# Lecture 1: Introduction to zeroth-order optimization (ZOO)

## Motivation

**Table 1:** Workers $W_{i,j}$

| Shift | Skill levels | | |
|---|---|---|---|
| | High | Med | Low |
| S1 | 1 | 3 | 7 |
| S2 | 0 | 5 | 2 |
| S3 | 3 | 1 | 2 |

**Table 2:** SLA targets $\gamma_{i,j}$

| Priority | Customers | |
|---|---|---|
| | Bossy Corp | Cool Inc |
| $P_1$ | 4h | 5h |
| $P_2$ | 8h | 12h |
| $P_3$ | 24h | 48h |
| $P_4$ | 18h | 144h |

$$G \equiv$$

$$x \rightarrow \boxed{\text{Gradient Oracle}} \rightarrow \nabla f(x) + \epsilon \rightarrow N(0,1)$$

$$x^* \in \arg\min_x f(x)$$

$$x_{k+1} = x_k - a_k G_k$$

Aim: Find the optimal number of workers for each shift and of each skill level

- that minimizes the labor cost and

- satisfies SLA requirements
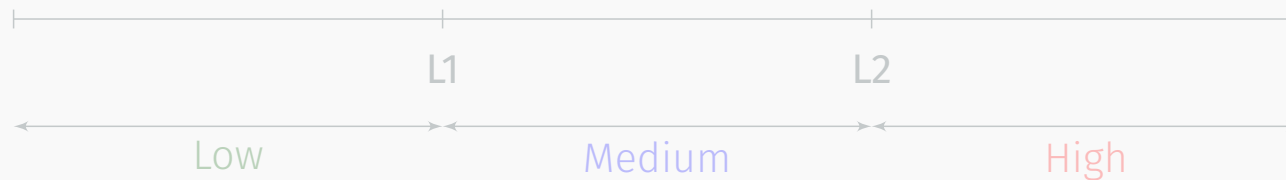
On a good day, the traffic is . . .

**Aim:** Maximize traffic flow

*Input:*
Coarse congestion estimates

*Output:*
Policy for switching traffic lights

*Input: Coarse congestion estimates*
Sensor loops at two points along the road

L1     L2

Low     Medium     High

How to switch traffic lights given L1 and L2?

How to choose L1 and L2 for a given policy and road network?

8

**Aim:** Maximize traffic flow

*Input:*
Coarse congestion estimates

*Output:*
Policy for switching traffic lights

*Input: Coarse congestion estimates*
Sensor loops at two points along the road

L1     L2

Low          Medium          High

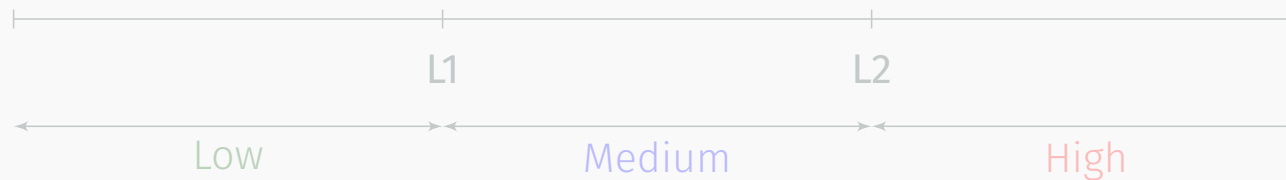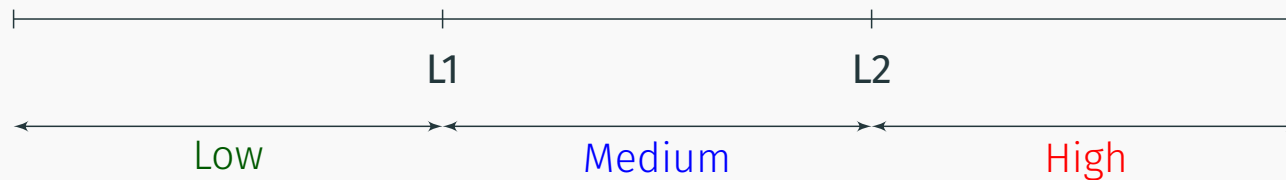How to switch traffic lights given L1 and L2?

How to choose L1 and L2 for a given policy and road network?

**Aim:** Maximize traffic flow

*Input:*
Coarse congestion estimates

*Output:*
Policy for switching traffic lights

*Input: Coarse congestion estimates*
Sensor loops at two points along the road

L1      L2

Low     Medium     High

How to switch traffic lights given L1 and L2?

How to choose L1 and L2 for a given policy and road network?

# Application III: Intrusion detection using sensor networks



● Sensor    —— Intruder

Aim:

- minimize the energy consumption of the sensors, while

- keeping tracking error to a minimum

# Common application traits

**Stochastic:**
noisy observations

**Model-free:**
sample access to objective
* gradients unavailable

**High-dimensional:**
brute-force search infeasible

**Solution:**
Simultaneous perturbation
methods

# The framework

# Basic optimization problem

noisy observation

$$F(\theta_n, \xi_n)$$



$F(\theta_n, \xi_n)$

Observation

Environment    Agent

Query

$\theta_n$

$\theta_n$

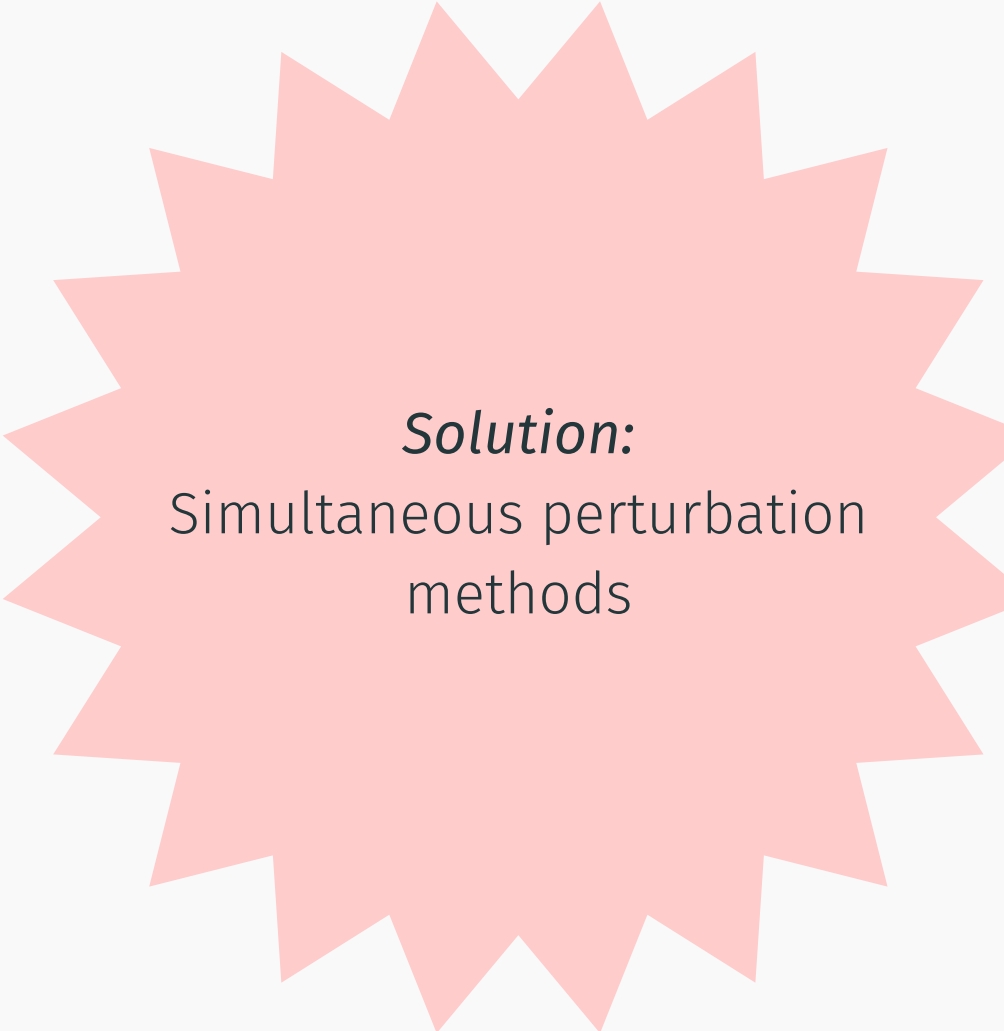Aim: $\theta^* = \arg\min_{\theta \in \Theta} \left\{ f(\theta) \triangleq \mathbb{E}[F(\theta, \xi)] \right\},$

- $f \colon \mathbb{R}^d \to \mathbb{R}$ is the performance measure

  - $f$ *not* assumed to be convex

- $F(\theta, \xi)$ is the sample performance

- $\xi$ is the noise factor that captures stochastic nature of the problem

- $\theta$ is the (vector) parameter of interest

- $\Theta \subseteq \mathbb{R}^d$ is the feasible region in which $\theta$ takes values.

Stochastic optimization deals with highly nonlinear and high dimensional systems. The challenges with these problems are:

- Too complex to solve analytically.

- Many simplifying assumptions are required.

A good alternative of modeling and analysis is "Simulation"

Zero mean

$$\theta_n \longrightarrow \boxed{\text{Simulator}} \longrightarrow f(\theta_n) + \xi_n$$

Figure 1: Simulation optimization

Stochastic optimization deals with highly nonlinear and high dimensional systems. The challenges with these problems are:

- Too complex to solve analytically.

- Many simplifying assumptions are required.

A good alternative of modeling and analysis is "Simulation"



**Figure 1:** Simulation optimization

Recall: $f(\theta) = \mathbb{E}\left[F(\theta, \xi)\right]$.

**Two settings for** noise:

Controlled noise $\xi$ can be kept fixed between queries to obtain $F(\theta_1, \xi)$ and $F(\theta_2, \xi)$

Uncontrolled noise $F(\theta, \xi)$ can be obtained at any point, but $\xi$ is not controllable

Recall: $f(\theta) = \mathbb{E}\left[F(\theta, \xi)\right]$.

Two settings for noise:

Controlled noise: $\xi$ can be kept fixed between queries to obtain $F(\theta_1, \xi)$ and $F(\theta_2, \xi)$

Uncontrolled noise: $F(\theta, \xi)$ can be obtained at any point, but $\xi$ is not controllable

**Deterministic optimization problem**

- focus is on search for better solutions

- Complete information about objective function $f$, esp. gradients

Stochastic optimization problem

- $f$ cannot be obtained directly, but we are given sample access, i.e.,

$$f(\theta) \equiv E_\xi[F(\theta, \xi)]$$

- Each sample $F(\theta, \xi)$ is obtained from an expensive simulation experiment or a (real) field test

- focus is on both search and evaluation

  - Tradeoff between evaluating better vs. finding more candidate solutions

Challenge: to find $\theta^* = \arg\min_{\theta \in \Theta} f(\theta)$, given only noisy function evaluations.

**Deterministic optimization problem**

- focus is on search for better solutions

- Complete information about objective function $f$, esp. gradients

**Stochastic optimization problem**

- $f$ cannot be obtained directly, but we are given sample access, i.e.,

$$f(\theta) \equiv E_\xi[F(\theta, \xi)]$$

- Each sample $F(\theta, \xi)$ is obtained from an expensive simulation experiment or a (real) field test

- focus is on both search and evaluation

  - Tradeoff between evaluating better vs. finding more candidate solutions

Challenge: to find $\theta^* = \arg\min_{\theta \in \Theta} f(\theta)$, given only noisy function evaluations.

# Some more applications

## Energy Demand management

- Consumer demand, energy generation are uncertain.

- Objective is to minimize the difference.

## Transportation

- Car-following model

- route choice

- traffic assignment model

Service systems (banks, restaurants, call centers, amusement parks)

# and some more..

Transportation systems (airports: air space, runways, baggage, roads, queues)

Manufacturing    Semiconductor fab    Supply chains

Networks    Finance    Insurance

Education    Healthcare    Banking

Mining    Oil & Gas    Call centers

Automotive OEM    Aerospace    Retirement planning

# Some vendors...

aGPSS    Analytic solver    Analytica

AnyLogic    FlexSim    ExtendSim Pro

Arena    MedModel Opt Suite    Oracle Crystal Ball

Pedestrian dynamics    Polaris    ProModel Opt Suite

SLIM    Solver SDK Platform    Vanguard

Tecnomatix    Simio    DiscoverSim

[1] James J. Swain, "Simulation Software Survey — Simulation Takes Over: Reality is for Sissies," *OR/MS Today*, Oct 2017.

18

# Success stories...

- **Kroger (Edelman 2013 finalist, gradient-based)** Kroger Uses Simulation-Optimization to Improve Pharmacy Inventory Management
    - `www.youtube.com/watch?v=BNyDbBy-KYY` (start at 0:45)
    - `https://www.informs.org/About-INFORMS/News-Room/Press-Releases/Edelman-2013-Announcement`

      *The Franz Edelman Award recognizes outstanding examples of innovative operations research and analytics that improves organizations and often change people's lives.*

- Financial engineering
    - Monte Carlo simulation used widely on Wall Street.
    - Gradient estimates needed for hedging.
    - Hot research area: several research papers continue to be published

# The Matrix has you..

# First-order methods

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{1}$$

Suppose that

- $G_n$ is an noisy estimate of the gradient $\nabla f(\theta_n)$, i.e.,
  $\mathbb{E}(G_n) = \nabla f(\theta_n).$

- $\{a_n\}$ are pre-determined step-sizes satisfying:
  $$\sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty$$

- iterates are stable: $\sup_n \|\theta_n\| < \infty.$

Theorem (Variant of Robbins Monro stochastic approximation)

Letting $K := \{\theta \mid \nabla f(\theta) = 0\}$, we have

$$\theta_n \to K \ a.s. \ as \ n \to \infty.$$

# Stochastic analog of gradient descent

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{1}$$

Suppose that

- $G_n$ is an noisy estimate of the gradient $\nabla f(\theta_n)$, i.e., $\mathbb{E}(G_n) = \nabla f(\theta_n)$.

