

# Random subspaces and expected decrease in derivative-free optimization

Clément W. Royer (Université Paris Dauphine-PSL)

Workshop Bayesian Optimization and related applications

June 20, 2024

**Dauphine** | PSL   
UNIVERSITÉ PARIS

**PR[AI]RIE**  
PaRis Artificial Intelligence Research InstitutE

**ffCRfCRf**  
FONDS FRANCE CANADA POUR LA RECHERCHE  
FRANCE CANADA RESEARCH FUND

*Joint work with Warren Hare (UBC) & Lindon Roberts (U. of Sydney)*



- *Direct search based on probabilistic descent in reduced spaces*  
L. Roberts and C. W. Royer, SIAM J. Optimization, 2023.
- *Expected decrease for derivative-free algorithms using random subspaces*  
W. Hare, L. Roberts and C. W. Royer, under review, 2024.

- 1 Derivative-free algorithm
- 2 Reduced subspace approach
- 3 Numerics with subspaces
- 4 Subspace dimensions

- 1 Derivative-free algorithm
- 2 Reduced subspace approach
- 3 Numerics with subspaces
- 4 Subspace dimensions

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}).$$

## Assumptions

- $f$  bounded below;
- $f$  continuously differentiable (nonconvex).

## Blackbox/Derivative-free optimization

- **Derivatives unavailable for algorithmic use.**
- Only access to values of  $f$ .

## My goal

Develop algorithms with controlled

- Number of calls to  $f$ ;
- Dependency on  $n$ .

## My goal

Develop algorithms with controlled

- Number of calls to  $f$ ;
- Dependency on  $n$ .

## Complexity bound

Given  $\epsilon \in (0, 1)$  and, bound the number of **function evaluations** needed by a method to reach  $\mathbf{x}$  such that

$$\|\nabla f(\mathbf{x})\| \leq \epsilon,$$

deterministically or **in expectation/probability**.

## My goal

Develop algorithms with controlled

- Number of calls to  $f$ ;
- Dependency on  $n$ .

## Complexity bound

Given  $\epsilon \in (0, 1)$  and, **bound** the number of **function evaluations** needed by a method to reach  $\mathbf{x}$  such that

$$\|\nabla f(\mathbf{x})\| \leq \epsilon,$$

deterministically or **in expectation/probability**.

**Focus: dependency w.r.t.  $n$ .**



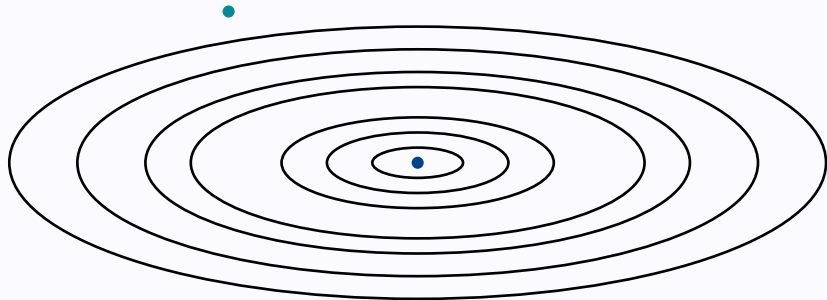
## Main algorithmic families

- Direct search: Explore the space through selected directions.
- Model based: Build a surrogate for the objective function.

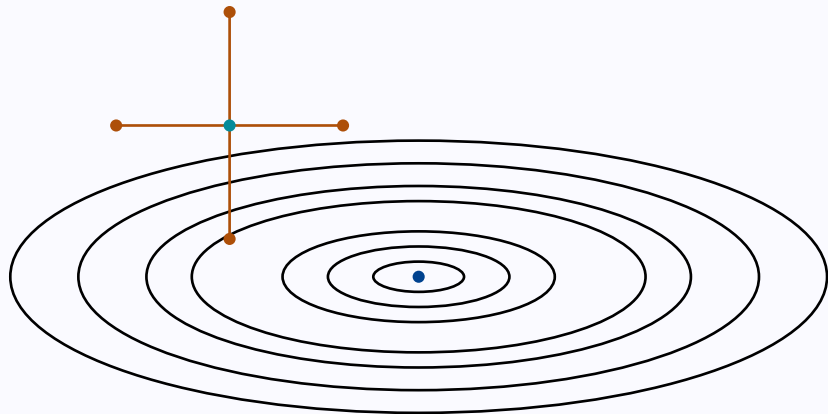
## Choosing a family for a 2pm talk

- Direct search simpler to explain.
- **All results have a model-based counterpart.**

