

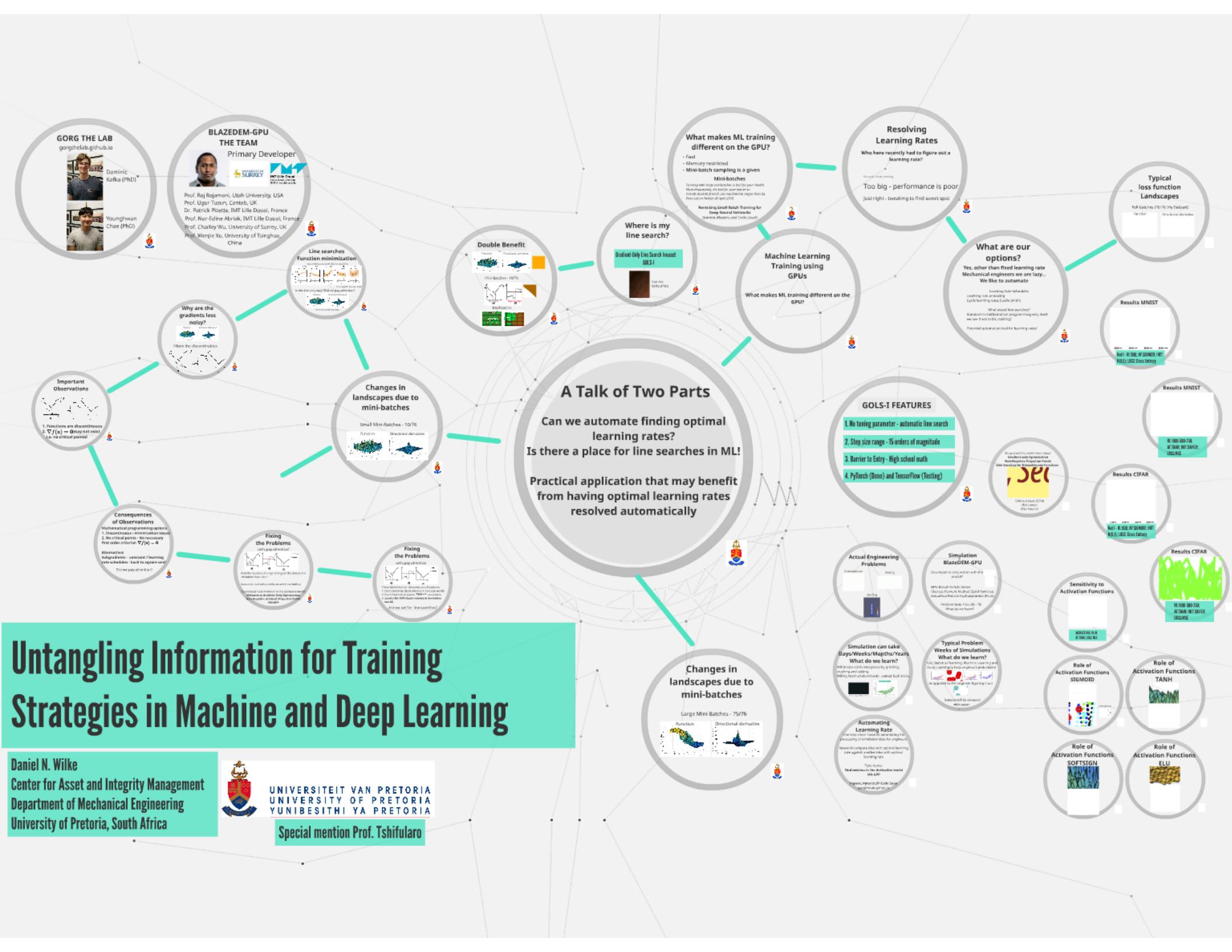
Untangling Information for Training Strategies in Machine and Deep Learning

Daniel N. Wilke
Center for Asset and Integrity Management
Department of Mechanical Engineering
University of Pretoria, South Africa



UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

Special mention Prof. Tshifularo



Untangling Information for Training Strategies in Machine and Deep Learning

Daniel N. Wilke
Center for Asset and Integrity Management
Department of Mechanical Engineering
University of Pretoria, South Africa



UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

Special mention Prof. Tshifularo

learning
Is there a place for
Practical applications
from having optimised
resolved au

2. $\nabla f(x) = 0$ may not exist
i.e. no critical points!

Consequences of Observations

Mathematical programming options
1. Discontinuous - minimization issues
2. No critical points - No necessary first order criterion $\nabla f(x) = 0$

Alternatives
Subgradients - constant / learning rate schedules - back to square one!
Did we pay attention?

Fixing the Problems

Let's pay attention!
Find the location of a sign change in the directional derivative from - to +
Descent (- derivative) while ascent (+ derivative)

Guaranteed local minimum in the derivative world!
Welcome to Gradient Descent Optimization
Non-Negative Gradient Projection Points
NN-GPP

Fixing the Problems

Let's pay attention!
Three Guidelines for Discontinuous Functions:
1. Don't minimize (local minima in function world)
2. Locate and critical points $\nabla f(x) \rightarrow$ candidates
3. Locate NN-GPP (local minima in derivative world)
Are we set for line searches?

Untangling Information for Training Strategies in Machine and Deep Learning

Daniel N. Wilke

Center for Asset and Integrity Management
Department of Mechanical Engineering
University of Pretoria, South Africa



UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

Special mention Prof. Tshifularo

A Talk of Two Parts

- Can we automate finding optimal learning rates?

Is there a place for line searches in ML!

Practical application that may benefit from having optimal learning rates resolved automatically



Machine Learning Training using GPUs

**What makes ML training different on the
GPU?**



What makes ML training different on the GPU?

- Fast
- Memory restricted
- **Mini-batch sampling is a given**

Mini-batches

Training with large minibatches is bad for your health.

More importantly, it's bad for your test error.

Friends dont let friends use minibatches larger than 32.

Yann LeCun Twitter 26 April 2018

Revisiting Small Batch Training for Deep Neural Networks

Dominic Masters and Carlo Luschi



Resolving Learning Rates

Who here recently had to figure out a learning rate?

Too small - it takes too long

Too big - performance is poor

Just right - tweaking to find sweet spot



What are our options?

**Yes, other than fixed learning rate
Mechanical engineers we are lazy...
We like to automate**

Learning Rate Schedules
Learning rate annealing
Cyclic learning rates (Leslie Smith)

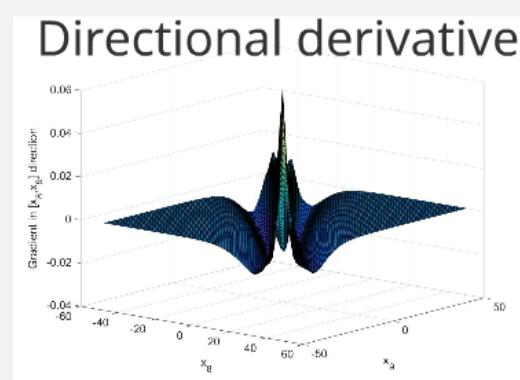
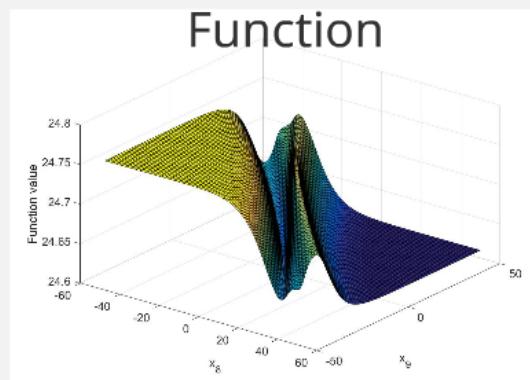
What about line searches?
Standard in mathematical programming why don't
we see them in ML training?

Potential automation tool for learning rates!

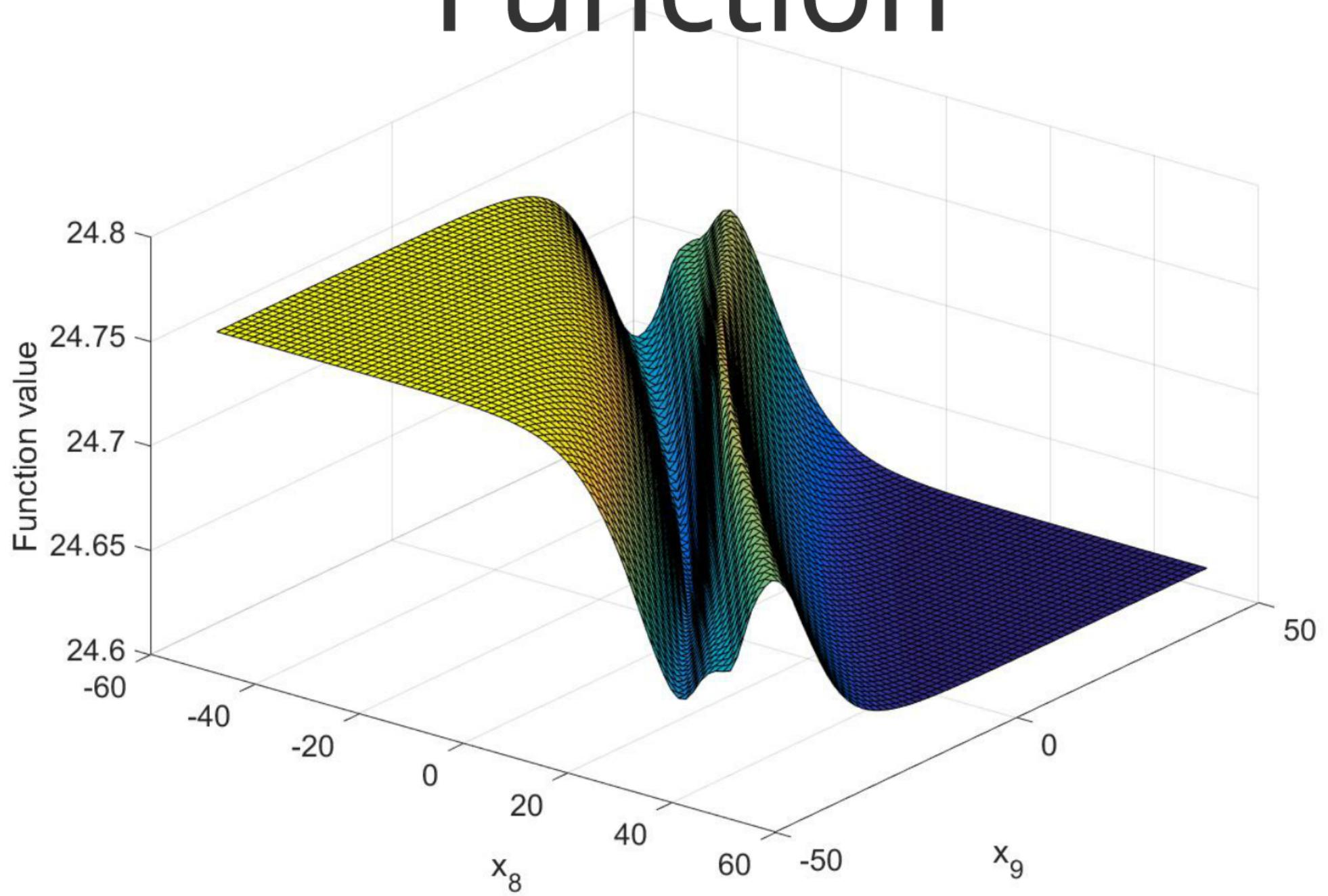


Typical loss function Landscapes

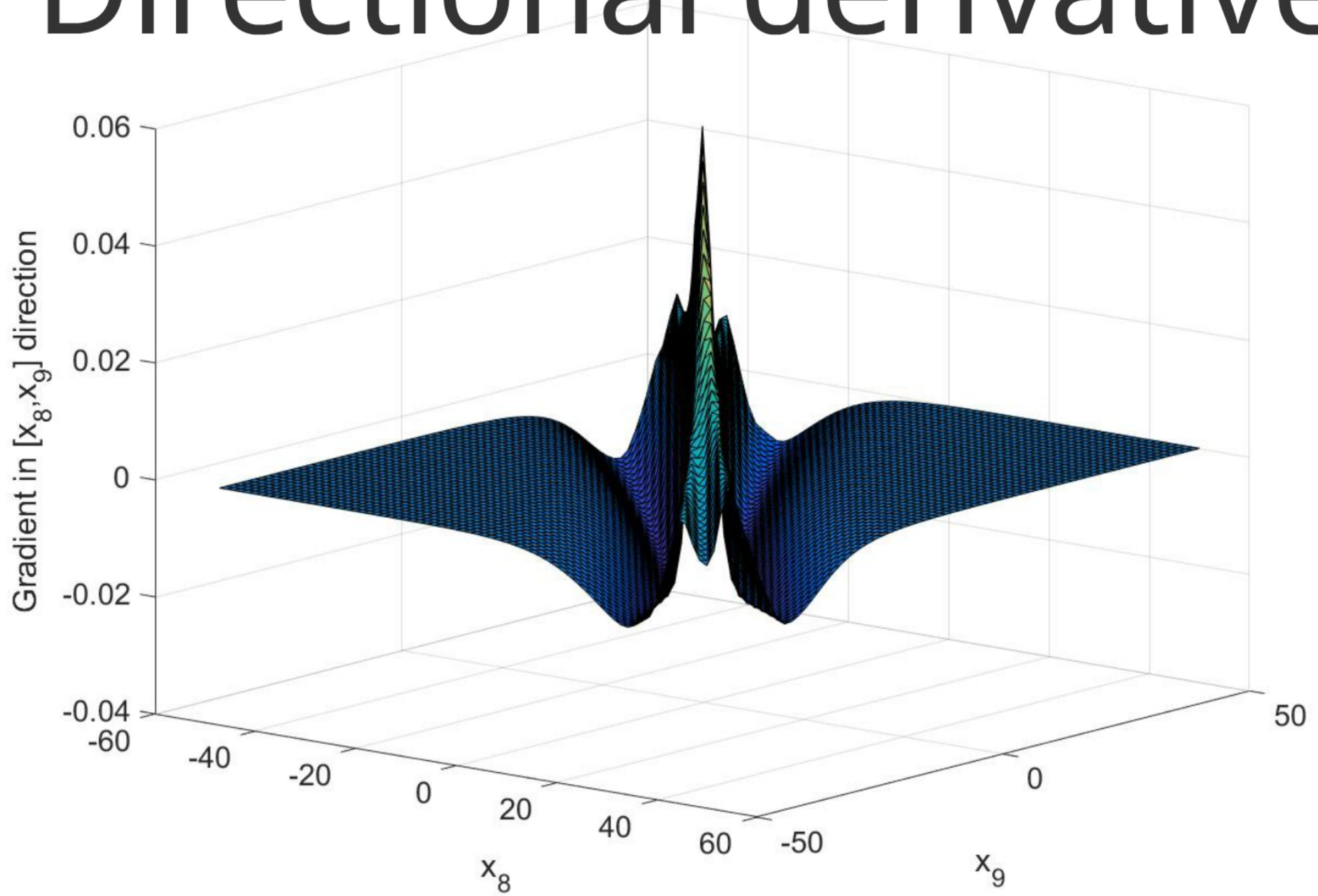
Full-batches (76/76 Iris Dataset)



Function

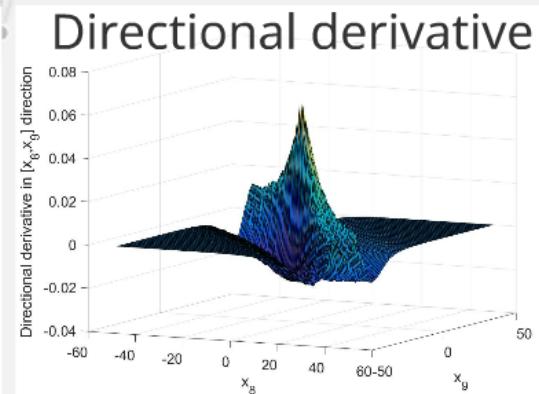
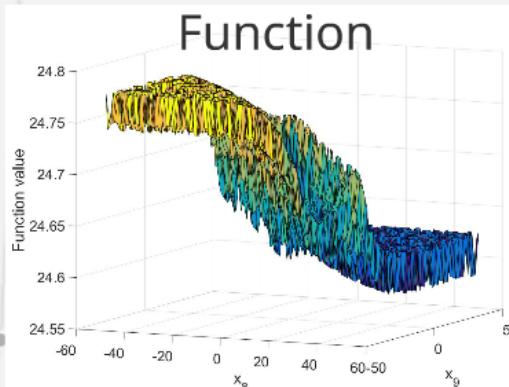


Directional derivative



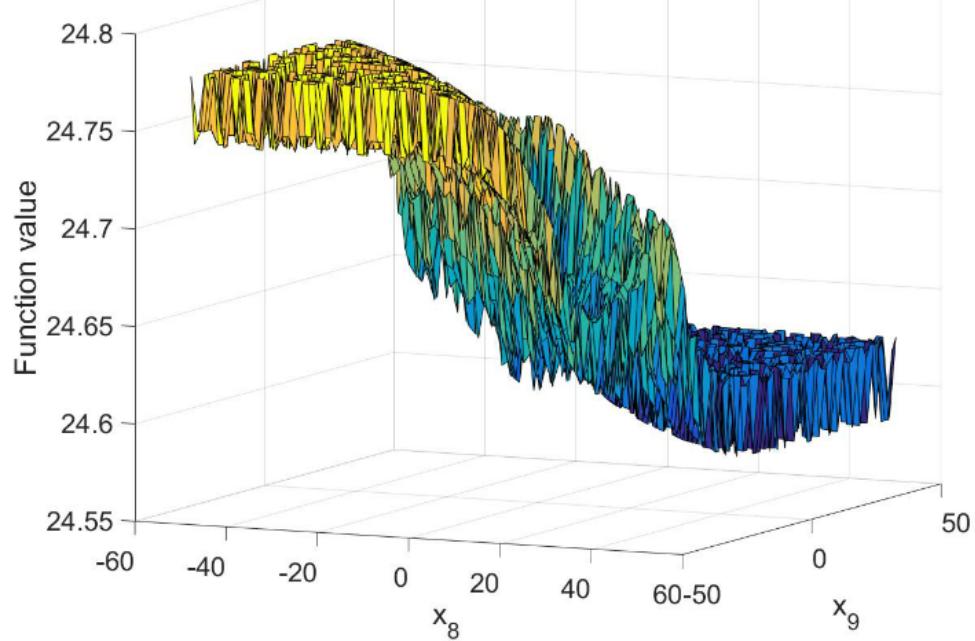
Changes in landscapes due to mini-batches