- $\{a_n\}$ are pre-determined step-sizes satisfying:

$$\sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty$$

- iterates are stable: $\sup_n \|\theta_n\| < \infty$.

Theorem (Variant of Robbins Monro stochastic approximation)

Letting $K := \{\theta \mid \nabla f(\theta) = 0\}$, we have

$$\theta_n \to K \text{ a.s. as } n \to \infty.$$

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{1}$$

Suppose that

- $G_n$ is an noisy estimate of the gradient $\nabla f(\theta_n)$, i.e.,
  $$\mathbb{E}(G_n) = \nabla f(\theta_n).$$

- $\{a_n\}$ are pre-determined step-sizes satisfying:
  $$\sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty$$

- iterates are stable: $\sup_n \|\theta_n\| < \infty$.

**Theorem (Variant of Robbins Monro stochastic approximation)**

*Letting $K := \{\theta \mid \nabla f(\theta) = 0\}$, we have*

$$\theta_n \to K \text{ a.s. as } n \to \infty.$$

# Stochastic analog of gradient descent

$$\theta_{n+1} = \theta_n - a_n G_n. \qquad (1)$$

*gradient estimate*

*step size*

Suppose that

- $G_n$ is an noisy estimate of the gradient $\nabla f(\theta_n)$, i.e.,
  $$\mathbb{E}(G_n) = \nabla f(\theta_n).$$ *unbiased gradients*

- $\{a_n\}$ are pre-determined step-sizes satisfying:
  $$\sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty$$
  *e.g. $a_n = \dfrac{1}{n}$*

- iterates are stable: $\sup_n \|\theta_n\| < \infty$.

- *f is smooth*

## Theorem (Variant of Robbins Monro stochastic approximation)

*Letting $K := \{\theta \mid \nabla f(\theta) = 0\}$, we have*

$$\theta_n \to K \text{ a.s. as } n \to \infty.$$

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{2}$$

How to keep iterates stable?

Project $\theta_n$ onto a compact and convex set $\Theta \leftarrow$ Projected stochastic approximation

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{2}$$

How to estimate the gradient of $f$ from samples?

$$\theta_n \longrightarrow \boxed{\text{Simulator}} \longrightarrow f(\theta_n) + \xi_n$$

*or*
*real world*
*measurement*

Simultaneous perturbation methods.

Stochastic approximation (SA) alphabet soup

FDSA   Finite difference stochastic approximation

SPSA   Simultaneous perturbation stochastic approximation

SFSA   Smoothed functional stochastic approximation

RDSA   Random direction stochastic approximation

# In the next few slides . . .

$$\theta_{n+1} = \theta_n - a_n G_n. \tag{3}$$

Q1) How to form $G_n$ from function samples so that $G_n \approx \nabla f(\theta_n)$

Q2) Such a $G_n$ - is it unbiased?

Q3) Does $\theta_n$ converge to $\theta^*$ with such a $G_n$?

Q4) If answer is yes to above, what is the convergence rate?

# Outline

$$\Theta \longrightarrow \boxed{\begin{array}{c} \text{Zeroth-order} \\ \text{Oracle} \end{array}} \longrightarrow \begin{array}{c} f(\theta) \\ + \epsilon \end{array}$$

For $\delta$ small,

$$f'(\theta) \approx \frac{f(\theta + \delta) - f(\theta)}{\delta}$$

Finite-differencing

24

Finite-difference stochastic approximation (FDSA) (Kiefer and Wolfowitz, 1952):

*(handwritten: One-sided gradient estimate →)*

*(handwritten: $\theta \rightarrow \boxed{\phantom{x}} \rightarrow f(\theta)$)*

$$g^i = \frac{1}{\delta}\left(f(\theta + \delta e_i) - f(\theta)\right), \quad i = 1, \ldots, N\,d$$

Assume $f \in \mathcal{C}^3$

Taylor-series expansion:

*(handwritten: Assume Hessian's norm bounded above)*

$$f(\theta + \delta e_i) = f(\theta) + \delta \nabla f(\theta)^\top e_i + \frac{\delta^2}{2} e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

*(handwritten under term: $\nabla_i f(\theta)$)*

Accuracy: $\|g - \nabla f(\theta)\|_2 = O(\delta)$.

*(handwritten: Suppose $\exists \tilde{\delta}$ s.t. $\|\tilde{g} - \nabla f(\theta)\| = O(\delta^2)$)*

Needs $N + 1$ queries.

Finite-difference stochastic approximation (FDSA) (Kiefer and Wolfowitz, 1952):

$$g^i = \frac{1}{\delta}\left(f(\theta + \delta e_i) - f(\theta)\right), \quad i = 1, \ldots, d.$$

Assume $f \in \mathcal{C}^3$ ( three times continuously differentiable )
Taylor-series expansion:

$$f(\theta + \delta e_i) = f(\theta) + \delta\,\nabla f(\theta)e_i + \frac{\delta^2}{2}\,e_i^\top \nabla^2 f(\theta)e_i + O(\delta^3).$$

Accuracy: $\|g - \nabla f(\theta)\|_2 = O(\delta)$.

for one gradient descent update Fterction

Needs $d + 1$ queries.

25

Improved estimate:

*Balanced estimator* →

$$g^i = \frac{1}{2\delta}\left(f(\theta + \delta e_i) - f(\theta - \delta e_i)\right), \quad i = 1, \ldots, N$$

Taylor-series expansions:

① $\quad f(\theta + \delta e_i) = f(\theta) + \delta\,\nabla f(\theta)e_i + \frac{\delta^2}{2}\,e_i^\top \nabla^2 f(\theta)e_i + O(\delta^3).$

② $\quad f(\theta - \delta e_i) = f(\theta) - \delta\,\nabla f(\theta)e_i + \frac{\delta^2}{2}\,e_i^\top \nabla^2 f(\theta)e_i + O(\delta^3).$

① − ② $= \quad 2\delta\ \nabla f(\theta)e_i\ +\ O(\delta^3)$

$$\text{Accuracy: } \|g - \nabla f(\theta)\|_2 = O(\delta^2).$$

Needs $2N$ queries.

