

Revisiting derivative-free optimization algorithms in the context of black-box adversarial attacks

Clément Royer (Université Paris Dauphine-PSL)

Séminaire Statistiques et Optimisation

Toulouse, 10 mars 2026





What is MILES?

- Team-project within Dauphine.
- Focus: **Trustworthy AI**.
- Expertise: Machine learning, Game theory, Optimization,...

We have an MCF opening for 2026!



What is MILES?

- Team-project within Dauphine.
- Focus: **Trustworthy AI**.
- Expertise: Machine learning, Game theory, Optimization,...

We have an MCF opening for 2026!

One of our strengths: Adversarial attacks

Perturb data to fool models

- Computing attacks (Meunier et al '19, Timoz et al '26).
- Game theoretical perspective (Meunier et al '21).
- Robustify models via randomization (Pinot et al. '19 '20 '22, Gnecco Heredia et al '24).

COMPROMIS project (PEPR Cybersécurité)

- Thesis of Bastien Cavarretta.
- Goal: Develop structured adversarial attacks in a black-box setting using derivative-free optimization.

COMPROMIS project (PEPR Cybersécurité)

- Thesis of Bastien Cavarretta.
- Goal: Develop structured adversarial attacks in a black-box setting using derivative-free optimization.

What we started with

- Zoo of black-box adversarial attacks.
- Derivative-free methods somewhat disregarded because of scalability issues.

COMPROMIS project (PEPR Cybersécurité)

- Thesis of Bastien Cavarretta.
- Goal: Develop structured adversarial attacks in a black-box setting using derivative-free optimization.

What we started with

- Zoo of black-box adversarial attacks.
- Derivative-free methods somewhat disregarded because of scalability issues.

What we have now

- An algorithmic framework with enhanced scalability (theory).
- Preliminary numerical validation.

- 1 Adversarial attacks
- 2 Black-box adversarial attacks
- 3 Adversarial attacks via direct search

- 1 Adversarial attacks
- 2 Black-box adversarial attacks
- 3 Adversarial attacks via direct search

Basic classification paradigm

- Data points $\{\mathbf{x}_i, y_i\}_{i=1, \dots, N}$ with
 - \mathbf{x}_i input (e.g. image) viewed as a vector in $\mathcal{X} \subset \mathbb{R}^n$.
 - $y_i \in \{1, \dots, C\}$ discrete (class) label.
- Goal: Predict y_i given \mathbf{x}_i through a model h .
- Training: Compute best h within a model class \mathcal{H} .

Basic classification paradigm

- Data points $\{\mathbf{x}_i, y_i\}_{i=1, \dots, N}$ with
 - \mathbf{x}_i input (e.g. image) viewed as a vector in $\mathcal{X} \subset \mathbb{R}^n$.
 - $y_i \in \{1, \dots, C\}$ discrete (class) label.
- Goal: Predict y_i given \mathbf{x}_i through a model h .
- Training: Compute best h within a model class \mathcal{H} .

Empirical risk minimization

$$\underset{h \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{N} \sum_{i=1}^N \mathcal{L}(h(\mathbf{x}_i), y_i).$$

- Typical model: Logit function $h(\mathbf{x}) \in \mathbb{R}^C$ output of a neural network.
- Typical loss: Cross entropy $\mathcal{L}(h(\mathbf{x}), y) = -\log \left(\frac{\exp(h(\mathbf{x})^y)}{\sum_{c=1}^C \exp(h(\mathbf{x})^c)} \right)$.

What could go wrong in classification?

Suppose

- Data $\{(\mathbf{x}_i, y_i)\}$ (e.g. ImageNet).
- Model h obtained by training on this data?

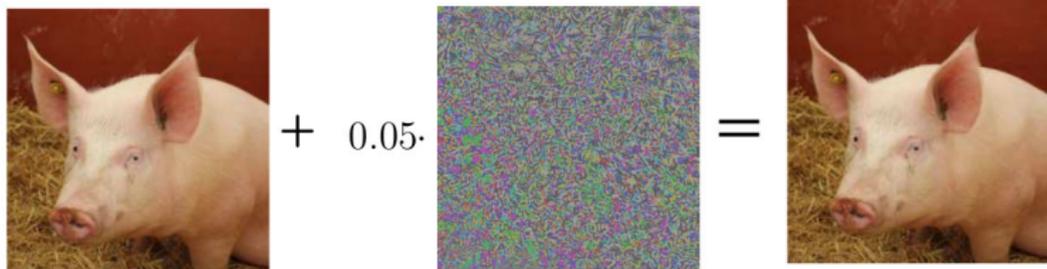
Can the model be fooled by small perturbations of data points?

What could go wrong in classification?

Suppose

- Data $\{(\mathbf{x}_i, y_i)\}$ (e.g. ImageNet).
- Model h obtained by training on this data?

Can the model be fooled by small perturbations of data points?



A hog misclassified as an airliner by a ResNet50 with high confidence
<https://adversarial-ml-tutorial.org/introduction/>

Adversarial examples/attacks

- Perturbation of data points that fool models (typically not humans).
- Observed in spam detection (Dalvi et al '04), speech recognition, language processing, etc.

Adversarial examples/attacks

- Perturbation of data points that fool models (typically not humans).
- Observed in spam detection (Dalvi et al '04), speech recognition, language processing, etc.



Perturbed data points from ImageNet classified as Ostrich by AlexNet (Szegedy et al '14).

How are adversarial attacks computed?

How are adversarial attacks computed?

Adversarial attack from data point (\mathbf{x}, y)

$$\underset{\mathbf{z} \in \mathcal{B}(\mathbf{x}, \epsilon) \cap \mathcal{X}}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{z}), \mathbf{y}).$$

- $\mathcal{B}(\mathbf{x}, \epsilon) = \{\mathbf{z} \mid \|\mathbf{z} - \mathbf{x}\| \leq \epsilon\}$, $\|\cdot\|$ norm on \mathbb{R}^n .
- h either un-trained (data poisoning) or **trained (evasion attack) model**.

→ Focus for this talk: Evasion attacks.

How are adversarial attacks computed?

Adversarial attack from data point (\mathbf{x}, y)

$$\underset{\mathbf{z} \in \mathcal{B}(\mathbf{x}, \epsilon) \cap \mathcal{X}}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{z}), \mathbf{y}).$$

- $\mathcal{B}(\mathbf{x}, \epsilon) = \{\mathbf{z} \mid \|\mathbf{z} - \mathbf{x}\| \leq \epsilon\}$, $\|\cdot\|$ norm on \mathbb{R}^n .
- h either un-trained (data poisoning) or **trained (evasion attack) model**.

→ Focus for this talk: Evasion attacks.

Original attack (Szegedy et al '14)

$$\underset{\mathbf{z}}{\text{minimize}} \quad \|\mathbf{z} - \mathbf{x}\|_2 \quad \text{s.t.} \quad h(\mathbf{z}) = y_{adv}, \mathbf{z} \in \mathcal{X}.$$

→ Targeted attack (label y_{adv})

→ Nonlinear optimization problem (solved with L-BFGS).

Gradient-based attacks

Idea Solve the attack problem inexactly using gradient-based methods.

Gradient-based attacks

Idea Solve the attack problem inexactly using gradient-based methods.

Fast Gradient Method (Goodfellow et al '15)

$$\underset{\mathbf{z} \in \mathcal{B}(\mathbf{x}, \epsilon)}{\text{maximize}} \quad (\mathbf{z} - \mathbf{x})^T \nabla_{\mathbf{x}} \mathcal{L}(h(\mathbf{x}), y).$$

- First-order approximation of $\mathcal{L}(h(\mathbf{z}), y) - \mathcal{L}(h(\mathbf{x}), y)$.
- Compute $\nabla_{\mathbf{x}} \mathcal{L}(h(\mathbf{x}), y)$ through automatic differentiation.
- Closed-form expressions when using $\|\cdot\|_2$ and $\|\cdot\|_{\infty}$.

Gradient-based attacks

Idea Solve the attack problem inexactly using gradient-based methods.

