

# Small errors in random zeroth-order optimization are imaginary

Sunday Oct. 15.

SC38. Optimization Society's Award Session II.

---

Based on: “*Small errors in random zeroth-order optimization are imaginary*”

arXiv: <https://arxiv.org/abs/2103.05478>.

by **Wouter Jongeneel** (EPFL), Man-Chung Yue (HKU) and Daniel Kuhn (EPFL).

This work was supported by the (SNSF) *NCCR*

*Automation*: <https://nccr-automation.ch>.

mail: [wouter.jongeneel@epfl.ch](mailto:wouter.jongeneel@epfl.ch),

web: [wjongeneel.nl](http://wjongeneel.nl).

# First-order optimization 101

For  $f: C^1(X; \mathbb{R}^n; \mathbb{R})$ , how to find

$$x \in \operatorname{argmin}_X f(x)?$$

---

<sup>1</sup>Popularity measure: in the last year (May 2022 - May 2023), searching for “SGD” was on average just a factor 1/25 as popular as searching for “Covid” (worldwide) <https://trends.google.com/trends/explore?q=SGD,Covid>.

<sup>2</sup>Still a very active research topic, see <https://www.quantamagazine.org/risky-giant-steps-can-solve-optimization-problems-faster-20230811/>.

<sup>3</sup>Nesterov 2003, § 2.1.5.

# First-order optimization 101

For  $f \in C^1(X \subseteq \mathbb{R}^n; \mathbb{R})$ , how to find

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

Common<sup>1</sup> approach: **gradient descent**

$$x_{k+1} = x_k - \mu_k \nabla f(x_k), \quad k = 1, 2, \dots \quad (1)$$

Let  $f$  be convex with a  $L$ -Lipschitz **gradient**, i.e.,

$$\|f(x) - f(y)\| \leq L \|x - y\|, \quad x, y \in X,$$

then, for<sup>2</sup>  $\mu_k = 1/L$  and  $x_1, x_2, \dots, x_K$  generated by (1) one obtains<sup>3</sup>

$$\|f(x_K) - f(x^*)\| = O\left(\frac{L \cdot \|x_1 - x^*\|^2}{K}\right).$$

---

<sup>1</sup>Popularity measure: in the last year (May 2022 - May 2023), searching for "SGD" was on average just a factor 1/25 as popular as searching for "Covid" (worldwide) <https://trends.google.com/trends/explore?q=SGD,Covid>.

<sup>2</sup>Still a very active research topic, see <https://www.quantamagazine.org/risky-giant-steps-can-solve-optimization-problems-faster-20230811/>.

<sup>3</sup>Nesterov 2003, § 2.1.5.

**If a gradient exists, does it mean we always *have* a gradient?**

## Example: DE constrained problems

Energy efficiency of transportation systems becomes increasingly important; must be optimized<sup>4</sup>. Good news: regularity is understood/studied.

Let  $f(x)$  represent *aerodynamic performance* for  $x$  a set of *design parameters*, do we have an expression for  $f(x)$ ?

---

<sup>4</sup>Images from: <https://predatorcycling.com/>, <https://www.3ds.com/> and [https://www.youtube.com/watch?v=FGmYpo-gkpU&ab\\_channel=EdwinLinders](https://www.youtube.com/watch?v=FGmYpo-gkpU&ab_channel=EdwinLinders).

## Example: DE constrained problems

Energy efficiency of transportation systems becomes increasingly important; must be optimized<sup>4</sup>. Good news: regularity is understood/studied.

Let  $f(x)$  represent *aerodynamic performance* for  $x$  a set of *design parameters*, do we have an expression for  $f(x)$ ?

Idea: we *can evaluate*  $f(x)$  for some design choice  $x$ , *i.e.*, by simulation, and subsequently use  $x_1, x_2, \dots, f(x_1), f(x_2), \dots$ , (**zeroth-order information**) for optimization.

---

<sup>4</sup>Images from: <https://predatorcycling.com/>, <https://www.3ds.com/> and [https://www.youtube.com/watch?v=FGmYpo-gkpU&ab\\_channel=EdwinLinders](https://www.youtube.com/watch?v=FGmYpo-gkpU&ab_channel=EdwinLinders).

# Zeroth-order optimization

Obtain (approximate)

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

via function evaluations  $f(x_1), f(x_2), \dots, f(x_K)$  for some set of *selected* points  $x_1, x_2, \dots, x_K$ . (For simplicity, we omit noise for now.)

---

<sup>5</sup>For references, consult the recent survey articles: Larson, Menickelly, and Wild 2019; Liu et al. 2020.

<sup>6</sup>See the books by Conn, Scheinberg, and Vicente 2009 and Audet and Hare 2017.

<sup>7</sup>Kiefer and Wolfowitz 1952; Nemirovsky and Yudin 1983; Flaxman, Kalai, and McMahan 2004; Spall 2005; Nesterov and Spokoiny 2017.

# Zeroth-order optimization

Obtain (approximate)

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

via function evaluations  $f(x_1), f(x_2), \dots, f(x_K)$  for some set of *selected* points  $x_1, x_2, \dots, x_K$ . (For simplicity, we omit noise for now.)

Two common paths<sup>5</sup>:

- (i) *Approximate a model*: construct a local model of  $f$ , optimize using that model, e.g., using a trust region method<sup>6</sup>.
- (ii) *Approximate an algorithm*: e.g., approximate  $f$  directly and apply some form of gradient descent<sup>7</sup>.

---

<sup>5</sup>For references, consult the recent survey articles: Larson, Menickelly, and Wild 2019; Liu et al. 2020.

<sup>6</sup>See the books by Conn, Scheinberg, and Vicente 2009 and Audet and Hare 2017.

<sup>7</sup>Kiefer and Wolfowitz 1952; Nemirovsky and Yudin 1983; Flaxman, Kalai, and McMahan 2004; Spall 2005; Nesterov and Spokoyny 2017.



## A gradient-based approach

For any smooth  $f : \mathbb{R} \rightarrow \mathbb{R}$

$$f(x + \delta) = f(x) + \delta f'(x) + O(\delta^2).$$

---

<sup>8</sup>d'Aspremont 2008; Devolder, Glineur, and Nesterov 2014.