Large Mini-Batches - 75/76

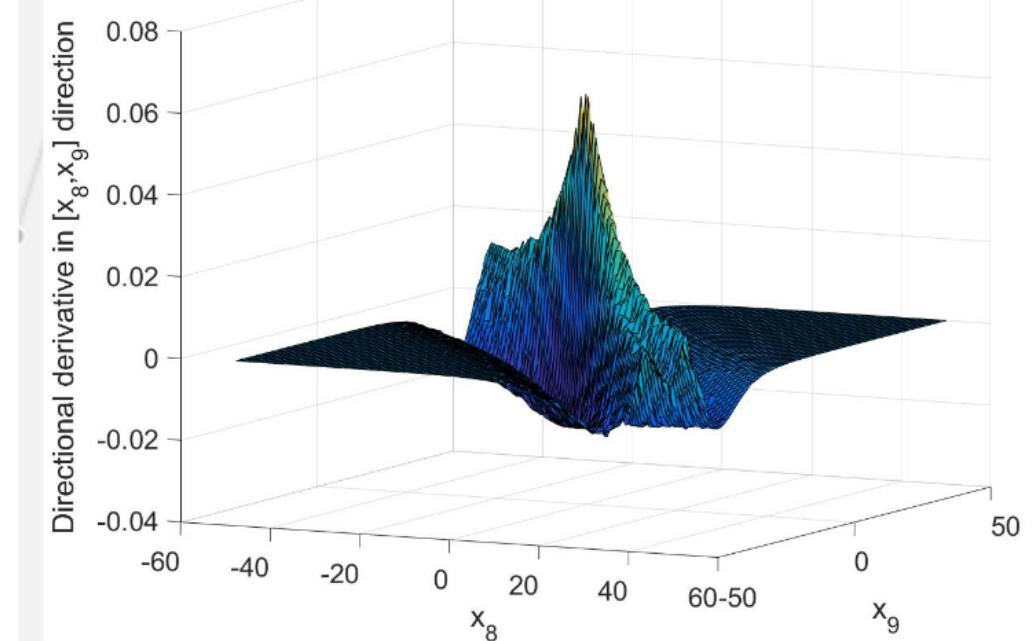


Large Mini-Batches - 75/76

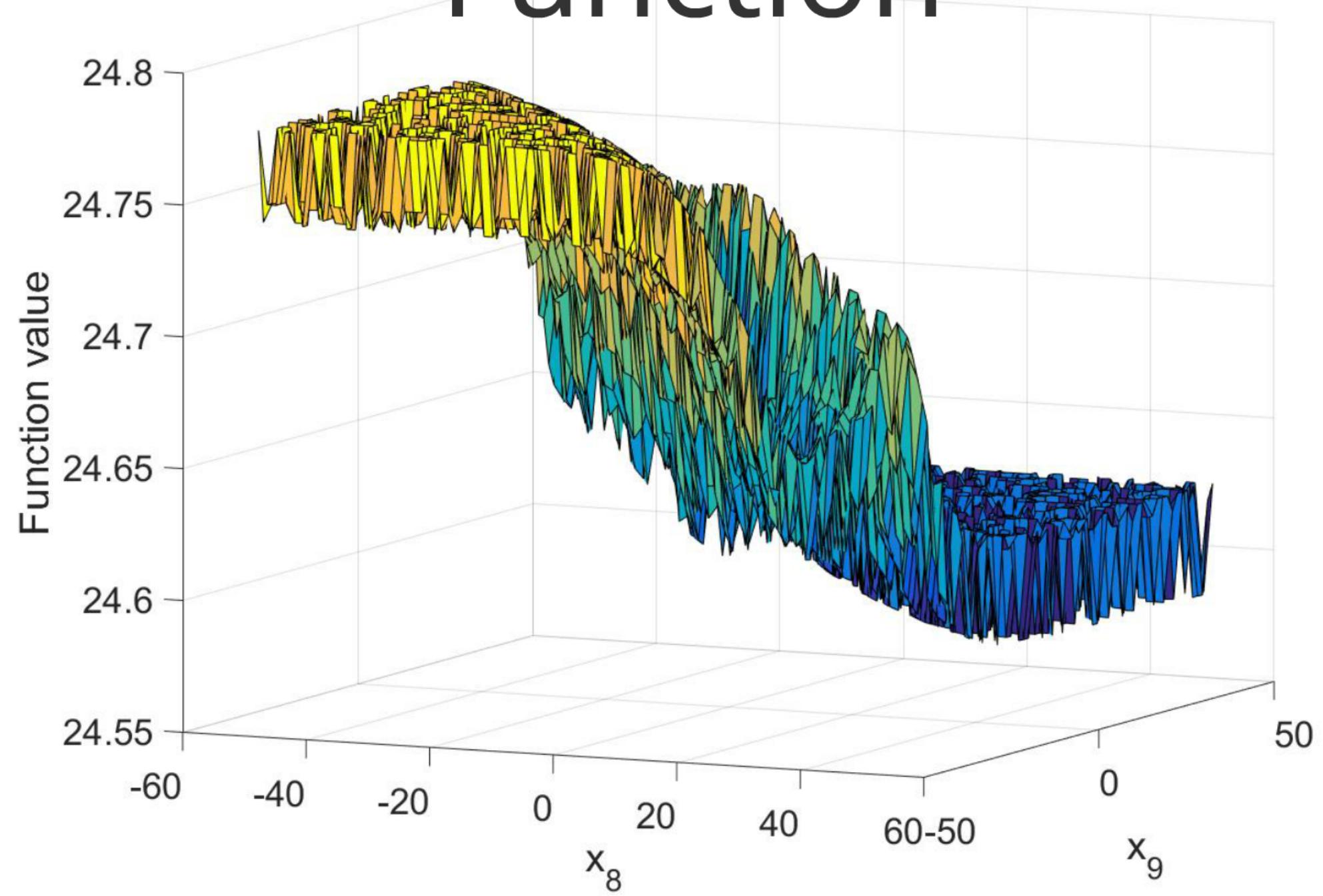
Function



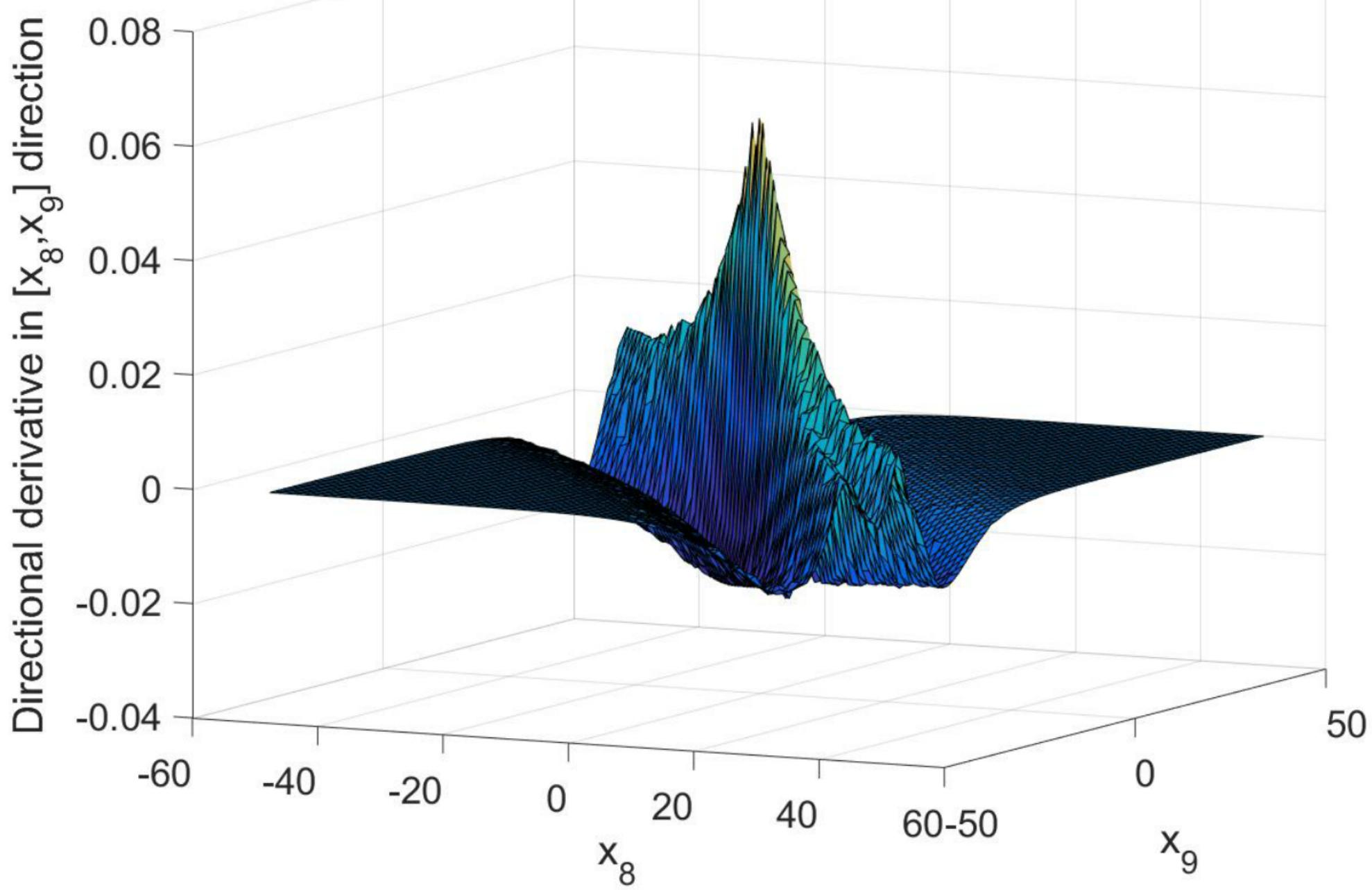
Directional derivative



Function

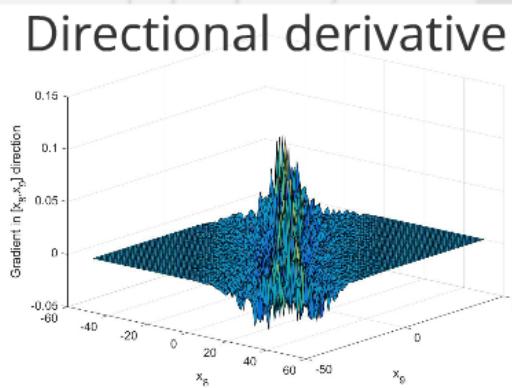
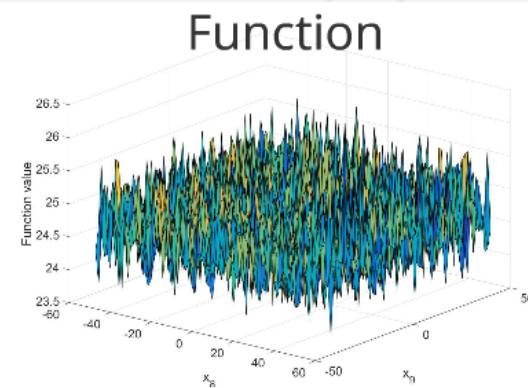


Directional derivative

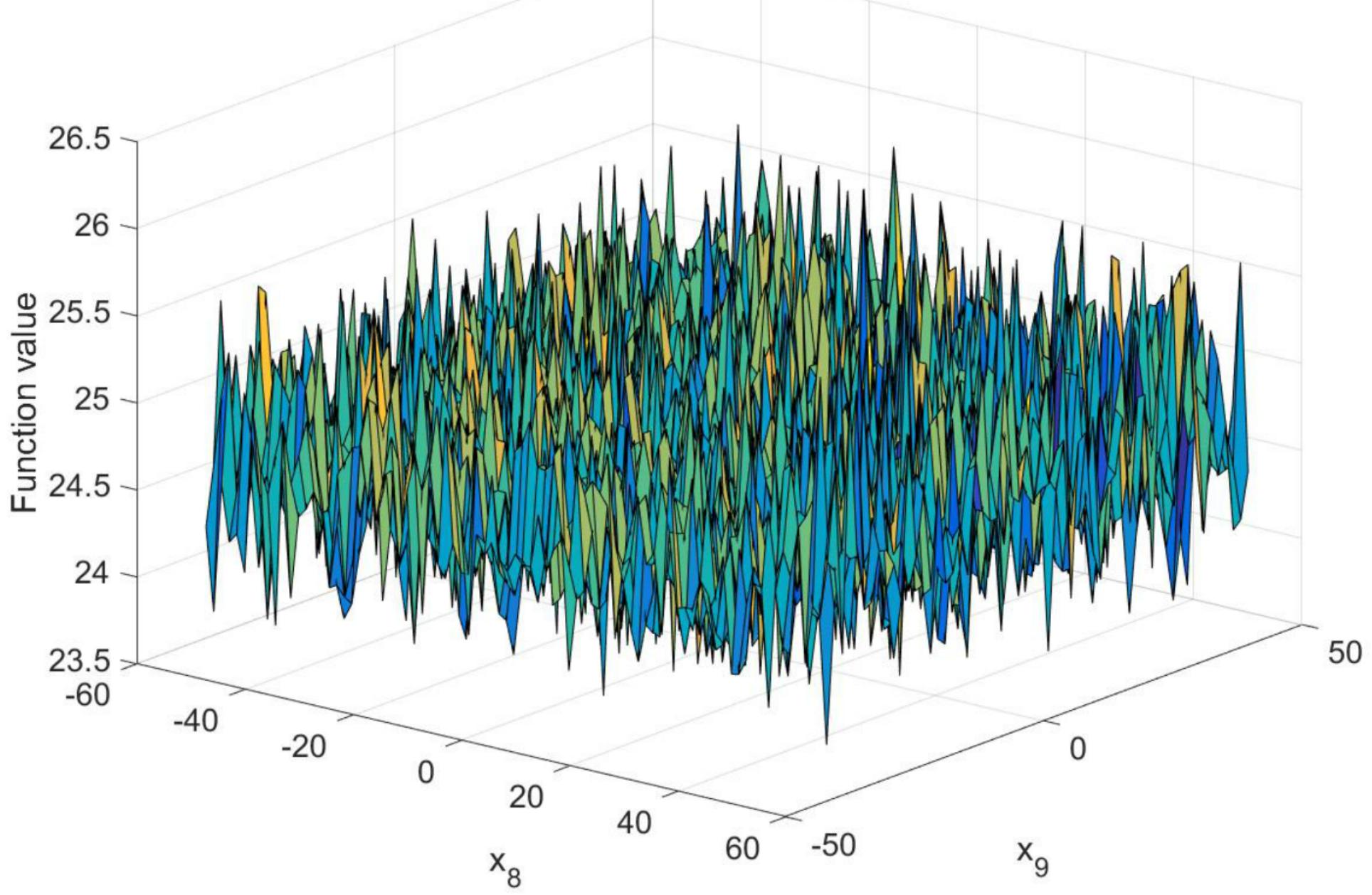


Changes in landscapes due to mini-batches

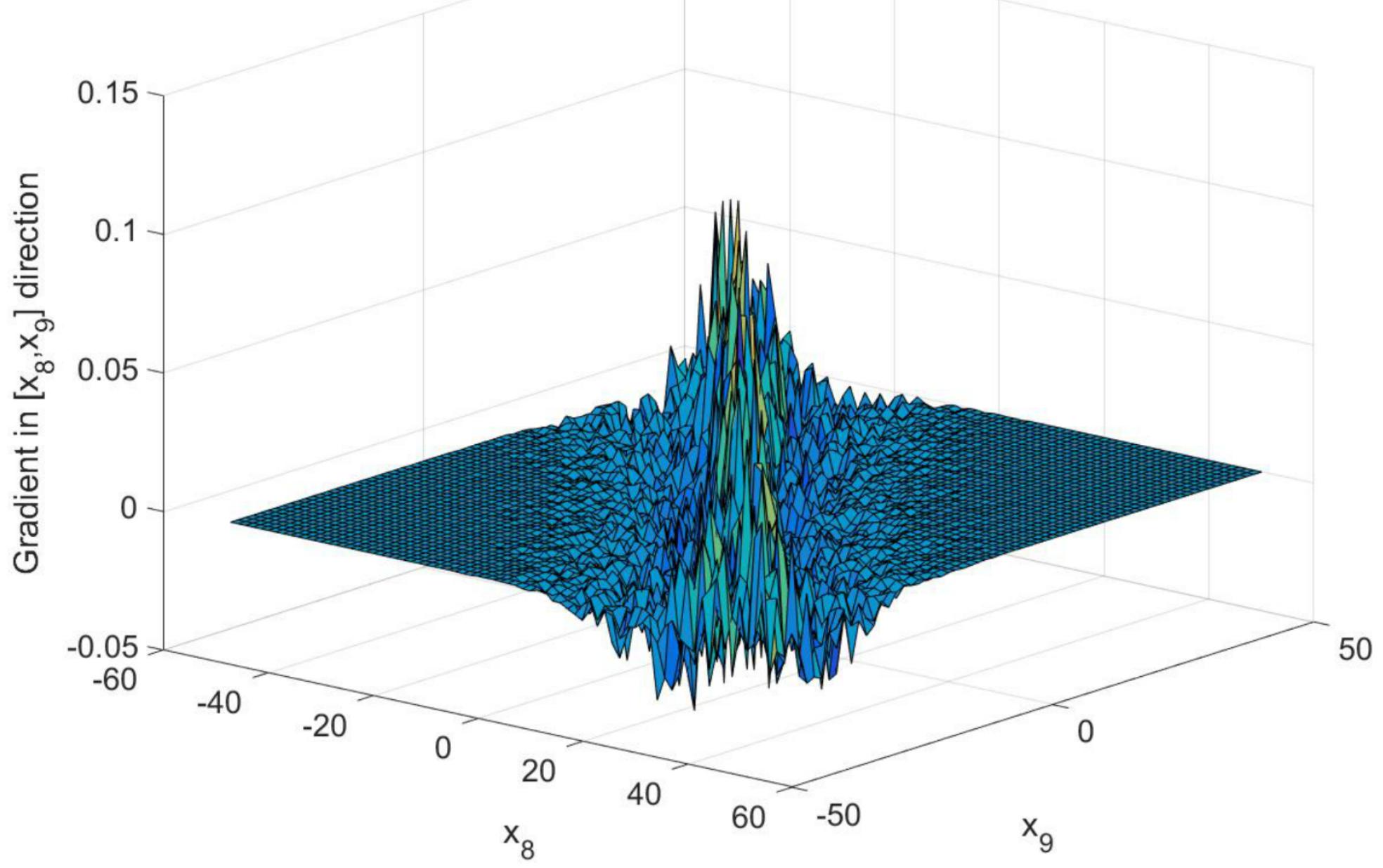
Small Mini-Batches - 10/76



Function



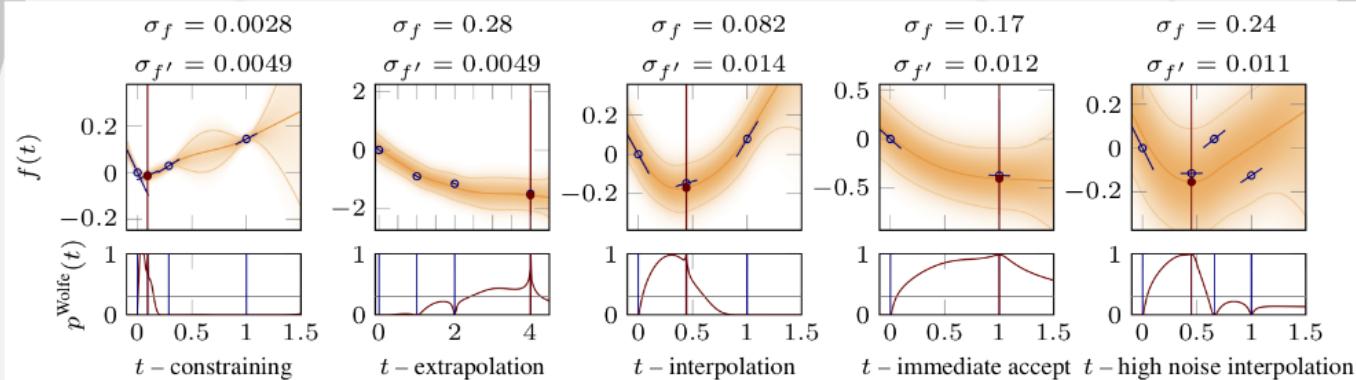
Directional derivative



Line searches

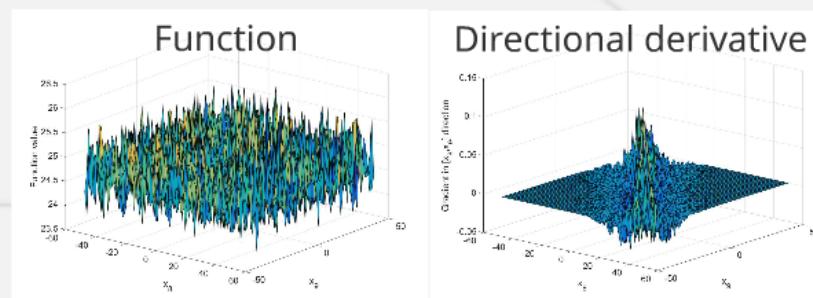
Function minimization

Maren Mahsereci and Phillip Hennig (2015)



arxiv.org/pdf/1502.02846.pdf

Is this the only way? Did we pay attention?



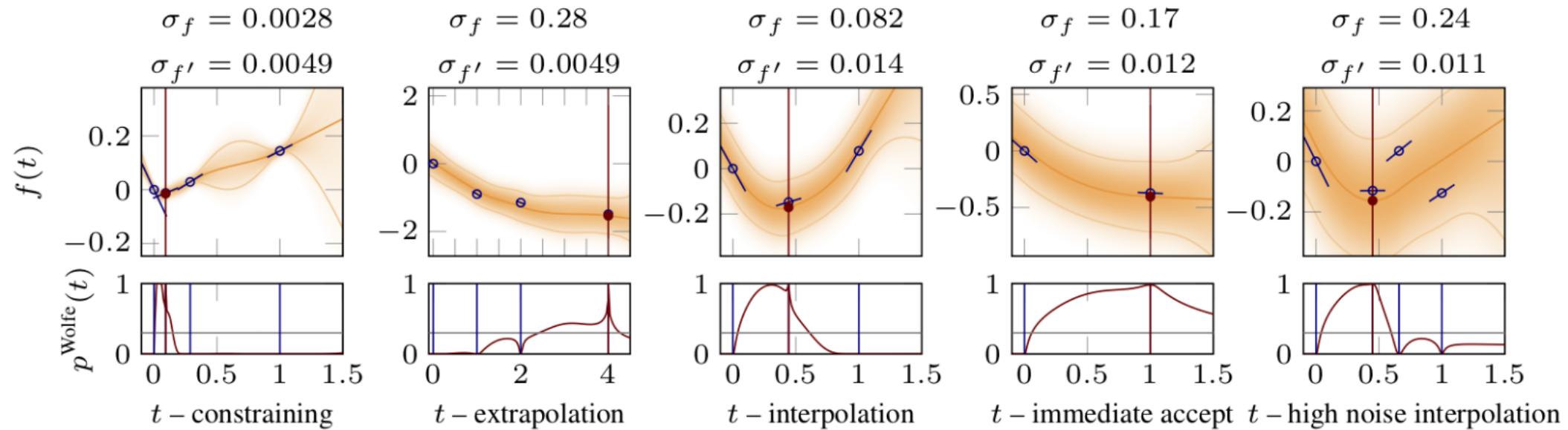
Which side looks less noisy?