# FDSA with two-sided Differences

Improved estimate:

$$g^i = \frac{1}{2\delta}\left(f(\theta + \delta e_i) - f(\theta - \delta e_i)\right), \quad i = 1, \ldots, N.$$

Taylor-series expansions:

$$f(\theta + \delta e_i) = f(\theta) + \delta\,\nabla f(\theta)e_i + \frac{\delta^2}{2}\,e_i^\top \nabla^2 f(\theta)e_i + O(\delta^3).$$

$$f(\theta - \delta e_i) = f(\theta) - \delta\,\nabla f(\theta)e_i + \frac{\delta^2}{2}\,e_i^\top \nabla^2 f(\theta)e_i + O(\delta^3).$$

$$\text{Accuracy: } \|g - \nabla f(\theta)\|_2 = O(\delta^2).$$

Needs 2N queries.

Improved estimate:

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N.$$

*(handwritten annotations: θ → [box] → 2f(θ)+ξ ; 2 r\cdot v. near noise)*

Taylor-series expansions:

$$f(\theta + \delta e_i) = f(\theta) + \delta \, \nabla f(\theta)^\top e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

$$f(\theta - \delta e_i) = f(\theta) - \delta \, \nabla f(\theta)^\top e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

Assumption: $\mathbb{E}\left[\xi^\pm\right] = 0$, $\mathbb{E}\left[(\xi^\pm)\right] \leq \sigma^2 < +\infty$.

$\mathbb{E}\left[G^i\right] = g^i$.   Hence

$$\left\| \mathbb{E}\left[G\right] - \nabla f(\theta) \right\|_2 = O(\delta^2). \longleftarrow \text{bias}$$

Improved estimate:

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N$$

Taylor-series expansions:

$$f(\theta + \delta e_i) = f(\theta) + \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

$$f(\theta - \delta e_i) = f(\theta) - \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

Assumption: $\mathbb{E}\left[\xi^\pm\right] = 0, \mathbb{E}\left[(\xi^\pm)\right] \leq \sigma^2 < +\infty$.

$\mathbb{E}\left[G^i\right] = g^i$. Hence

$$\left\| \mathbb{E}\left[G\right] - \nabla f(\theta) \right\|_2 = O(\delta^2). \longleftarrow \text{bias}$$

Improved estimate:

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N.$$

Taylor-series expansions:

$$f(\theta + \delta e_i) = f(\theta) + \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

$$f(\theta - \delta e_i) = f(\theta) - \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

Assumption: $\mathbb{E}\left[\xi^\pm\right] = 0$, $\mathbb{E}\left[(\xi^\pm)\right] \leq \sigma^2 < +\infty$.

$\mathbb{E}\left[G^i\right] = g^i$. Hence

$$\left\| \mathbb{E}\left[G\right] - \nabla f(\theta) \right\|_2 = O(\delta^2). \longleftarrow \text{bias}$$

Improved estimate:

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N.$$
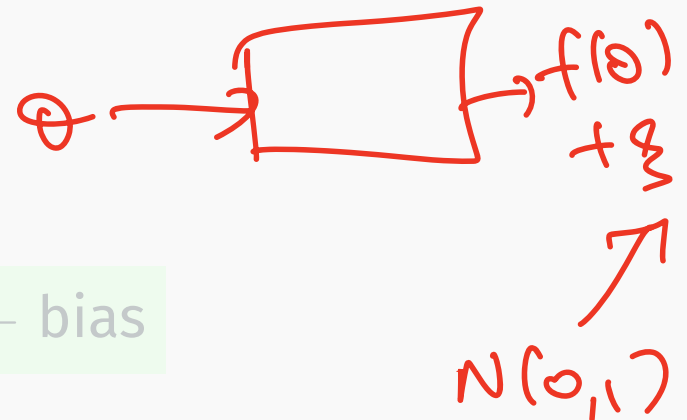
Taylor-series expansions:

$$f(\theta + \delta e_i) = f(\theta) + \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

$$f(\theta - \delta e_i) = f(\theta) - \delta \, \nabla f(\theta) e_i + \frac{\delta^2}{2} \, e_i^\top \nabla^2 f(\theta) e_i + O(\delta^3).$$

Assumption: $\mathbb{E}\left[\xi^\pm\right] = 0, \mathbb{E}\left[(\xi^\pm)\right] \leq \sigma^2 < +\infty.$

$\mathbb{E}\left[G^i\right] = g^i.$  Hence

$$\left\| \mathbb{E}\left[G\right] - \nabla f(\theta) \right\|_2 = O(\delta^2) . \longleftarrow \text{bias}$$

So far: with FDSA, we can get a gradient estimate

$$G^i = \frac{1}{2\delta}\left\{f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-)\right\}, \quad i = 1, \ldots, N \quad \text{with}$$

bias $O(\delta^2)$

what is second moment: $\mathbb{E}\left[\|G\|_2^2\right] = ?$

$G_i = g_i + \dfrac{\xi_i^+ - \xi_i^-}{2\delta}$ , hence $\mathbb{E}\left[G_i^2\right] = g_i^2 + \dfrac{2\sigma^2}{4\delta^2} = g_i^2 + \dfrac{\sigma^2}{2\delta^2}$ and

$$\mathbb{E}\left[\|G\|_2^2\right] = \|g\|_2^2 + O\left(\frac{N}{\delta^2}\right).$$

So far: with FDSA, we can get a gradient estimate

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N \quad \text{with}$$

bias $O(\delta^2)$

what is second moment: $\mathbb{E}\left[\|G\|_2^2\right] = ?$

$$G_i = g_i + \frac{\xi_i^+ - \xi_i^-}{2\delta}, \quad \text{hence } \mathbb{E}\left[G_i^2\right] = g_i^2 + \frac{2\sigma^2}{4\delta^2} = g_i^2 + \frac{\sigma^2}{2\delta^2} \quad \text{and}$$

$$\mathbb{E}\left[\|G\|_2^2\right] = \|g\|_2^2 + O\left(\frac{N}{\delta^2}\right).$$

So far: with FDSA, we can get a gradient estimate

$$G^i = \frac{1}{2\delta} \left\{ f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-) \right\}, \quad i = 1, \ldots, N \quad \text{with}$$

bias $O(\delta^2)$

what is second moment: $\mathbb{E}\left[ \|G\|_2^2 \right] =$?