Fast Gradient Method (Goodfellow et al '15)

$$\underset{\mathbf{z} \in \mathcal{B}(\mathbf{x}, \epsilon)}{\text{maximize}} \quad (\mathbf{z} - \mathbf{x})^T \nabla_{\mathbf{x}} \mathcal{L}(h(\mathbf{x}), y).$$

- First-order approximation of $\mathcal{L}(h(\mathbf{z}), y) - \mathcal{L}(h(\mathbf{x}), y)$.
- Compute $\nabla_{\mathbf{x}} \mathcal{L}(h(\mathbf{x}), y)$ through automatic differentiation.
- Closed-form expressions when using $\|\cdot\|_2$ and $\|\cdot\|_{\infty}$.

Projected Gradient Descent/PGD (Kurakin et al '17, Madry et al '18)

$$\mathbf{z}_{t+1} = \Pi_{\mathcal{B}(\mathbf{x}, \epsilon)} (\mathbf{z}_t + \alpha \nabla_{\mathbf{x}} \mathcal{L}(h(\mathbf{z}_t), y))$$

- $\Pi_{\mathcal{B}(\mathbf{x}, \epsilon)}$: Projection.
- Multiple steps of gradient ascent for the problem.

More attacks, and defenses (sample!)

- **CW Attack** (Carlini & Wagner '17b): Modified version of (Szegedy et al '14) using changes of variables.
→ (Carlini & Wagner '17a): Circumvent detection of adversarial examples.

More attacks, and defenses (sample!)

- **CW Attack** (Carlini & Wagner '17b): Modified version of (Szegedy et al '14) using changes of variables.
→ (Carlini & Wagner '17a): Circumvent detection of adversarial examples.
- **Adversarial training** (Goodfellow et al '15, Madry et al '18): Train classifier with adversarial examples+training data.

More attacks, and defenses (sample!)

- **CW Attack** (Carlini & Wagner '17b): Modified version of (Szegedy et al '14) using changes of variables.
→ (Carlini & Wagner '17a): Circumvent detection of adversarial examples.
- **Adversarial training** (Goodfellow et al '15, Madry et al '18): Train classifier with adversarial examples+training data.
- **Circumventing existing defenses**
 - (Athalye et al '18) 7/9 defenses from ICLR 2018.
 - (Tramer et al '20) 13 defenses from ICML/ICLR/NeurIPS.
 - (Sitawarin et al '22) breaks (Raff et al '19).

More attacks, and defenses (sample!)

- **CW Attack** (Carlini & Wagner '17b): Modified version of (Szegedy et al '14) using changes of variables.
→ (Carlini & Wagner '17a): Circumvent detection of adversarial examples.
- **Adversarial training** (Goodfellow et al '15, Madry et al '18): Train classifier with adversarial examples+training data.
- **Circumventing existing defenses**
 - (Athalye et al '18) 7/9 defenses from ICLR 2018.
 - (Tramer et al '20) 13 defenses from ICML/ICLR/NeurIPS.
 - (Sitawarin et al '22) breaks (Raff et al '19).
- **Methodology for evaluating defenses**
 - (Carlini et al '19): Good practice.
 - **RobustBench** (Croce et al '21): Standardized benchmarks.

More attacks, and defenses (sample!)

- **CW Attack** (Carlini & Wagner '17b): Modified version of (Szegedy et al '14) using changes of variables.
→ (Carlini & Wagner '17a): Circumvent detection of adversarial examples.
 - **Adversarial training** (Goodfellow et al '15, Madry et al '18): Train classifier with adversarial examples+training data.
 - **Circumventing existing defenses**
 - (Athalye et al '18) 7/9 defenses from ICLR 2018.
 - (Tramer et al '20) 13 defenses from ICML/ICLR/NeurIPS.
 - (Sitawarin et al '22) breaks (Raff et al '19).
 - **Methodology for evaluating defenses**
 - (Carlini et al '19): Good practice.
 - **RobustBench** (Croce et al '21): Standardized benchmarks.
- A never-ending game of cat and mouse...
 - Randomized defenses (popular topic in MILES) not necessarily better (Lucas et al '23).

Adversarial examples

- Built through an optimization problem.
- Attacks help build defenses/robust models.
- Classical tools are gradient-based.

Adversarial examples

- Built through an optimization problem.
- Attacks help build defenses/robust models.
- Classical tools are gradient-based.

What if the model is not available?

- 1 Adversarial attacks
- 2 Black-box adversarial attacks
- 3 Adversarial attacks via direct search

Recall: Adversarial attack problem

$$\underset{z \in \mathcal{B}(\mathbf{x}, \epsilon)}{\text{maximize}} \mathcal{L}(h(\mathbf{z}), \mathbf{y})$$

Recall: Adversarial attack problem

$$\underset{z \in \mathcal{B}(x, \epsilon)}{\text{maximize}} \mathcal{L}(h(z), y)$$

Black-box setup

- h model (e.g. neural net) only available through queries.
- Different from previous *white-box* models.

Proxy-based attacks (Papernot et al. '17, Tu et al '18)

- Replace h with a proxy network \tilde{h} .
 - Build white-box attacks for such a network!
- Efficient for computing many adversarial examples.

Proxy-based attacks (Papernot et al. '17, Tu et al '18)

- Replace h with a proxy network \tilde{h} .
- Build white-box attacks for such a network!

→ Efficient for computing many adversarial examples.

ZOO/Zeroth-Order Optimization (Chen et al '17)

- Use finite-difference schemes to approximate gradient ascent.
- Efficient for computing a small number of attacks.

Proxy-based attacks (Papernot et al. '17, Tu et al '18)

- Replace h with a proxy network \tilde{h} .
- Build white-box attacks for such a network!

→ Efficient for computing many adversarial examples.

ZOO/Zeroth-Order Optimization (Chen et al '17)

- Use finite-difference schemes to approximate gradient ascent.
- Efficient for computing a small number of attacks.

→ A derivative-free optimization technique!

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

Derivative-free algorithms

(Audet & Hare '25, Conn et al '09, Larson et al '21)

- **Derivatives unavailable for algorithmic use.**
- Function calls typically considered expensive.

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

Derivative-free algorithms

(Audet & Hare '25, Conn et al '09, Larson et al '21)

- **Derivatives unavailable for algorithmic use.**
- Function calls typically considered expensive.

Main algorithmic families

- **Zeroth-order** Estimate gradient via finite differences.

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

Derivative-free algorithms

(Audet & Hare '25, Conn et al '09, Larson et al '21)

- **Derivatives unavailable for algorithmic use.**
- Function calls typically considered expensive.

Main algorithmic families

- **Zeroth-order** Estimate gradient via finite differences.
- **Genetic algorithms** Population-based methods that sample.

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

Derivative-free algorithms

(Audet & Hare '25, Conn et al '09, Larson et al '21)

- **Derivatives unavailable for algorithmic use.**
- Function calls typically considered expensive.

Main algorithmic families

- **Zeroth-order** Estimate gradient via finite differences.
- **Genetic algorithms** Population-based methods that sample.
- **Direct search** Explore the space through directions.

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

Derivative-free algorithms

(Audet & Hare '25, Conn et al '09, Larson et al '21)

- **Derivatives unavailable for algorithmic use.**
- Function calls typically considered expensive.

Main algorithmic families

- **Zeroth-order** Estimate gradient via finite differences.
- **Genetic algorithms** Population-based methods that sample.
- **Direct search** Explore the space through directions.
- **Model based** Use surrogate for f to compute new points (similar to proxy but iterative).

Goal l_∞ attacks

$$\underset{\mathbf{z} \in \mathcal{B}_\infty(\mathbf{x}, \epsilon_\infty)}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{x}), y).$$

- $\mathcal{B}_\infty(\mathbf{x}, \epsilon_\infty) = \{\mathbf{z} \mid \max_i |z_i - x_i| \leq \epsilon\}$.
- h : ResNet50 with or without adversarial training (Engstrom et al '19).
- Dataset: CIFAR-10, $32 \times 32 \times 3$ images, 10 classes.
- 900 attacks (Assign all other labels to one of 100 images).