## A gradient-based approach

For any smooth  $f : \mathbb{R} \rightarrow \mathbb{R}$

$$f(x + \delta) = f(x) + \delta f'(x) + O(\delta^2).$$

Then, run *inexact* ( $\epsilon > 0$  fixed) gradient descent

$$x_{k+1} = x_k - \mu_k \frac{f(x_k + \delta) - f(x_k)}{\delta}.$$

When does  $f(x_k) \leq f(x^*) + \epsilon$ ?

---

<sup>8</sup>d'Aspremont 2008; Devolder, Glineur, and Nesterov 2014.

## A gradient-based approach

For any smooth  $f : \mathbb{R} \rightarrow \mathbb{R}$

$$\Delta f(x) = \frac{f(x + \delta) - f(x)}{\delta} + O(\delta):$$

Then, run inexact ( $\delta > 0$  fixed) gradient descent

$$x_{k+1} = x_k - \delta \frac{f(x_k + \delta) - f(x_k)}{\delta}.$$

When does  $\|x_k - x^*\| \leq \epsilon$ ?

For fixed  $\delta > 0$ , a bias prevails,  $\|x_k - x^*\| \leq \epsilon + O(\delta)$ , e.g., for  $f(x) = x^2$  we effectively compute the gradient of  $f(x) + \delta x$ , shifting  $x^* = 0$  to  $-\frac{\delta}{2}$ .

Similarly, for  $f \in C^1(\mathbb{R}^n; \mathbb{R})$ , one should not naively use

$$\sum_{i=1}^n \frac{f(x + e_i) - f(x)}{e_i} e_i \quad \text{for } (e_1; e_2; \dots; e_n) = I_n:$$

---

<sup>8</sup>d'Aspremont 2008; Devolder, Glineur, and Nesterov 2014.

## A gradient-based approach cont.

(i) For appropriate (adaptive)  $\epsilon > 0$ , apply line-search<sup>9</sup> using

$$\chi^n \frac{f(x + b_i) - f(x)}{b_i} \approx \nabla f(x); \quad \text{for } \det(b_1; b_2; \dots; b_n) \neq 0;$$

i=1

---

<sup>9</sup>Berahas, Cao, and Scheinberg 2021.

<sup>10</sup>Randomization can be optimal Duchi et al. 2015, but no uniformly superior method exists yet randomized Monte Carlo schemes can be implemented to be  $n$  times cheaper [than deterministic Monte Carlo]; but an algorithm based on them has to take at least  $n$  times more steps. Scheinberg 2022, see also Berahas et al. 2022.

## A gradient-based approach cont.

(i) For appropriate (adaptive)  $\epsilon > 0$ , apply line-search<sup>9</sup> using

$$x^{n+1} = x^n + \frac{f(x^n + b_i) - f(x^n)}{\langle \nabla f(x^n), b_i \rangle} b_i; \quad \text{for } \det(b_1; b_2; \dots; b_n) \neq 0:$$

(ii) Suppose we find a random variable  $\epsilon \in \mathbb{R}^n$  such that

$$\mathbb{E} \left[ \frac{f(x + \epsilon) - f(x)}{\langle \nabla f(x), \epsilon \rangle} \right] = \nabla f(x):$$

Consider the randomized algorithm

$$x_{k+1} = x_k + \epsilon_k \frac{f(x_k + \epsilon_k) - f(x_k)}{\langle \nabla f(x_k), \epsilon_k \rangle}; \quad \epsilon_k \in \mathbb{R}^n:$$

<sup>9</sup>Berahas, Cao, and Scheinberg 2021.

<sup>10</sup>Randomization can be optimal Duchi et al. 2015, but no uniformly superior method exists yet randomized finite difference schemes can be implemented to be  $n$  times cheaper [than deterministic finite difference]; but an algorithm based on them has to take at least  $n$  times more steps. Scheinberg 2022, see also Berahas et al. 2022.

## A gradient-based approach cont.

(i) For appropriate (adaptive)  $\alpha > 0$ , apply line-search<sup>9</sup> using

$$x_{k+1} = x_k + \alpha \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}; \quad \text{for } \det(b_1; b_2; \dots; b_n) \neq 0:$$

(ii) Suppose we find a random variable  $h \in \mathbb{R}^n$  such that

$$\mathbb{E} \left[ \frac{h^T \nabla f(x_k)}{\|h\|} \right] = \|\nabla f(x_k)\|:$$

Consider the randomized algorithm

$$x_{k+1} = x_k + \alpha \frac{h^T \nabla f(x_k)}{\|h\|}; \quad h \sim \mathcal{D}$$

(!) Active topic of research<sup>10</sup>.

<sup>9</sup>Berahas, Cao, and Scheinberg 2021.

<sup>10</sup>Randomization can be optimal Duchi et al. 2015, but no uniformly superior method exists yet randomized nite di erence schemes can be implemented to be  $n$  times cheaper [than deterministic nite di erence]; but an algorithm based on them has to take at least  $n$  times more steps. Scheinberg 2022, see also Berahas et al. 2022.

## Highly-in uential exercise by Nemirovski and Yudin

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , Nemirovski and Yudin<sup>11</sup> consider: -smoothing

$$f(x) = \mathbb{E}_y \int_{B^n} [f(x+y)] = \text{vol}(B^n)^{-1} \int_{B^n} f(x+y) dy; \quad (2a)$$

$$r f(x) = \frac{n}{n-1} \mathbb{E}_y \int_{S^{n-1}} [f(x+y)] y = \frac{n}{n-1} \int_{S^{n-1}} f(x+y) y(dy); \quad (2b)$$

---

<sup>11</sup>Nemirovsky and Yudin 1983, credits usually given to Flaxman, Kalai, and McMahan 2004.

<sup>12</sup>Agarwal, Dekel, and Xiao 2010; Nesterov 2011.

# Highly-in uential exercise by Nemirovski and Yudin

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , Nemirovski and Yudin<sup>11</sup> consider:  $\epsilon$ -smoothing

$$f_\epsilon(x) = \mathbb{E}_y \int_{B^n} [f(x+y)] = \frac{1}{\text{vol}(B^n)} \int_{B^n} f(x+y) dy; \quad (2a)$$

$$r f_\epsilon(x) = \frac{1}{\text{vol}(S^{n-1})} \mathbb{E}_y \int_{S^{n-1}} [f(x+y)] y = \frac{1}{\text{vol}(S^{n-1})} \int_{S^{n-1}} f(x+y) y (dy); \quad (2b)$$

Natural single-point candidate to approximate  $r f_\epsilon$ :

$$g_\epsilon(x) = \frac{1}{\text{vol}(S^{n-1})} \int_{S^{n-1}} f(x+y) y; \quad (3a)$$

---

<sup>11</sup>Nemirovsky and Yudin 1983, credits usually given to Flaxman, Kalai, and McMahan 2004.