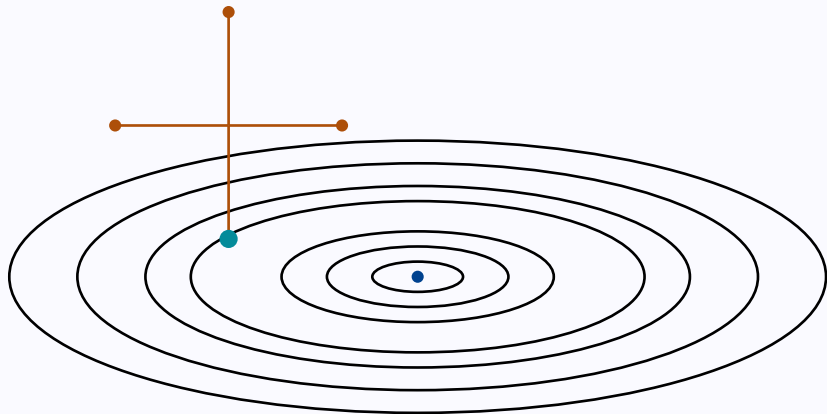
# Example: Coordinate search



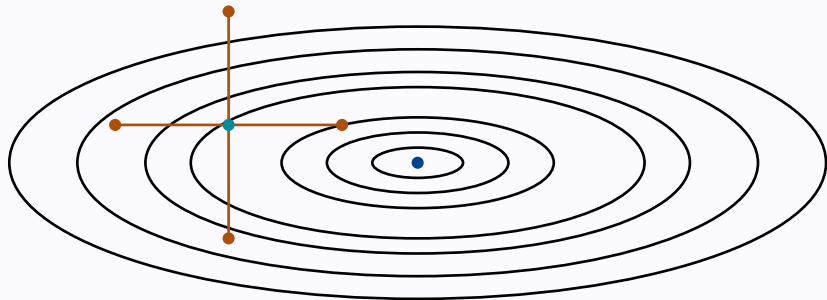
# Example: Coordinate search



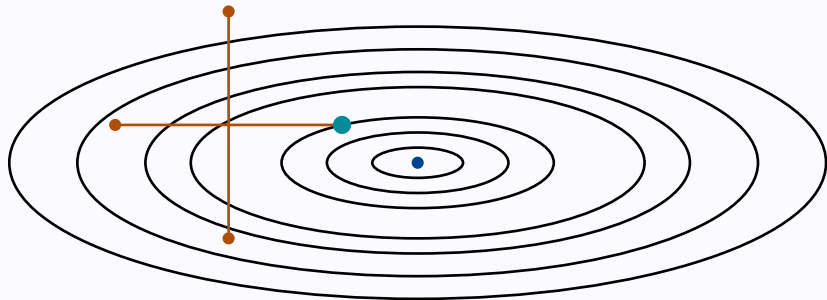
# Example: Coordinate search



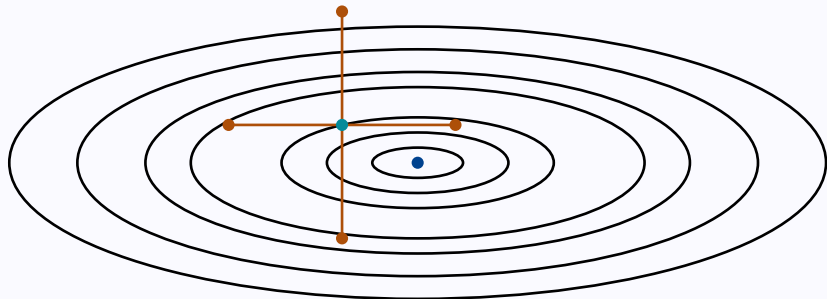
# Example: Coordinate search



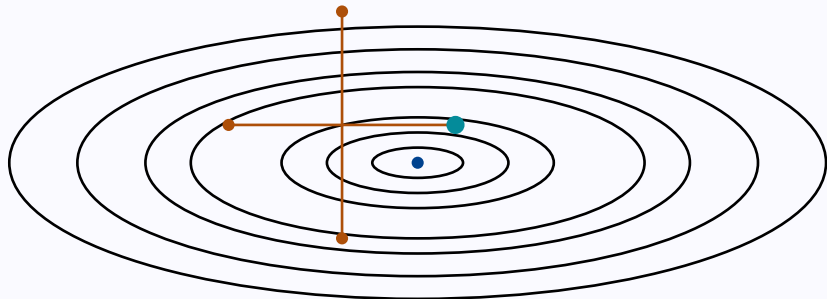
# Example: Coordinate search



# Example: Coordinate search

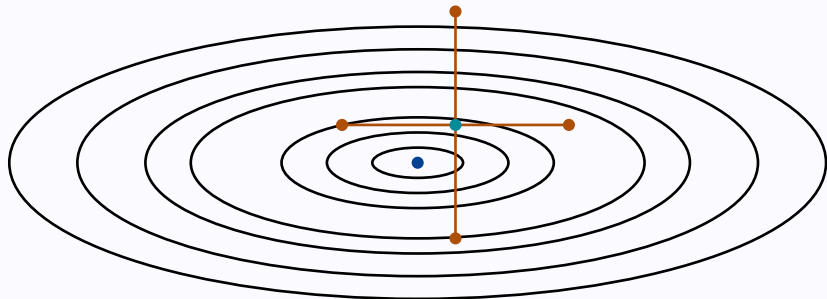


# Example: Coordinate search

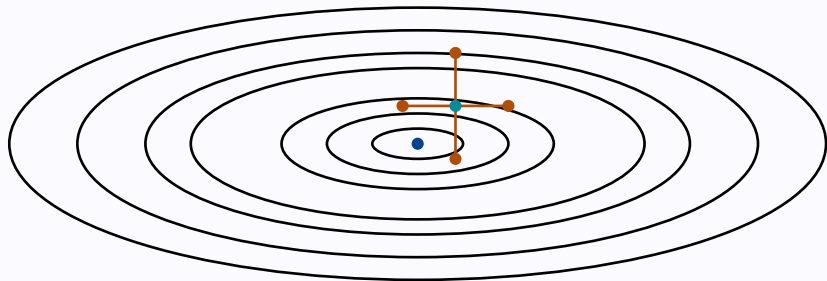




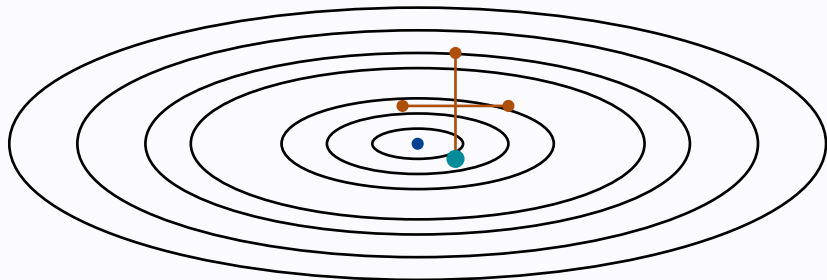
# Example: Coordinate search



# Example: Coordinate search



# Example: Coordinate search



# A (simplified) direct-search framework

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- Choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors.

# A (simplified) direct-search framework

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- Choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors.
- If  $\exists \mathbf{d}_k \in \mathcal{D}_k$  such that

$$f(\mathbf{x}_k + \delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \delta_k^2 \|\mathbf{d}_k\|^2$$

set  $\mathbf{x}_{k+1} := \mathbf{x}_k + \delta_k \mathbf{d}_k$ ,  $\delta_{k+1} := 2\delta_k$ .

# A (simplified) direct-search framework

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- Choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors.
- If  $\exists \mathbf{d}_k \in \mathcal{D}_k$  such that

$$f(\mathbf{x}_k + \delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \delta_k^2 \|\mathbf{d}_k\|^2$$

set  $\mathbf{x}_{k+1} := \mathbf{x}_k + \delta_k \mathbf{d}_k$ ,  $\delta_{k+1} := 2\delta_k$ .

- Otherwise, set  $\mathbf{x}_{k+1} := \mathbf{x}_k$ ,  $\delta_{k+1} := \delta_k/2$ .

# A (simplified) direct-search framework

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- Choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors.
- If  $\exists \mathbf{d}_k \in \mathcal{D}_k$  such that

$$f(\mathbf{x}_k + \delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \delta_k^2 \|\mathbf{d}_k\|^2$$

set  $\mathbf{x}_{k+1} := \mathbf{x}_k + \delta_k \mathbf{d}_k$ ,  $\delta_{k+1} := 2\delta_k$ .

- Otherwise, set  $\mathbf{x}_{k+1} := \mathbf{x}_k$ ,  $\delta_{k+1} := \delta_k/2$ .

# A (simplified) direct-search framework

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- **Choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors.**