$$G_i = g_i + \frac{\xi_i^+ - \xi_i^-}{2\delta}, \text{ hence } \mathbb{E}\left[ G_i^2 \right] = g_i^2 + \frac{2\sigma^2}{4\delta^2} = g_i^2 + \frac{\sigma^2}{2\delta^2} \quad \text{and}$$

$$\mathbb{E}\left[ \|G\|_2^2 \right] = \|g\|_2^2 + O\left( \frac{N}{\delta^2} \right).$$

$$\|\mathbb{E}G - \nabla f\| = O(\delta^2)$$

So far: with FDSA, we can get a gradient estimate

$$G^i = \frac{1}{2\delta}\left\{f(\theta + \delta e_i) + \xi_i^+ - (f(\theta - \delta e_i) + \xi_i^-)\right\}, \quad i = 1, \dots, N. \quad \text{with}$$

bias $O(\delta^2)$

what is second moment: $\mathbb{E}\left[\|G\|_2^2\right] = ?$

$$G = \begin{pmatrix} G_1 \\ \vdots \\ G_N \end{pmatrix}$$

$$G_i = g_i + \frac{\xi_i^+ - \xi_i^-}{2\delta}, \text{ hence } \mathbb{E}\left[G_i^2\right] = g_i^2 + \frac{2\sigma^2}{4\delta^2} = g_i^2 + \frac{\sigma^2}{2\delta^2} \quad \text{and}$$

$$\mathbb{E}\left[\|G\|_2^2\right] = \|g\|_2^2 + O\left(\frac{N}{\delta^2}\right). \quad \Rightarrow \text{Var}(G) = O\left(\frac{d}{\delta^2}\right)$$

because $\mathbb{E}\left(\xi_i^+ - \xi_i^-\right) = 0$

$$\mathbb{E}G_i^2 = \mathbb{E}g_i^2 + \mathbb{E}\left(\frac{(\xi_i^+ - \xi_i^-)^2}{4\delta^2}\right) + 2\frac{g_i}{2\delta}\mathbb{E}\left(\xi_i^+ - \xi_i^-\right)$$

28

FDSA perturbed dimensions one-at-a-time, leading to $2N$ queries.
Can we reduce the number of queries?

Idea: Simultaneously randomly perturb all dimensions! (Spall, 1992)

## Function measurements

$$y_n^+ = f(\ \theta_n + \delta_n d_n\ ) + \xi_n^+, \quad y_n^- = f(\ \theta_n - \delta_n d_n\ ) + \xi_n^-$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

How to choose $d_n^i, i = 1, \dots, N$?

$-1$ — $1$

w.p. $\dfrac{1}{2}$     w.p. $\dfrac{1}{2}$

Only 2-queries, regardless of $N$!

$\mathbb{E}\left[G^i\right] = g^i!$  Hence, $\|\mathbb{E}[G] - \nabla f(\theta)\|_2 = O(\delta^2)$.

29

FDSA perturbed dimensions one-at-a-time, leading to 2*N* queries.
Can we reduce the number of queries?

Idea: Simultaneously randomly perturb all dimensions! (Spall, 1992)

## Function measurements

$$y_n^+ = f(\boxed{\theta_n + \delta_n d_n}) + \xi_n^+, \quad y_n^- = f(\boxed{\theta_n - \delta_n d_n}) + \xi_n^-$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

*→ Rademacher r.v.*

How to choose $d_n^i, i = 1, \ldots, N$?

-1 ——— 1

w.p. $\dfrac{1}{2}$      w.p. $\dfrac{1}{2}$

Only 2-queries, regardless of *N*!

$\mathbb{E}\left[G^i\right] = g^i!$  Hence, $\|\mathbb{E}[G] - \nabla f(\theta)\|_2 = O(\delta^2)$.

FDSA perturbed dimensions one-at-a-time, leading to $2N$ queries.
Can we reduce the number of queries?

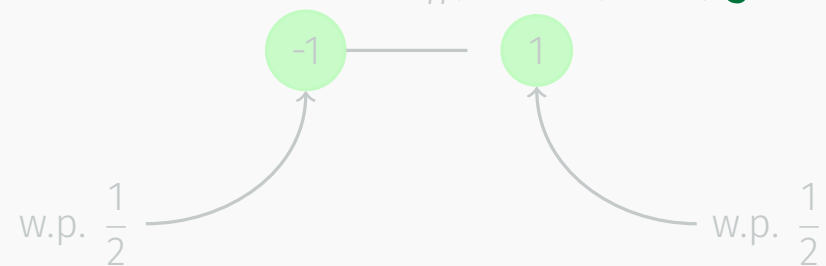Idea:  Simultaneously randomly perturb all dimensions! (Spall, 1992)

$\rightarrow d_n \rightarrow$ random vector

## Function measurements

$$y_n^+ = f(\,\theta_n + \delta_n d_n\,) + \xi_n^+, \quad y_n^- = f(\,\theta_n - \delta_n d_n\,) + \xi_n^-$$

## Gradient estimate

$$G_n^i = \left[ \frac{y_n^+ - y_n^-}{2\delta_n d_n^i} \right].$$

How to choose $d_n^i, i = 1, \ldots, N$?

Rademacher r.v.

-1 — 1

w.p. $\frac{1}{2}$      w.p. $\frac{1}{2}$

$\theta_{n+1} = \theta_n - a_n G_n$

## Only 2-queries, regardless of $d$

$$\mathbb{E}\left[G^i\right] = g^i! \quad \text{Hence,} \ \|\mathbb{E}[G] - \nabla f(\theta)\|_2 = O(\delta^2).$$

FDSA perturbed dimensions one-at-a-time, leading to $2N$ queries.
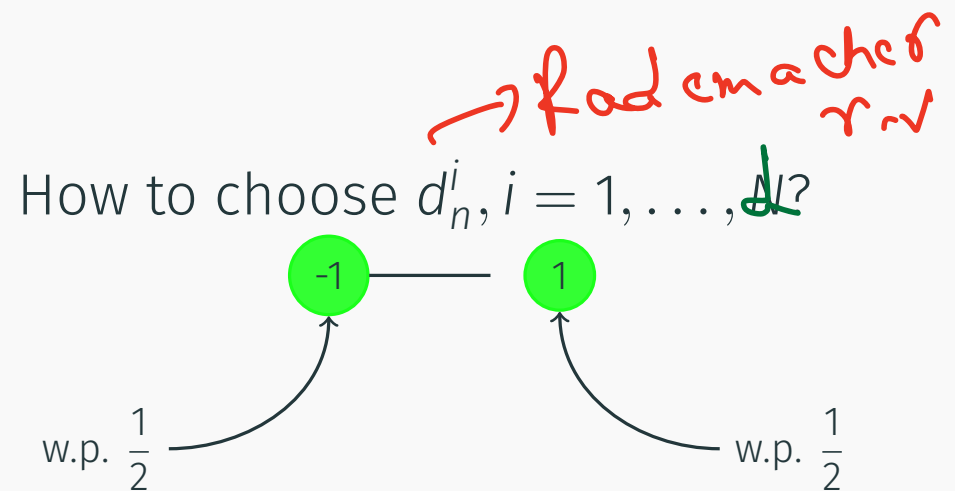Can we reduce the number of queries?

