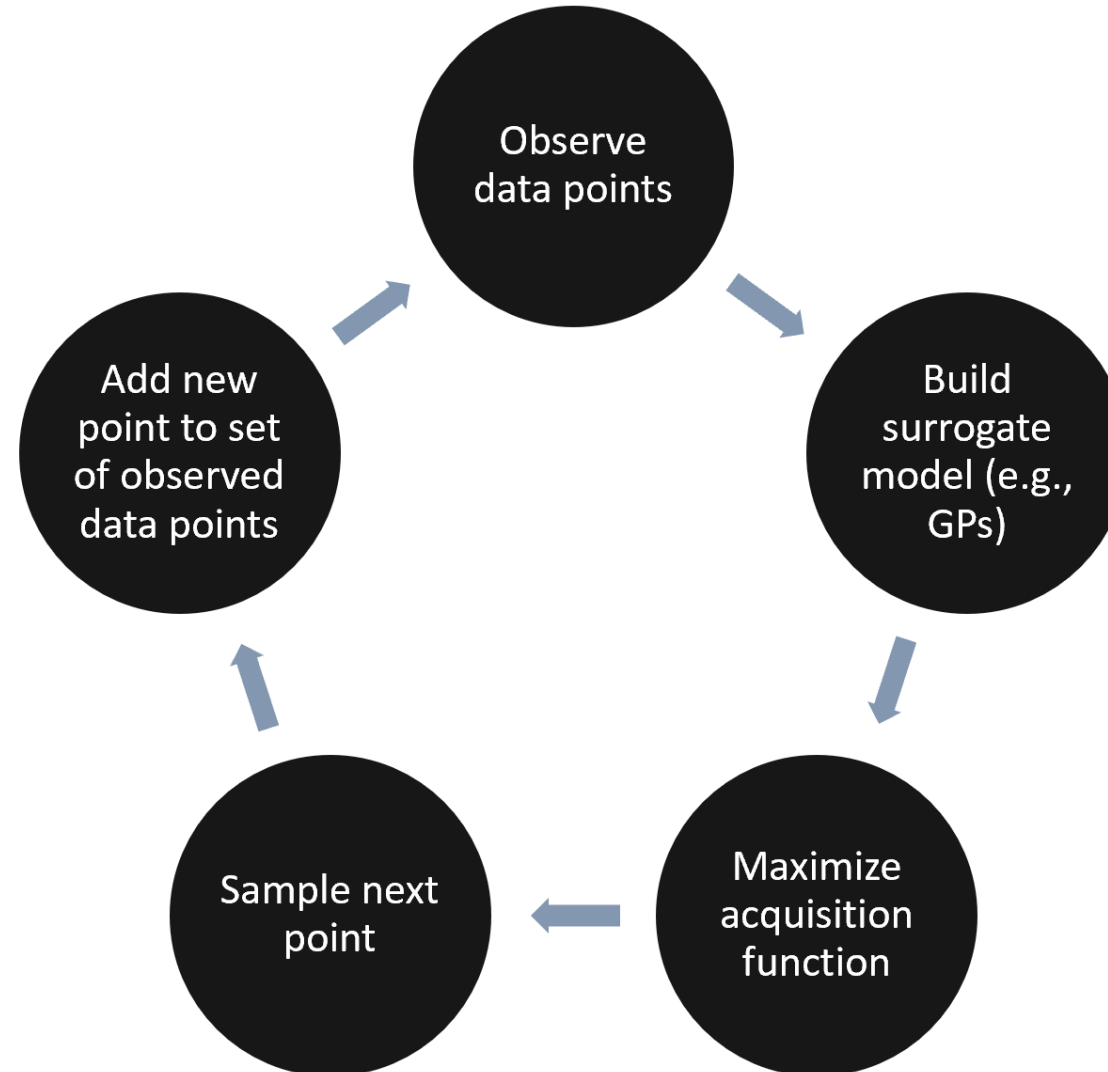
Expected Improvement



Expected Improvement



Algoritmo



Algoritmo

Algorithm 1: Single-objective Bayesian optimization

- 1 $\mathcal{D} \leftarrow$ if available: $\{\mathbf{x}, f(\mathbf{x})\}$
 - 2 Prior \leftarrow if available: Prior of the response surface
 - 3 **repeat**
 - 4 Train response surface \hat{f} from \mathcal{D}
 - 5 Find \mathbf{x}^* maximizing the acquisition surface $\alpha(\mathbf{x})$
 - 6 Evaluate \mathbf{x}^* on the real system
 - 7 Add $\{\mathbf{x}^*, f(\mathbf{x}^*)\}$ to \mathcal{D}
 - 8 **until** *stopCriteria*
-

Limitaciones de la optimización Bayesiana

Cons:

- No escala bien para espacio de parámetros de alta dimensionalidad (~30 dimensiones).
- No garantiza la convergencia.
- Si la función original es difícil de modelar, se obtienen desempeños muy malos (funciones discontinuas, o con muchos mínimos locales).

Pro:

- Sorpresivamente eficiente en muchas aplicaciones.
- Fácilmente se pueden incluir estructura o evaluaciones pasadas para un problema en particular.
- Modelos interpretables pueden proveer información.
- Funciona en el mundo real!

Referencias

- Basada en la presentación “Introduction to Bayesian Optimization” por el profesor Roberto Calandra. <https://slides.com/rcalandra/introduction-to-bayesian-optimization>



Institución
Universitaria
Reacreditada en Alta Calidad

¡Gracias!

Somos Innovación Tecnológica con *Sentido Humano*



Alcaldía de Medellín