<sup>12</sup>Agarwal, Dekel, and Xiao 2010; Nesterov 2011.



## Highly-in uential exercise by Nemirovski and Yudin

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , Nemirovski and Yudin<sup>11</sup> consider:  $\epsilon$ -smoothing

$$f_\epsilon(x) = \mathbb{E}_y \int_{B^n} [f(x+y)] = \frac{1}{\text{vol}(B^n)} \int_{B^n} f(x+y) dy; \quad (2a)$$

$$r f_\epsilon(x) = \frac{1}{\text{vol}(S^{n-1})} \mathbb{E}_y \int_{S^{n-1}} [f(x+y)] y = \frac{1}{\text{vol}(S^{n-1})} \int_{S^{n-1}} f(x+y) y (dy); \quad (2b)$$

Natural single-point candidate to approximate  $r f$ :

$$g(x) = \frac{1}{\text{vol}(S^{n-1})} f(x+y) y; \quad y \in S^{n-1}; \quad (3a)$$

Observation<sup>12</sup>: avoid high-variance for  $\epsilon \neq 0$  and give (3a) again the interpretation of a directional derivative and use a multi-point oracle like:

$$g^0(x) = \frac{1}{\text{vol}(S^{n-1})} (f(x+y) - f(x)) y; \quad y \in S^{n-1}; \quad (3b)$$

<sup>11</sup>Nemirovsky and Yudin 1983, credits usually given to Flaxman, Kalai, and McMahan 2004.

<sup>12</sup>Agarwal, Dekel, and Xiao 2010; Nesterov 2011.

## Early algorithmic analysis by Nesterov and Spokoiny

For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  (locally convex), Gaussiansmoothing<sup>13</sup>

$$f_r(x) = \frac{1}{Z} \int_{\mathbb{R}^n} f(x + y) e^{-\frac{1}{2}kyk_2^2} dy \quad (4a)$$

$$r f(x) = \frac{1}{Z} \int_{\mathbb{R}^n} \frac{f(x + y) - f(x - y)}{2} e^{-\frac{1}{2}kyk_2^2} dy \quad (4b)$$

with  $kr f_r f k = O(n^{-2})$ .

---

<sup>13</sup>Nesterov 2011; Nesterov and Spokoiny 2017.

## Early algorithmic analysis by Nesterov and Spokoiny

For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  (locally convex), Gaussiansmoothing<sup>13</sup>

$$f(x) = \frac{1}{Z} \int_{\mathbb{R}^n} f(x+y) e^{-\frac{1}{2}kyk_2^2} dy \quad (4a)$$

$$r f(x) = \frac{1}{Z} \int_{\mathbb{R}^n} \frac{f(x+y) - f(x-y)}{2} e^{-\frac{1}{2}kyk_2^2} dy \quad (4b)$$

with  $kr f r f k = O(n^{-2})$ .

$$\text{Oracle (cd): } g(x) = \frac{f(x+y) - f(x-y)}{2} y; \quad y \sim N(0; I_n)$$

with  $E_{u \sim N(0; I_n)} kg(x)k_2^2 = O(n^{-2} + nkr f(x)k_2^2)$ .

Algorithm:  $x_{k+1} = x_k - \eta g_k(x_k); \quad \eta = O(1/nL)$

Performance: for  $k \leq K$  and  $x_K := x_1 - \sum_{k=1}^K \eta g_k$

$$E[f(x_K)] - f(x^*) = O\left(\frac{nL \eta^2 \sum_{k=1}^K \|x_k - x^*\|_2^2}{K}\right) = O(n) \quad \text{gradient descent}$$

<sup>13</sup>Nesterov 2011; Nesterov and Spokoiny 2017.

## Numerical considerations

Analysis continued after 2011-2017, still, all common<sup>14</sup> oracles of the form

$$\text{(finite difference): } \frac{f(x+y) - f(x)}{y} = \nabla f(x) + O(y)$$

$$\text{(central difference): } \frac{f(x+y) - f(x-y)}{2y} = \nabla f(x) + O(y^2);$$

$$\vdots = \nabla f(x) + O(y^{p-1})$$

As such, many algorithms require  $\epsilon_k = O(\epsilon^{1/k^q})$ , with  $q > 0$  for  $k = 1; 2; \dots$ .

---

<sup>14</sup>Hazan and Levy 2014; Duchi et al. 2015; Nesterov and Spokoiny 2017; Gasnikov et al. 2017; Shamir 2017; Akhavan, Pontil, and Tsybakov 2020; Lam, Li, and Zhang 2021; Novitskii and Gasnikov 2021.

<sup>15</sup>Generally,  $M = 2^{52} \cdot 10^{16}$ .

## Numerical considerations

Analysis continued after 2011-2017, still, all common<sup>14</sup> oracles of the form

$$\text{(finite difference): } \frac{f(x+y) - f(x)}{y} = \nabla f(x) + O(y)$$

$$\text{(central difference): } \frac{f(x+y) - f(x-y)}{2y} = \nabla f(x) + O(y^2);$$

$$\vdots = \nabla f(x) + O(y^{p-1})$$

As such, many algorithms require  $\epsilon_k = O(1/k^q)$ , with  $q > 0$  for  $k = 1; 2; \dots$

However, can we practically select  $k \rightarrow 0$  for  $k \rightarrow +1$ ?

For sufficiently small  $\epsilon$ ,  $f(x+y) - f(x)$  machine precision<sup>15</sup>

=) cancellation error, i.e., oracle output is nonsense.

Not that frequently discussed, does it matter?

---

<sup>14</sup>Hazan and Levy 2014; Duchi et al. 2015; Nesterov and Spokoiny 2017; Gasnikov et al. 2017; Shamir 2017; Akhavan, Pontil, and Tsybakov 2020; Lam, Li, and Zhang 2021; Novitskii and Gasnikov 2021.