Idea: Simultaneously randomly perturb all dimensions! (Spall, 1992)
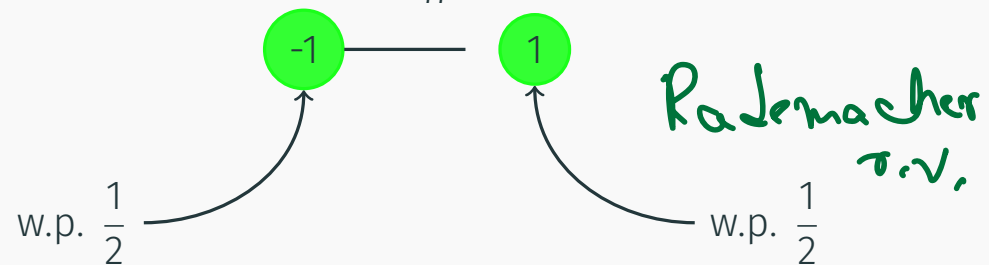
## Function measurements

$$y_n^+ = f(\boxed{\theta_n + \delta_n d_n}) + \xi_n^+, \quad y_n^- = f(\boxed{\theta_n - \delta_n d_n}) + \xi_n^-$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

How to choose $d_n^i, i = 1, \ldots, \cancel{N}$?

-1 — 1

w.p. $\frac{1}{2}$       w.p. $\frac{1}{2}$

## Only 2-queries, regardless of $\cancel{N}$

$$\mathbb{E}\left[G^i\right] = g^i! \quad \text{Hence, } \|\mathbb{E}[G] - \nabla f(\theta)\|_2 = O(\delta^2).$$

FDSA perturbed dimensions one-at-a-time, leading to $2N$ queries.
Can we reduce the number of queries?

Idea:  Simultaneously randomly perturb all dimensions! (Spall, 1992)
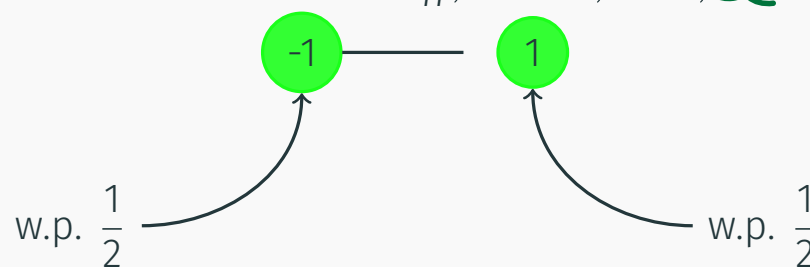
## Function measurements

$$y_n^+ = f(\ \theta_n + \delta_n d_n\ ) + \xi_n^+, \quad y_n^- = f(\ \theta_n - \delta_n d_n\ ) + \xi_n^-$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

How to choose $d_n^i, i = 1, \ldots, N$?

-1 —— 1

w.p. $\frac{1}{2}$          w.p. $\frac{1}{2}$

## Only 2-queries, regardless of $N$!

$$\mathbb{E}\left[G^i\right] = g^i! \ \ \text{Hence,} \ \|\mathbb{E}\left[G\right] - \nabla f(\theta)\|_2 = O(\delta^2).$$

29

## Function measurements

$$y_n^+ = f(\boxed{\theta_n + \delta_n d_n}) + \xi_n^+, \quad y_n^- = f(\boxed{\theta_n - \delta_n d_n}) + \xi_n^-.$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

## Taylor series expansions

$$f(\theta_n \pm \delta_n d_n) = f(\theta_n) \pm \delta_n d_n^\top \nabla f(\theta_n) + \frac{\delta_n^2}{2} d_n^\top \nabla^2 f(\theta_n) d_n + O(\delta_n^3)$$

$$\frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2\delta_n d_n^i} = \nabla_i f(\theta_n) + \boxed{\sum_{j=1, j\neq i}^{N} \frac{d_n^j}{d_n^i} \nabla_j f(\theta_n)} + O(\delta_n^2)$$

zero-mean since $d_n$ symmetric Bernoulli $\pm 1$ r.v.s

Hence, $\left\|\mathbb{E}\left[G^i\right] - \nabla f(\theta_n)\right\|_2 = O(\delta_n^2)$.
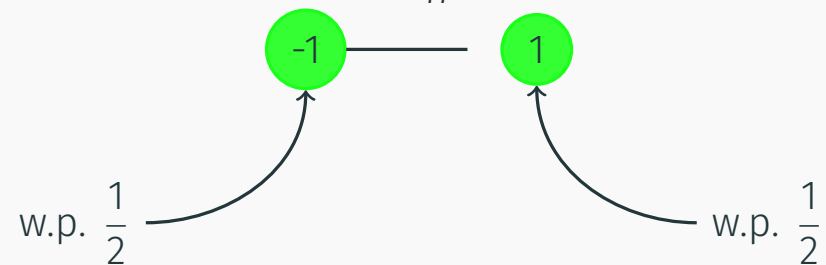
## Function measurements

$$y_n^+ = f(\boxed{\theta_n + \delta_n d_n}) + \xi_n^+, \quad y_n^- = f(\boxed{\theta_n - \delta_n d_n}) + \xi_n^-.$$

## Gradient estimate

$$G^i = \left[\frac{y_n^+ - y_n^-}{2\delta_n d_n^i}\right].$$

$$\mathbb{E}\, G^i = \mathbb{E}\left[\frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2\delta_n d_n^i}\right]$$

## Taylor series expansions

$$f(\theta_n \pm \delta_n d_n) = f(\theta_n) \pm \delta_n d_n^\top \nabla f(\theta_n) + \frac{\delta_n^2}{2} d_n^\top \nabla^2 f(\theta_n) d_n + O(\delta_n^3)$$

$$\frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2\delta_n d_n^i} = \nabla_i f(\theta_n) + \boxed{\sum_{j=1, j\neq i}^N \frac{d_n^j}{d_n^i} \nabla_j f(\theta_n)} + O(\delta_n^2)$$

zero-mean since $d_n$ symmetric Bernoulli $\pm 1$ r.v.s

Hence, $\left\| \mathbb{E}\left[G^i\right] - \nabla f(\theta_n) \right\|_2 = O(\delta_n^2).$

$$f(\theta_n \pm \delta_n d_n) = f(\theta_n) \pm \delta_n d_n^\top \nabla f(\theta_n) + \frac{\delta_n^2}{2} d_n^\top \nabla^2 f(\theta_n) d_n + O(\delta_n^3)$$

$$f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n) = 2\delta_n \, d_n^\top \nabla f(\theta_n) + O(\delta_n^3)$$

$$\frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2 \delta_n d_n^i}$$

$$= \sum_{\partial = 1}^{d} \frac{d_n^\partial \, \nabla_\partial f(\theta_n)}{d_n^i} + O(\delta_n^2)$$

$$= \nabla_i f(\theta_n) + \sum_{\partial \neq i} \frac{d_n^\partial}{d_n^i} \nabla_\partial f(\theta_n) + O(\delta_n^2)$$