Goal l_∞ attacks

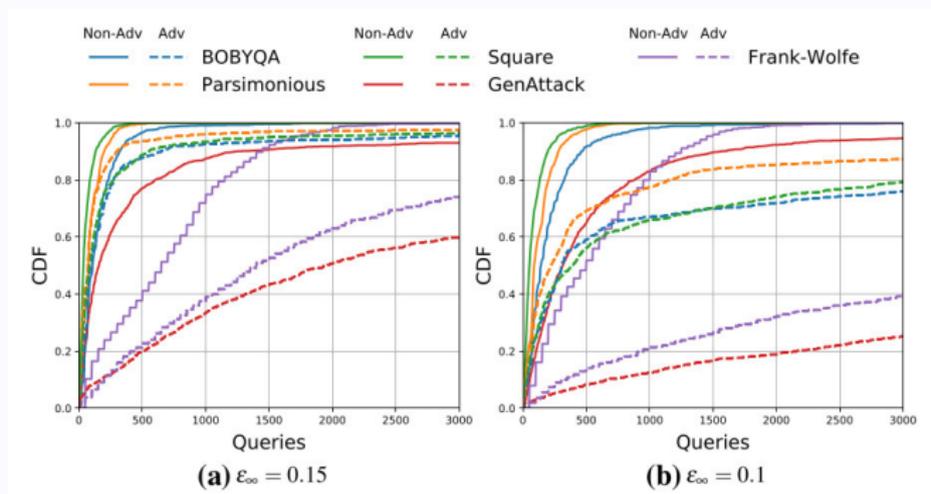
$$\underset{\mathbf{z} \in \mathcal{B}_\infty(\mathbf{x}, \epsilon_\infty)}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{x}), y).$$

- $\mathcal{B}_\infty(\mathbf{x}, \epsilon_\infty) = \{\mathbf{z} \mid \max_i |z_i - x_i| \leq \epsilon\}$.
- h : ResNet50 with or without adversarial training (Engstrom et al '19).
- Dataset: CIFAR-10, $32 \times 32 \times 3$ images, 10 classes.
- 900 attacks (Assign all other labels to one of 100 images).

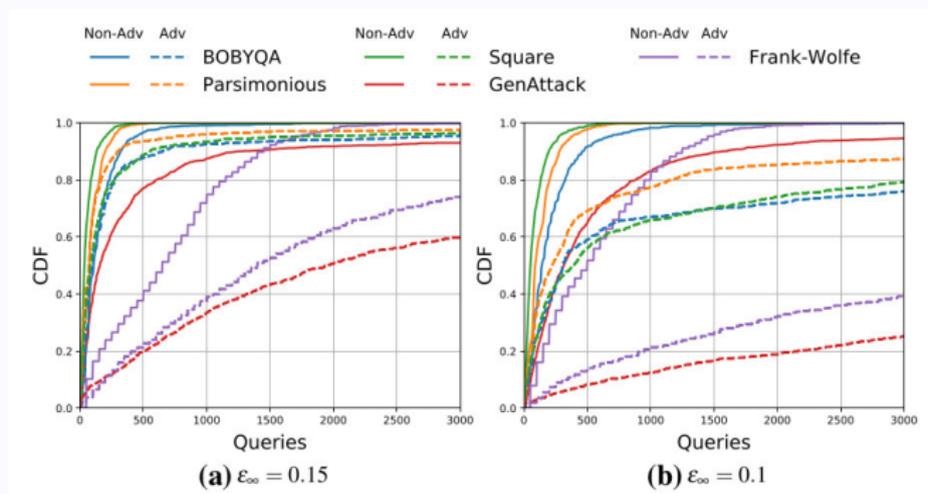
Compared methods

- *Frank-Wolfe* (Chen et al '20): Randomized zeroth-order method.
- *GenAttack* (Alzantot et al '19): Genetic algorithm.
- *Parsimonious* (Moon et al '19), *Square* (Andriushchenko et al '20): Direct search over combinatorial domain.
- BOBYQA (Powell '09) general-purpose DFO solver (in Fortran!)

A numerical comparison (from Ughi et al '22)



A numerical comparison (from Ughi et al '22)



Good performance for off-the-shelf BOBYQA?

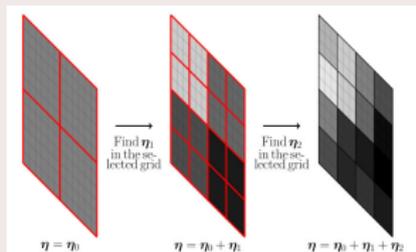
BOBYQA solver (Original, Powell '09)

- Typically used up to a few hundred variables.
- Requires to build n -dimensional models.
- Here $n = 3072!$

BOBYQA solver (Original, Powell '09)

- Typically used up to a few hundred variables.
- Requires to build n -dimensional models.
- Here $n = 3072!$

BOBYQA for adversarial (Ughi et al '19)



- **Random sub-sampling** of variables (also used in ZOO/Frank-Wolfe).
- **Hierarchical lifting**, built through randomness (also used in GenAttack).

Other derivative-free techniques

True label:	brambling	cannon	pug-dog	fly	armadillo	balloon
Perturbed image (ZO-PSGD):						
Prediction label:	goldfinch	plow	bucket	longicorn	croquet ball	parachute
ℓ_2 Distortion:	35.2137	87.7304	72.3397	111.068	172.719	35.9330
# of queries:	7640	150	130	50	210	5830
Perturbed image (ZO-SMD):						
Prediction label:	goldfinch	plow	bucket	cicada	croquet ball	parachute
ℓ_2 distortion:	8.9708	26.0126	20.9504	30.0968	45.097	10.9023
# of queries:	29980	350	260	140	570	15830
Perturbed image (ZO-AdaMM):						
Prediction label:	goldfinch	plow	bucket	cicada	croquet ball	parachute
ℓ_2 distortion:	8.0502	5.7359	4.5753	4.4456	6.3149	7.7405
# of queries:	53300	1710	790	540	2780	32900
Perturbed image (ZO-NES):						
Prediction label:	goldfinch	plow	bucket	cicada	croquet ball	parachute
ℓ_2 distortion:	54.9956	34.5852	28.7035	28.9158	40.1483	51.7116
# of queries:	22110	1080	830	430	2280	15340

Zeroth-order with ImageNet on InceptionV3 (Liu et al '20).

→ Mostly randomized variants with convergence rates.

Riemannian Zeroth-Order (Li et al '23)

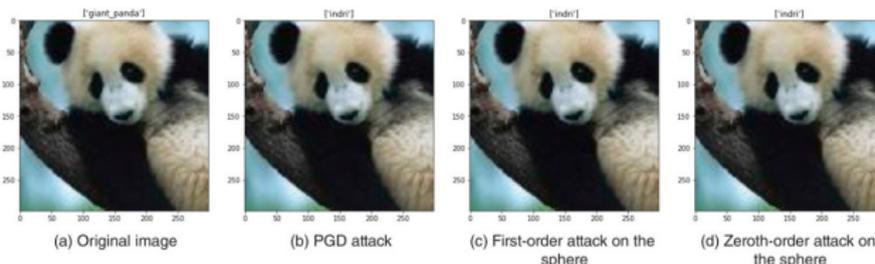
- Compute attacks on the sphere

$$\underset{\|z-x\|=\epsilon}{\text{maximize}} \mathcal{L}(h(z), y).$$

→ Riemannian structure!

- Randomized zeroth-order scheme

→ Aggregate finite-difference estimates.



Attack on ImageNet, $\|\cdot\| = \|\cdot\|_2$, $\epsilon = 0.05$.

Recap: Black-box adversarial attacks

- Query-based methods, akin to DFO.
- Numerous variants and benchmarking tasks.
- Mostly of zeroth-order type.

Recap: Black-box adversarial attacks

- Query-based methods, akin to DFO.
- Numerous variants and benchmarking tasks.
- Mostly of zeroth-order type.

Patterns

- Randomization plays a key role for scalability.
- Exploiting structure is also important.
- Challenging to apply DFO solvers other than zeroth-order.

Recap: Black-box adversarial attacks

- Query-based methods, akin to DFO.
- Numerous variants and benchmarking tasks.
- Mostly of zeroth-order type.