- If  $\exists \mathbf{d}_k \in \mathcal{D}_k$  such that

$$f(\mathbf{x}_k + \delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \delta_k^2 \|\mathbf{d}_k\|^2$$

set  $\mathbf{x}_{k+1} := \mathbf{x}_k + \delta_k \mathbf{d}_k$ ,  $\delta_{k+1} := 2\delta_k$ .

- Otherwise, set  $\mathbf{x}_{k+1} := \mathbf{x}_k$ ,  $\delta_{k+1} := \delta_k/2$ .

**Which vectors should we use?**



## A measure of set quality

The set  $\mathcal{D}_k$  is called  $\kappa$ -descent for  $f$  at  $\mathbf{x}_k$  if

$$\max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^T \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\| \|\nabla f(\mathbf{x}_k)\|} \geq \kappa \in (0, 1].$$

## A measure of set quality

The set  $\mathcal{D}_k$  is called  $\kappa$ -descent for  $f$  at  $\mathbf{x}_k$  if

$$\max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^T \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\| \|\nabla f(\mathbf{x}_k)\|} \geq \kappa \in (0, 1].$$

- Guaranteed when  $\mathcal{D}_k$  is a Positive Spanning Set (PSS);
- $\mathcal{D}_k$  PSS  $\Rightarrow |\mathcal{D}_k| \geq n + 1$ ;
- Ex)  $\mathcal{D}_\oplus := \{\mathbf{e}_1, \dots, \mathbf{e}_n, -\mathbf{e}_1, \dots, -\mathbf{e}_n\}$  is always  $\frac{1}{\sqrt{n}}$ -descent.

**Assumption:** For every  $k$ ,  $\mathcal{D}_k$  is  $\kappa$ -descent and contains  $m$  unit directions.

## Theorem (Vicente '12)

Let  $\epsilon \in (0, 1)$  and  $N_\epsilon$  be the number of function evaluations needed to reach  $\mathbf{x}_k$  such that  $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ . Then,

$$N_\epsilon \leq \mathcal{O}(m \kappa^{-2} \epsilon^{-2}).$$

# Complexity of deterministic direct search

**Assumption:** For every  $k$ ,  $\mathcal{D}_k$  is  $\kappa$ -descent and contains  $m$  unit directions.

## Theorem (Vicente '12)

Let  $\epsilon \in (0, 1)$  and  $N_\epsilon$  be the number of function evaluations needed to reach  $\mathbf{x}_k$  such that  $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ . Then,

$$N_\epsilon \leq \mathcal{O}(m \kappa^{-2} \epsilon^{-2}).$$

- Unit norm can be replaced by bounded norm.
- Choosing  $\mathcal{D}_k = \mathcal{D}_\oplus$ , one has  $\kappa = \frac{1}{\sqrt{n}}$ ,  $m = 2n$ , and the bound becomes

$$N_\epsilon \leq \mathcal{O}(n^2 \epsilon^{-2}).$$

$\Rightarrow$  **Best possible dependency** w.r.t.  $n$  for **deterministic** direct-search algorithms.