<sup>15</sup>Generally,  $M = 2^{52} \approx 10^{16}$ .

## Intermezzo: a beautiful insight from complex analysis

As pioneered in the 60s<sup>16</sup>, let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be real analytic ( $C^\infty$ ) and consider

$$f(x + i\epsilon) = f(x) + \epsilon f'(x) + \frac{1}{2} \epsilon^2 f''(x) + \frac{1}{6} \epsilon^3 f'''(x) + O(\epsilon^4); \quad \epsilon^2 = 1:$$

such that (for  $z \in \mathbb{C}$ ,  $z = \langle z \rangle + i \Im(z)$ ):

$$\Im(f(x + i\epsilon)) = \epsilon f'(x) + \frac{1}{6} \epsilon^3 f'''(x) + O(\epsilon^5)$$

---

<sup>16</sup>Lyness and Moler 1967; Squire and Trapp 1998; Martins, Sturdza, and Alonso 2003; Abreu et al. 2018.

<sup>17</sup>A value of  $\epsilon = 10^{-100}$  (!) is successfully used in National Physical Laboratory software Cox and Harris 2004, Page 44.

## Intermezzo: a beautiful insight from complex analysis

As pioneered in the 60s<sup>16</sup>, let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be real analytic ( $C^\infty$ ) and consider

$$f(x + i) = f(x) + if'(x) - \frac{1}{2}f''(x) - \frac{1}{6}f'''(x)i^3 + O(|i|^4); \quad |i|^2 = 1:$$

such that (for  $z \in \mathbb{C}$ ,  $z = \text{Re}(z) + i\text{Im}(z)$ ):

$$f(x + i) = f(x) + if'(x) - \frac{1}{2}f''(x) + O(|i|^5)$$

and thus

$$if'(x) = \frac{f(x + i) - f(x)}{i} + O(|i|^2); \quad f'(x) = \text{Re}(f(x + i)) + O(|i|^2):$$

Hence, consider using

$$\frac{f(x + i) - f(x)}{i} = f'(x):$$

Cancellation errors are impossible<sup>17</sup>. Again, does it matter?

<sup>16</sup>Lyness and Moler 1967; Squire and Trapp 1998; Martins, Sturdza, and Alonso 2003; Abreu et al. 2018.

<sup>17</sup>A value of  $\epsilon = 10^{-100}$  (!) is successfully used in National Physical Laboratory software Cox and Harris 2004, Page 44.

## Numerical considerations cont.: an example

For  $f(x) = x^3$ , approximate  $f'(x)$  at  $x = 2$  using  $h = 1; 0.1; 0.01$  using

$$\text{(forward difference): } f'_{fd}(x; h) = \frac{f(x+h) - f(x)}{h}; \quad (5a)$$

$$\text{(central difference): } f'_{cd}(x; h) = \frac{f(x+h) - f(x-h)}{2h}; \quad (5b)$$

$$\text{(complex-step): } f'_{cs}(x; h) = \frac{f(x + ih) - f(x)}{ih} \quad (5c)$$

and compare the error for  $h \neq 0$ :



## Numerical considerations cont.: an example

For  $f(x) = x^3$ , approximate  $f'(x)$  at  $x = 2$  using  $h = 1; 0.1; 0.01$  using

$$\text{(forward difference): } f'_{fd}(x; h) = \frac{f(x+h) - f(x)}{h}; \quad (5a)$$

$$\text{(central difference): } f'_{cd}(x; h) = \frac{f(x+h) - f(x-h)}{2h}; \quad (5b)$$

$$\text{(complex-step): } f'_{cs}(x; h) = \frac{f(x+ih) - f(x)}{ih} \quad (5c)$$

and compare the error for  $h \neq 0$ :

$$(a) \quad x = 2$$

$$(b) \quad x = 0$$

$$(c) \quad x = 10$$

Failures well before  $h = 10^{-16}$ , so, it does matter.

## On the necessity of leaving $\mathbb{R}$

Although single-point estimators exist<sup>18</sup>, variance blows up for  $\epsilon \neq 0$ . Is this complex-lifting business really needed? Is there not a real analogue of

$$\text{Re} f(x) = \frac{f(x + i\epsilon) + f(x - i\epsilon)}{2} + O(\epsilon^2)? \quad (6)$$

---

<sup>18</sup>Flaxman, Kalai, and McMahan 2004.

<sup>19</sup>Jongeneel 2021.

## On the necessity of leaving $\mathbb{R}$

Although single-point estimators exist<sup>18</sup>, variance blows up for  $\epsilon \neq 0$ . Is this complex-lifting business really needed? Is there not a real analogue of

$$\mathbb{E} f(x) = \frac{f(x + i\epsilon)}{\epsilon} + O(\epsilon^2)? \quad (6)$$

Partial answer<sup>19</sup>: no.

Consider some non-empty open, convex set  $D \subset \mathbb{R}^n$  then, there does not exist a continuous map  $G : \mathbb{R} \rightarrow \mathbb{R}$  such that for all real-analytic functions  $f : D \rightarrow \mathbb{R}$

$$G(f(x + iy)) = \int_0^1 f(x + iy) dy + o(\epsilon) \quad \forall x \in D; \epsilon > 0; y \in S^{n-1} \quad (7)$$

not surprising, but provides motivation.

---

<sup>18</sup>Flaxman, Kalai, and McMahan 2004.

<sup>19</sup>Jongeneel 2021.

# Comment on Algorithmic Differentiation (AD)

Why bother with approximations?

---

<sup>20</sup>The Deep Learning Toolbox in MATLAB and AD tools in Julia (See Bezanson et al. 2017; Revels, Lubin, and Papamarkou 2016; Innes 2018; Moses and Churavy 2020) e.g., ForwardDiff.jl, Zygote.jl and Enzyme.jl or in Python, e.g., JAX Bradbury et al. 2018

## Comment on Algorithmic Differentiation (AD)

Why bother with approximations?

Dual numbers:  $a + b\epsilon$  with  $a, b \in \mathbb{R}$  and  $\epsilon \neq 0$ , yet,  $\epsilon^2 = 0$ , i.e., elements of the quotient ring  $\mathbb{R}[\epsilon]/\langle \epsilon^2 \rangle$ , not a field  $\Rightarrow$ , e.g.,  $\epsilon^2 = 0$  and  $\frac{1}{\epsilon^2}$  not defined.