Patterns

- Randomization plays a key role for scalability.
- Exploiting structure is also important.
- Challenging to apply DFO solvers other than zeroth-order.

Our objectives

- Adapt general-purpose direct search to adversarial attacks.
- Develop structured attacks with randomness.

- 1 Adversarial attacks
- 2 Black-box adversarial attacks
- 3 Adversarial attacks via direct search

$$\underset{z \in \mathbb{R}^n}{\text{maximize}} \quad f(z)$$

Blackbox/Derivative-free optimization

- **Derivatives unavailable for algorithmic use.**
- Only access to values of f .
- **Cost** Number of **function evaluations**.

$$\underset{z \in \mathbb{R}^n}{\text{maximize}} \quad f(z)$$

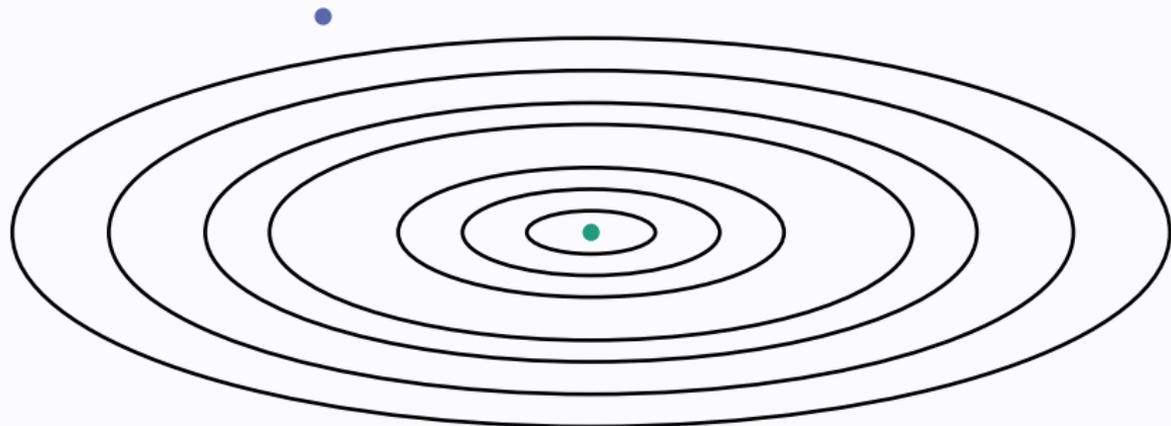
Blackbox/Derivative-free optimization

- **Derivatives unavailable for algorithmic use.**
- Only access to values of f .
- **Cost** Number of **function evaluations**.

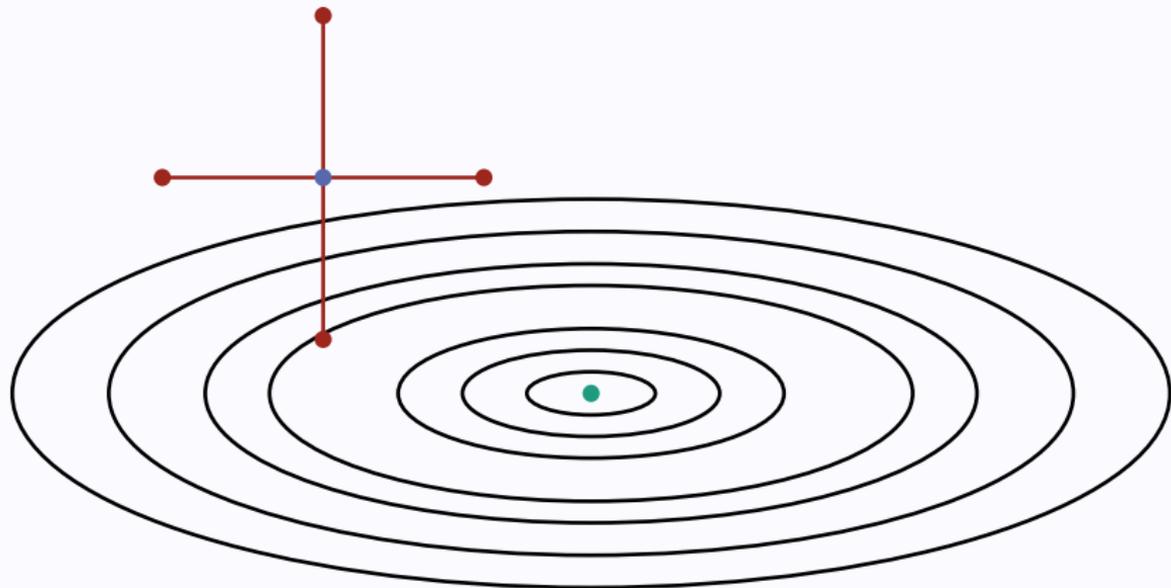
Direct search

- Proceed by polling the space along suitable directions.
- Rich convergence theory (Kolda et al '03, Dzhahini et al. '25).