## Classical direct search

- Set  $\mathcal{D}_k \subset \mathbb{R}^n$ ,  $|\mathcal{D}_k| = m$ ,  $\text{cm}(\mathcal{D}_k) \geq \kappa$ ;
- Complexity:

$$\mathcal{O}(m\kappa^{-2}\epsilon^{-2}).$$

- $m$  depends on  $n$  ( $m \geq n + 1$ ).
- $\kappa$  depends on  $n$  (approximate  $\nabla f(\mathbf{x}_k) \in \mathbb{R}^n$ ).

## Classical direct search

- Set  $\mathcal{D}_k \subset \mathbb{R}^n$ ,  $|\mathcal{D}_k| = m$ ,  $\text{cm}(\mathcal{D}_k) \geq \kappa$ ;
- Complexity:

$$\mathcal{O}(m\kappa^{-2}\epsilon^{-2}).$$

- $m$  depends on  $n$  ( $m \geq n + 1$ ).
- $\kappa$  depends on  $n$  (approximate  $\nabla f(\mathbf{x}_k) \in \mathbb{R}^n$ ).

## My original thought

- Generate directions in random subspaces of  $\mathbb{R}^n$ ;
- Use results from dimensionality reduction;
- Remove all dependencies on  $n$ !

# Randomizing direct search

## Classical direct search

- Set  $\mathcal{D}_k \subset \mathbb{R}^n$ ,  $|\mathcal{D}_k| = m$ ,  $\text{cm}(\mathcal{D}_k) \geq \kappa$ ;
- Complexity:

$$\mathcal{O}(m\kappa^{-2} \epsilon^{-2}).$$

- $m$  depends on  $n$  ( $m \geq n + 1$ ).
- $\kappa$  depends on  $n$  (approximate  $\nabla f(\mathbf{x}_k) \in \mathbb{R}^n$ ).

## My original thought

- Generate directions in random subspaces of  $\mathbb{R}^n$ ;
- Use results from dimensionality reduction;
- Remove all dependencies on  $n$ !

**Spoiler alert:** You can only *reduce* the dependency on  $n$ .

# What can you do?

## Our approach

- Consider a random subspace of dimension  $r \leq n$ ;
- Use a PSS to approximate the projected gradient in the subspace;
- Guarantee sufficient gradient information **in probability**.

## What it brings us

- Use random directions.
- Possibly less than  $n$ .
- Possibly **unbounded**.



## Probabilistic descent (Gratton et al '15)

- Use directions  $[\mathbf{d} - \mathbf{d}^*]$  with  $\mathbf{d} \sim \mathcal{U}(\mathbb{S}^{n-1})$ .
- Complexity improves from  $\mathcal{O}(n^2\epsilon^{-2})$  to  $\mathcal{O}(n\epsilon^{-2})$  ( $m = 2$ ).
- Limited to one distribution.

# Not the only game in town (1/2)

## Probabilistic descent (Gratton et al '15)

- Use directions  $[\mathbf{d} - \mathbf{d}]$  with  $\mathbf{d} \sim \mathcal{U}(\mathbb{S}^{n-1})$ .
- Complexity improves from  $\mathcal{O}(n^2\epsilon^{-2})$  to  $\mathcal{O}(n\epsilon^{-2})$  ( $m = 2$ ).
- Limited to one distribution.

**Gaussian smoothing approach:** Draw  $\mathbf{d} \sim \mathcal{N}(0, \mathbf{I})$  and use

$$\frac{f(\mathbf{x} + \delta\mathbf{d}) - f(\mathbf{x})}{\delta}\mathbf{d} \quad \text{or} \quad \frac{f(\mathbf{x} + \delta\mathbf{d}) - f(\mathbf{x} - \delta\mathbf{d})}{\delta}\mathbf{d}.$$

*Random gradient-free method (Nesterov and Spokoiny 2017),  
Stochastic three-point method (Bergou et al, 2020).*

- Also achieve  $\mathcal{O}(n\epsilon^{-2})$  bound.
- Use one-dimensional subspace based on Gaussian vectors.
- Use fixed or decreasing stepsizes.

## Zeroth-order (Kozak et al '21, '22)

- Estimate directional derivatives directly.
- Use orthogonal random directions  $\mathbf{Q} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .
- Complexity results for convex/PL functions.

## Zeroth-order (Kozak et al '21, '22)

- Estimate directional derivatives directly.
- Use orthogonal random directions  $\mathbf{Q} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .
- Complexity results for convex/PL functions.

## Our approach

- General, **subspace-based** framework.
- Inspiration: Model-based methods (Cartis and Roberts '23, Dzhahini and Wild '22a).