$\mathbb{C}$  is an algebraically closed field.

AD: for  $f : \mathbb{R} \rightarrow \mathbb{R}$  is sufficiently regular, e.g.,  $f \in C^1(\mathbb{R})$ , then,

$f(x + \epsilon) = f(x) + \epsilon f'(x)$ , i.e.,  $f(x + \epsilon)$  provides us with the pair  $(f(x); \epsilon f'(x))$ .

---

<sup>20</sup>The Deep Learning Toolbox in MATLAB and AD tools in Julia (See Bezanson et al. 2017; Revels, Lubin, and Papamarkou 2016; Innes 2018; Moses and Churavy 2020) e.g., ForwardDiff.jl, Zygote.jl and Enzyme.jl or in Python, e.g., JAX Bradbury et al. 2018

## Comment on Algorithmic Differentiation (AD)

Why bother with approximations?

Dual numbers:  $a + b\epsilon$  with  $a, b \in \mathbb{R}$  and  $\epsilon \neq 0$ , yet,  $\epsilon^2 = 0$ , i.e., elements of the quotient ring  $\mathbb{R}[\epsilon]/\langle \epsilon^2 \rangle$ , not a field  $\Rightarrow$ , e.g.,  $\epsilon^2 = 0$  and  $\frac{1}{\epsilon^2}$  not defined.

$\mathbb{C}$  is an algebraically closed field.

AD: for  $f : \mathbb{R} \rightarrow \mathbb{R}$  is sufficiently regular, e.g.,  $f \in C^1(\mathbb{R})$ , then,

$f(x + \epsilon) = f(x) + \epsilon f'(x)$ , i.e.,  $f(x + \epsilon)$  provides us with the pair  $(f(x); \epsilon f'(x))$ .

Consider  $\epsilon f'(x)|_{x=0}$  for the following  $C^1$  functions:

$$f(x) = x^2; \quad f(x) = \sin(x); \quad f(x) = e^{-\frac{1}{x^2}}$$

No free lunch: most popular AD tools<sup>20</sup> evaluate to NaN whereas the complex-step derivative correctly approximates  $\epsilon f'(x)|_{x=0} = 0$ .

Theoretical solution: Levi-Civita field  $\sum_{q \in \mathbb{Q}} a_q \epsilon^q$  with  $a_q \in \mathbb{R}$   $\forall q \in \mathbb{Q}$  (inf. dim).

<sup>20</sup>The Deep Learning Toolbox in MATLAB and AD tools in Julia (See Bezanson et al. 2017; Revels, Lubin, and Papamarkou 2016; Innes 2018; Moses and Churavy 2020) e.g., ForwardDiff.jl, Zygote.jl and Enzyme.jl or in Python, e.g., JAX Bradbury et al. 2018

## A solution: the complex-step oracle <sup>22</sup> (exercise)

Let  $f \in C^1(X \subseteq \mathbb{R}^n; \mathbb{R})$ , using Cauchy-Riemann/Stokes show that:

$$\begin{aligned} \text{Re } f'(x) &= \mathbb{E}_y \mathbb{E}_{B^n} [\langle f'(x + iy) \rangle] \\ \text{Im } f'(x) &= \frac{1}{|y|} \mathbb{E}_y \mathbb{E}_{S^{n-1}} [f'(x + iy) \cdot y] \end{aligned}$$

with  $\|f'(x)\| \leq L \|x\| = O(n^2)$ .

---

<sup>21</sup>The paper provides similar results for strong-convex and non-convex functions. This approach recently surfaced in the optimization community Nikolovski and Stojkovska 2018; Hare and Srivastava 2023 with the first complete deterministic non-asymptotic analysis appearing in Jongeneel, Yue, and Kuhn 2021. The first applications of the complex-step derivative to Reinforcement Learning appeared in Wang and Spall 2021; Wang, Zhu, and Spall 2021.

<sup>22</sup>Jongeneel, Yue, and Kuhn 2021.

## A solution: the complex-step oracle <sup>22</sup> (exercise)

Let  $f \in C^1(X \subseteq \mathbb{R}^n; \mathbb{R})$ , using Cauchy-Riemann/Stokes show that:

$$f'(x) = E_y \in \mathbb{B}^n [\langle f(x + iy) \rangle]$$

$$r f'(x) = \frac{1}{n} E_y \in S^{n-1} [f(x + iy) y]$$

with  $\|r f'(x)\|_2 = O(n^{-2})$ .

Oracle (cs):  $g(x) = \frac{1}{n} \sum_{y \in S^{n-1}} (f(x + iy) - f(x)) y$

with  $E_{y \in S^{n-1}} \|g(x)\|_2^2 = O(n^{-4} + n^{-2} \|r f'(x)\|_2 + n \|r f'(x)\|_2^2)$ .

Algorithm:  $x_{k+1} = x_k - \alpha g_k(x_k); \quad \alpha = O(1/nL)$

Performance: for  $f$  convex  $\alpha = O(1/nL)$  and  $x_K := x_1 - \sum_{k=1}^K \alpha g_k(x_k)$

$$E[f(x_K)] - f(x^*) = O\left(\frac{nL \|x_1 - x^*\|_2^2}{K}\right) = O(n) \text{ gradient descent}^{21}$$

<sup>21</sup>The paper provides similar results for strong-convex and non-convex functions. This approach recently surfaced in the optimization community Nikolovski and Stojkovska 2018; Hare and Srivastava 2023 with the first complete deterministic non-asymptotic analysis appearing in Jongeneel, Yue, and Kuhn 2021. The first applications of the complex-step derivative to Reinforcement Learning appeared in Wang and Spall 2021; Wang, Zhu, and Spall 2021.

<sup>22</sup>Jongeneel, Yue, and Kuhn 2021.



## Solution:

$f \in C^1$  convex  $\Leftrightarrow f$  convex, e.g.,

for  $f(x) = x^4$  we have  $\Re(f(x + iy)) = x^4 - 6x^2y^2 + y^4$ .

Hence, look beyond typical convex proofs<sup>23</sup>.

---

<sup>23</sup>We frequently appeal to Schmidt, Roux, and Bach 2011, Lem. 1.

<sup>24</sup>Although the general form is largely due to Cartan (Élie).