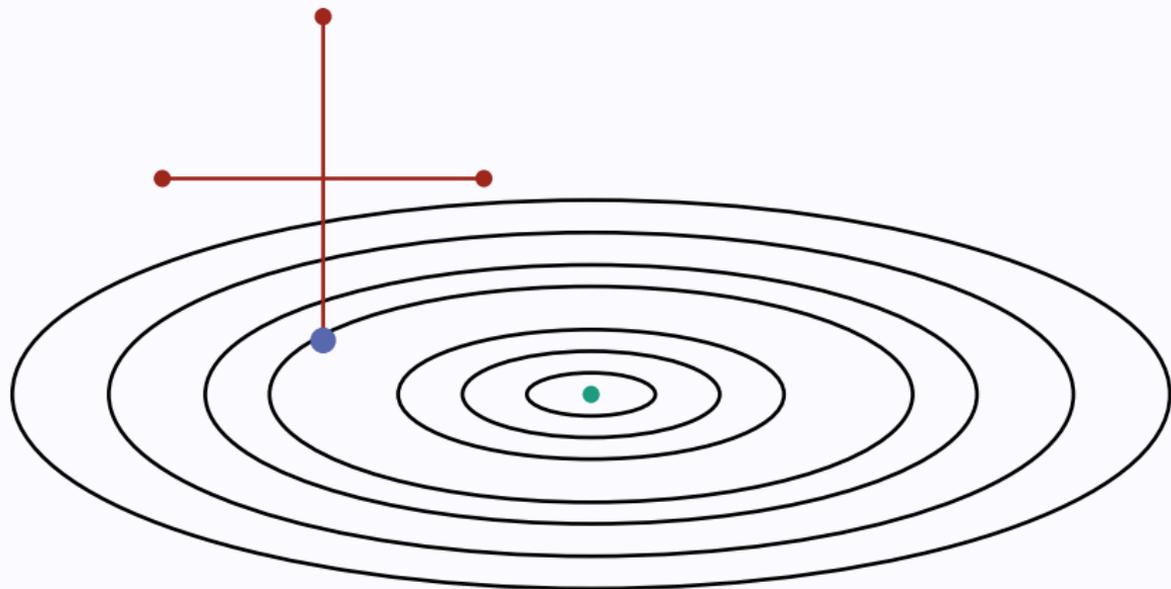
Example: Coordinate search



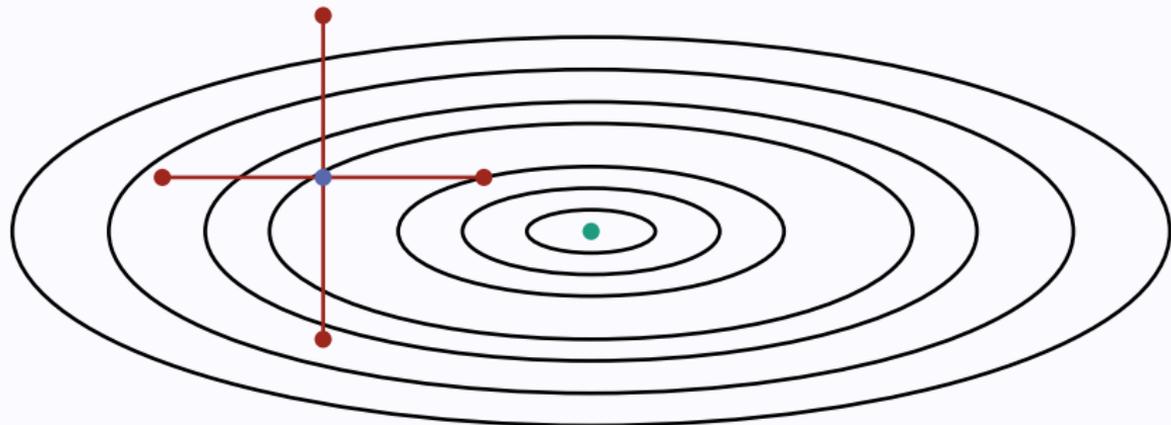
Example: Coordinate search



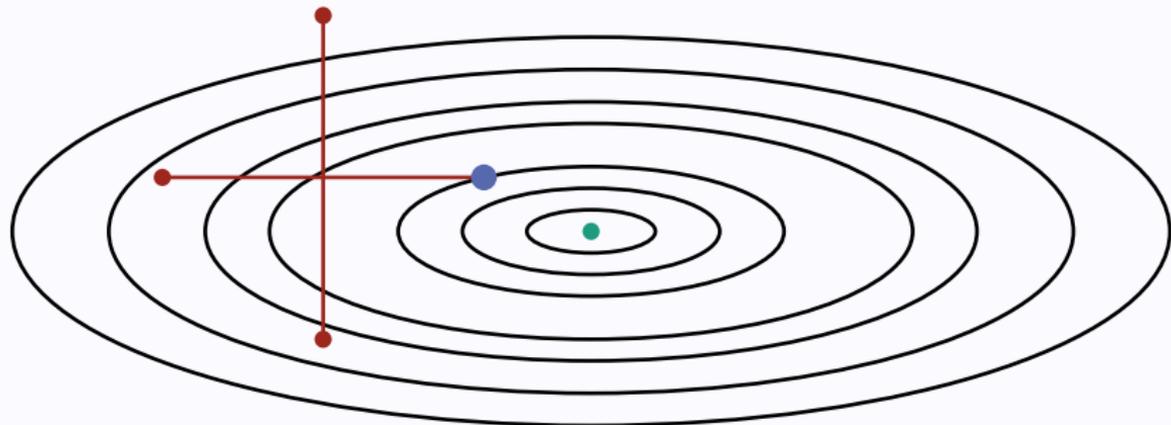
Example: Coordinate search



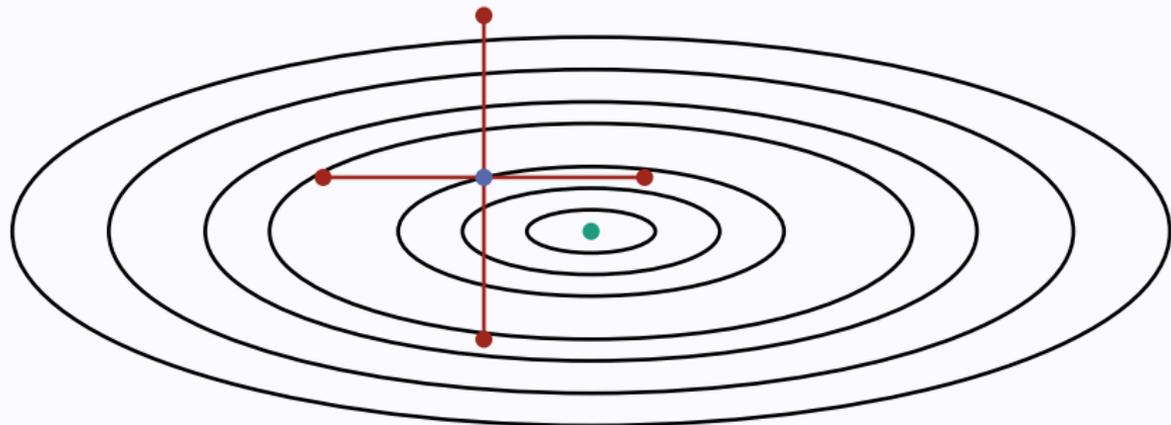
Example: Coordinate search



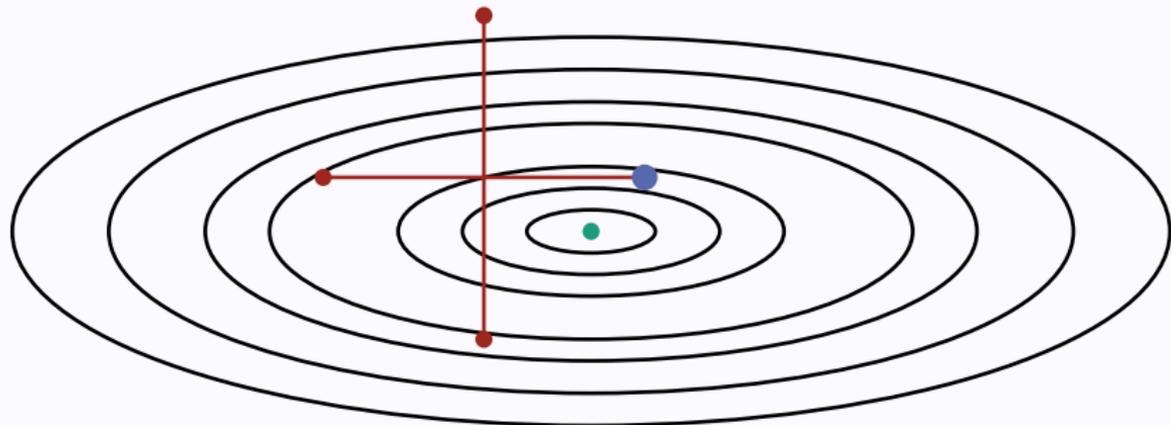
Example: Coordinate search



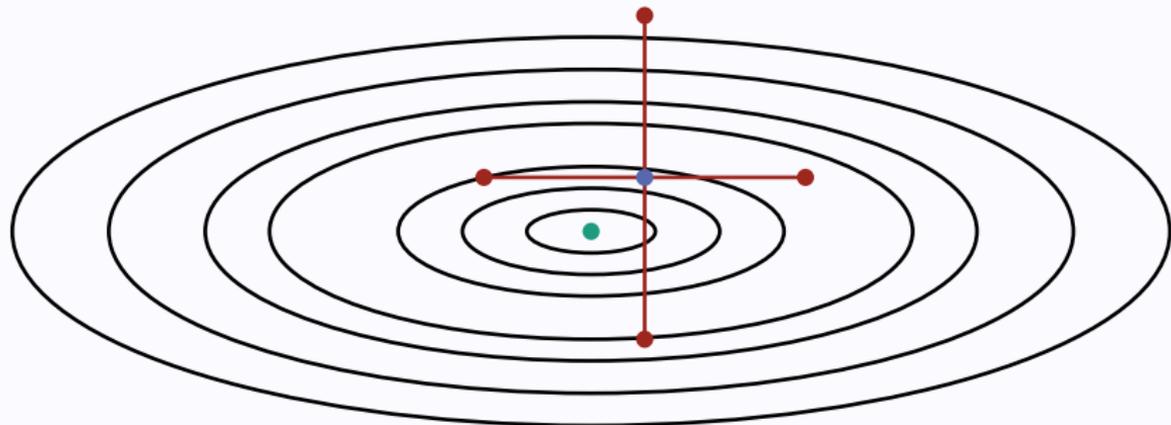
Example: Coordinate search



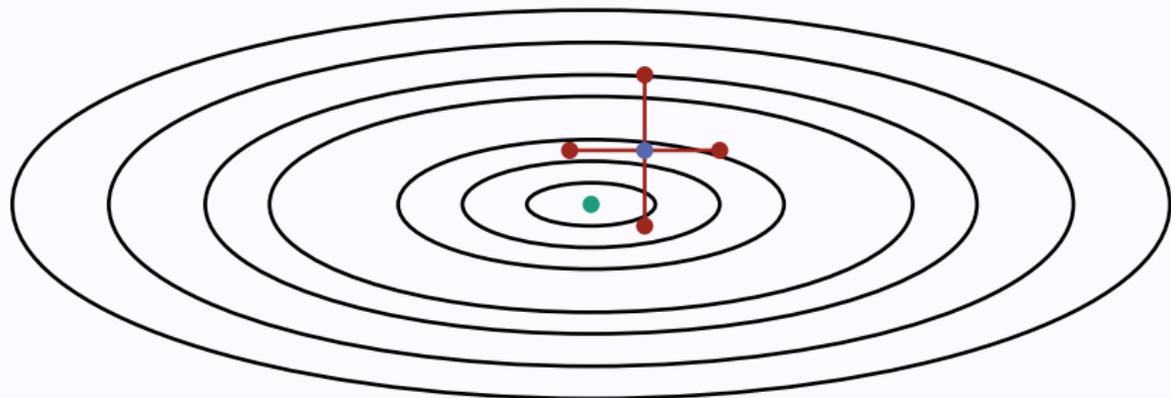
Example: Coordinate search



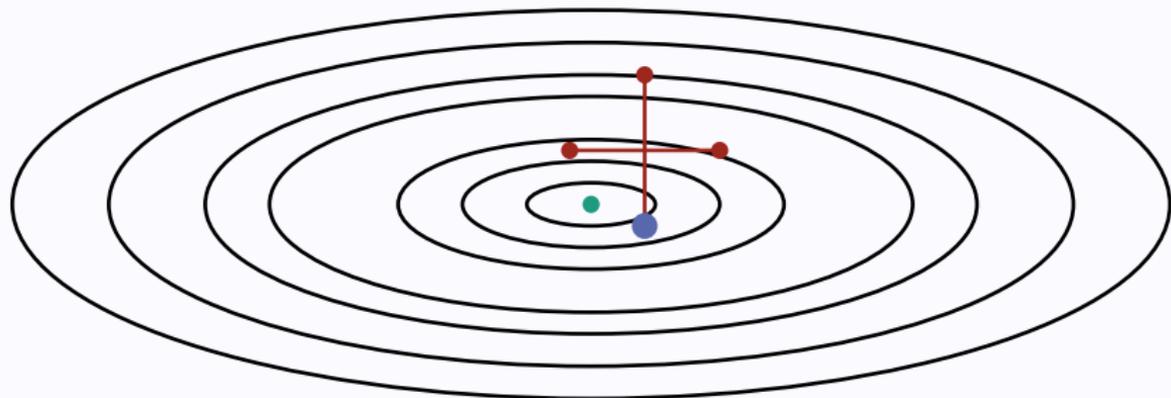
Example: Coordinate search



Example: Coordinate search



Example: Coordinate search



Basic direct-search framework

Problem maximize $_{\mathbf{z} \in \mathbb{R}^n}$ $f(\mathbf{z})$.

Inputs $\mathbf{z}_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

- Choose $\mathcal{D}_k \subset \mathbb{R}^n$.

Basic direct-search framework

Problem maximize $_{\mathbf{z} \in \mathbb{R}^n}$ $f(\mathbf{z})$.

Inputs $\mathbf{z}_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

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

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) > f(\mathbf{x}_k) + \eta \alpha_k^2 \|\mathbf{d}_k\|^2$$

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

Basic direct-search framework

Problem maximize $_{\mathbf{z} \in \mathbb{R}^n}$ $f(\mathbf{z})$.

Inputs $\mathbf{z}_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

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

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) > f(\mathbf{x}_k) + \eta \alpha_k^2 \|\mathbf{d}_k\|^2$$

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

- Otherwise, set $\mathbf{x}_{k+1} := \mathbf{x}_k$, $\alpha_{k+1} := 0.5\alpha_k$.

Basic direct-search framework

Problem maximize $\mathbf{z} \in \mathbb{R}^n$ $f(\mathbf{z})$.

Inputs $\mathbf{z}_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

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

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) > f(\mathbf{x}_k) + \eta \alpha_k^2 \|\mathbf{d}_k\|^2$$

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

- Otherwise, set $\mathbf{x}_{k+1} := \mathbf{x}_k$, $\alpha_{k+1} := 0.5\alpha_k$.

Positive Spanning Set (PSS)

$$\mathcal{D} \subset \mathbb{R}^n \text{ PSS} \iff \text{cm}(\mathcal{D}) = \max_{\mathbf{d} \in \mathcal{D}} \min_{\|\mathbf{v}\|=1} \frac{\mathbf{d}^T \mathbf{v}}{\|\mathbf{d}\|} > 0.$$

$\text{cm}(\mathcal{D})$: Cosine measure.

Goal $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$.

Theorem (Vicente '12)

If $\mathcal{D}_k = \mathcal{D}$ PSS $\forall k$, the method takes at most