- 1 Derivative-free algorithm
- 2 Reduced subspace approach**
- 3 Numerics with subspaces
- 4 Subspace dimensions

**Inputs:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\delta_0 > 0$ .

**Iteration  $k$ :** Given  $(\mathbf{x}_k, \delta_k)$ ,

- Choose  $\mathbf{P}_k \in \mathbb{R}^{r \times n}$  **at random**.
- Choose  $\mathcal{D}_k \subset \mathbb{R}^r$  having  $m$  vectors.
- If  $\exists \mathbf{d}_k \in \mathcal{D}_k$  such that

$$f(\mathbf{x}_k + \delta_k \mathbf{P}_k^T \mathbf{d}_k) < f(\mathbf{x}_k) - \delta_k^2 \|\mathbf{P}_k^T \mathbf{d}_k\|^2,$$

set  $\mathbf{x}_{k+1} := \mathbf{x}_k + \delta_k \mathbf{P}_k^T \mathbf{d}_k$ ,  $\delta_{k+1} := 2\delta_k$ .

- Otherwise, set  $\mathbf{x}_{k+1} := \mathbf{x}_k$ ,  $\delta_{k+1} := \delta_k/2$ .

## New polling sets

$$\{\mathbf{P}_k^T \mathbf{d} \mid \mathbf{d} \in \mathcal{D}_k\} \subset \mathbb{R}^n.$$

- $\mathbf{P}_k \in \mathbb{R}^{r \times n}$ : Maps onto  $r$ -dimensional subspace;
- $\mathcal{D}_k$ : Direction set in  $\mathbb{R}^r$ .

## What do we want?

- Preserve information while applying  $\mathbf{P}_k / \mathbf{P}_k^T$ .
- Approximate  $-\mathbf{P}_k \nabla f(\mathbf{x}_k)$  using  $\mathcal{D}_k$ .

$\mathbf{P}_k$  is  $(\eta, \sigma, P_{\max})$ -well aligned for  $(f, \mathbf{x}_k)$  if

$$\left\{ \begin{array}{l} \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\| \geq \eta \|\nabla f(\mathbf{x}_k)\|, \\ \sigma_{\min}(\mathbf{P}_k) \geq \sigma, \\ \sigma_{\max}(\mathbf{P}_k) \leq P_{\max}. \end{array} \right.$$



$\mathbf{P}_k$  is  $(\eta, \sigma, P_{\max})$ -well aligned for  $(f, \mathbf{x}_k)$  if

$$\begin{cases} \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\| \geq \eta \|\nabla f(\mathbf{x}_k)\|, \\ \sigma_{\min}(\mathbf{P}_k) \geq \sigma, \\ \sigma_{\max}(\mathbf{P}_k) \leq P_{\max}. \end{cases}$$

Ex)  $\mathbf{P}_k = \mathbf{I}_n \in \mathbb{R}^{n \times n}$  is  $(1, 1, 1)$ -well aligned.

$\mathbf{P}_k$  is  $(\eta, \sigma, P_{\max})$ -well aligned for  $(f, \mathbf{x}_k)$  if

$$\begin{cases} \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\| \geq \eta \|\nabla f(\mathbf{x}_k)\|, \\ \sigma_{\min}(\mathbf{P}_k) \geq \sigma, \\ \sigma_{\max}(\mathbf{P}_k) \leq P_{\max}. \end{cases}$$

Ex)  $\mathbf{P}_k = \mathbf{I}_n \in \mathbb{R}^{n \times n}$  is  $(1, 1, 1)$ -well aligned.

## Probabilistic version

$\{\mathbf{P}_k\}$  is  $(q, \eta, \sigma, P_{\max})$ -well aligned if:

$$\begin{aligned} & \mathbb{P}(\mathbf{P}_0 \text{ } (q, \eta, \sigma, P_{\max})\text{-well aligned}) \geq q \\ \forall k \geq 1, & \quad \mathbb{P}((q, \eta, \sigma, P_{\max})\text{-well aligned} \mid \mathbf{P}_0, \mathcal{D}_0, \dots, \mathbf{P}_{k-1}, \mathcal{D}_{k-1}) \geq q, \end{aligned}$$

## Deterministic descent

The set  $\mathcal{D}_k$  is  $(\kappa, d_{\max})$ -descent for  $(f, \mathbf{x}_k)$  if

$$\left\{ \begin{array}{l} \max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^\top \mathbf{P}_k \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\| \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\|} \geq \kappa, \\ \forall \mathbf{d} \in \mathcal{D}_k, \quad d_{\max}^{-1} \leq \|\mathbf{d}\| \leq d_{\max}. \end{array} \right.$$

## Deterministic descent

The set  $\mathcal{D}_k$  is  $(\kappa, d_{\max})$ -descent for  $(f, \mathbf{x}_k)$  if

$$\left\{ \begin{array}{l} \max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^\top \mathbf{P}_k \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\| \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\|} \geq \kappa, \\ \forall \mathbf{d} \in \mathcal{D}_k, \quad d_{\max}^{-1} \leq \|\mathbf{d}\| \leq d_{\max}. \end{array} \right.$$

Ex)  $D_{\oplus} = \{\mathbf{e}_1, \dots, \mathbf{e}_n, -\mathbf{e}_1, \dots, -\mathbf{e}_n\}$  is  $(\frac{1}{\sqrt{n}}, 1)$ -descent.

## Deterministic descent

The set  $\mathcal{D}_k$  is  $(\kappa, d_{\max})$ -descent for  $(f, \mathbf{x}_k)$  if

$$\left\{ \begin{array}{l} \max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^\top \mathbf{P}_k \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\| \|\mathbf{P}_k \nabla f(\mathbf{x}_k)\|} \geq \kappa, \\ \forall \mathbf{d} \in \mathcal{D}_k, \quad d_{\max}^{-1} \leq \|\mathbf{d}\| \leq d_{\max}. \end{array} \right.$$

Ex)  $\mathcal{D}_\oplus = \{\mathbf{e}_1, \dots, \mathbf{e}_n, -\mathbf{e}_1, \dots, -\mathbf{e}_n\}$  is  $(\frac{1}{\sqrt{n}}, 1)$ -descent.