Line searches

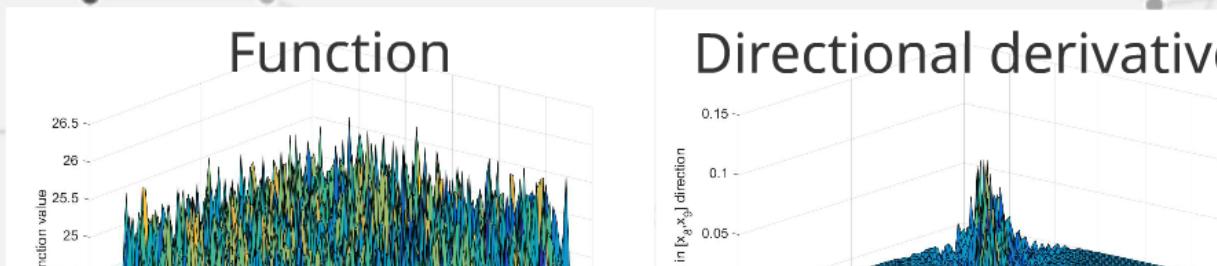
Function minimization

Maren Mahsereci and Phillip Hennig (2015)



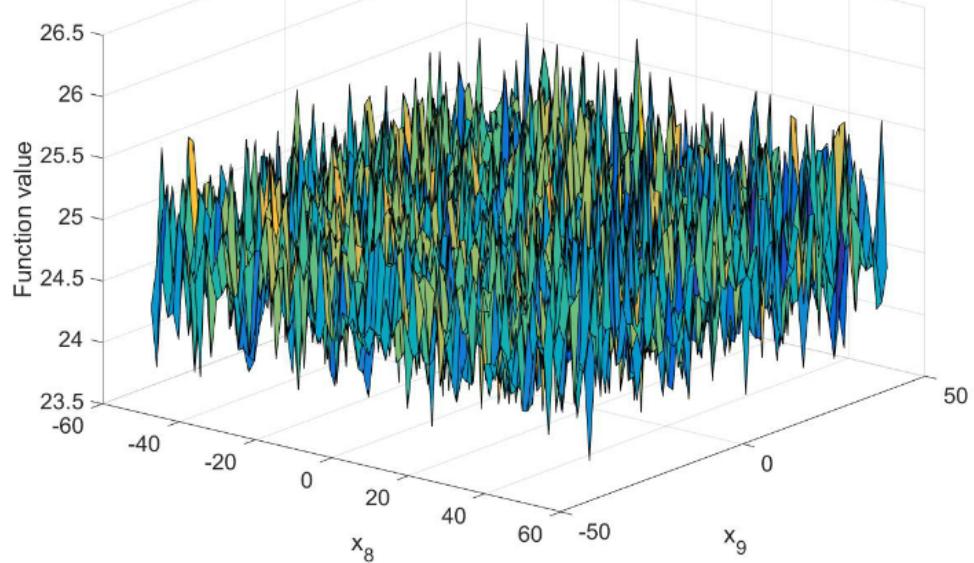
arxiv.org/pdf/1502.02846.pdf

Is this the only way? Did we pay attention?

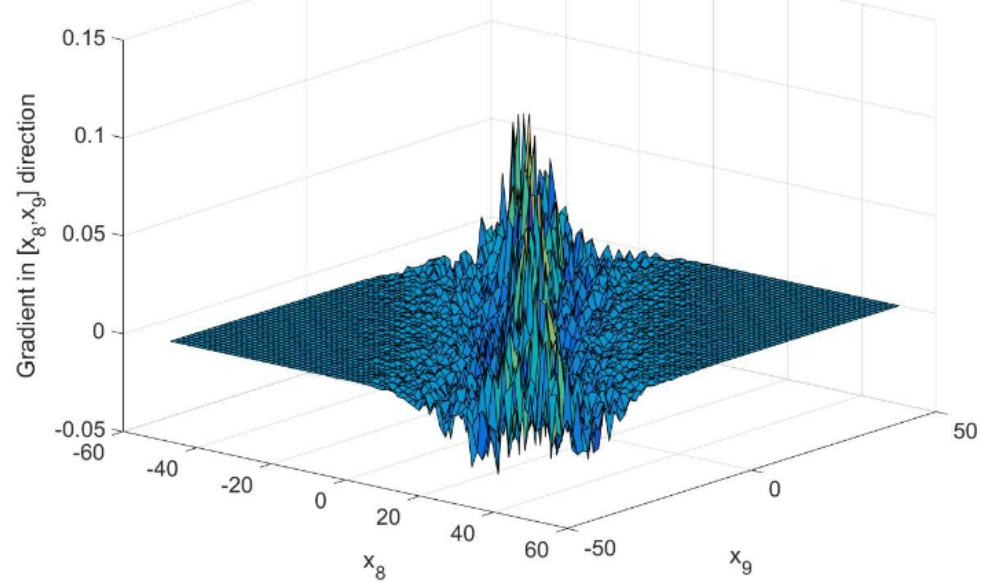


the only way? Did we pay attention?

Function

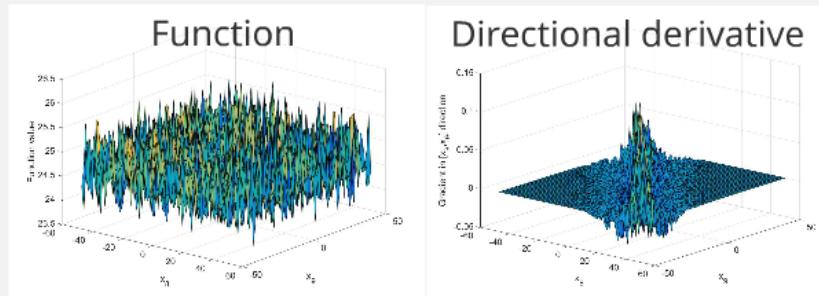


Directional derivative

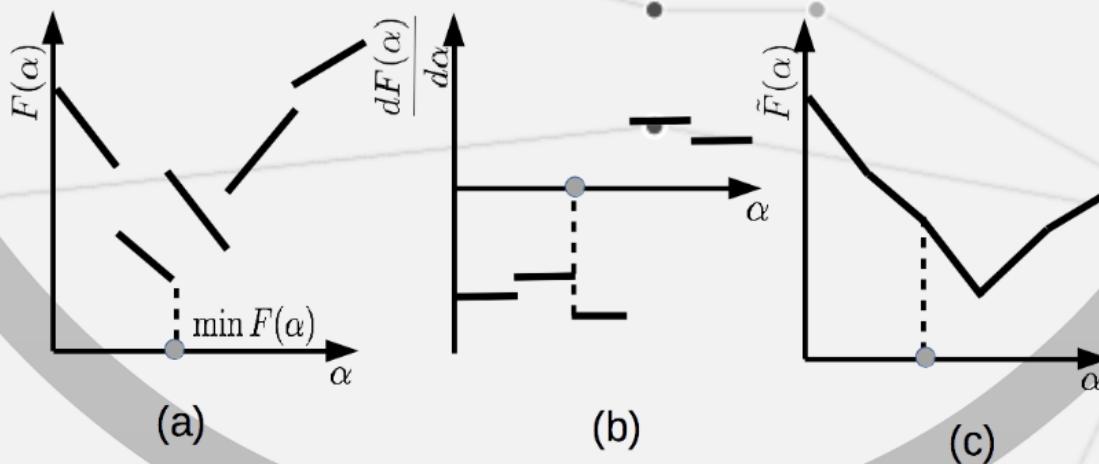


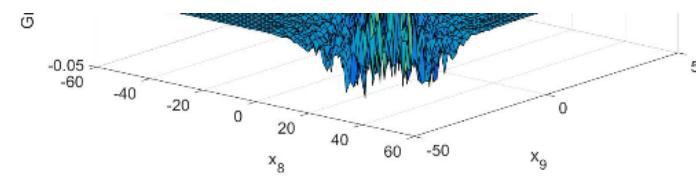
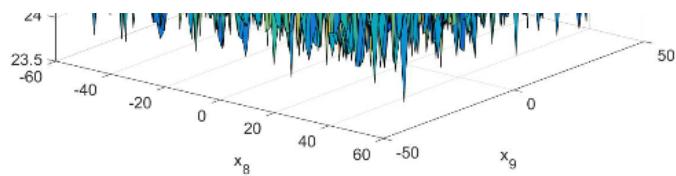
which side looks less noisy?

Why are the gradients less noisy?

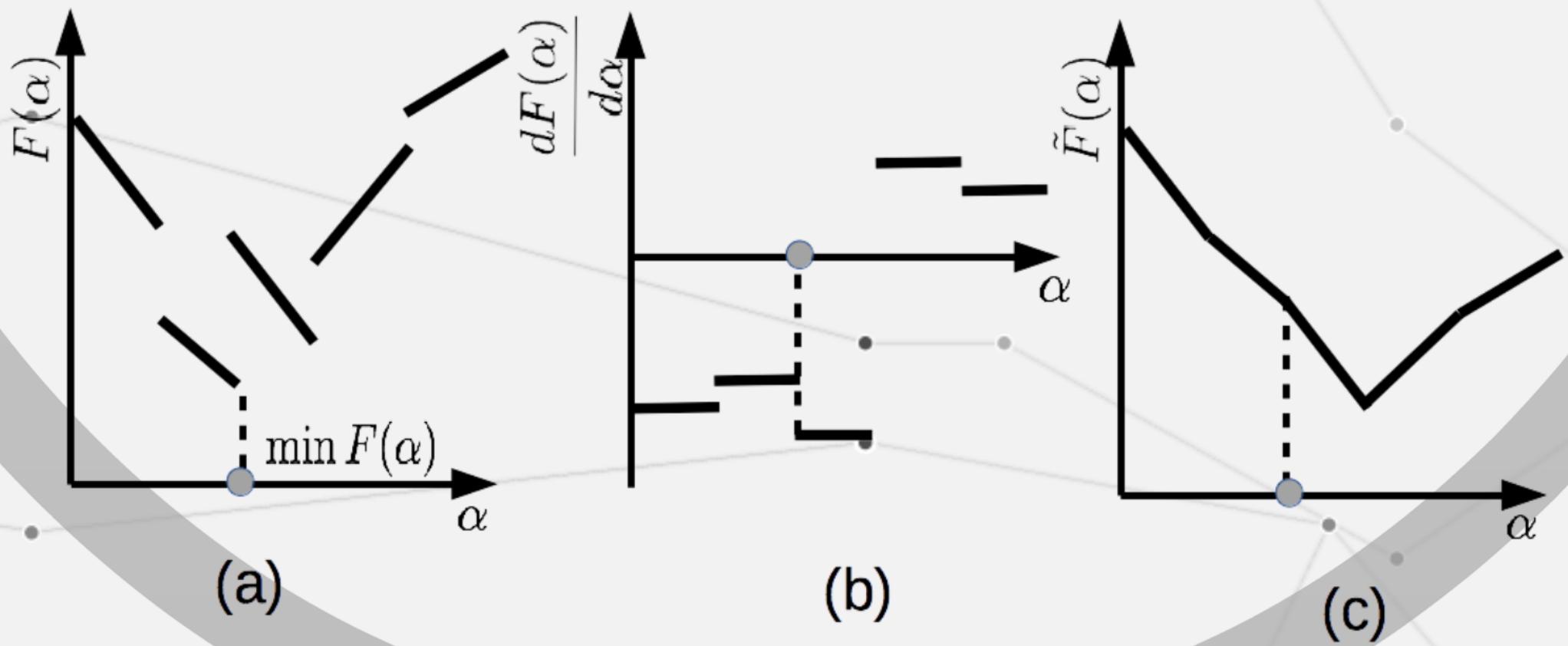


Filters the discontinuities

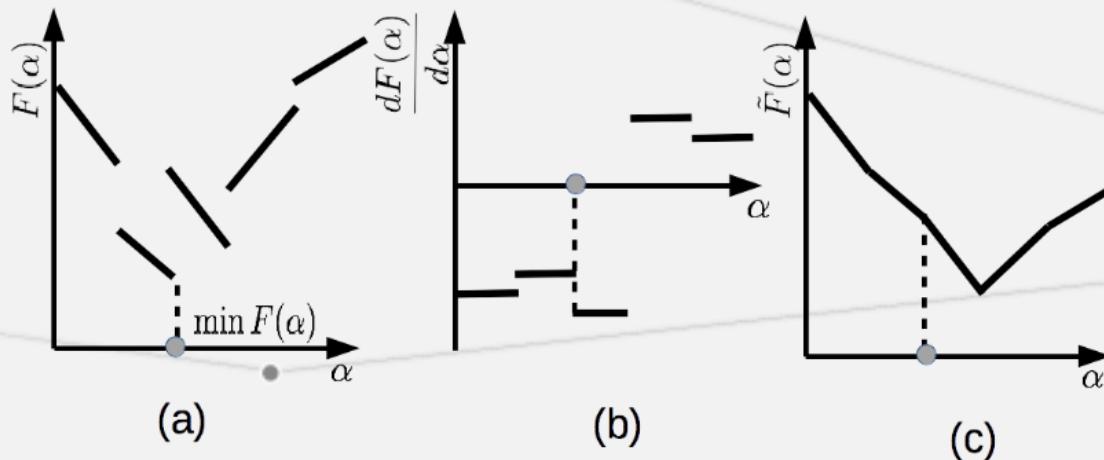




Filters the discontinuities



Important Observations



1. Functions are discontinuous
2. $\nabla f(x) = 0$ may not exist
i.e. no critical points!



Consequences of Observations

Mathematical programming options

1. Discontinuous - minimization issues
2. No critical points - No necessary first order criterion $\nabla f(\mathbf{x}) = \mathbf{0}$

Alternatives

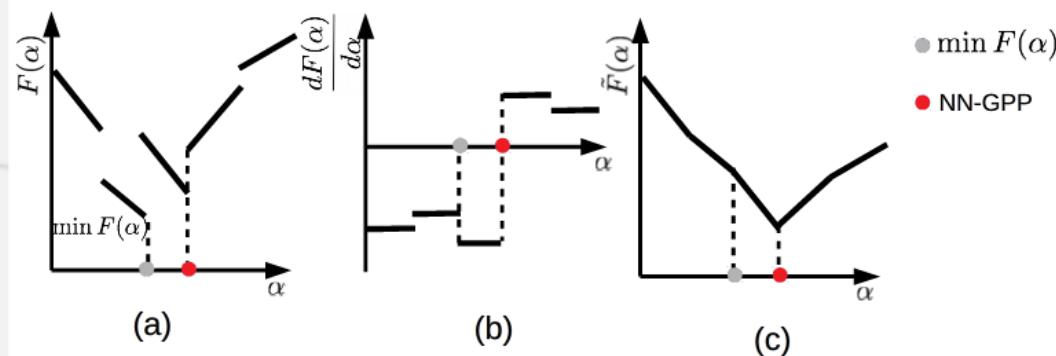
Subgradients - constant / learning rate schedules - back to square one!

Did we pay attention?



Fixing the Problems

Let's pay attention!



Find the location of a sign change in the directional derivative from - to +

Descent (- derivative) while ascent (+ derivative)

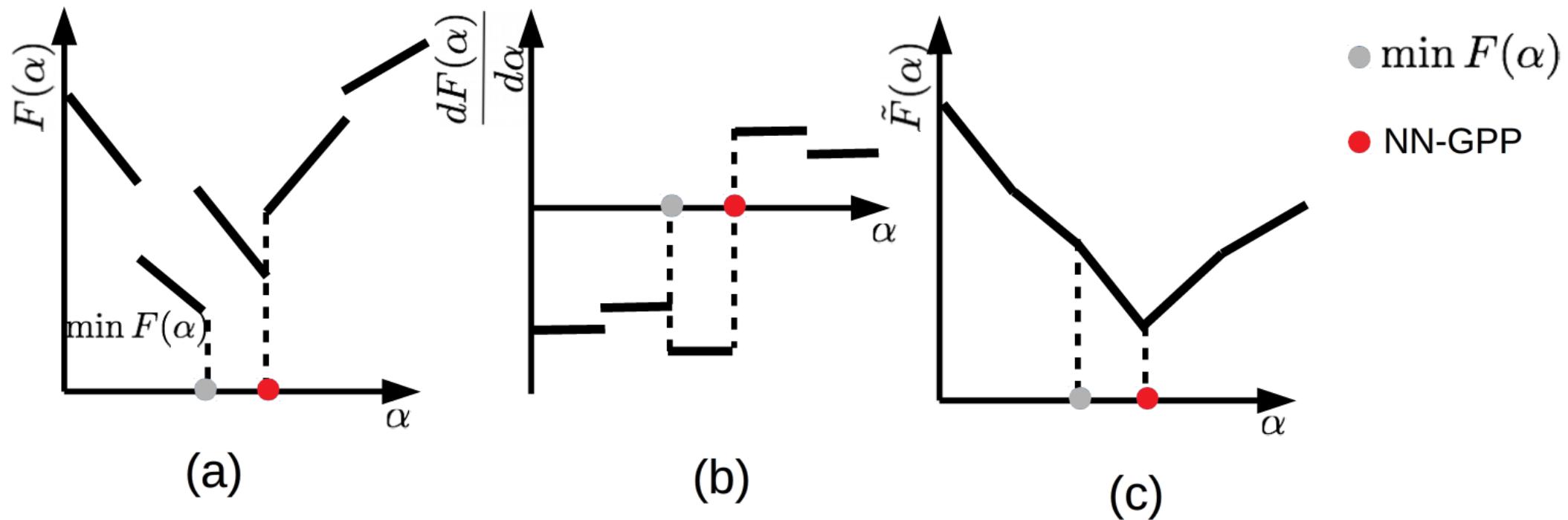
Guaranteed local minimum in the derivative world!

Welcome to Gradient-Only Optimization
Non-Negative Gradient Projection Points
NN-GPP



the Problems

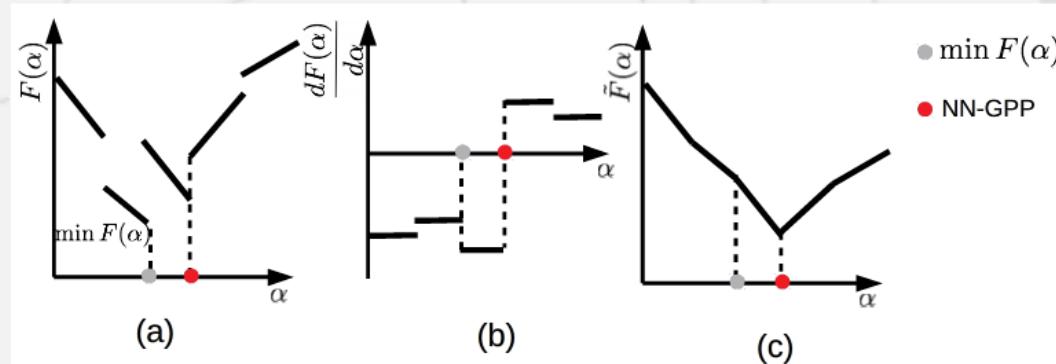
Let's pay attention!