<sup>25</sup> $C^1$  is sufficient, but not necessary.

## Solution:

$f \in C^1$  convex  $\Leftrightarrow f$  convex, e.g.,  
 for  $f(x) = x^4$  we have  $\langle f(x+iy) \rangle = x^4 - 6x^2(y)^2 + (y)^4$ .  
 Hence, look beyond typical convex proofs<sup>23</sup>.

Cauchy, Riemann and Stokes<sup>24</sup> meet:

for  $f \in C^1(\mathbb{R}^n; \mathbb{R})$  with  $f(x+iy) = u(x,y) + iv(x,y)$ , then

$$\partial_i u = \partial_i v; \partial_i u = -\partial_i v \quad \text{CR}$$

and for orientable  $\mathbb{R}^n$  we have that  $d_{\mathbb{R}} = \partial$  (Stokes), implication: the divergence theorem  $\text{div}(X) \text{dvol} = \partial \langle X \rangle \text{dvol}$ . Hence:

$$\begin{aligned} \int_{\mathbb{R}^n} f(x) \text{dvol} &\stackrel{(\text{def.}; \text{DCT})}{=} \int_{B^n} \langle f(x+iy) \rangle \text{dvol} \\ &\stackrel{(\text{CR})}{=} (\text{vol}(B^n))^{-1} \int_{B^n} \langle f(x+iy) \rangle \text{dvol} \\ &\stackrel{(\text{Stokes})}{=} \text{vol}(S^{n-1}) = (\text{vol}(B^n))^{-1} \int_{S^{n-1}} \langle f(x+iy) \rangle \text{dvol} \\ &\stackrel{(\text{vol}(S^{n-1}) = (\text{vol}(B^n))^{-1})}{=} (n) E_y \quad [= \langle f(x+iy) \rangle y]: \end{aligned}$$

<sup>23</sup>We frequently appeal to Schmidt, Roux, and Bach 2011, Lem. 1.

<sup>24</sup>Although the general form is largely due to Cartan (Élie).

<sup>25</sup> $C^1$  is sufficient, but not necessary.

## Example: worst function in the world

Consider the test function from Nesterov 2003, § 2.1.2

$$f_n(x) = L \frac{1}{2} (x^{(1)})^2 + \prod_{i=1}^{n-1} (x^{(i+1)} - x^{(i)})^2 + (x^{(n)})^2 \quad x^{(1)} \quad (9)$$

for  $x_1 = 0$ ,  $L = 10^8$ ,  $L_1(f) = 4L$  and  $(x^*)^{(i)} = 1 \quad i=(n+1)$  with  $x^{(i)}$ .

(ia) Suboptimality gap  
 $f(x_K) - f^*$  for the test  
 function (9).

(ib) Suboptimality gap  
 $f(x_K) - f^*$  for the test  
 function (9).

Figure: The single-point Complex-smoothing (CS) compared to the multi-point Gaussian smoothing (GS) (fd) method from Nesterov and Spokoiny 2017.

Example: strong convexity  $f(x) = \frac{1}{2}\|x\|_2^2$

(a)  $f(x_K) - f^*, n = 10^0$ .      (b)  $f(x_K) - f^*, n = 10^2$ .      (c)  $f(x_K) - f^*, n = 10^4$ .

Figure: The single-point Complex-smoothing (CS) compared to the multi-point Gaussian smoothing (GS) (fd and cd) method from Nesterov and Spokoiny 2017, Eq. (55). The rate is for (GS) (cd).

## Example: non-convex optimization

Consider a Rosenbrock optimization problem

$$\min_{x \in \mathbb{R}^2} (1 - x^{(1)})^2 + 100(x^{(2)} - (x^{(1)})^2)^2 \quad (10)$$

with  $x^* = (1; 1)$ .

(a) Suboptimality gap  
 $f(x_K) - f^*$  for (10).

(b) Paths taken corresponding to  
Figure 4a.

Figure: The single-point Complex-smoothing (CS) method versus Gaussian-smoothing Nesterov and Spokoiny 2017.

## What about noise?

We can handle<sup>26</sup> **simulation noise**, e.g.,  $= (f(z)) + \epsilon$ ,  $z \in \mathbb{C}^n$ ,  $\epsilon \in (0; \epsilon^2)$ .

$$\text{Oracle (cs, noisy): } g(x) = \frac{1}{n} \sum_{i=1}^n (f(x + i y)) y + \frac{1}{n} \sum_{i=1}^n \epsilon_i y; \quad y \in \mathbb{S}^{n-1}; \quad (11)$$

Handling ( $n = \infty$ ) non-trivial. In general, we need  $k = O(1/\epsilon)$  and  $\epsilon_k = o(\epsilon)$ .

---

<sup>26</sup>Jongeneel 2021.

<sup>27</sup>Building upon Hazan, Rakhlin, and Bartlett 2008; Akhavan, Pontil, and Tsybakov 2020.

<sup>28</sup>Shown by building upon Shamir 2013.

## What about noise?

We can handle<sup>26</sup> **simulation noise**, e.g.,  $= (f(z)) + \epsilon$ ,  $z \in \mathbb{C}^n$ ,  $\epsilon \in (0; \sigma^2)$ .

$$\text{Oracle (cs, noisy): } g(x) = \frac{1}{n} \sum_{i=1}^n (f(x + i y)) y + \frac{1}{n} \sum_{i=1}^n \epsilon_i y; \quad y \in \mathbb{S}^{n-1}; \quad (11)$$

Handling  $(n = \infty)$  non-trivial. In general, we need  $k = O(1/\epsilon)$  and  $\epsilon_k = o(\epsilon)$ .

Non-asymptotic results<sup>27</sup> for: constrained/unconstrained strongly convex functions and some non-convex functions (locally).

The algorithm is rate-optimal in the quadratic setting<sup>28</sup> (not surprising).

---

<sup>26</sup>Jongeneel 2021.

<sup>27</sup>Building upon Hazan, Rakhlin, and Bartlett 2008; Akhavan, Pontil, and Tsybakov 2020.

<sup>28</sup>Shown by building upon Shamir 2013.