## Probabilistic descent sets

$\{\mathcal{D}_k\}$  is  $(p, \kappa, d_{\max})$ -descent if:

$$\begin{aligned} & \mathbb{P}(\mathcal{D}_0 \text{ } (\kappa, d_{\max})\text{-descent} \mid \mathbf{P}_0) \geq p \\ \forall k \geq 1, & \quad \mathbb{P}(\mathcal{D}_k \text{ } (\kappa, d_{\max})\text{-descent} \mid \mathbf{P}_0, \mathcal{D}_0, \dots, \mathbf{P}_{k-1}, \mathcal{D}_{k-1}, \mathbf{P}_k) \geq p, \end{aligned}$$

## Theorem (Roberts, R. '23)

Assume:

- $\{\mathcal{D}_k\}$   $(p, \kappa, d_{\max})$ -descent,  $|\mathcal{D}_k| = m$ ;
- $\{\mathbf{P}_k\}$   $(q, \eta, \sigma, P_{\max})$ -well aligned,  $pq > \frac{1}{2}$ .

Let  $N_\epsilon$  the number of function evaluations needed to have  $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ .

$$\mathbb{P}\left(N_\epsilon \leq \mathcal{O}\left(\frac{m\phi\epsilon^{-2}}{2pq-1}\right)\right) \geq 1 - \exp\left(-\mathcal{O}\left(\frac{2pq-1}{pq}\phi\epsilon^{-2}\right)\right).$$

where  $\phi = d_{\max}^8 \kappa^{-2} \eta^{-2} \sigma^{-2} P_{\max}^4$ .

## Theorem (Roberts, R. '23)

Assume:

- $\{\mathcal{D}_k\}$   $(p, \kappa, d_{\max})$ -descent,  $|\mathcal{D}_k| = m$ ;
- $\{\mathbf{P}_k\}$   $(q, \eta, \sigma, P_{\max})$ -well aligned,  $pq > \frac{1}{2}$ .

Let  $N_\epsilon$  the number of function evaluations needed to have  $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ .

$$\mathbb{P}\left(N_\epsilon \leq \mathcal{O}\left(\frac{m\phi\epsilon^{-2}}{2pq-1}\right)\right) \geq 1 - \exp\left(-\mathcal{O}\left(\frac{2pq-1}{pq}\phi\epsilon^{-2}\right)\right).$$

where  $\phi = d_{\max}^8 \kappa^{-2} \eta^{-2} \sigma^{-2} P_{\max}^4$ .

Does this bound depend on  $n$ ?

$$m\phi\epsilon^{-2} = m d_{\max}^8 \kappa^{-2} \eta^{-2} \sigma^{-2} P_{\max}^4 \epsilon^{-2}.$$



$$m\phi\epsilon^{-2} = m d_{\max}^8 \kappa^{-2} \eta^{-2} \sigma^{-2} P_{\max}^4 \epsilon^{-2}.$$

## Best directions in subspaces

- $\mathcal{D}_k = \{\mathbf{e}_1, \dots, \mathbf{e}_r, -\mathbf{e}_1, \dots, -\mathbf{e}_r\}$  in  $\mathbb{R}^r$ ;
- $\kappa = \frac{1}{\sqrt{r}}$ ,  $m = 2r$ ,  $d_{\max} = 1$ .

$\Rightarrow$  With  $r = \mathcal{O}(1)$ ,  $m d_{\max}^8 \kappa^{-2} = \mathcal{O}(1)$ !

# Complexity and dimension dependencies

$$m\phi\epsilon^{-2} = \mathcal{O}(1)\eta^{-2}\sigma^{-2}P_{\max}^4\epsilon^{-2}.$$

## Best directions in subspaces

- $\mathcal{D}_k = \{\mathbf{e}_1, \dots, \mathbf{e}_r, -\mathbf{e}_1, \dots, -\mathbf{e}_r\}$  in  $\mathbb{R}^r$ ;
- $\kappa = \frac{1}{\sqrt{r}}$ ,  $m = 2r$ ,  $d_{\max} = 1$ .

$\Rightarrow$  With  $r = \mathcal{O}(1)$ ,  $m d_{\max}^8 \kappa^{-2} = \mathcal{O}(1)$ !

## Best subspaces?

$P_k$	$\sigma$	$P_{\max}$
Gaussian	$\Theta(\sqrt{n/r})$	$\Theta(\sqrt{n/r})$
Hashing	$\Theta(\sqrt{n/r})$ (Dzahini & Wild '22b)	$\sqrt{n}$
Orthogonal	$\sqrt{n/r}$	$\sqrt{n/r}$ .

$\Rightarrow$  Even with  $r = \mathcal{O}(1)$  and  $\eta = \mathcal{O}(1)$ ,  $\eta^{-2}\sigma^{-2}P_{\max}^4 = \mathcal{O}(n)$ !