Find the location of a sign change in the derivative from - to +

Fixing the Problems

Let's pay attention



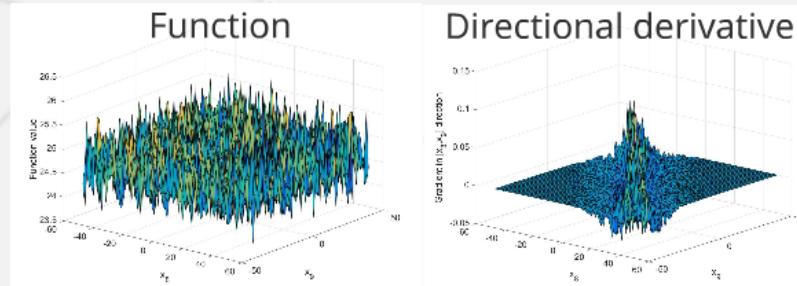
Three Guidelines for Discontinuous Functions:

1. Don't minimize (local minima in function world)
2. Don't find critical points $\nabla f(x) = 0$ candidates
3. **Locate NN-GPP (local minima in derivative world)**

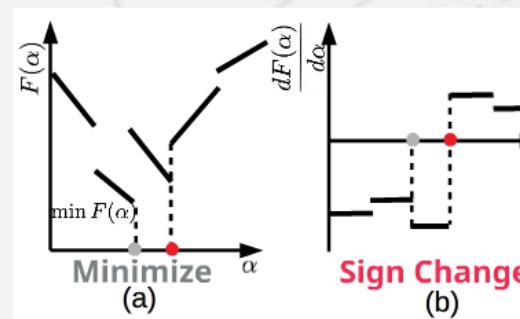
Are we set for line searches?



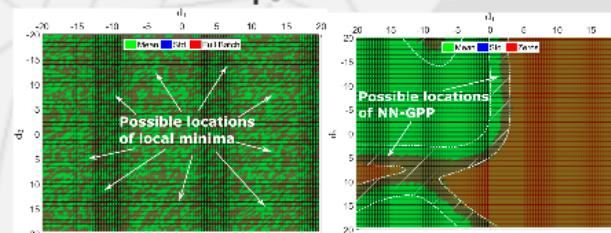
Double Benefit



Mini-Batches - 10/76



Implication

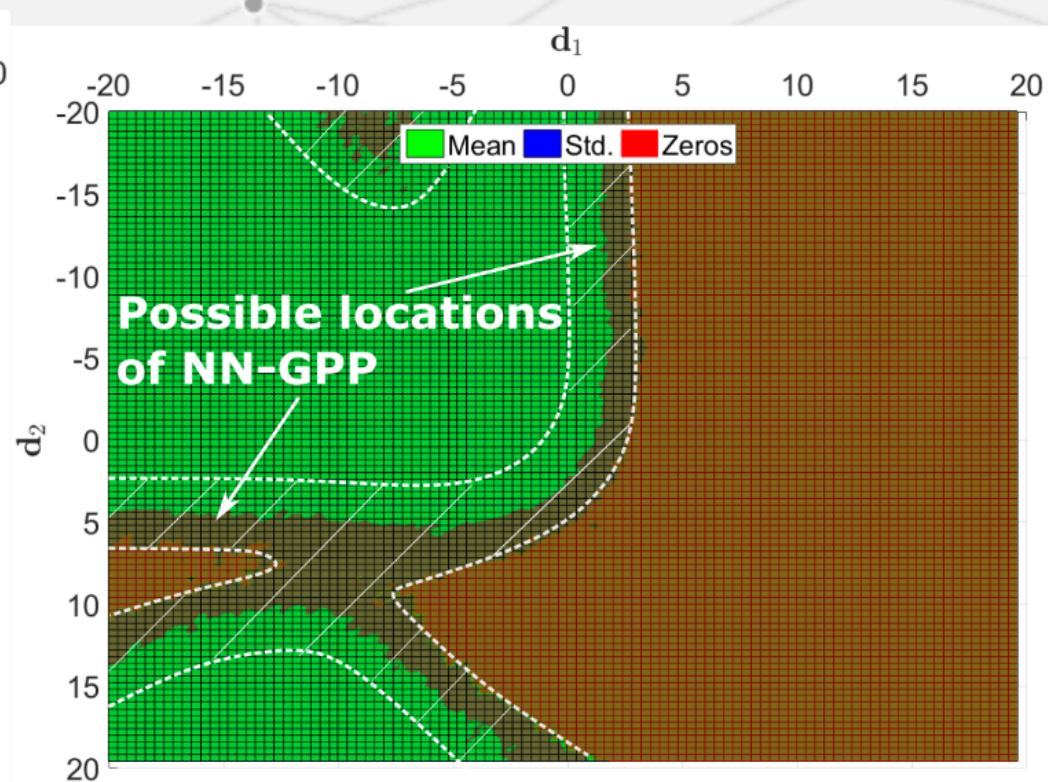
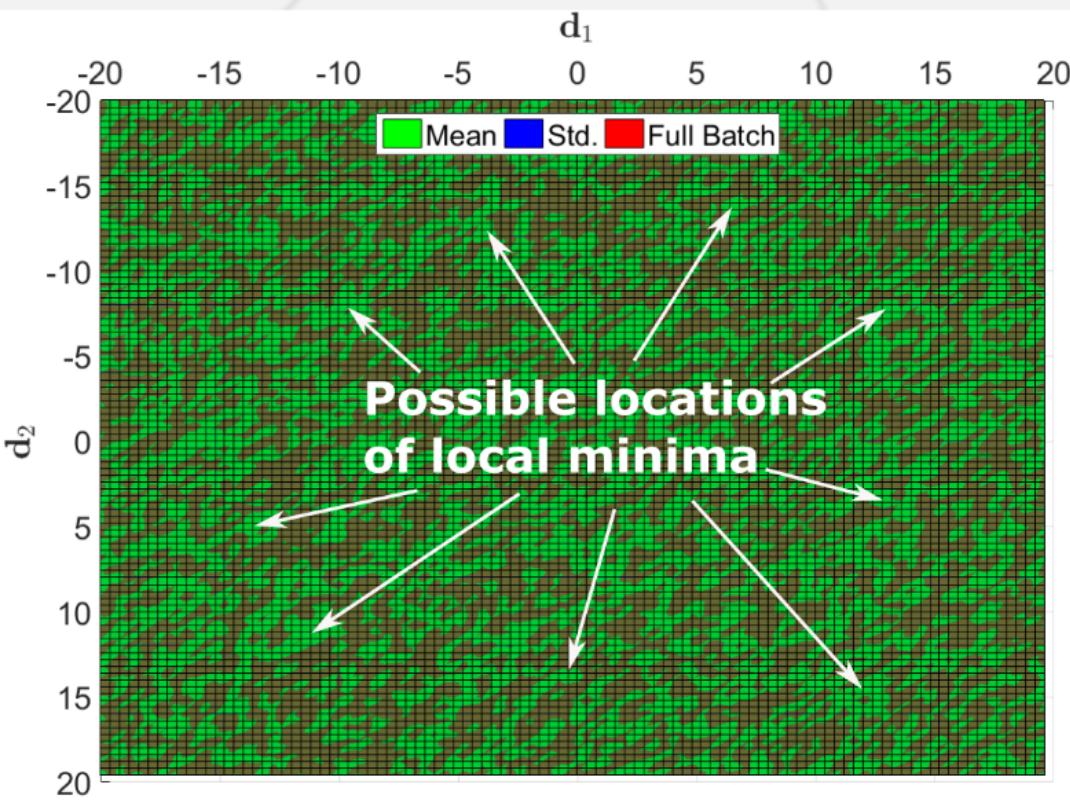


Minimize
(a)

α

Sign Change
(b)

Implication



Where is my line search?

Gradient-Only Line Search Inexact
GOLS-I



Dominic
Kafka (PhD)

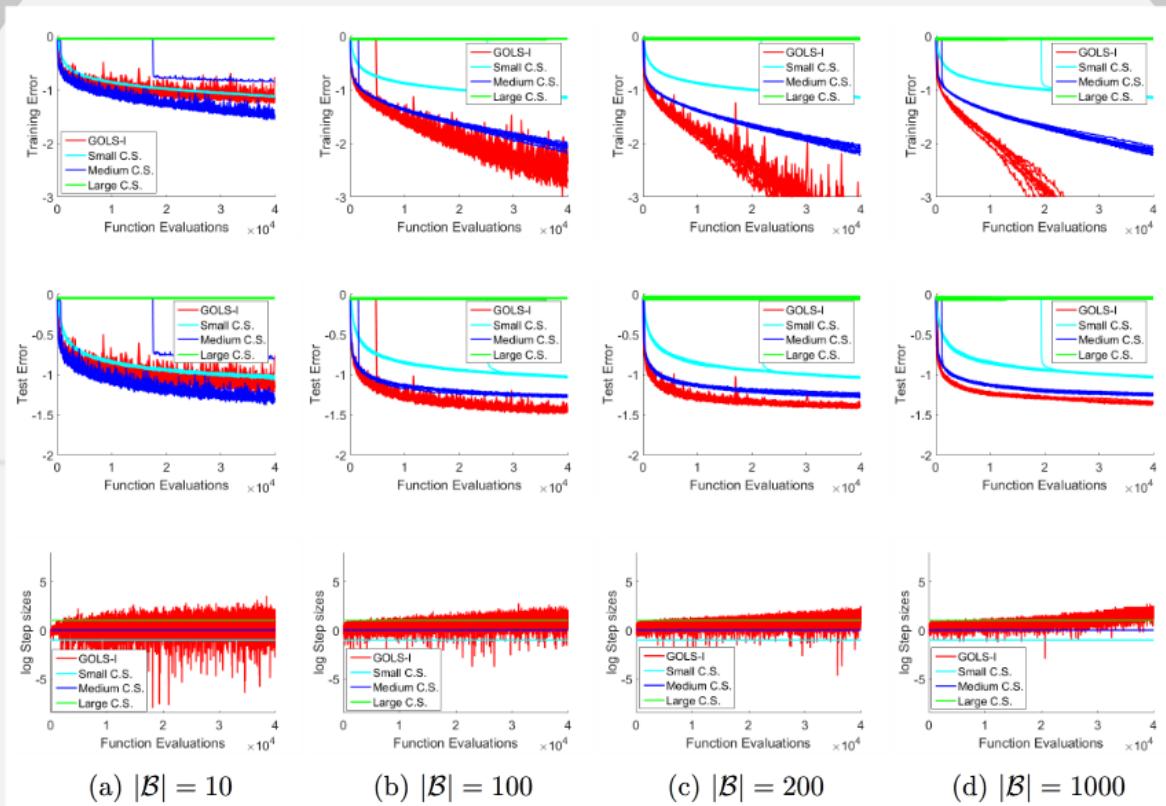


GOLS-I FEATURES

- 1. No tuning parameter - automatic line search**
- 2. Step size range - 15 orders of magnitude**
- 3. Barrier to Entry - High school math**
- 4. PyTorch (Done) and TensorFlow (Testing)**

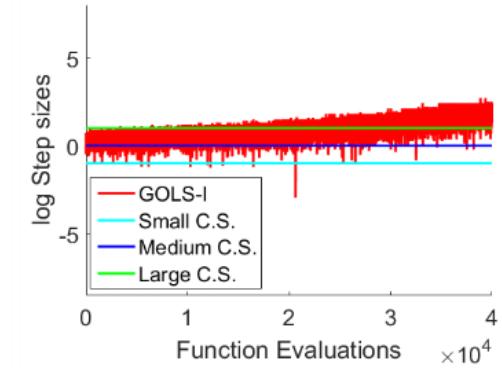
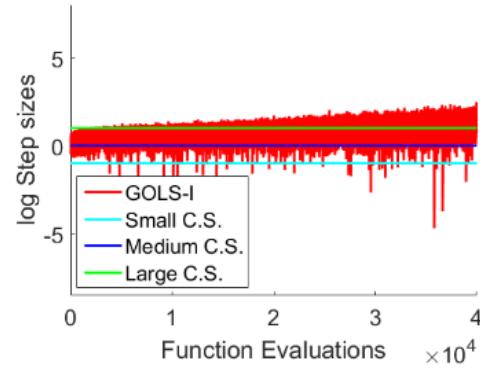
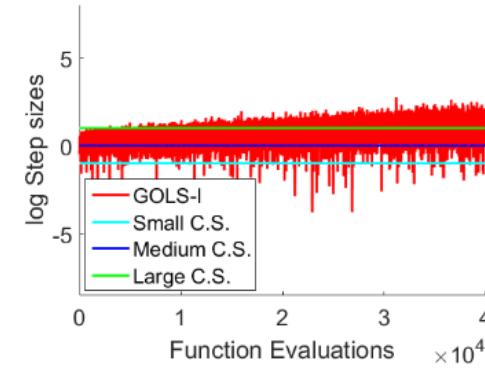
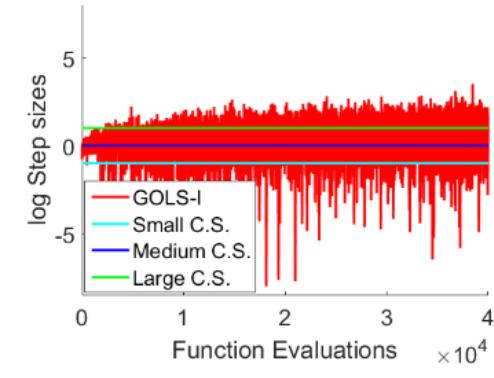
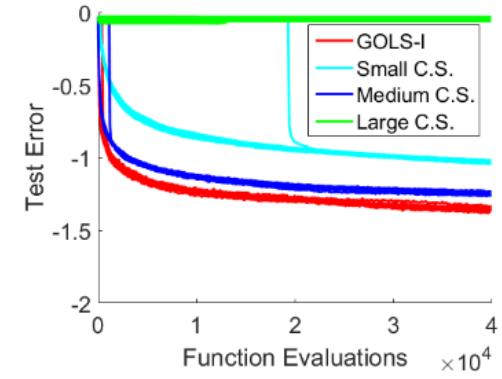
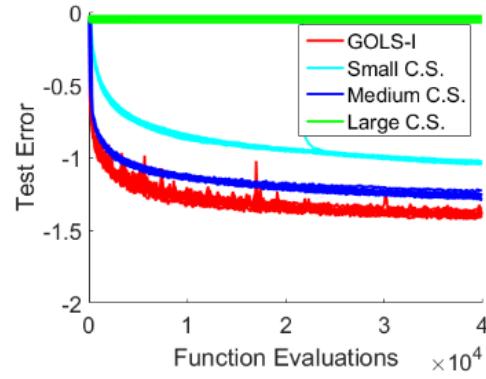
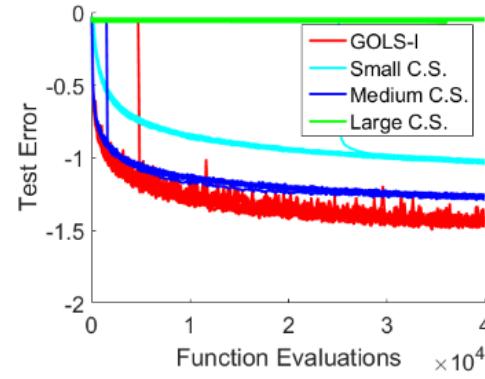
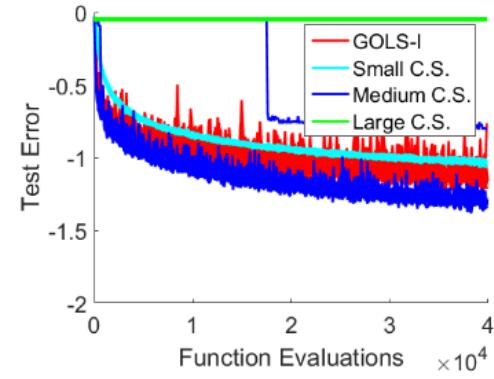
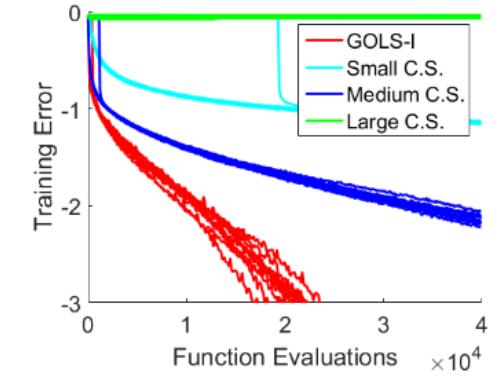
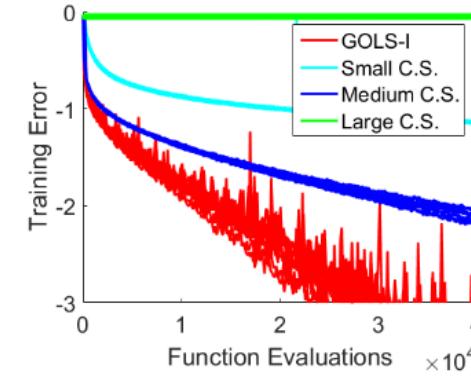
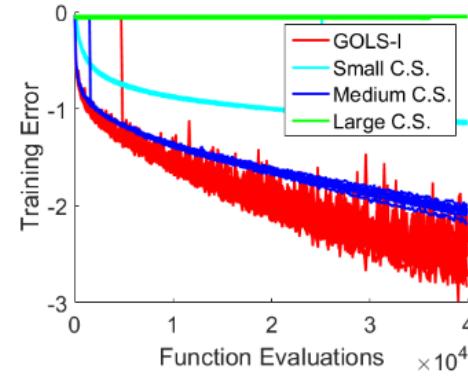
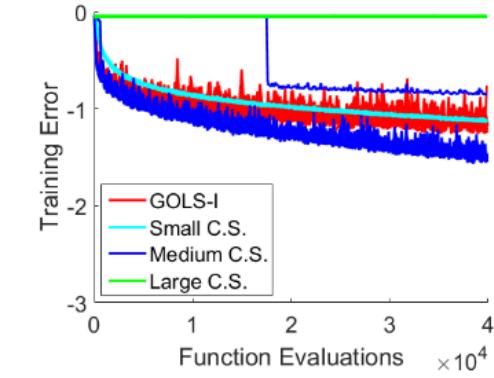