<span style="color:blue">} treat $\theta_n$ as a constant in this computation</span>

<span style="color:red">$$\mathbb{E}\left( \frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2 \delta_n d_n^i} \right)$$</span>

<span style="color:red">$$= \nabla_i f(\theta_n) + \sum_{\partial \neq i} \mathbb{E}\left( \frac{d_n^\partial}{d_n^i} \right) \nabla_\partial f(\theta_n) + O(\delta_n^2)$$</span>

<span style="color:blue">zero in expectation</span>

<span style="color:red">$$= \nabla_i f(\theta_n) + O(\delta_n^2)$$</span>

<span style="color:blue">Note!
$$d_n = (d_n^1 \cdots d_n^d)$$
independent Rademacher</span>

## Function measurements

$$y_n^+ = f(\;\theta_n + \delta_n d_n\;) + \xi_n^+, \quad y_n^- = f(\;\theta_n - \delta_n d_n\;) + \xi_n^-.$$

## Gradient estimate

$$G^i = \left[ \frac{y_n^+ - y_n^-}{2\delta_n d_n^i} \right].$$

## Taylor series expansions

$$f(\theta_n \pm \delta_n d_n) = f(\theta_n) \pm \delta_n d_n^\top \nabla f(\theta_n) + \frac{\delta_n^2}{2} d_n^\top \nabla^2 f(\theta_n) d_n + O(\delta_n^3)$$

$$\frac{f(\theta_n + \delta_n d_n) - f(\theta_n - \delta_n d_n)}{2\delta_n d_n^i} = \nabla_i f(\theta_n) + \boxed{\sum_{j=1, j\neq i}^{N} \frac{d_n^j}{d_n^i} \nabla_j f(\theta_n)} + O(\delta_n^2)$$

zero-mean since $d_n$ symmetric Bernoulli $\pm 1$ r.v.s

Hence, $\left\| \mathbb{E}\left[ G \right] - \nabla f(\theta_n) \right\|_2 = O(\delta_n^2).$

Consider the following estimator:

$$G = \Delta \frac{\left[\left(f(\theta + \delta\Delta) + \xi^+\right) - \left(f(\theta - \delta\Delta) + \xi^-\right)\right]}{2\delta}$$

$$\Delta = (\Delta_1, --- \Delta_d)^T \text{ independent}$$

$$\Delta_i \sim N(0,1)$$

Is $G$ a good enough estimator of $\nabla f(\theta)$?

Yes. Gaussian Smoothed functional aka Gaussian Smoothing

$$G = \frac{(f(\theta + U) + \xi^+) - (f(\theta - U) + \xi^-)}{2\delta} V.$$

Choose $U, V$ such that $\mathbb{E}\left[VU^\top\right] = I$, $\mathbb{E}[V] = 0$.

One-point estimate!

$$G = \frac{(f(\theta + U) + \xi^+)}{\delta} V.$$

Choose $U, V$ such that $\mathbb{E}\left[VU^\top\right] = I$, $\mathbb{E}[V] = 0$. Works??

$$\mathbb{E}[G] = \mathbb{E}\left[G - \frac{f(\theta)}{\delta} V\right] = \mathbb{E}\left[\frac{(f(\theta + U) + \xi^+) - f(\theta)}{\delta} V\right].$$

# "Mother" of all two-point sim-pert estimates

$$G = \frac{(f(\theta + U) + \xi^+) - (f(\theta - U) + \xi^-)}{2\delta} V.$$

Choose $U, V$ such that $\mathbb{E}\left[VU^\top\right] = I$, $\mathbb{E}\left[V\right] = 0$.

One-point estimate!

$$G = \frac{(f(\theta + U) + \xi^+)}{\delta} V.$$

Choose $U, V$ such that $\mathbb{E}\left[VU^\top\right] = I$, $\mathbb{E}\left[V\right] = 0$. Works??

$$\mathbb{E}\left[G\right] = \mathbb{E}\left[G - \frac{f(\theta)}{\delta} V\right] = \mathbb{E}\left[\frac{(f(\theta + U) + \xi^+) - f(\theta)}{\delta} V\right].$$

- $U \sim \delta \mathcal{N}(0, I)$, $V = \delta^{-1} U$
  - Smoothed functional by Katkovnik and Kulchitsky (1972);
  - Refined by Polyak and Tsybakov (1990); also studied by Dippon (2003); Nesterov and Spokoiny (2011).
- $U \sim \delta \operatorname{Unif}(\mathbb{S}_N)$, $V = N\delta^{-1} U$
  - RDSA by Kushner and Clark (1978); Enhanced by Prashanth et al. (2017)
  - Rediscovered by Flaxman et al. (2005)
- $U_i \sim \delta \operatorname{Rademacher}(\pm 1)$, $V = \delta^{-1} U$
  - SPSA by Spall (1992).
- Deterministic perturbations by Bhatnagar et al. (2003)
- ...

Does it matter which of these we select? Not really:
Bias is always $O(\delta^2)$, while variance is $O(1)$ or $O(\delta^{-2})$ (noise controlled or not)

- $U \sim \delta \mathcal{N}(0, I)$, $V = \delta^{-1} U$
  - Smoothed functional by Katkovnik and Kulchitsky (1972);
  - Refined by Polyak and Tsybakov (1990); also studied by Dippon (2003); Nesterov and Spokoiny (2011).
- $U \sim \delta \operatorname{Unif}(\mathbb{S}_N)$, $V = N\delta^{-1} U$
  - RDSA by Kushner and Clark (1978); Enhanced by Prashanth et al. (2017)
  - Rediscovered by Flaxman et al. (2005)
- $U_i \sim \delta \operatorname{Rademacher}(\pm 1)$, $V = \delta^{-1} U$
  - SPSA by Spall (1992).
- Deterministic perturbations by Bhatnagar et al. (2003)
- ...

Does it matter which of these we select? Not really:
Bias is always $O(\delta^2)$, while variance is $O(1)$ or $O(\delta^{-2})$ (noise controlled or not)

- $U \sim \delta \mathcal{N}(0, I)$, $V = \delta^{-1} U$
  - Smoothed functional by Katkovnik and Kulchitsky (1972);
  - Refined by Polyak and Tsybakov (1990); also studied by Dippon (2003); Nesterov and Spokoiny (2011).
- $U \sim \delta \operatorname{Unif}(\mathbb{S}_N)$, $V = N\delta^{-1} U$
  - RDSA by Kushner and Clark (1978); Enhanced by Prashanth et al. (2017)
  - Rediscovered by Flaxman et al. (2005)
- $U_i \sim \delta \operatorname{Rademacher}(\pm 1)$, $V = \delta^{-1} U$
  - SPSA by Spall (1992).
- Deterministic perturbations by Bhatnagar et al. (2003)
- ...