# Complexity and dimension dependencies

$$m\phi\epsilon^{-2} = \mathcal{O}(1) \mathcal{O}(n)\epsilon^{-2}.$$

## Best directions in subspaces

- $\mathcal{D}_k = \{\mathbf{e}_1, \dots, \mathbf{e}_r, -\mathbf{e}_1, \dots, -\mathbf{e}_r\}$  in  $\mathbb{R}^r$ ;
- $\kappa = \frac{1}{\sqrt{r}}$ ,  $m = 2r$ ,  $d_{\max} = 1$ .

$\Rightarrow$  With  $r = \mathcal{O}(1)$ ,  $m d_{\max}^8 \kappa^{-2} = \mathcal{O}(1)$ !

## Best subspaces?

$P_k$	$\sigma$	$P_{\max}$
Gaussian	$\Theta(\sqrt{n/r})$	$\Theta(\sqrt{n/r})$
Hashing	$\Theta(\sqrt{n/r})$ (Dzahini & Wild '22b)	$\sqrt{n}$
Orthogonal	$\sqrt{n/r}$	$\sqrt{n/r}$ .

$\Rightarrow$  Even with  $r = \mathcal{O}(1)$  and  $\eta = \mathcal{O}(1)$ ,  $\eta^{-2}\sigma^{-2}P_{\max}^4 = \mathcal{O}(n)$ !

- Can compute steps in  $r$ -dim. subspaces,  $r = \mathcal{O}(1)$ .
- Reduced evaluation cost per iteration.
- Complexity:  $\mathcal{O}(n^2) \Rightarrow \mathcal{O}(n)$ !

- 1 Derivative-free algorithm
- 2 Reduced subspace approach
- 3 Numerics with subspaces**
- 4 Subspace dimensions

## Benchmark:

- Medium-scale test set (90 CUTEst problems of dimension  $\approx 100$ );
- Large-scale test set (28 CUTEst problems of dimension  $\approx 1000$ ).

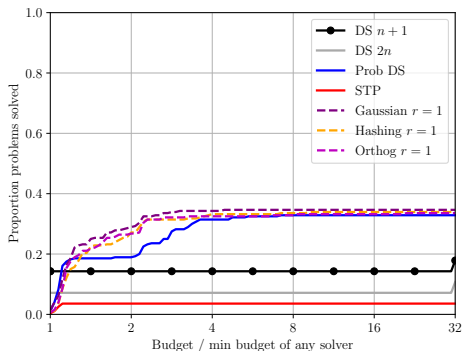
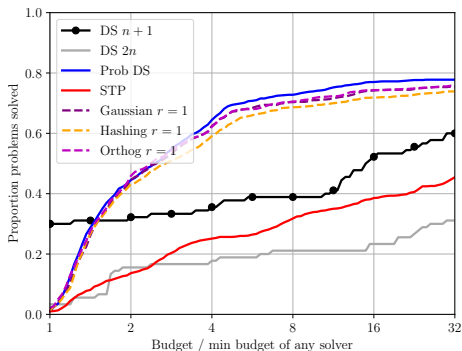
**Budget:**  $200(n + 1)$  evaluations.

## Comparison:

- Deterministic DS with  $\mathcal{D}_k = \mathcal{D}_{\oplus}$  or  $\mathcal{D}_k = \{\mathbf{e}_1, \dots, \mathbf{e}_n, -\sum_{i=1}^n \mathbf{e}_i\}$ ;
- Probabilistic direct search with 2 uniform directions;
- Stochastic Three Point;
- Probabilistic direct search with Gaussian/Hashing/Orthogonal  $\mathbf{P}_k$  matrices + 2 directions in the subspace.

**Goal:** Satisfy  $f(\mathbf{x}_k) - f_{opt} \leq 0.1(f(\mathbf{x}_0) - f_{opt})$ .

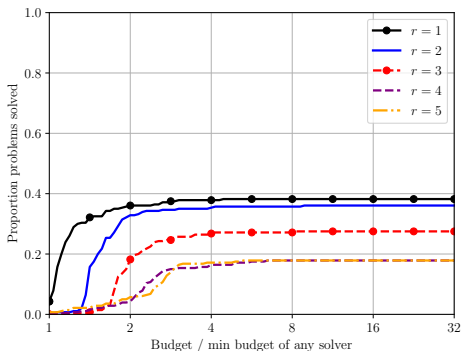
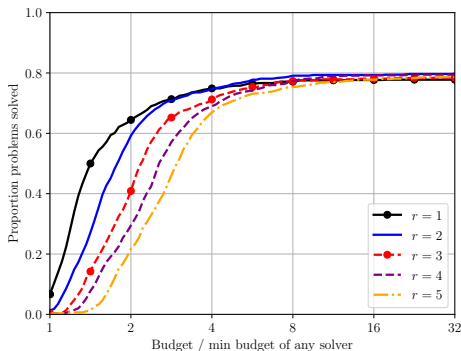
# Comparison of all methods



Left: Medium scale; Right: Large scale.

- Operating in random subspaces works!
- But always a (hidden) dependency on  $n$ !

# Gaussian matrices and subspace dimensions



Left: Medium scale; Right: Large scale.

## Numerically

- Sketches of dimension  $> 1$  may improve things...
- ...but in general opposite (Gaussian) directions work best!



## The package

- `https://github.com/lindonroberts/directsearch`
- Python code + paper experiments.
- `pip install directsearch`

## The package

- `https://github.com/lindonroberts/directsearch`
- Python code + paper experiments.
- `pip install directsearch`

## Recent use at Meta:



**Olivier Teytaud**

Admin · 23 janvier · 🌐



In progress: adding <https://github.com/lindonroberts/directsearch> inside Nevergrad.