Results MNIST



**Net I - HL:800; AF:SIGMOID; INIT:
N(0,1); LOSS: Cross Entropy**





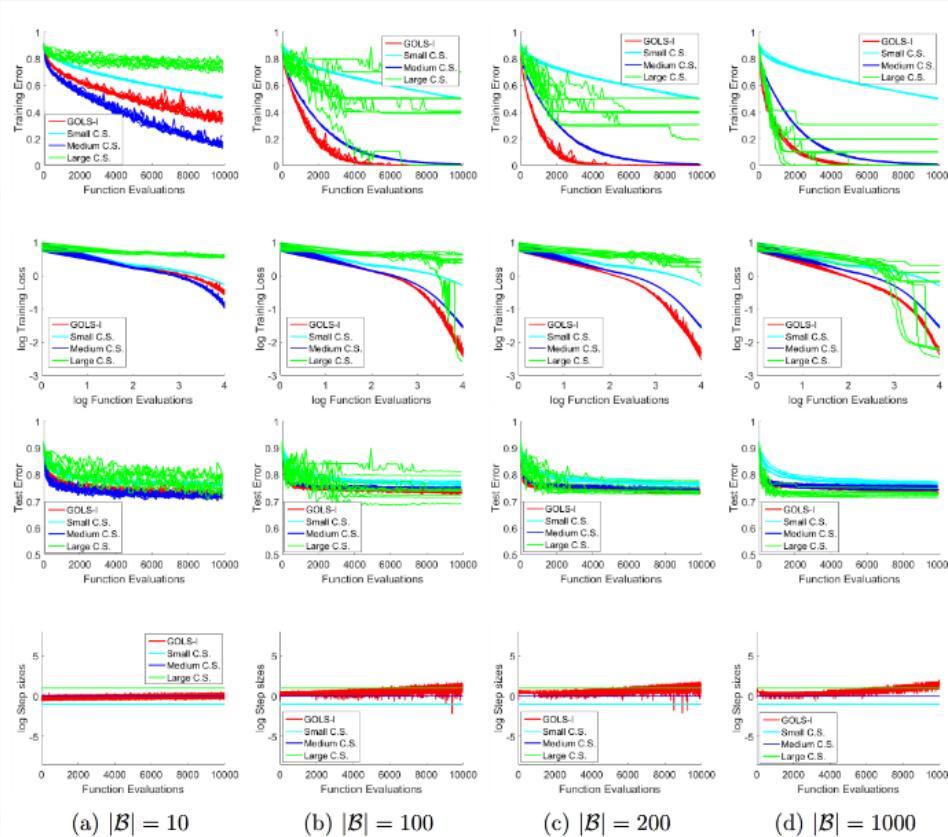
(a) $|\mathcal{B}| = 10$

(b) $|\mathcal{B}| = 100$

(c) $|\mathcal{B}| = 200$

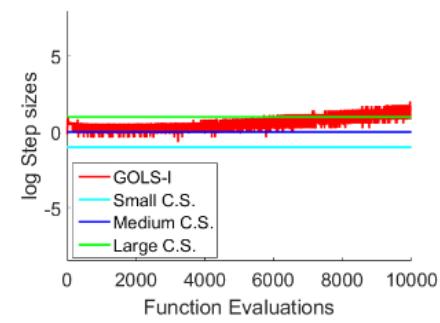
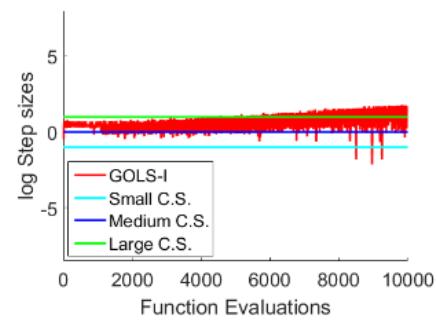
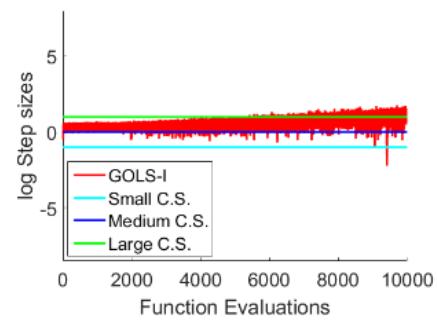
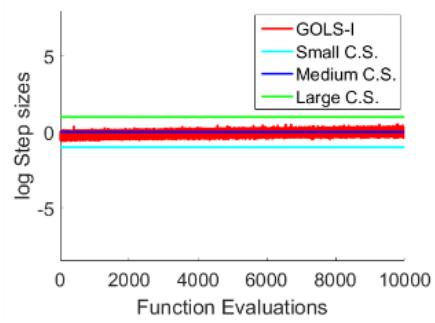
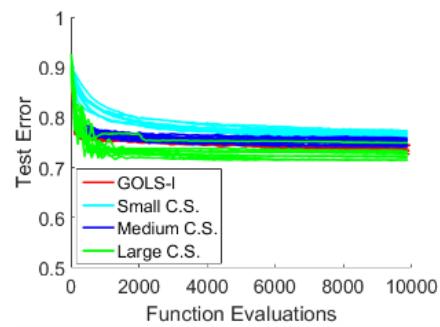
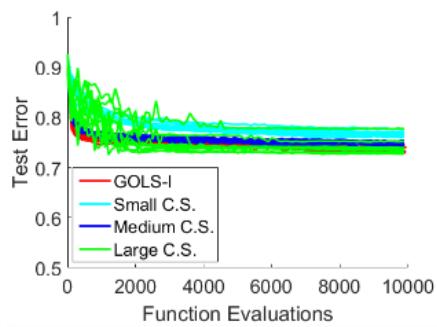
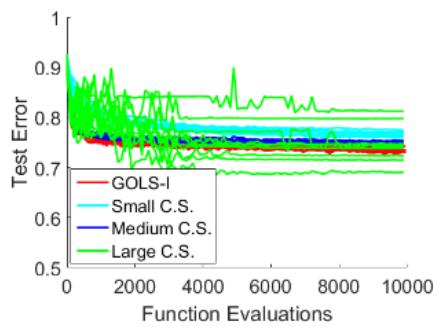
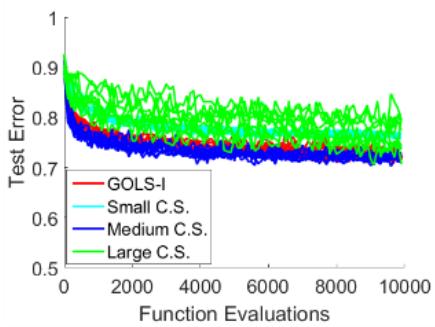
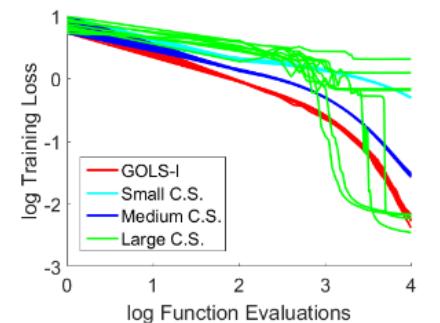
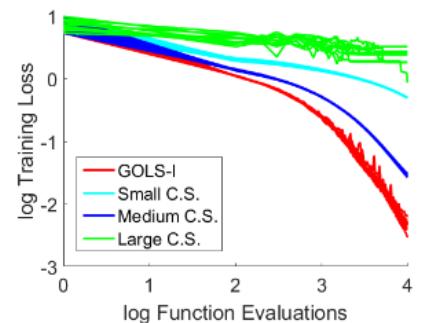
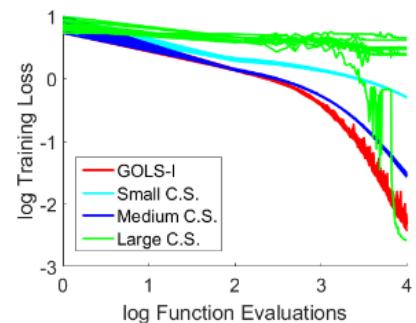
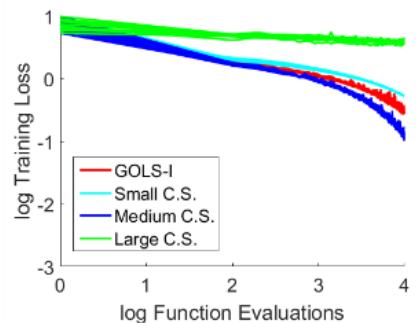
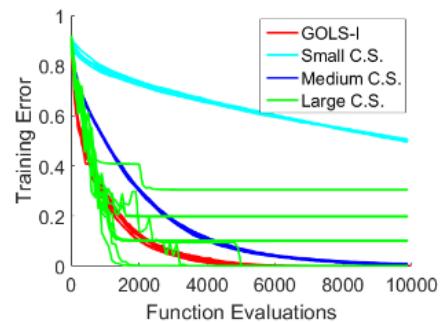
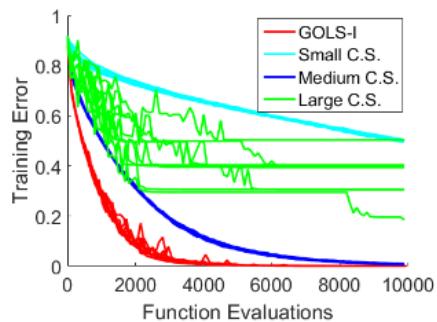
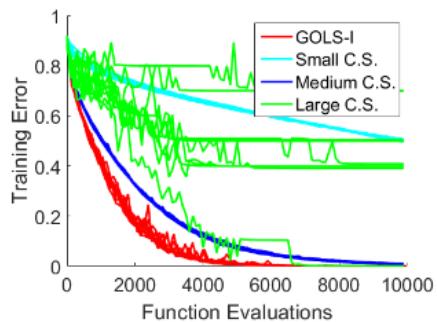
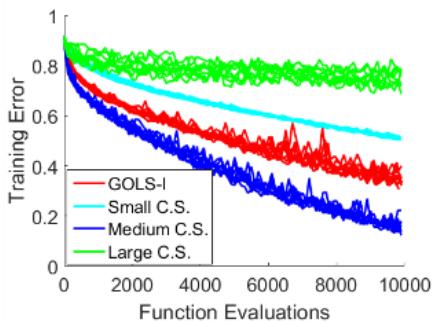
(d) $|\mathcal{B}| = 1000$

Results CIFAR



Net I - HL:800; AF:SIGMOID; INIT:
 $N(0,1)$; LOSS: Cross Entropy





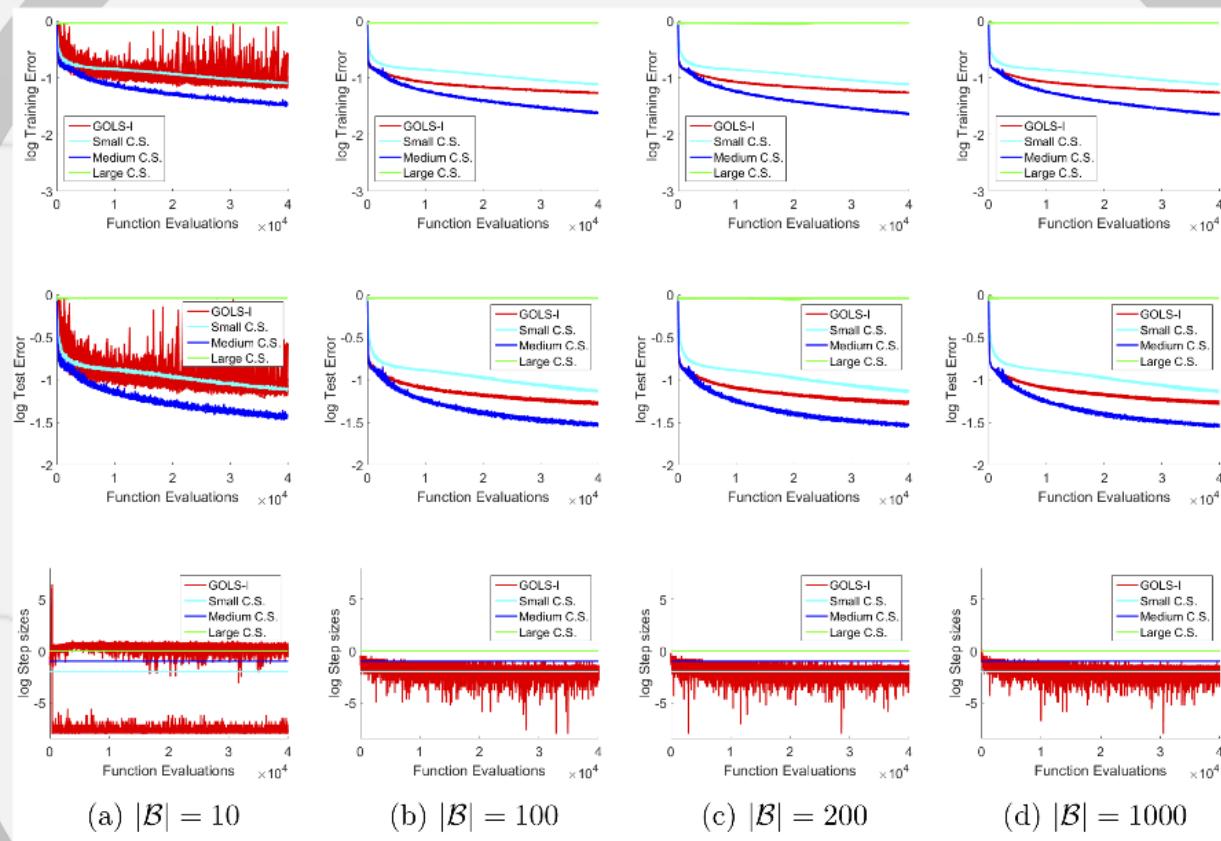
(a) $|\mathcal{B}| = 10$

(b) $|\mathcal{B}| = 100$

(c) $|\mathcal{B}| = 200$

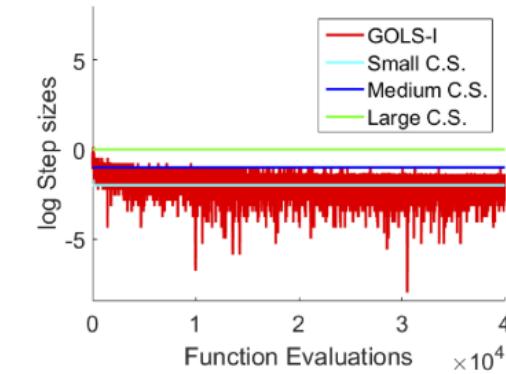
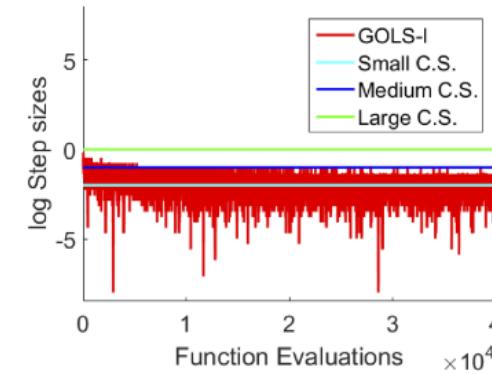
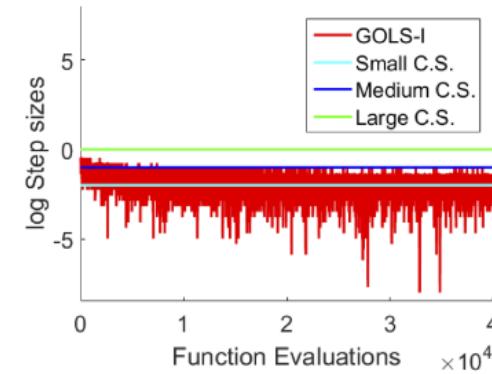
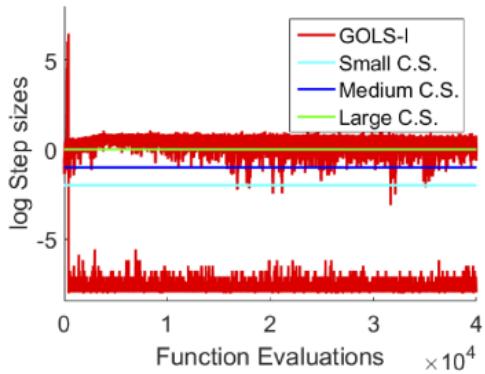
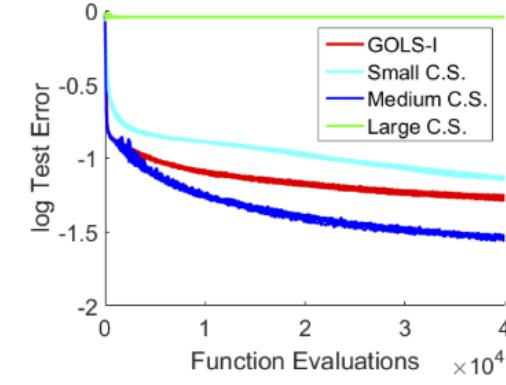
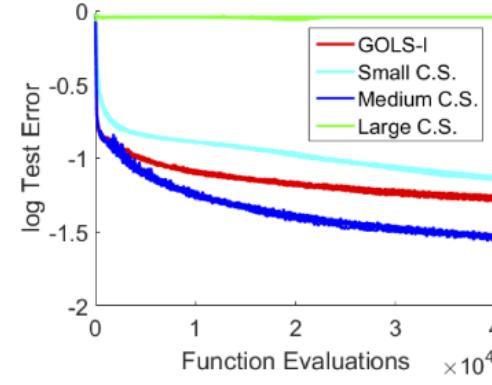
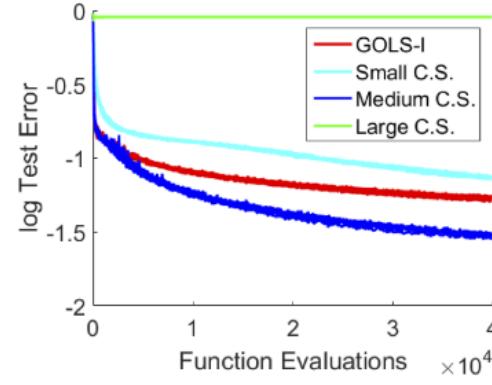
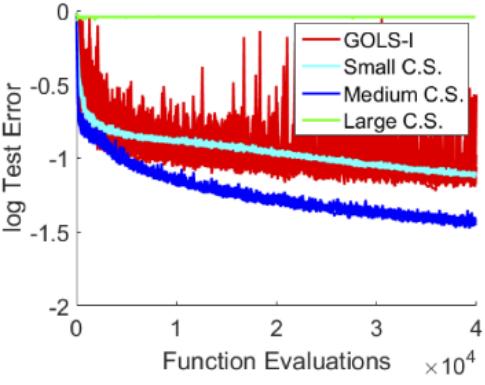
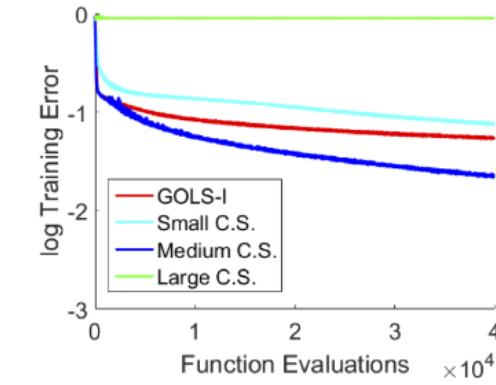
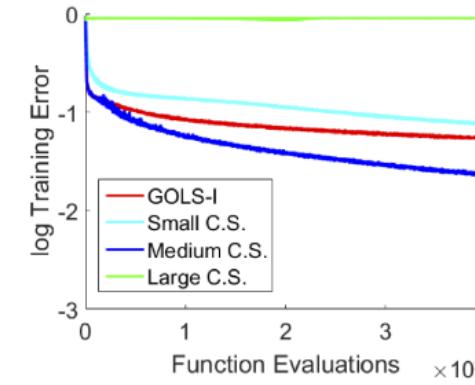
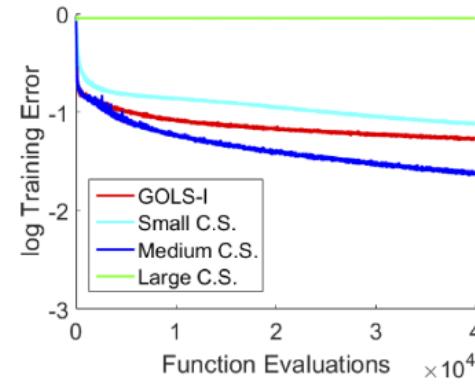
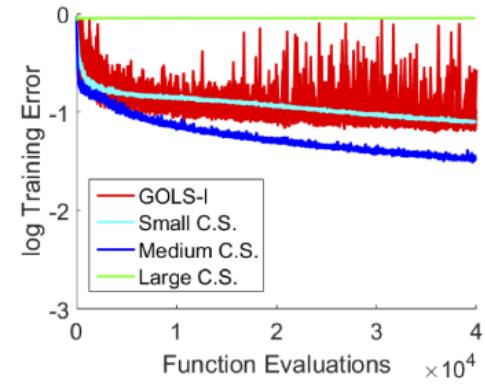
(d) $|\mathcal{B}| = 1000$