Does it matter which of these we select? Not really:
Bias is always $O(\delta^2)$, while variance is $O(1)$ or $O(\delta^{-2})$ (noise controlled or not)

$$y^+ = f(\theta + \delta U) + \xi^+, \text{ and } y^- = f(\theta - \delta U) + \xi^-$$

$$G = \frac{(y^+ - y^-) V}{2\delta}$$

## Assumption:

**A2.1.** Let $U, V$ be random $N$-vectors satisfying $\mathbb{E}\left[VU^\top\right] = I, \mathbb{E}[V] = 0,$ and $\mathbb{E}\left[\|V\| \|U\|^3\right] < \infty.$

**A2.2.** The noise factors $\xi^\pm$ in (2.6) satisfy

$$\mathbb{E}[\xi^+ - \xi^- | U, V] = 0, \quad \text{and} \quad \mathbb{E}[(\xi^+ - \xi^-)^2 | U, V] \leq \sigma^2 < \infty. \quad (2.7)$$

**A2.3.** The objective $f$ satisfies

$$\sup_{\theta \in \mathbb{R}^{dk}} \mathbb{E}[f(\theta \pm \delta U)^2] \leq B < \infty. \quad (2.8)$$

A2.4   $f$ is three-times continuously differentiable with $\left| \nabla^3_{i_1 i_2 i_3} f(\theta) \right| < \bar{B} < \infty$  $\forall i_1, i_2, i_3.$

Claim: Assume A2.1 − 2.4. Then,

$$\| \mathbb{E}\, G - \nabla f(\theta) \| \leq C_1 \delta^2 \text{ and}$$

$$\mathbb{E}\left[\| G - \mathbb{E}\, G \|^2\right] \leq \frac{C_2}{\delta^2}$$

**Proof:**

$$\mathbb{E}\, G = \mathbb{E}\left[ V\left( \frac{f(\theta + \delta v) - f(\theta - \delta v)}{2\delta} \right) \right]$$

$$\left(\text{since} \quad \mathbb{E}\left( V\left( \frac{\xi^+ - \xi^-}{2\delta} \right) \right) = \frac{1}{2\delta} \mathbb{E}\left[ \mathbb{E}\left[ V(\xi^+ - \xi^-) \mid v \right] \right] = \mathbb{E}\left[ v\, \mathbb{E}\left( \frac{\xi^+ - \xi^-}{2\delta} \mid v \right) \right]\right.$$

$$= 0$$

Using Taylor series expansions,

$$f(\theta \pm \delta v) = f(\theta) \pm \delta\, v^T \nabla f(\theta) + \frac{\delta^2}{2} v^T \nabla^2 f(\theta) v + O(\delta^3)$$

$$V\left( \frac{f(\theta + \delta v) - f(\theta - \delta v)}{2\delta} \right) = V v^T \nabla f(\theta) + O(\delta^2)$$

Taking expectations,

using $\mathbb{E}(V V^T) = I$, we get

$$\| \mathbb{E} G - \nabla f(\theta) \| \leq C_1 \delta^2$$
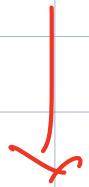
**Second claim proof:**

$$E\|G - EG\|^2 \le 4\, E\|G\|^2$$

$$= 4\, E\left[\|V\|^2 \left\{\left(\frac{\xi^+ - \xi^-}{2\delta}\right)^2 + 2\underbrace{\left(\frac{\xi^+ - \xi^-}{2\delta}\right)}_{\text{near zero}}\left(\frac{f(\theta+\delta U)-f(\theta-\delta U)}{2\delta}\right)\right.\right.$$

$$\left.\left. + \left(\frac{f(\theta+\delta U)-f(\theta-\delta U)}{2\delta}\right)^2\right\}\right]$$

$$= 4\, E\left[\frac{\|V\|^2 (\xi^+ - \xi^-)^2}{4\delta^2}\right] + 4\, E\left[\frac{\|V\|^2 \left(f(\theta+\delta U)-f(\theta-\delta U)\right)^2}{4\delta^2}\right]$$

$$\le \frac{C_2}{\delta^2}$$

$\longrightarrow$

$$\text{Why (i)}\quad E\, f(\theta \pm \delta U)^2 \le B_1 < \infty$$

$$\text{(ii)}\quad E\|V\|^2 \le B_2 < \infty$$

For performing gradient descent:

$$\theta_{n+1} = \theta_n - a_n G_n,$$

we can construct nearly unbiased gradient estimate $G_n$ using simultaneous perturbation trick

*to be handled later.*

| Noise → Gradient estimate ↓ | Controlled | Uncontrolled |
|---|---|---|
| Bias | $C_1 \delta$ | $C_1 \delta^2$ |
| Variance | $C_2$ | $\dfrac{C_2}{\delta^2}$ |

This assumed $f \in \mathcal{C}^3$. Holds also for $f$ convex, smooth.

# A few answers so far...

**Q1)** How to form $G_n$ from function samples so that
$G_n \approx \nabla f(\theta_n)$

Use simultaneous perturbation trick

**Q2)** Such a $G_n$ - is it unbiased?

Almost ... what we get is an asymptotically unbiased estimate?

**Q3)** Does $\theta_{n+1} = \theta_n - a_n G_n$ converge to $\theta^*$ with such a $G_n$?

??

**Q4)** If answer is yes to above, what is the convergence rate?

??

# A few answers so far...

**Q1)** How to form $G_n$ from function samples so that
$G_n \approx \nabla f(\theta_n)$

Use simultaneous perturbation trick

**Q2)** Such a $G_n$ - is it unbiased?

Almost ... what we get is an asymptotically
unbiased estimate? *if we let $\delta_n \to 0$ since bias is $O(\delta_n^2)$*

**Q3)** Does $\theta_{n+1} = \theta_n - a_n G_n$ converge to $\theta^*$ with such a
$G_n$?

??

**Q4)** If answer is yes to above, what is the convergence
rate?

??

**Q1)** How to form $G_n$ from function samples so that $G_n \approx \nabla f(\theta_n)$

Use simultaneous perturbation trick

**Q2)** Such a $G_n$ - is it unbiased?

Almost ... what we get is an asymptotically unbiased estimate?

**Q3)** Does $\theta_{n+1} = \theta_n - a_n G_n$ converge to $\theta^*$ with such a $G_n$?

$\rightarrow$ Stationary point $\nabla f(\theta^*) = 0$

??

**Q4)** If answer is yes to above, what is the convergence rate?

??

Can I bound

$E \, \|\theta_n - \theta^*\|^2$ or

$E \, (f(\theta_n) - f(\theta^*))$ or $E \, \|\nabla f(\theta_n)\|^2$

34

$\rightarrow$ To be covered at a later point in the course after introducing the necessary background on Stochastic approximation