

Optimization Algorithms

Derivative-Free (Black-Box) Optimization

Marc Toussaint
Technical University of Berlin
Winter 2024/25

Derivative-Free (Black-Box) Optimization

- Let $x \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, find

$$\operatorname{argmin}_x f(x)$$

- Derivative-Free/Blackbox optimization:
 - No access to ∇f or $\nabla^2 f$, sometimes also noisy evaluations $f(x)$

Derivative-Free (Black-Box) Optimization

- Let $x \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, find

$$\operatorname{argmin}_x f(x)$$

- Derivative-Free/Blackbox optimization:
 - No access to ∇f or $\nabla^2 f$, sometimes also noisy evaluations $f(x)$
- Algorithms needs to collect *data* D about f , and decide on next queries
- Many variants:
 - Classical derivative-free, implicit filtering, model-based optimization
 - Heuristics: Nelder-Mead, Coordinate search, Twiddle, Pattern Search
 - Stochastic Search, evolution strategies, EDAs, other EAs
 - Bayesian Optimization, Global Optimization
 - others?



Implicit Filtering

- Estimates the local gradient using finite differencing

$$\nabla_{\epsilon} f(x) \approx \left[\frac{1}{2\epsilon} (f(x + \epsilon e_i) - f(x - \epsilon e_i)) \right]_{i=1, \dots, n}$$

- Lines search along the gradient; if not successful, decrease ϵ
- Can be extended by using $\nabla_{\epsilon} f(x)$ to update an approximation of the Hessian (as in BFGS)

Model-based optimization

following Nodocal et al. “Derivative-free optimization”



Model-based optimization

- The previous stochastic search methods are heuristics to update θ

Why not store the previous data directly?

- Model-based optimization takes the approach

- Store a data set $\theta = D = \{(x_i, y_i)\}_{i=1}^n$ of previously explored points (let \hat{x} be the current minimum in D)

- Compute a (quadratic) model $D \mapsto \hat{f}(x) = \phi_2(x)^\top \beta$

- Choose the next point as

$$x^+ = \underset{x}{\operatorname{argmin}} \hat{f}(x) \quad \text{s.t.} \quad |x - \hat{x}| < \alpha$$

- Update D and α depending on $f(x^+)$

- The argmin is solved with constrained optimization methods

Model-based optimization

- 1: Initialize D with at least $\frac{1}{2}(n+1)(n+2)$ data points
- 2: **repeat**
- 3: Compute a regression $\hat{f}(x) = \phi_2(x)^\top \beta$ on D
- 4: Compute $x^+ = \operatorname{argmin}_x \hat{f}(x)$ s.t. $|x - \hat{x}| < \alpha$
- 5: Compute the improvement ratio $\rho = \frac{f(\hat{x}) - f(x^+)}{\hat{f}(\hat{x}) - \hat{f}(x^+)}$
- 6: **if** $\rho > \epsilon$ **then**
- 7: Increase the stepsize α
- 8: Accept $\hat{x} \leftarrow x^+$
- 9: Add to data, $D \leftarrow D \cup \{(x^+, f(x^+))\}$
- 10: **else**
- 11: **if** $\det(D)$ is too small **then** *// Data improvement*
- 12: Compute $x^+ = \operatorname{argmax}_x \det(D \cup \{x\})$ s.t. $|x - \hat{x}| < \alpha$
- 13: Add to data, $D \leftarrow D \cup \{(x^+, f(x^+))\}$
- 14: **else**
- 15: Decrease the stepsize α
- 16: **end if**
- 17: **end if**
- 18: Prune the data, e.g., remove $\operatorname{argmax}_{x \in \Delta} \det(D \setminus \{x\})$

19: **until** converges



Model-based optimization

- Optimal parameters (with data matrix $X \in \mathbb{R}^{n \times \dim(\beta)}$)

$$\hat{\beta}^{\text{ls}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$$

The determinant $\det(\mathbf{X}^T \mathbf{X})$ or $\det(\mathbf{X})$ (denoted $\det(D)$ on the previous slide) is a measure for well the data supports the regression. The data improvement explicitly selects a next evaluation point to increase $\det(D)$.

- Nocedal describes in more detail a geometry-improving procedure to update D .
- Model-based optimization is closely related to Bayesian approaches. But
 - Should we really prune data to have only a minimal set D (of size $\dim(\beta)$)?
 - Is there another way to think about the “data improvement” selection of x^+ ? (\rightarrow maximizing uncertainty/information gain)

Nelder-Mead method – Downhill Simplex Method

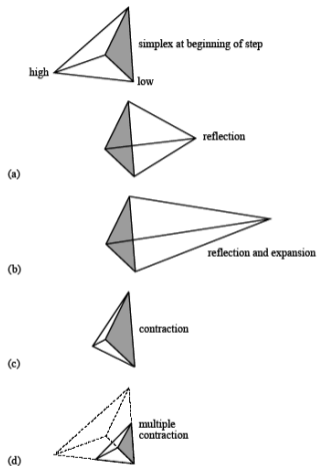


Figure 10.4.1. Possible outcomes for a step in the downhill simplex method. The simplex at the beginning of the step, here a tetrahedron, is shown, top. The simplex at the end of the step can be any one of (a) a reflection away from the high point, (b) a reflection and expansion away from the high point, (c) a contraction along one dimension from the high point, or (d) a contraction along all dimensions towards the low point. An appropriate sequence of such steps will always converge to a minimum of the function.

Nelder-Mead method – Downhill Simplex Method

- Let $x \in \mathbb{R}^n$
- Maintain $n + 1$ points x_0, \dots, x_n , sorted by $f(x_0) < \dots < f(x_n)$
- Compute center c of points
- Reflect: $y = c + \alpha(c - x_n)$
- If $f(y) < f(x_0)$: **Expand**: $y = c + \gamma(c - x_n)$
- If $f(y) > f(x_{n-1})$: **Contract**: $y = c + \varrho(c - x_n)$
- If still $f(y) > f(x_n)$: **Shrink** $\forall_{i=1, \dots, n} x_i \leftarrow x_0 + \sigma(x_i - x_0)$

- Typical parameters: $\alpha = 1, \gamma = 2, \varrho = -\frac{1}{2}, \sigma = \frac{1}{2}$

Coordinate Search

Input: Initial $x \in \mathbb{R}^n$

1: **repeat**

2: **for** $i = 1, \dots, n$ **do**

3: $\alpha^* = \operatorname{argmin}_{\alpha} f(x + \alpha e_i)$

// Line Search

4: $x \leftarrow x + \alpha^* e_i$

5: **end for**

6: **until** x converges

- The LineSearch must be approximated
 - E.g. abort on any improvement, when $f(x + \alpha e_i) < f(x)$
 - Remember the last successful stepsize α_i for each coordinate

Twiddle

Input: Initial $x \in \mathbb{R}^n$, initial stepsizes α_i for all $i = 1 : n$

1: **repeat**

2: **for** $i = 1, \dots, n$ **do**

3: $x \leftarrow \operatorname{argmin}_{y \in \{x - \alpha_i e_i, x, x + \alpha_i e_i\}} f(y)$

// twiddle x_i

4: Increase α_i if x changed; decrease α_i otherwise

5: **end for**

6: **until** x converges

Pattern Search

- In each iteration k , have a (new) set of search directions $D_k = \{d_{ki}\}$ and test steps of length α_k in these directions
- In each iteration, adapt the search directions D_k and step length α_k

Details: See Nocedal et al.