Results MNIST



HL:1000-500-250;
AF:TANH; INIT:XAVIER;
LOSS:MSE





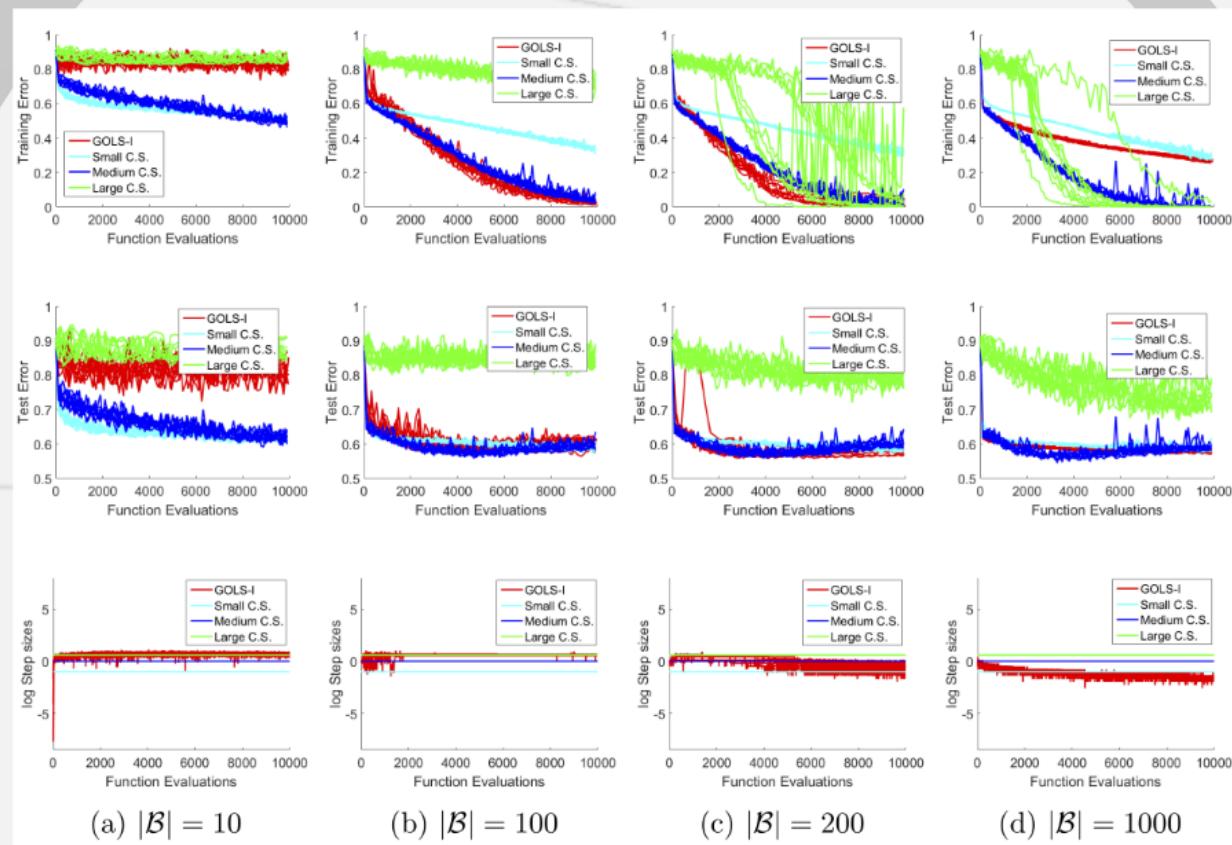
(a) $|\mathcal{B}| = 10$

(b) $|\mathcal{B}| = 100$

(c) $|\mathcal{B}| = 200$

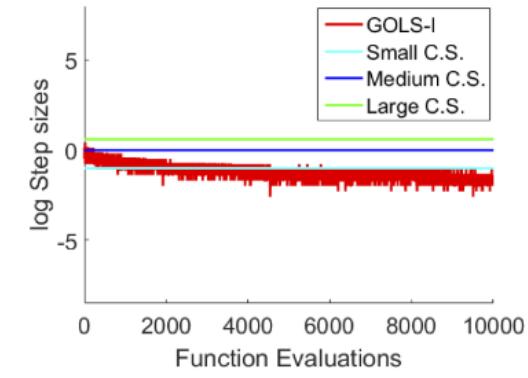
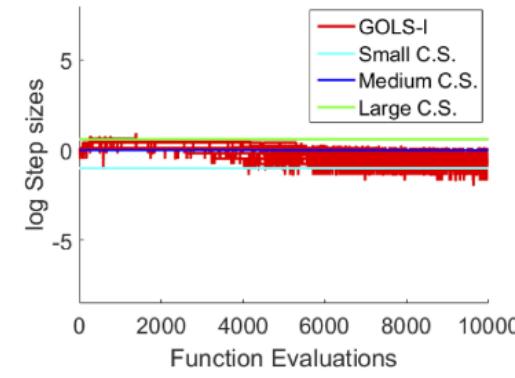
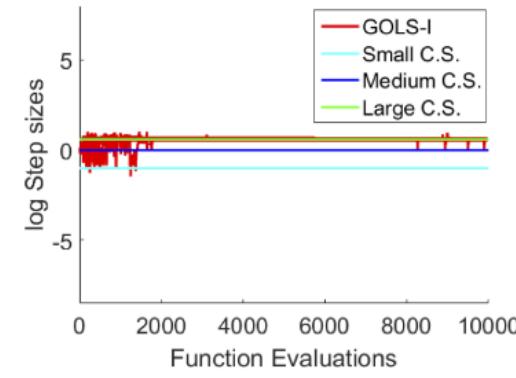
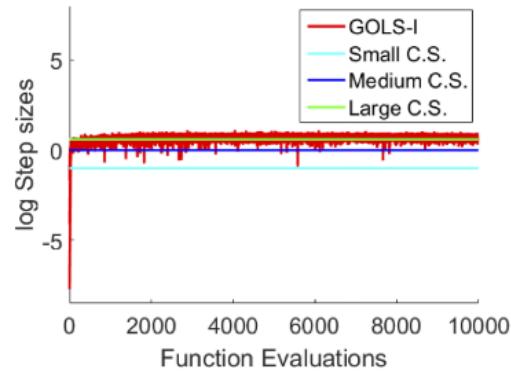
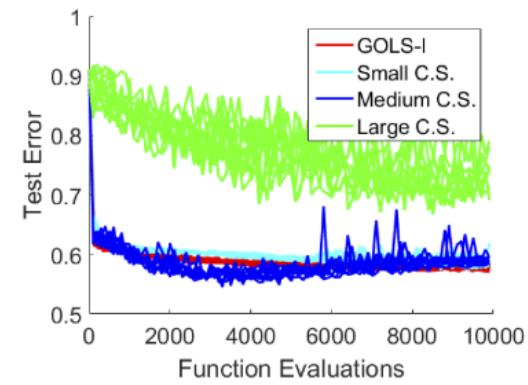
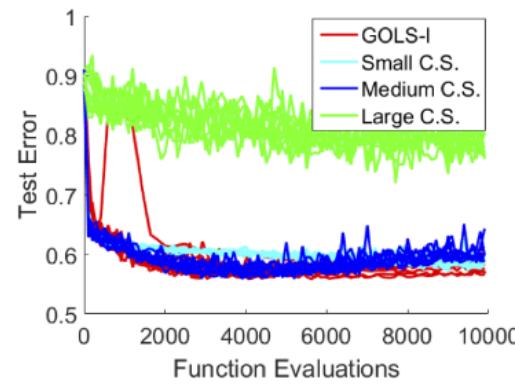
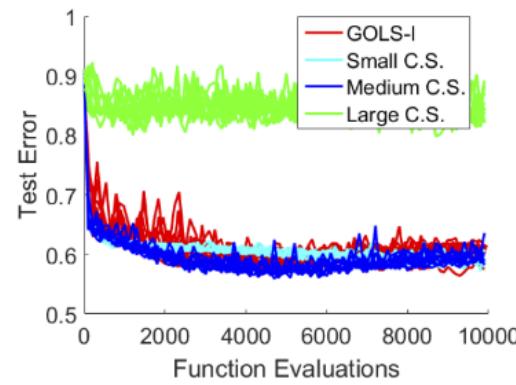
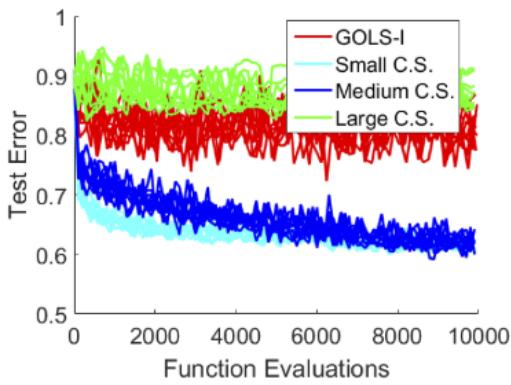
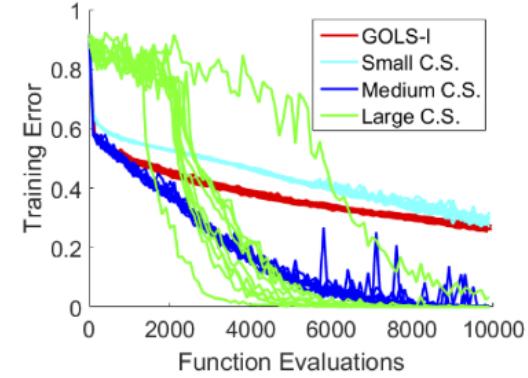
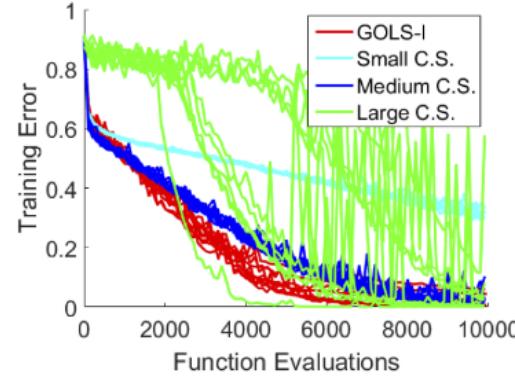
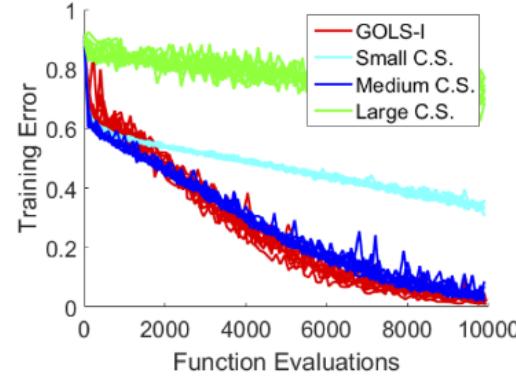
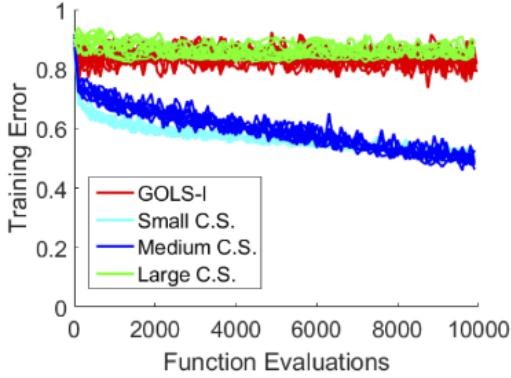
(d) $|\mathcal{B}| = 1000$