## What about noise?

We can handle<sup>26</sup> **simulation noise**, e.g.,  $= (f(z)) + \epsilon$ ,  $z \in \mathbb{C}^n$ ,  $\epsilon \in (0; \epsilon^2)$ .

$$\text{Oracle (cs, noisy): } g(x) = \frac{1}{n} \sum_{i=1}^n (f(x + i y)) y + \frac{1}{n} \sum_{i=1}^n \epsilon_i y; \quad y \in S^{n-1}; \quad (11)$$

Handling  $(n = \infty)$  non-trivial. In general, we need  $k = O(1/\epsilon)$  and  $k = o(k)$ .

Non-asymptotic results<sup>27</sup> for: constrained/unconstrained strongly convex functions and some non-convex functions (locally).

The algorithm is rate-optimal in the quadratic setting<sup>28</sup> (not surprising).

Why the ball  $B^n$  and not some other geometry

$M \subseteq \mathbb{R}^n = f(M) \subseteq [1; 1]^n$ :  $M$  diffeomorphic to  $B^n$ ? Optimal in the sense that

$$\min_{M \subseteq \mathbb{R}^n} \frac{\text{vol}(f(M))}{\text{vol}(M)} = \frac{1}{\text{vol}(B^n)}; \quad B^n = \operatorname{argmin}_{M \subseteq \mathbb{R}^n} \frac{\text{vol}(f(M))}{\text{vol}(M)}; \quad (12)$$

which follows from the isoperimetric inequality in  $\mathbb{R}^n$  Osserman 1978.

---

<sup>26</sup>Jongeneel 2021.

<sup>27</sup>Building upon Hazan, Rakhlin, and Bartlett 2008; Akhavan, Pontil, and Tsybakov 2020.

<sup>28</sup>Shown by building upon Shamir 2013.



## Example: non-convex optimization (outlook)

Regularity of ODE/PDE constrained optimization problems can often be understood. We apply our zeroth-order algorithm to a ODE problem<sup>29</sup>.

- (a) Decay of  $f(x_K)$  with the total number  $K$  of iterations for 10 independent simulation runs.
- (b) Trajectories starting from  $x(0)$  (unknown),  $x_0$  (initial) and  $x_K$  for  $K = 10^5$  (optimized).

Figure: Estimating the initial state  $x(0)$  of a Lorenz system from a noisy measurement of the state  $x(2) = x^2(x(0))$  (grey circle in 5b) at time 2. Even though the initial estimate  $x_0$  is close to the optimized estimate  $x_K$ ,  $x^2(x_0)$  is far from  $x^2(x(0))$ .

---

<sup>29</sup>The complex-step derivative is implemented in an airfoil optimization package. Their underlying algorithm relies on sequential quadratic programming Nocedal and Wright 2006, Ch. 18, as such, the guarantees one can provide are different, see <https://mdolab-cmplxfoil.readthedocs-hosted.com/en/latest/index.html>, our work aims at providing rigorous guarantees with respect to the optimization algorithm itself.

# The end

Main take away: single-point estimator where  $k = O(1/\epsilon)$  can be safely implemented.

Many open problems remain.

contact: wjongeneel.nl (slides will appear there).

## Key references:

- (a) Arkadi Semenovich Nemirovsky and David Borisovich Yudin (1983). Problem complexity and method efficiency in optimization . Wiley-Interscience
- (b) Boris Teodorovich Polyak and Aleksandr Borisovich Tsybakov (1990). Optimal order of accuracy of search algorithms in stochastic optimization . In: Problemy Peredachi Informatsii 26.2, pp. 45 53
- (c) John C Duchi et al. (2015). Optimal rates for zero-order convex optimization: The power of two function evaluations . In: IEEE Transactions on Information Theory 61.5, pp. 2788 2806
- (d) Francis Bach and Vianney Perchet (2016). Highly-smooth zero-th order online optimization . In: Conference on Learning Theory, pp. 1 27
- (e) Yurii Nesterov and Vladimir Spokoiny (2017). Random gradient-free minimization of convex functions . In: Foundations of Computational Mathematics 17.2, pp. 527 566
- (f) Albert S Berahas et al. (2022). A theoretical and empirical comparison of gradient approximations in derivative-free optimization . In: Foundations of Computational Mathematics 22.2, pp. 507 560
- (g) Wouter Jongeneel, Man-Chung Yue, and Daniel Kuhn (2021). Small Errors in Random Zeroth Order Optimization are Imaginary . In: arXiv: 2103.05478

## Appendix.

## References:

Abreu, Rafael et al. (2018). On the accuracy of the Complex-Step-Finite-Difference method . In: Journal of Computational and Applied Mathematics 340, pp. 390–403.

Agarwal, Alekh, Ofer Dekel, and Lin Xiao (2010). Optimal Algorithms for Online Convex Optimization with Multi-Point Bandit Feedback. . In: COLT, pp. 28–40.

Akhavan, Arya, Massimiliano Pontil, and Alexandre B Tsybakov (2020). Exploiting higher order smoothness in derivative-free optimization and continuous bandits . In: arXiv preprint arXiv:2006.07862.

Audet, Charles and Warren Hare (2017). Derivative-free and Blackbox Optimization . Springer.

Bach, Francis and Vianney Perchet (2016). Highly-smooth zero-th order online optimization . In: Conference on Learning Theory pp. 1–27.

Berahas, Albert S, Liyuan Cao, and Katya Scheinberg (2021). Global convergence rate analysis of a generic line search algorithm with noise . InSIAM Journal on Optimization 31.2, pp. 1489–1518.

Berahas, Albert S et al. (2022). A theoretical and empirical comparison of gradient approximations in derivative-free optimization . In: Foundations of Computational Mathematics 22.2, pp. 507 560.

Bezanson, Je et al. (2017). Julia: A Fresh Approach to Numerical Computing . In: SIAM Review 59.1, pp. 65 98. doi : 10.1137/141000671 .

Bradbury, James et al. (2018). JAX: composable transformations of Python+NumPy programs . Version 0.3.13. url : <http://github.com/google/jax> .

Conn, Andrew R., Katya Scheinberg, and Luis N. Vicente (2009). Introduction to Derivative-Free Optimization . Society for Industrial and Applied Mathematics.

Cox, MG and PM Harris (2004). Software Support for Metrology Best Practice Guide No. 11 . In.

d'Aspremont, Alexandre (2008). Smooth optimization with approximate gradient . In: SIAM Journal on Optimization 19.3, pp. 1171 1183.