$$\mathcal{O} (|\mathcal{D}| \text{cm}(\mathcal{D})^{-2} \epsilon_g^{-2})$$

calls to f .

Goal $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$.

Theorem (Vicente '12)

If $\mathcal{D}_k = \mathcal{D}$ PSS $\forall k$, the method takes at most

$$\mathcal{O}(|\mathcal{D}| \text{cm}(\mathcal{D})^{-2} \epsilon_g^{-2})$$

calls to f .

- $|\mathcal{D}|$ and $\text{cm}(\mathcal{D})$ depend on dimension n !
- Best known bound:

$$|\mathcal{D}| = \mathcal{O}(n), \text{cm}(\mathcal{D}) = \mathcal{O}(n^{-1/2}) \Rightarrow \mathcal{O}(n^2 \epsilon_g^{-2}).$$

Goal $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$.

Theorem (Vicente '12)

If $\mathcal{D}_k = \mathcal{D}$ PSS $\forall k$, the method takes at most

$$\mathcal{O}(|\mathcal{D}| \text{cm}(\mathcal{D})^{-2} \epsilon_g^{-2})$$

calls to f .

- $|\mathcal{D}|$ and $\text{cm}(\mathcal{D})$ depend on dimension n !
- Best known bound:

$$|\mathcal{D}| = \mathcal{O}(n), \text{cm}(\mathcal{D}) = \mathcal{O}(n^{-1/2}) \Rightarrow \mathcal{O}(n^2 \epsilon_g^{-2}).$$

Can we improve using randomness?

Our approach (w/ Lindon Roberts)

- 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.
- Much less than ambient dimension n !

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.

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).

Direct-search algorithm (in subspaces)

Inputs: $\mathbf{z}_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $\eta > 0$.

Iteration k : Given (\mathbf{z}_k, α_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

$$f(\mathbf{z}_k + \alpha_k \mathbf{P}_k^T \mathbf{d}_k) > f(\mathbf{z}_k) + \eta \alpha_k^2 \|\mathbf{P}_k^T \mathbf{d}_k\|^2,$$

set $\mathbf{z}_{k+1} := \mathbf{z}_k + \alpha_k \mathbf{P}_k^T \mathbf{d}_k$, $\alpha_{k+1} := 2\alpha_k$.

- Otherwise, set $\mathbf{z}_{k+1} := \mathbf{z}_k$, $\alpha_{k+1} := \alpha_k/2$.

Goal $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$.

Theorem (Roberts, R. '23)

Assume:

- $\{\mathbf{P}_k\}$ (p, σ, P_{\max}) -well aligned with $p > 1/2$.
- \mathcal{D}_k best PSS in r -dimensional subspace.

Let N_{ϵ_g} the number of function evaluations needed to have $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$. Then,

$$\mathbb{P}(N_{\epsilon} \leq \mathcal{O}(r^2 \phi \epsilon_g^{-2})) \geq 1 - \exp(-\mathcal{O}(\phi \epsilon_g^{-2})).$$

where $\phi = \sigma^{-2} P_{\max}^4$.

Goal $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$.

Theorem (Roberts, R. '23)

Assume:

- $\{\mathbf{P}_k\}$ (p, σ, P_{\max}) -well aligned with $p > 1/2$.
- \mathcal{D}_k best PSS in r -dimensional subspace.