Results CIFAR



**HL:1000-500-250;
AF:TANH; INIT:XAVIER;
LOSS:MSE**





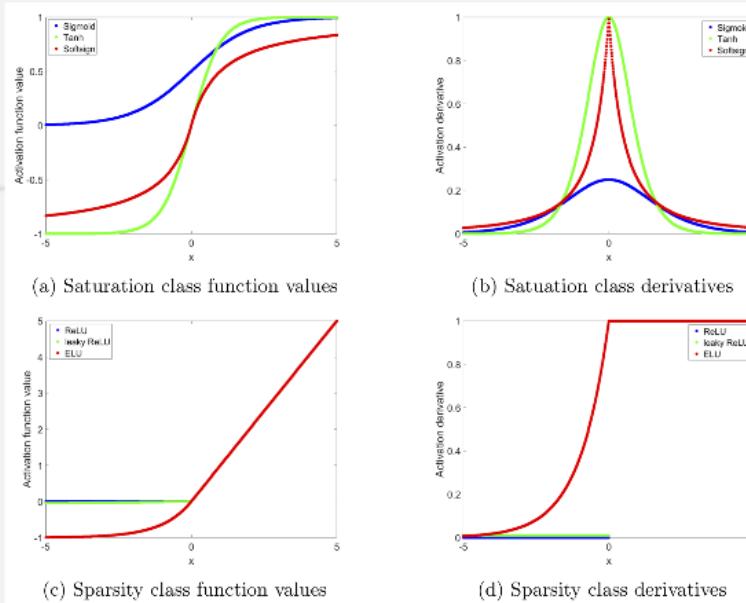
(a) $|\mathcal{B}| = 10$

(b) $|\mathcal{B}| = 100$

(c) $|\mathcal{B}| = 200$

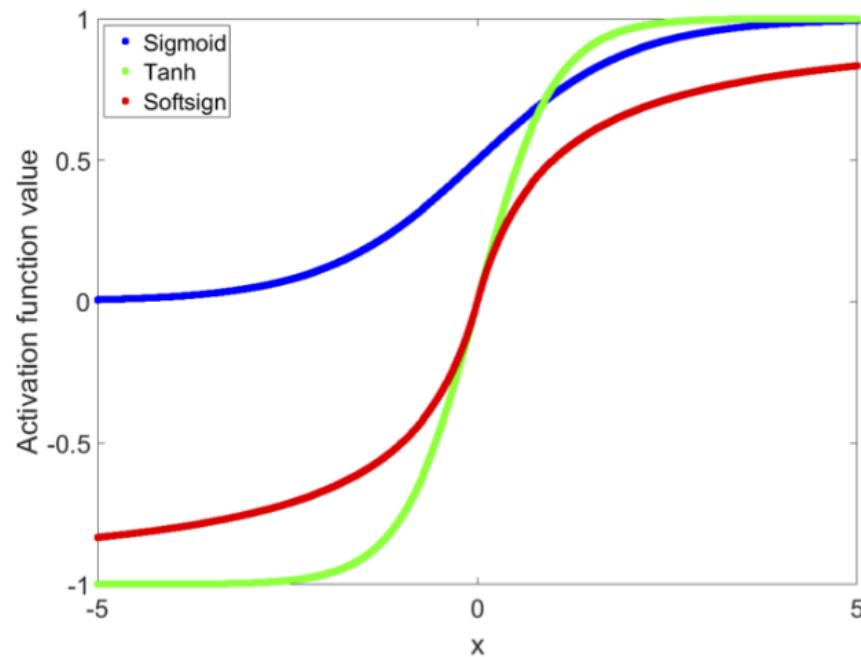
(d) $|\mathcal{B}| = 1000$

Sensitivity to Activation Functions

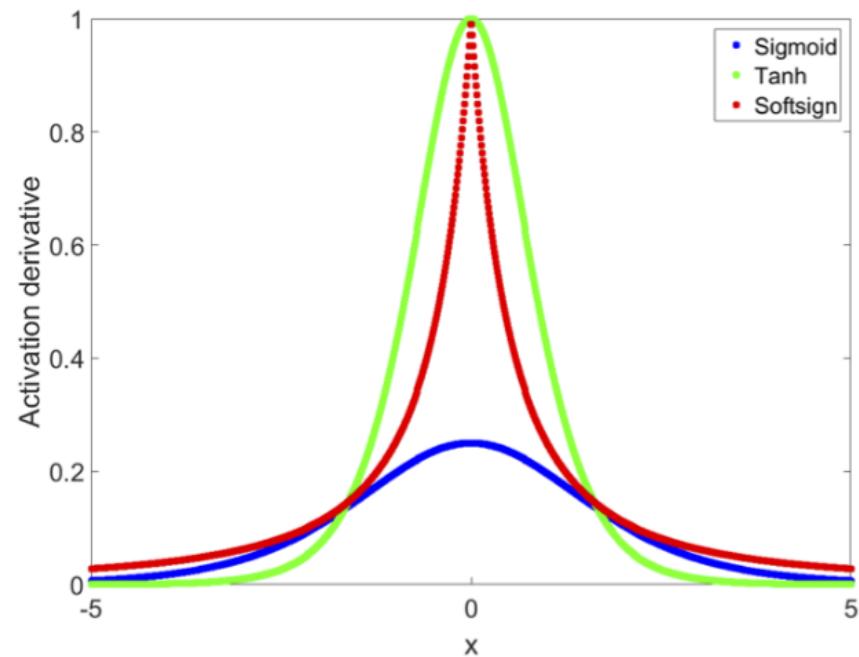


**DATASET:IRIS; HL:10
AF:TANH; LOSS: MSE**

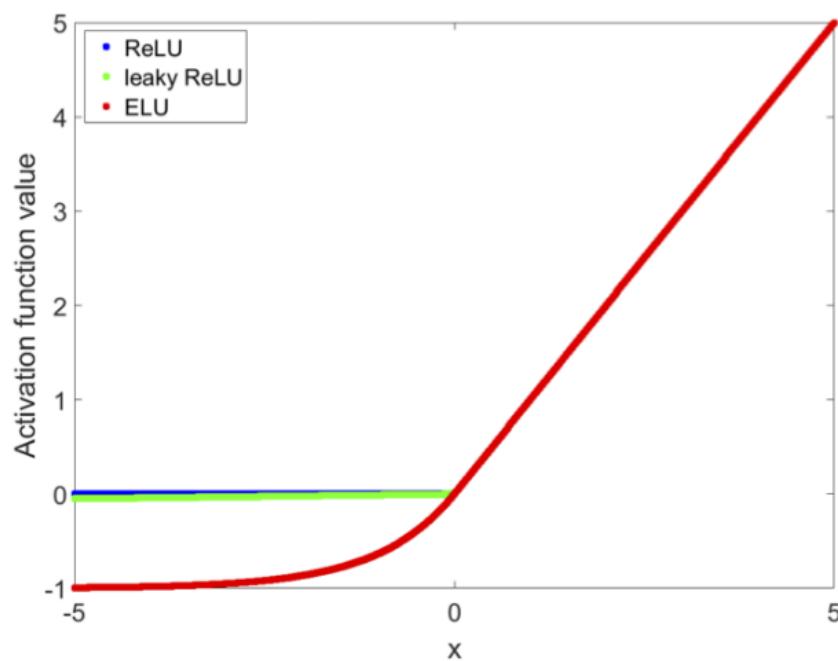




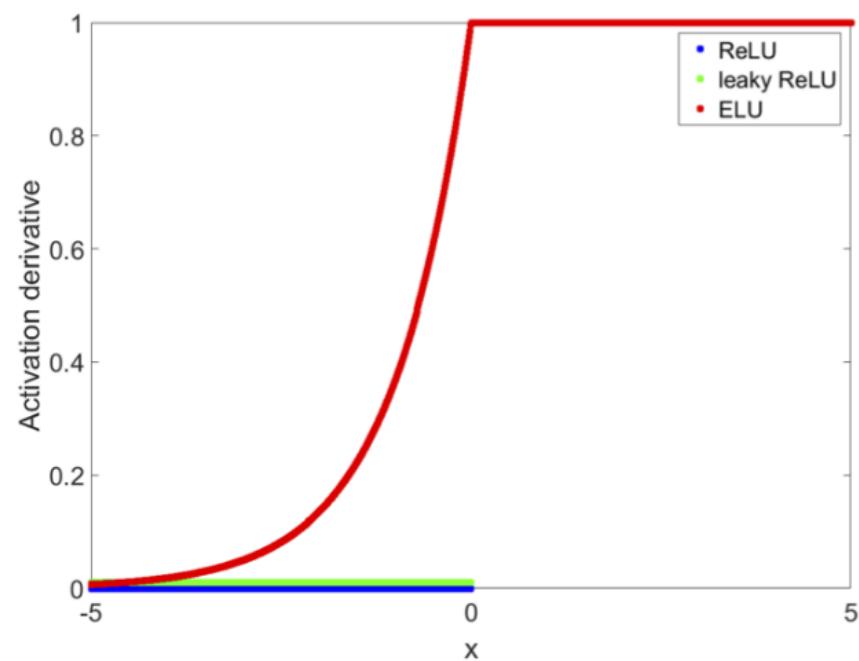
(a) Saturation class function values



(b) Satuation class derivatives

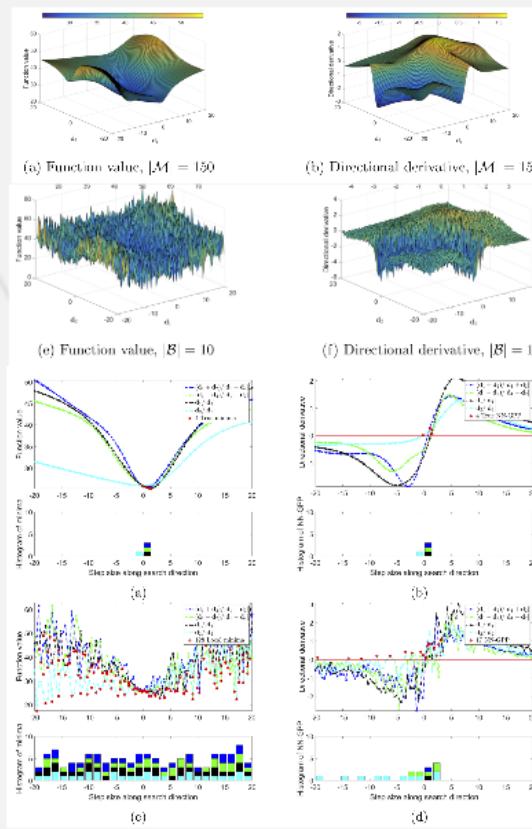


(c) Sparsity class function values

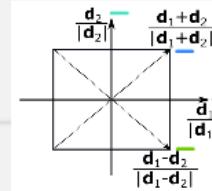


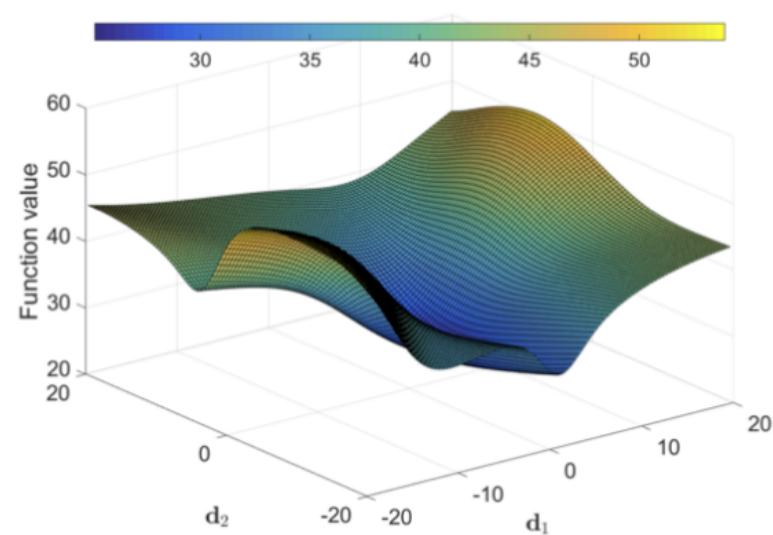
(d) Sparsity class derivatives

Role of Activation Functions SIGMOID

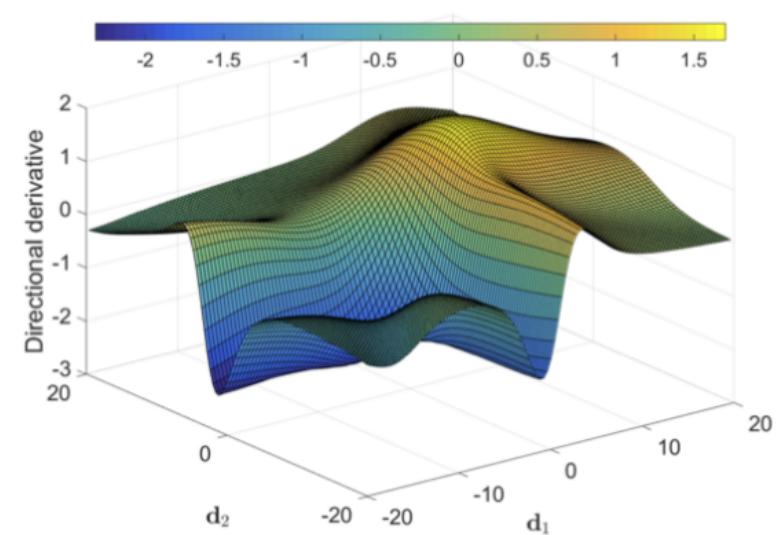


Directions

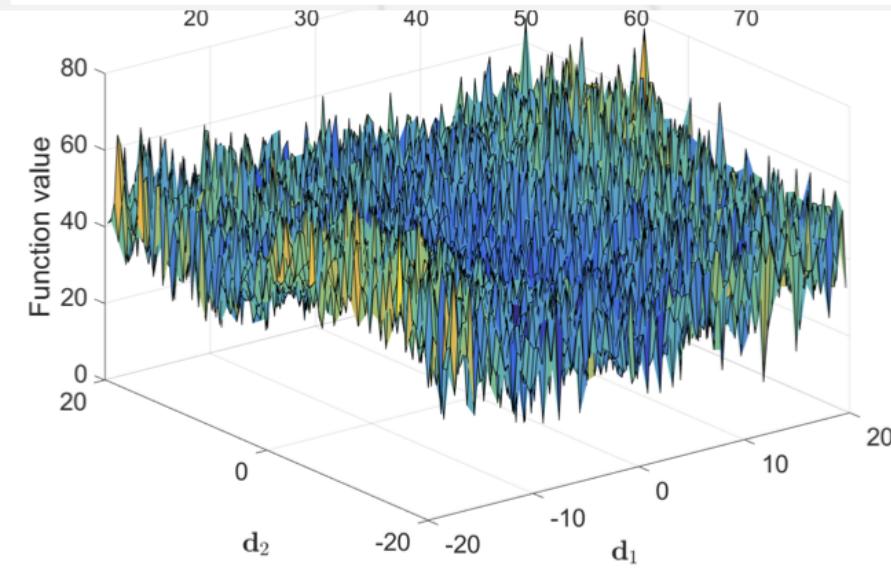




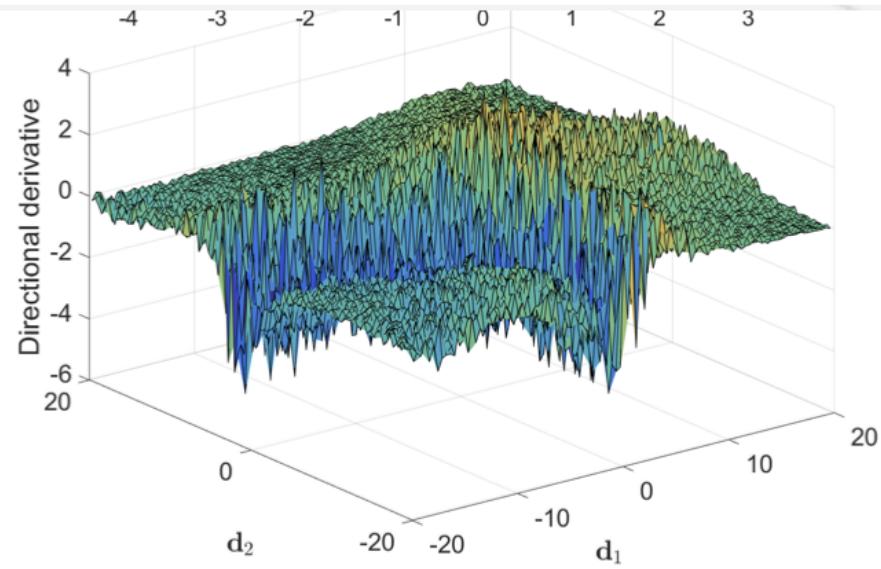
(a) Function value, $|\mathcal{M}| = 150$



(b) Directional derivative, $|\mathcal{M}| = 150$



(e) Function value, $|\mathcal{B}| = 10$

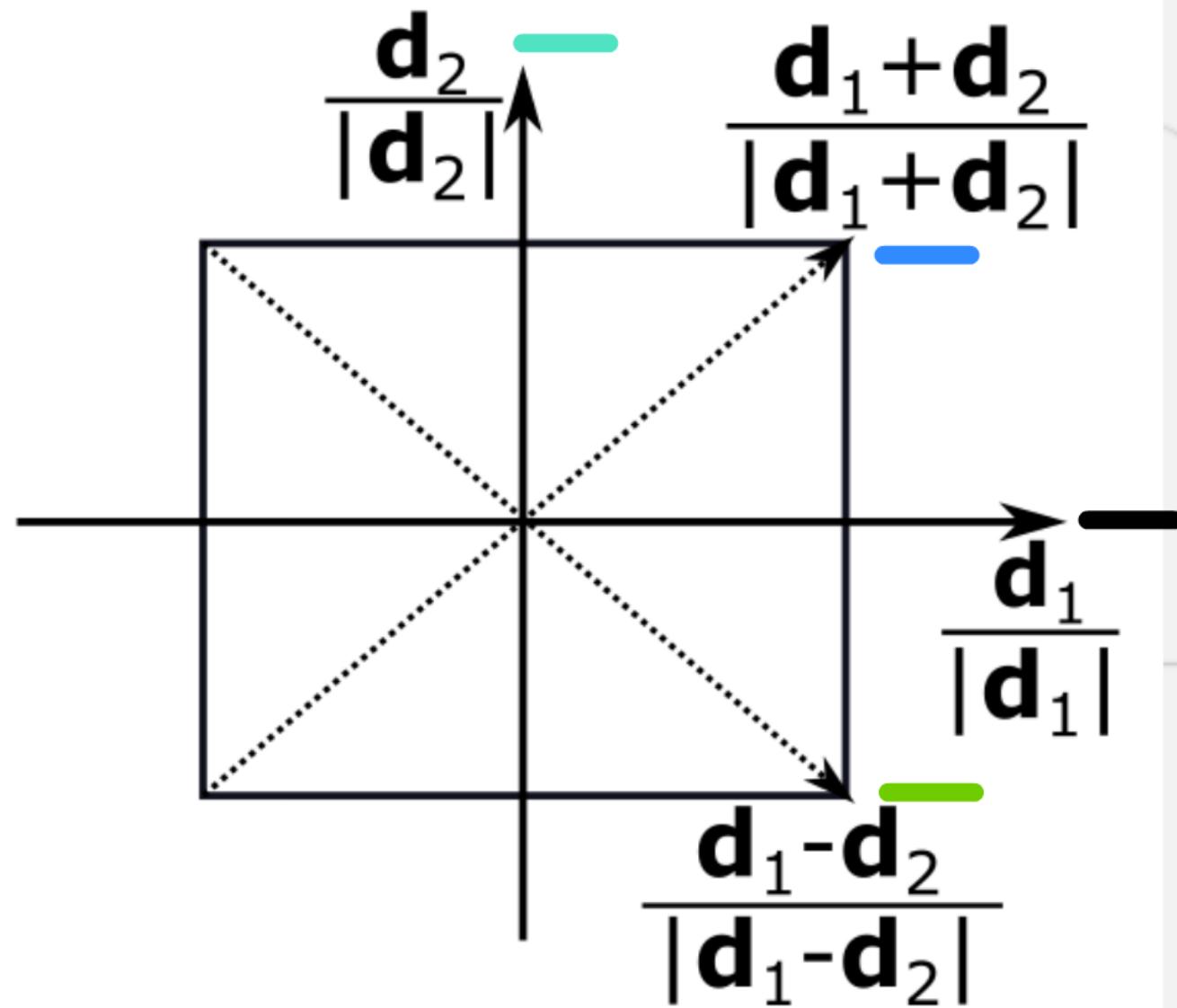


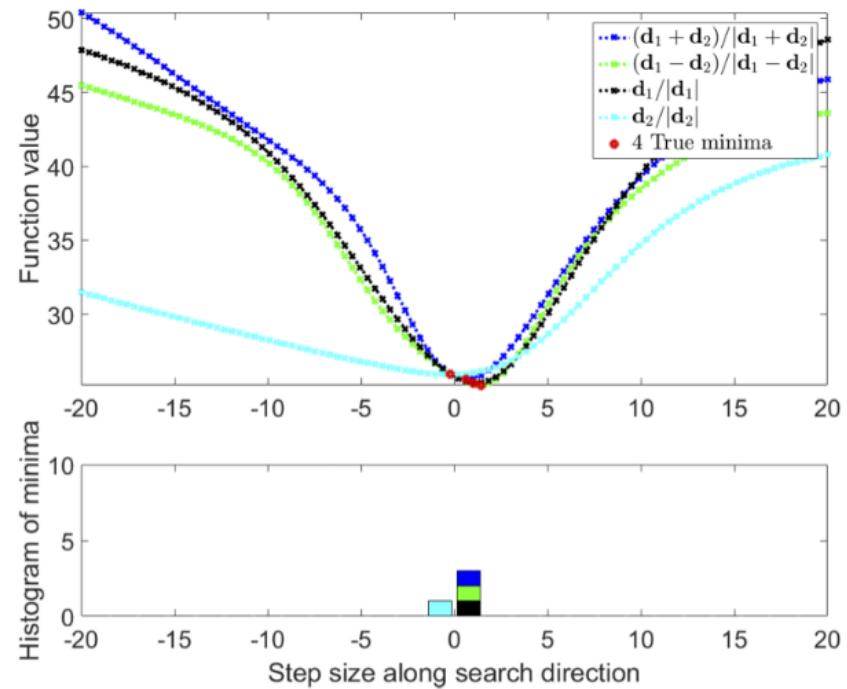
(f) Directional derivative, $|\mathcal{B}| = 10$

$\text{---} \bullet \bullet \bullet \quad (\mathbf{d}_1 + \mathbf{d}_2)/\|\mathbf{d}_1 + \mathbf{d}_2\|$

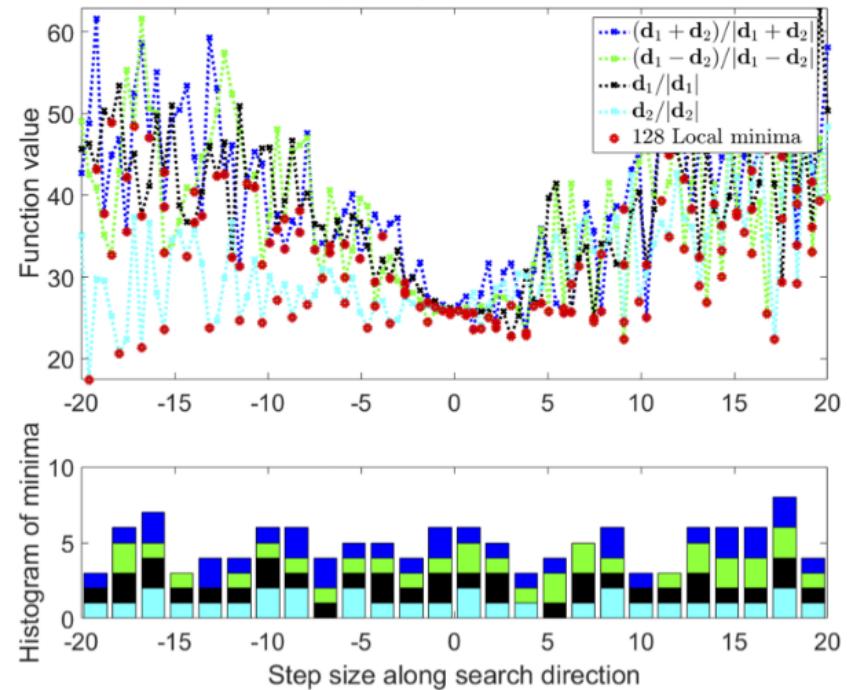
$\text{---} \bullet \bullet \bullet \quad (\mathbf{d}_1 + \mathbf{d}_2)/\|\mathbf{d}_1 + \mathbf{d}_2\|$

Directions

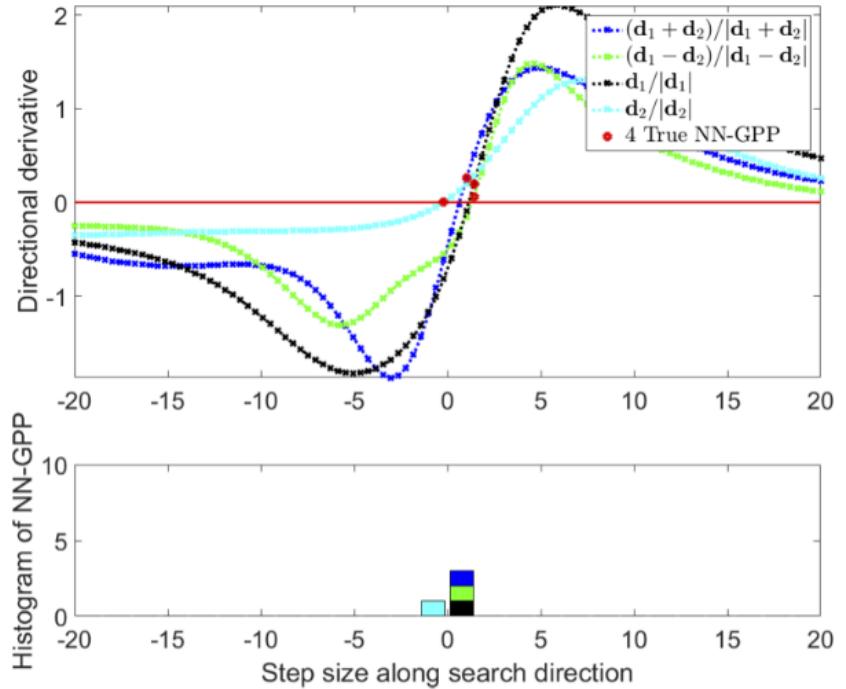




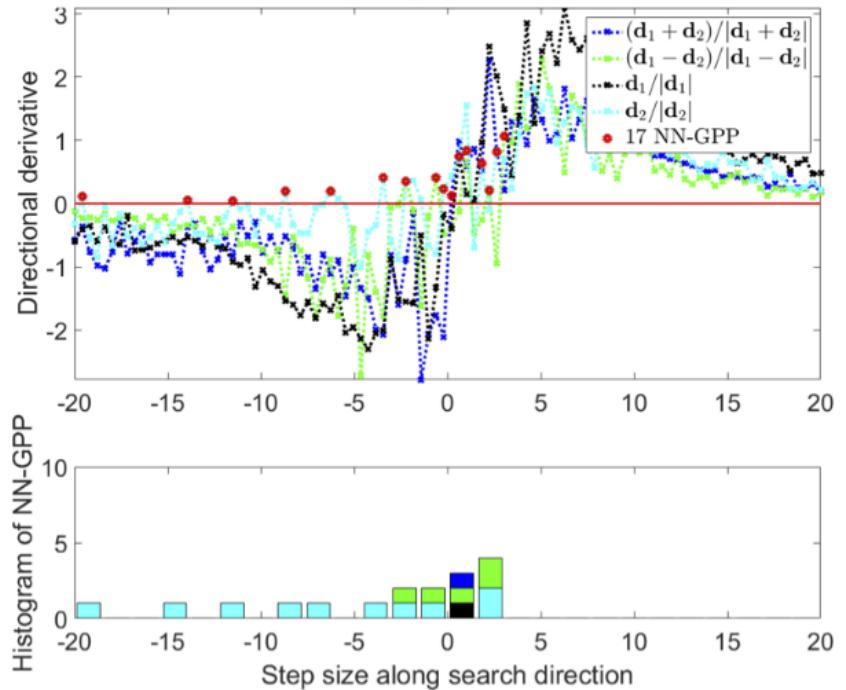
(a)



(c)



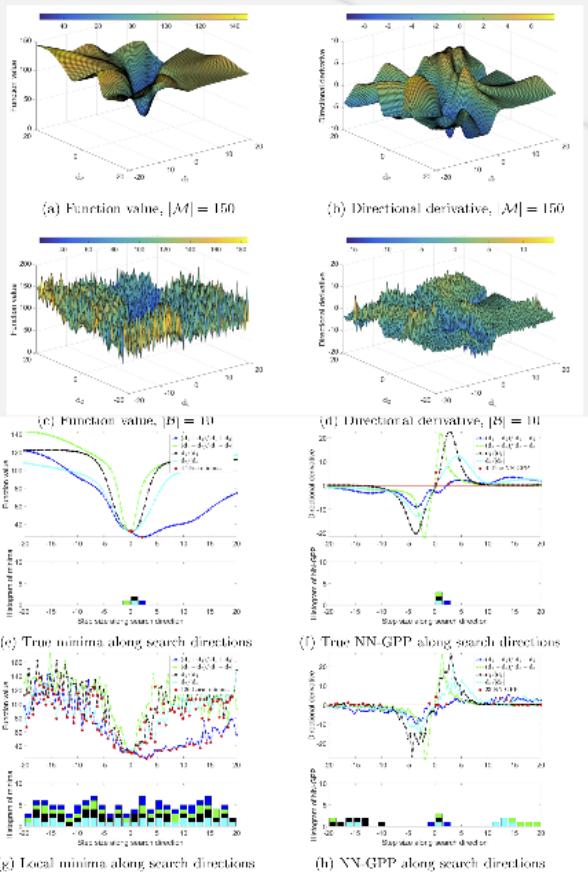
(b)

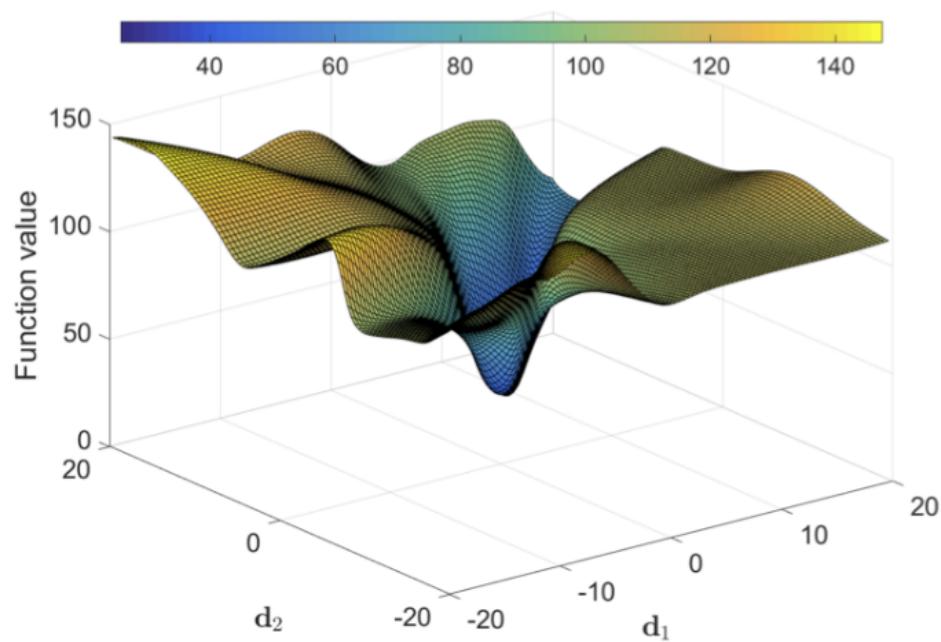


(d)