In particular there is an excellent stochastic direct search method. I don't know exactly the algorithm (yet). Thanks guys for this excellent code!

- 1 Derivative-free algorithm
- 2 Reduced subspace approach
- 3 Numerics with subspaces
- 4 Subspace dimensions**

## If you want to scale up...

- Can compute steps in  $r$ -dim. subspaces,  $r = \mathcal{O}(1)$ ;
- Reduced evaluation cost per iteration;
- Overall complexity:  $\mathcal{O}(n^2) \Rightarrow \mathcal{O}(n)$ !

## Numerically

- Subspaces of dimension  $r > 1$  may be good...
- ...but in general opposite Gaussian directions ( $r = 1$ ) are better!

Warren: “But *why* does this work?”

Why do 1-dim. subspaces give best performance?

## Why do 1-dim. subspaces give best performance?

### Our approach: Expected decrease guarantees

- Use Taylor approximation to focus on linear functions

$$f(\mathbf{x} + \mathbf{v}) - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^T \mathbf{v} + \frac{L}{2} \|\mathbf{v}\|^2.$$

- Generate  $\mathbf{v}$  in a random subspace.
- Analyze **expected value** of linear term:

$$\mathbb{E}_{\mathbf{v}} [\nabla f(\mathbf{x})^T \mathbf{v}].$$

## Why do 1-dim. subspaces give best performance?

### Our approach: Expected decrease guarantees

- Use Taylor approximation to focus on linear functions

$$f(\mathbf{x} + \mathbf{v}) - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^T \mathbf{v} + \frac{L}{2} \|\mathbf{v}\|^2.$$

- Generate  $\mathbf{v}$  in a random subspace.
- Analyze **expected value** of linear term:

$$\mathbb{E}_{\mathbf{v}} [\nabla f(\mathbf{x})^T \mathbf{v}].$$

- Equivalently, consider random  $\mathbf{g} \in \mathbb{R}^n$ , deterministic  $\mathbf{v}$ :

$$\mathbb{E}_{\mathbf{g}} [\mathbf{g}^T \mathbf{v}].$$

## Key result (Hare, Roberts, R. '22)

Let  $\mathbf{g} \in \mathbb{S}^{n-1}$ ,  $\mathbf{P} \in \mathbb{R}^{r \times n}$  and  $\mathcal{D} = \{\mathbf{e}_1, \dots, \mathbf{e}_r, -\mathbf{e}_1, \dots, -\mathbf{e}_r\}$ .

Then, the expected decrease ratio

$$\frac{\mathbb{E} [\min_{\mathbf{d} \in \mathcal{D}} \mathbf{g}^T \mathbf{P}^T \mathbf{d}]}{2r}$$

is minimized at  $r = 1$ .

## Side notes

- Key quantity:

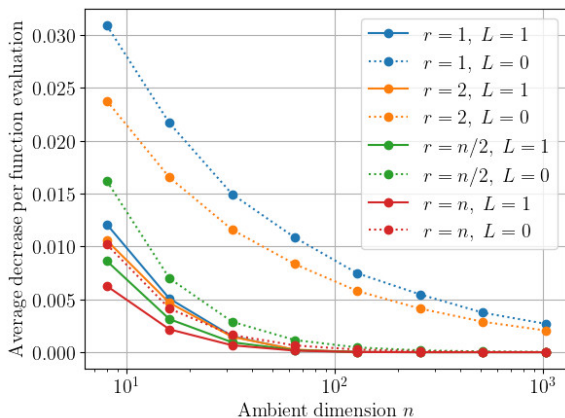
$$\mathbb{E}_{\mathbf{u} \sim \mathcal{U}(\mathbb{S}^{n-1})} \left[ \max_{1 \leq i \leq r} |[\mathbf{u}]_i| \right].$$

- Exact values hard to find in the literature!
- $r = 1$ : best “bang for your buck”.



## Setup

- Monte-Carlo approximations of expected decrease.
- Quadratic functions with a random linear term  $\mathbf{x} \mapsto \mathbf{g}^T \mathbf{x} + \frac{L}{2} \|\mathbf{x}\|^2$ .
- Normalization by the number of function evaluations.



## Our findings

- Probabilistic analysis/subspace viewpoint.
- Good complexity ( $\mathcal{O}(n)$ ).
- Low dimension provably better on average.

## Our findings

- Probabilistic analysis/subspace viewpoint.
- Good complexity ( $\mathcal{O}(n)$ ).
- Low dimension provably better on average.

## Going further

- Model-based algorithms (done for linear models).
- Stochastic/Noisy function values.

## References

- *Direct search based on probabilistic descent in reduced spaces*  
L. Roberts and C. W. Royer, SIAM J. Optim. 33(4):3057-3082, 2023.
- *Expected decrease for derivative-free algorithms using random subspaces*  
W. Hare, L. Roberts and C. W. Royer, Technical report  
arXiv:2308.04734v2, 2024.

## References

- *Direct search based on probabilistic descent in reduced spaces*  
L. Roberts and C. W. Royer, SIAM J. Optim. 33(4):3057-3082, 2023.
- *Expected decrease for derivative-free algorithms using random subspaces*  
W. Hare, L. Roberts and C. W. Royer, Technical report  
arXiv:2308.04734v2, 2024.

*Merci!* `clement.royer@lamsade.dauphine.fr`