Devolder, Olivier, François Glineur, and Yurii Nesterov (2014). First-order methods of smooth convex optimization with inexact oracle . In: Mathematical Programming 146.1, pp. 37 75.

Duchi, John C et al. (2015). Optimal rates for zero-order convex optimization: The power of two function evaluations . In: IEEE Transactions on Information Theory 61.5, pp. 2788 2806.

Flaxman, Abraham, Adam Tauman Kalai, and H. Brendan McMahan (2004). Online convex optimization in the bandit setting: gradient descent without a gradient . In: CoRR.

Gasnikov, Alexander V et al. (2017). Stochastic online optimization. Single-point and multi-point non-linear multi-armed bandits. Convex and strongly-convex case . In: Automation and remote control 78.2, pp. 224 234.

Hare, Warren and Kashvi Srivastava (2023). A numerical study of applying complex-step gradient and Hessian approximations in blackbox optimization . In: Pacific Journal of Optimization 19.3, pp. 391 410.

Hazan, Elad and K r Y Levy (2014). Bandit Convex Optimization: Towards Tight Bounds. . In: NIPS, pp. 784 792.

Hazan, Elad, Alexander Rakhlin, and Peter Bartlett (2008). Adaptive Online Gradient Descent . In: Advances in Neural Information Processing Systems pp. 65 72.

Innes, Michael (2018). Don't unroll adjoint: Differentiating SSA-form programs . In: arXiv preprint arXiv:1810.07951.

Jongeneel, Wouter (2021). Imaginary Zeroth-Order Optimization . In: arXiv: 2112.07488.








Jongeneel, Wouter, Man-Chung Yue, and Daniel Kuhn (2021). Small Errors in Random Zeroth Order Optimization are Imaginary . In: arXiv: 2103.05478.

Kiefer, Jack and Jacob Wolfowitz (1952). Stochastic estimation of the maximum of a regression function . In: The Annals of Mathematical Statistics , pp. 462-466.








Lam, Henry, Haidong Li, and Xuhui Zhang (2021). Minimally efficient finite-difference stochastic gradient estimators using black-box function evaluations . In: Operations Research Letters 49.1, pp. 40-47.








Larson, Jeffrey, Matt Menickelly, and Stefan M Wild (2019). Derivative-free optimization methods . In: arXiv preprint arXiv:1904.11585.

Liu, S. et al. (2020). A Primer on Zeroth-Order Optimization in Signal Processing and Machine Learning: Principles, Recent Advances, and Applications . In: IEEE Signal Processing Magazine 37.5, pp. 43-54.

-  Lyness, J. N. and C. B. Moler (1967). “Numerical Differentiation of Analytic Functions”. In: *SIAM Journal on Numerical Analysis* 4.2, pp. 202–210.
-  Martins, Joaquim R. R. A., Peter Sturdza, and Juan J. Alonso (Sept. 2003). “The Complex-Step Derivative Approximation”. In: *ACM Trans. Math. Softw.* 29.3, 245–262.
-  Moses, William and Valentin Churavy (2020). “Instead of Rewriting Foreign Code for Machine Learning, Automatically Synthesize Fast Gradients”. In: *Advances in Neural Information Processing Systems*. Vol. 33, pp. 12472–12485.
-  Nemirovsky, Arkadi Semenovitch and David Borisovich Yudin (1983). *Problem complexity and method efficiency in optimization*. Wiley-Interscience.
-  Nesterov, Yurii (2003). *Introductory lectures on convex optimization: A basic course*. Vol. 87. Springer Science & Business Media.
-  — (2011). *Random gradient-free minimization of convex functions*. CORE Discussion Papers 2011001.
-  Nesterov, Yurii and Vladimir Spokoiny (2017). “Random gradient-free minimization of convex functions”. In: *Foundations of Computational Mathematics* 17.2, pp. 527–566.



-  Nikolovski, Filip and Irena Stojkovska (2018). “Complex-step derivative approximation in noisy environment”. In: *Journal of Computational and Applied Mathematics* 327, pp. 64–78.
-  Nocedal, Jorge and Stephen J. Wright (2006). *Numerical Optimization*. Springer Series in Operations Research. Springer.
-  Novitskii, Vasilii and Alexander Gasnikov (2021). “Improved Exploiting Higher Order Smoothness in Derivative-free Optimization and Continuous Bandit”. In: *arXiv preprint arXiv:2101.03821*.
-  Osserman, Robert (1978). “The isoperimetric inequality”. In: *Bulletin of the American Mathematical Society* 84.6, pp. 1182–1238.
-  Polyak, Boris Teodorovich and Aleksandr Borisovich Tsybakov (1990). “Optimal order of accuracy of search algorithms in stochastic optimization”. In: *Problemy Peredachi Informatsii* 26.2, pp. 45–53.
-  Revels, Jarrett, Miles Lubin, and Theodore Papamarkou (2016). “Forward-mode automatic differentiation in Julia”. In: *arXiv preprint arXiv:1607.07892*.
-  Scheinberg, Katya (2022). “Finite Difference Gradient Approximation: To Randomize or Not?” In: *INFORMS Journal on Computing* 34.5, pp. 2384–2388.

-  Schmidt, Mark, Nicolas Roux, and Francis Bach (2011). “Convergence rates of inexact proximal-gradient methods for convex optimization”. In: *Neural Information Processing Systems*, pp. 1458–1466.
-  Shamir, Ohad (2013). “On the complexity of bandit and derivative-free stochastic convex optimization”. In: *Conference on Learning Theory*, pp. 3–24.
-  — (2017). “An optimal algorithm for bandit and zero-order convex optimization with two-point feedback”. In: *The Journal of Machine Learning Research* 18.1, pp. 1703–1713.
-  Spall, James C (2005). *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*. John Wiley & Sons.
-  Squire, William and George Trapp (1998). “Using Complex Variables to Estimate Derivatives of Real Functions”. In: *SIAM Review* 40.1, pp. 110–112.
-  Wang, Long and James C Spall (2021). “Improved SPSA using complex variables with applications in optimal control problems”. In: *American Control Conference*. IEEE, pp. 3519–3524.
-  Wang, Long, Jingyi Zhu, and James C Spall (2021). “Model-free optimal control using SPSA with complex variables”. In: *Annual Conference on Information Sciences and Systems*. IEEE, pp. 1–5.