Let N_{ϵ_g} the number of function evaluations needed to have $\|\nabla f(\mathbf{z}_K)\| \leq \epsilon_g$. Then,

$$\mathbb{P}(N_{\epsilon} \leq \mathcal{O}(r^2 \phi \epsilon_g^{-2})) \geq 1 - \exp(-\mathcal{O}(\phi \epsilon_g^{-2})).$$

where $\phi = \sigma^{-2} P_{\max}^4$.

How does $r^2 \phi$ depend on n ?

Setup

- Use random $\mathbf{P}_k \in \mathbb{R}^{r \times n}$.
- Complexity scales as $r^2 \sigma^{-2} P_{\max}^{-4}$.

Setup

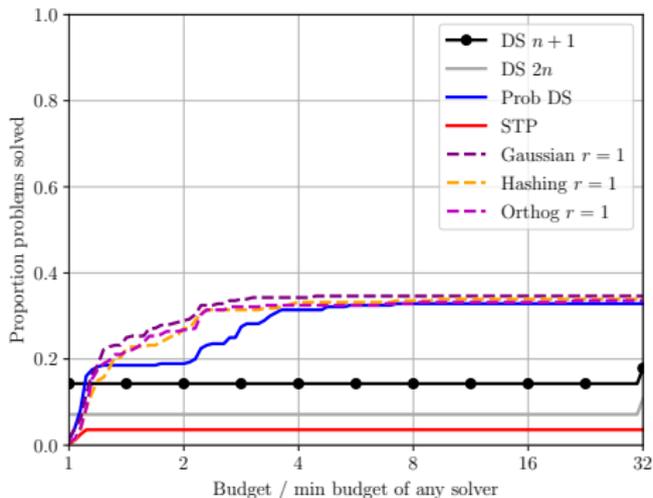
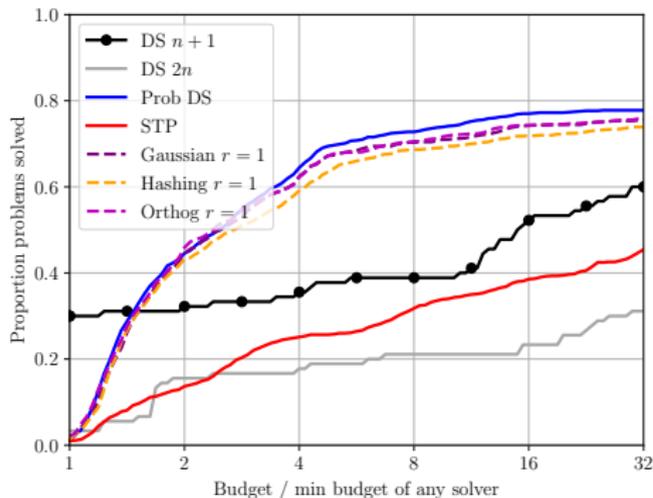
- Use random $\mathbf{P}_k \in \mathbb{R}^{r \times n}$.
- Complexity scales as $r^2 \sigma^{-2} P_{\max}^{-4}$.

\mathbf{P}_k	σ	P_{\max}
Gaussian	$\Theta(\sqrt{n/r})$	$\Theta(\sqrt{n/r})$
Hashing	$\Theta(\sqrt{n/r})$ (Dzahini & Wild '22b)	\sqrt{n}
Orthogonal	$\sqrt{n/r}$	$\sqrt{n/r}$.

Best bound

$$r = \mathcal{O}(1) \Rightarrow \mathcal{O}(n\epsilon_g^{-2}).$$

Numerical illustration



Left: Medium scale; Right: Large scale.

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

→ pip install directsearch

Direct search

- Simple to implement.
- Theoretical guarantees → Dependency on n .
- Competitor: Zeroth-order schemes!

Random subspaces

- Help scale things up (on standard DFO benchmarks).
- Still dependent on problem dimension n (not enough for adversarial attacks).

$$\underset{\mathbf{Z} \in \mathcal{M}}{\text{maximize}} f(\mathbf{Z}).$$

- f : Black-box objective.
- \mathbf{Z} : Matrix/Tensor object (e.g. $32 \times 32 \times 3$ images).
- \mathcal{M} : Known Riemannian manifold.

Motivating examples

- Spheres in Matrix/Tensor spaces: $\mathcal{M} = \{\mathbf{Z}, \|\mathbf{Z} - \mathbf{X}\| = \epsilon\}$.
- Symmetric Positive Definite Matrices (SPDNet).
- Low-rank matrix/tensors.

Problem: maximize $_{\mathbf{Z} \in \mathcal{M}}$ $f(\mathbf{Z})$.

Inputs $\mathbf{Z}_0 \in \mathcal{M}$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

- Choose a PSS $\mathcal{D}(\mathbf{Z}_k)$ for $\mathcal{T}_{\mathbf{Z}_k}\mathcal{M}$.
- If $\exists \mathbf{D}_k \in \mathcal{D}_k$ such that

$$f(\mathcal{R}_{\mathbf{Z}_k}(\alpha_k \mathbf{D}_k)) > f(\mathbf{Z}_k) + \eta \alpha_k^2 \|\mathbf{D}_k\|_{\mathbf{Z}_k}^2$$

set $\mathbf{Z}_{k+1} := \mathcal{R}_{\mathbf{Z}_k}(\alpha_k \mathbf{D}_k)$, $\alpha_{k+1} := 2\alpha_k$.

- Otherwise, set $\mathbf{Z}_{k+1} := \mathbf{Z}_k$, $\alpha_{k+1} := 0.5\alpha_k$.

Problem: maximize $_{\mathbf{Z} \in \mathcal{M}}$ $f(\mathbf{Z})$.

Inputs $\mathbf{Z}_0 \in \mathcal{M}$, $\alpha_0 > 0$, $\eta > 0$.

For $k = 0, 1, \dots$

- Choose a PSS $\mathcal{D}(\mathbf{Z}_k)$ for $\mathcal{T}_{\mathbf{Z}_k}\mathcal{M}$.
- If $\exists \mathbf{D}_k \in \mathcal{D}_k$ such that

$$f(\mathcal{R}_{\mathbf{Z}_k}(\alpha_k \mathbf{D}_k)) > f(\mathbf{Z}_k) + \eta \alpha_k^2 \|\mathbf{D}_k\|_{\mathbf{Z}_k}^2$$

set $\mathbf{Z}_{k+1} := \mathcal{R}_{\mathbf{Z}_k}(\alpha_k \mathbf{D}_k)$, $\alpha_{k+1} := 2\alpha_k$.

- Otherwise, set $\mathbf{Z}_{k+1} := \mathbf{Z}_k$, $\alpha_{k+1} := 0.5\alpha_k$.

- $\mathcal{T}_{\mathbf{Z}_k}\mathcal{M}$: Tangent space for \mathcal{M} at \mathbf{Z}_k with norm $\|\cdot\|_{\mathbf{Z}_k}$ \rightarrow Rigorous extension of PSS possible!
- $\mathcal{R}_{\mathbf{Z}_k}(\alpha_k \mathbf{D}_k)$: Retraction (brings back onto the manifold) \rightarrow Guarantees feasible iterates!

A theoretically sound approach

Goal $\|\text{grad}f(\mathbf{Z}_K)\| \leq \epsilon_g$ (Riemannian gradient).

Theorem (Cavarretta, Goyens, R., Yger '25)

The method takes at most

$$\mathcal{O}\left(\max_{\mathbf{Z} \in \mathcal{M}} |\mathcal{D}(\mathbf{Z})| \text{cm}(\mathcal{D}(\mathbf{Z}))^{-2} \epsilon_g^{-2}\right)$$

calls to f .

Goal $\|\text{grad}f(\mathbf{Z}_K)\| \leq \epsilon_g$ (Riemannian gradient).

Theorem (Cavarretta, Goyens, R., Yger '25)

The method takes at most

$$\mathcal{O}\left(\max_{\mathbf{Z} \in \mathcal{M}} |\mathcal{D}(\mathbf{Z})| \text{cm}(\mathcal{D}(\mathbf{Z}))^{-2} \epsilon_g^{-2}\right)$$

calls to f .

Extension: probabilistic variant

- Sample tangent subspaces (Gutman & Ho-Nguyen '22).
- Generalization of Euclidean case (Roberts, R. '23).

Test case: Low-rank attacks

Recall: Adversarial attack problem (matrix version)

$$\underset{\mathbf{Z} \in \mathcal{B}(\mathbf{X}, \epsilon)}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{Z}), y)$$

→ $\mathcal{B}(\cdot)$ ℓ_2 or nuclear norm ball.

Test case: Low-rank attacks

Recall: Adversarial attack problem (matrix version)

$$\underset{\mathbf{Z} \in \mathcal{B}(\mathbf{X}, \epsilon)}{\text{maximize}} \mathcal{L}(h(\mathbf{Z}), y)$$

→ $\mathcal{B}(\cdot)$ ℓ_2 or nuclear norm ball.

Low-rank PGD attacks (Savostianova et al '24)

- Observation: Projected Gradient Descent computes approximately low-rank attacks ($\mathbf{Z} - \mathbf{X}$ has many small singular values).

Test case: Low-rank attacks

Recall: Adversarial attack problem (matrix version)

$$\underset{\mathbf{Z} \in \mathcal{B}(\mathbf{X}, \epsilon)}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{Z}), y)$$

→ $\mathcal{B}(\cdot)$ ℓ_2 or nuclear norm ball.

Low-rank PGD attacks (Savostianova et al '24)

- Observation: Projected Gradient Descent computes approximately low-rank attacks ($\mathbf{Z} - \mathbf{X}$ has many small singular values).
- Approach: Enforce low-rank perturbation:

$$\underset{\mathbf{U}, \mathbf{V}}{\text{maximize}} \quad \mathcal{L}(h(\mathbf{X} + \mathbf{U} \otimes \mathbf{V}), y) \quad \text{s.t.} \quad \mathbf{X} + \mathbf{U} \otimes \mathbf{V} \in \mathcal{B}(\mathbf{X}, \epsilon).$$

and solve with PGD applied to \mathbf{U} and \mathbf{V} .

Test case: Low-rank attacks

Recall: Adversarial attack problem (matrix version)

$$\underset{\mathbf{Z} \in \mathcal{B}(\mathbf{X}, \epsilon)}{\text{maximize}} \mathcal{L}(h(\mathbf{Z}), y)$$

→ $\mathcal{B}(\cdot)$ ℓ_2 or nuclear norm ball.

Low-rank PGD attacks (Savostianova et al '24)

- Observation: Projected Gradient Descent computes approximately low-rank attacks ($\mathbf{Z} - \mathbf{X}$ has many small singular values).
- Approach: Enforce low-rank perturbation:

$$\underset{\mathbf{U}, \mathbf{V}}{\text{maximize}} \mathcal{L}(h(\mathbf{X} + \mathbf{U} \otimes \mathbf{V}), y) \quad \text{s.t.} \quad \mathbf{X} + \mathbf{U} \otimes \mathbf{V} \in \mathcal{B}(\mathbf{X}, \epsilon).$$

and solve with PGD applied to \mathbf{U} and \mathbf{V} .

What about low-rank black-box attacks?

Four methods (applied to factors \mathbf{U} , \mathbf{V})

- **DS** Deterministic, uses fixed directions.
→ $\mathcal{O}(\dim(\mathbf{U}) + \dim(\mathbf{V}))$ evaluations/iteration.
- **DS-PD(1)** Probabilistic, 1-dimensional subspace.
→ $\mathcal{O}(1)$ evaluations/iteration.
- **DS-PD(rank)** Probabilistic, r -dimensional subspace for rank- r attacks.
→ $\mathcal{O}(r)$ evaluations/iteration.
- **GFM** Gradient-Free Method, based on zeroth-order formula

$$\mathbf{W}_{k+1} = \mathbf{W}_k + \sigma_k \frac{f(\mathbf{W}_k + \mu_k \mathbf{D}_k) - f(\mathbf{W}_k - \mu_k \mathbf{D}_k)}{2\mu_k}$$

over 1-dimensional subspaces.
→ $\mathcal{O}(1)$ evaluations/iteration.

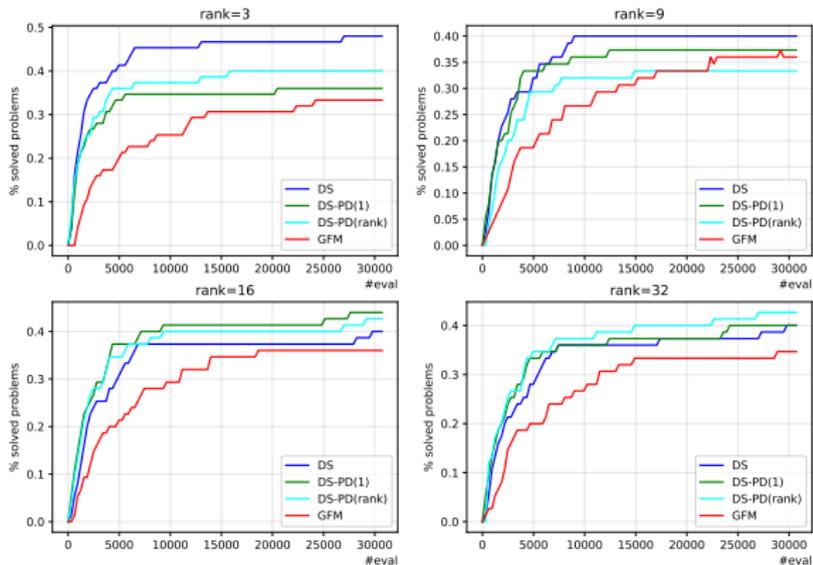
Model

- From RobustBench (Croce et al '21).
- WideResNet-28-10 trained on CIFAR-10.
- Problem dimension: $32 \times 32 \times 3 = 3072$.

Task

- Attack 75 images from CIFAR-10 (Target loss value).
- Budget: 10 PGD iterations with finite differences.
→ $32 \times 32 \times 3 \times 10 = 30720$ evaluations.

Black-box low-rank attacks: Results (fresh!)



Data profiles (top left is better)

- Small rank: Deterministic direct search is good.
- Increasing rank: Probabilistic variants for scalability!

Black-box adversarial attacks

- Hugely popular topic.
- Challenge: High dimensionality (but structure).

Black-box adversarial attacks

- Hugely popular topic.
- Challenge: High dimensionality (but structure).

Derivative-free optimization

- Classical techniques for low dimensions.
- Recent effort to scale up through randomization.

Black-box adversarial attacks

- Hugely popular topic.
- Challenge: High dimensionality (but structure).

Derivative-free optimization

- Classical techniques for low dimensions.
- Recent effort to scale up through randomization.

Our approach: Direct search

- Theory for Riemannian+Randomized versions.
- Promising results with low-rank perturbation of images.

- B. Cavarretta, F. Goyens, C. W. Royer and F. Yger, *Complexity guarantees and polling strategies for Riemannian direct-search methods*, arXiv:2511.15360, under review.
- K. J. Dzahini, F. Rinaldi, C. W. Royer and D. Zeffiro, *Direct search methods in the year 2025: Theoretical guarantees and algorithmic paradigms*, EURO Journal on Computational Optimization, 2025.
- L. Roberts and C. W. Royer, *Direct search based on probabilistic descent in reduced spaces*, SIAM Journal on Optimization, 2023.

→ Stay tuned for more!

- B. Cavarretta, F. Goyens, C. W. Royer and F. Yger, *Complexity guarantees and polling strategies for Riemannian direct-search methods*, arXiv:2511.15360, under review.
- K. J. Dzahini, F. Rinaldi, C. W. Royer and D. Zeffiro, *Direct search methods in the year 2025: Theoretical guarantees and algorithmic paradigms*, EURO Journal on Computational Optimization, 2025.
- L. Roberts and C. W. Royer, *Direct search based on probabilistic descent in reduced spaces*, SIAM Journal on Optimization, 2023.

→ Stay tuned for more!

Merci!

`clement.royer@lamsade.dauphine.fr`