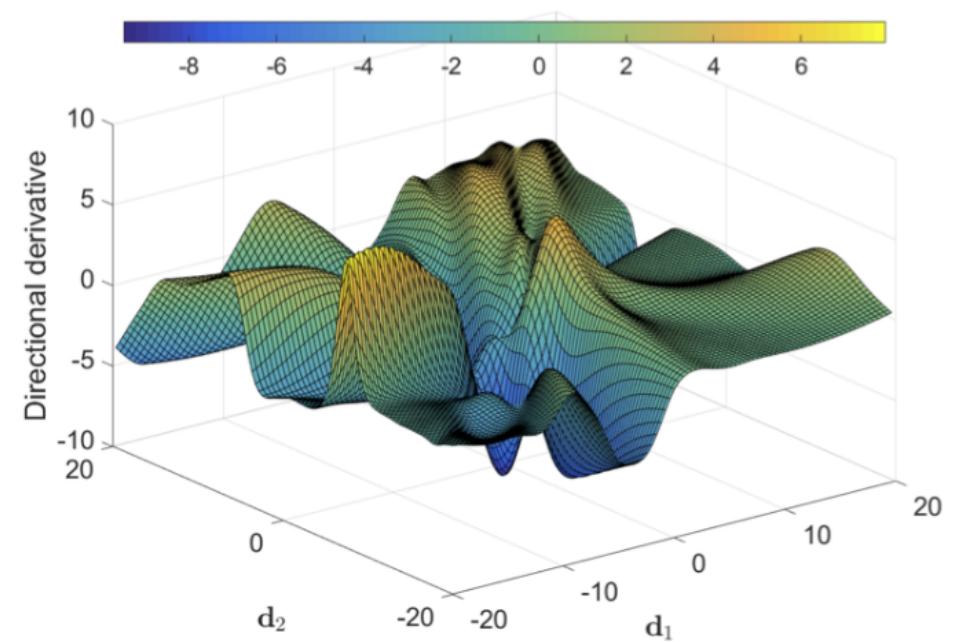
Role of Activation Functions

TANH

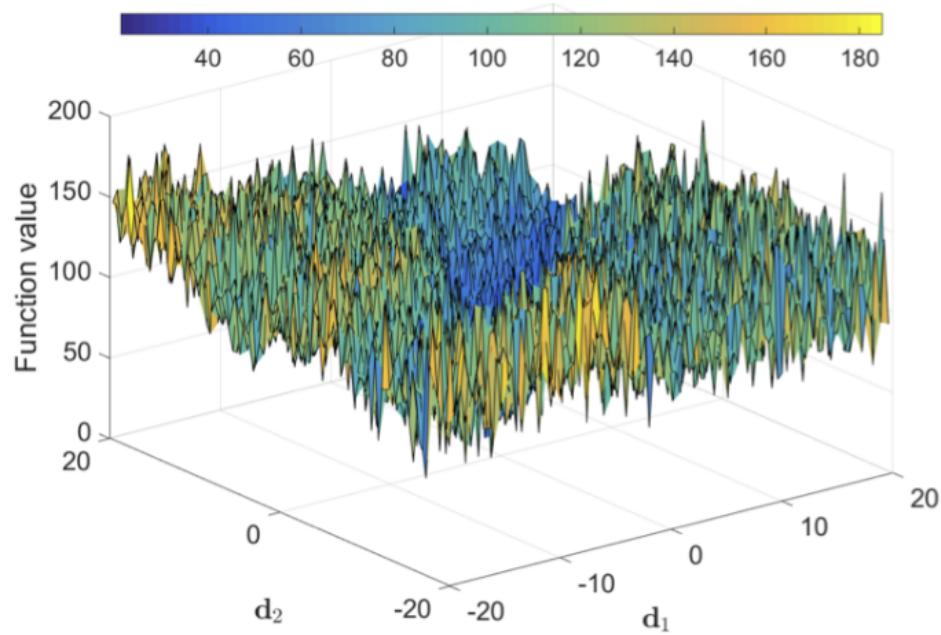




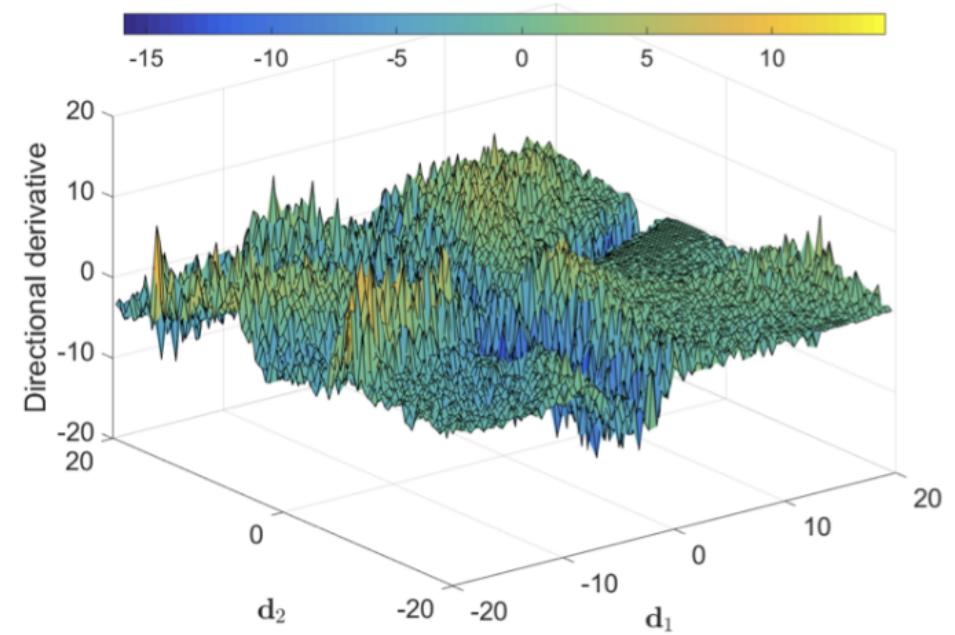
(a) Function value, $|\mathcal{M}| = 150$



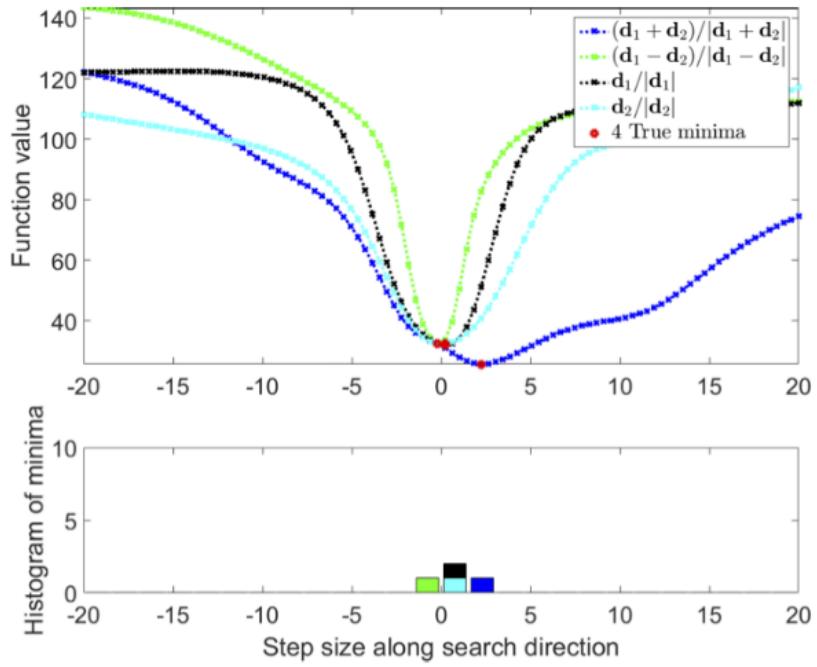
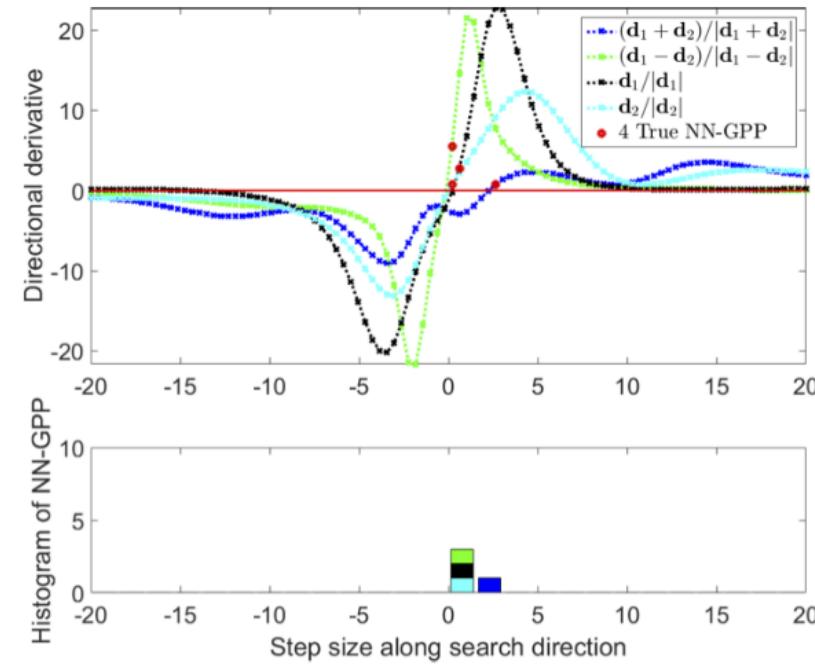
(b) Directional derivative, $|\mathcal{M}| = 150$



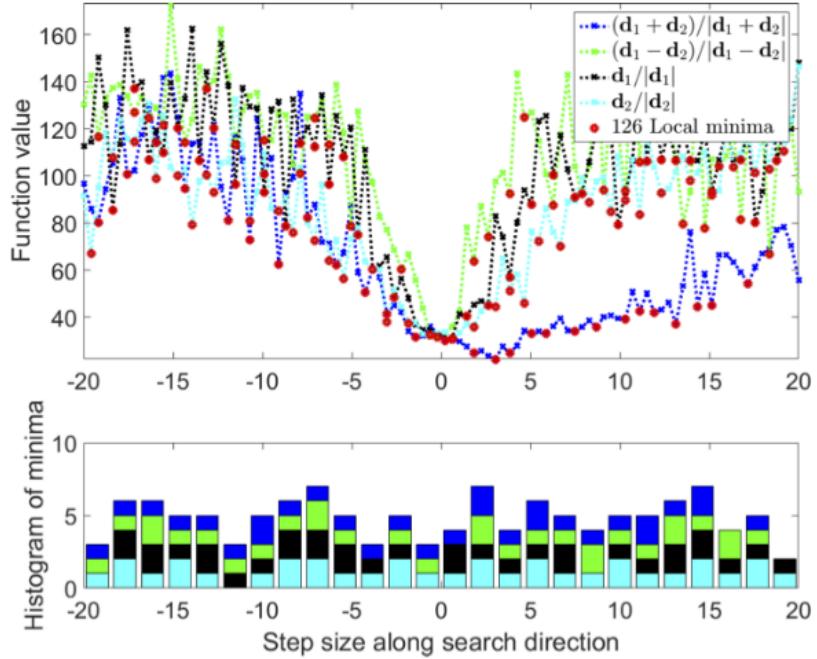
(c) Function value, $|\mathcal{B}| = 10$



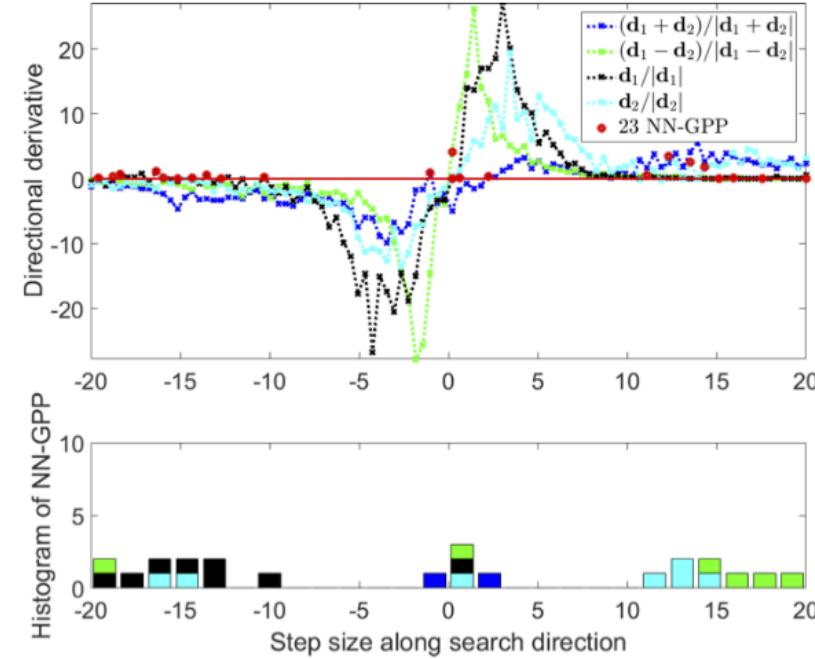
(d) Directional derivative, $|\mathcal{B}| = 10$

(c) Function value, $|B| = 10$ (d) Directional derivative, $|B| = 10$ 

(e) True minima along search directions



(f) True NN-GPP along search directions

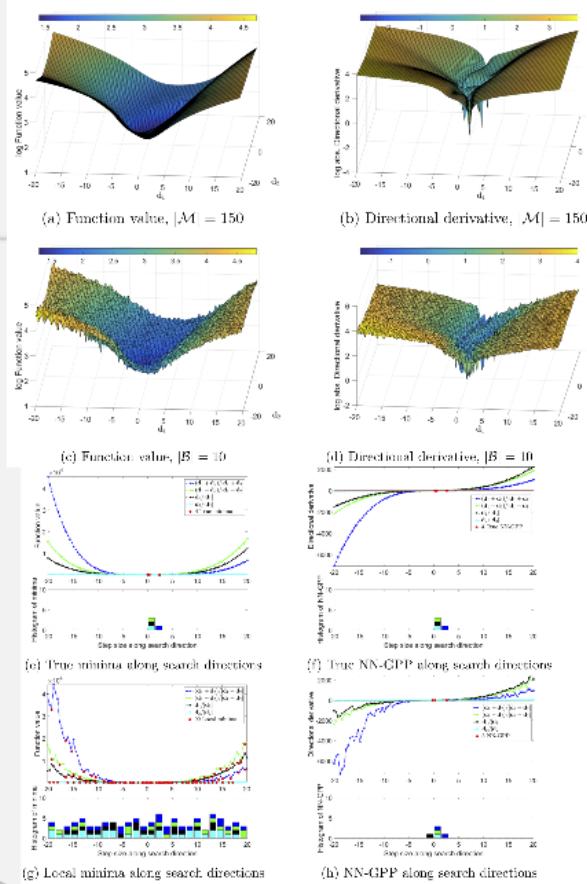


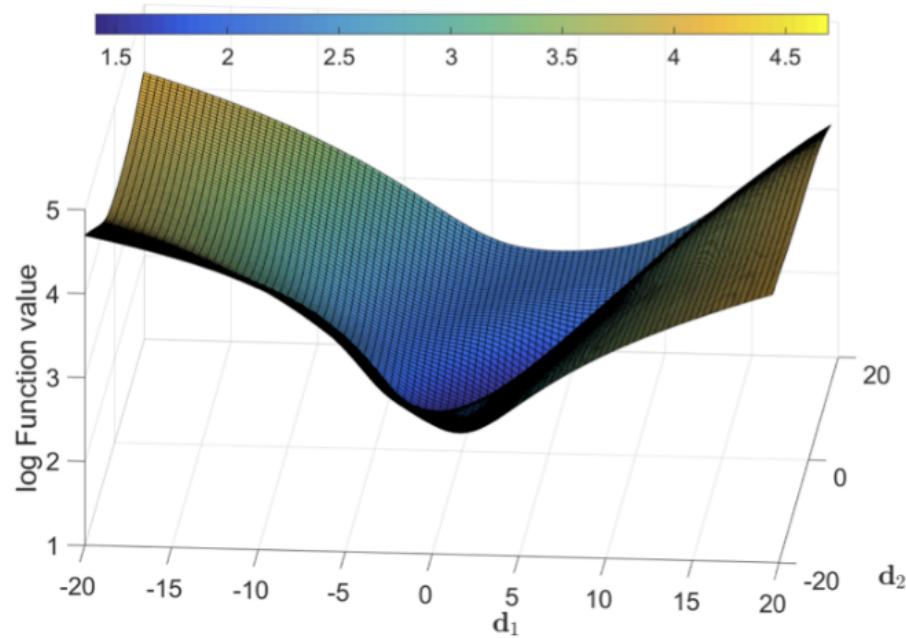
(g) Local minima along search directions

(h) NN-GPP along search directions

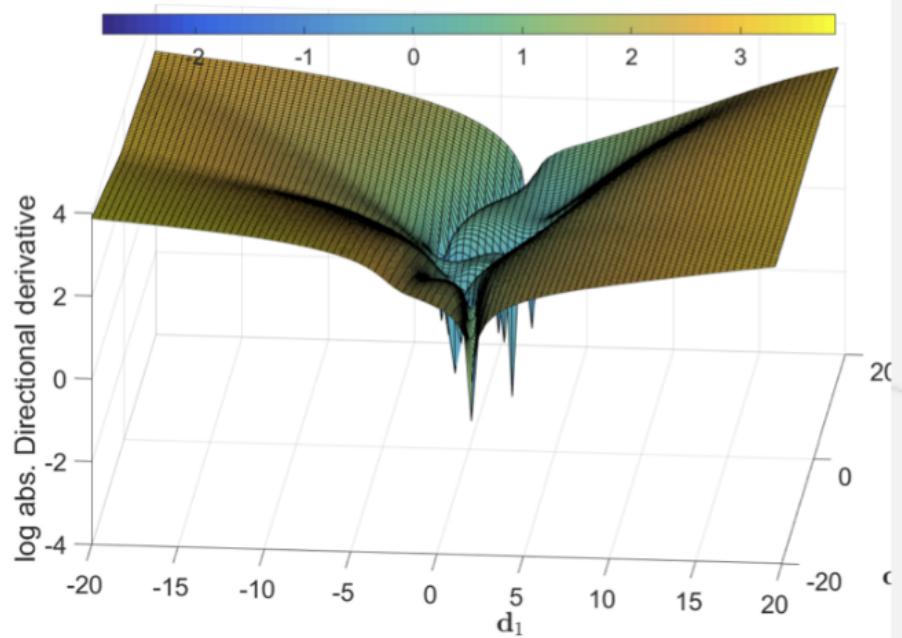
Role of Activation Functions

ELU

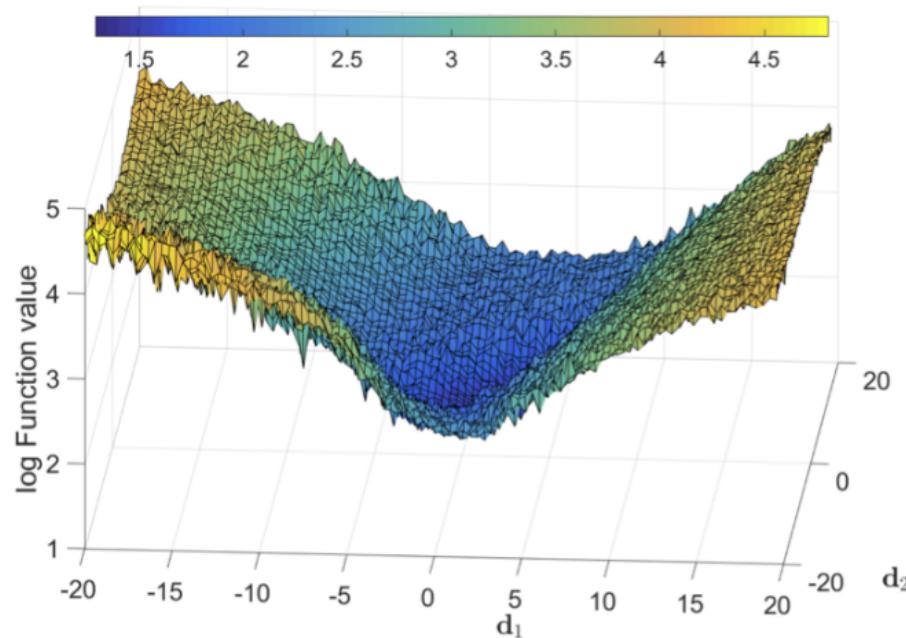




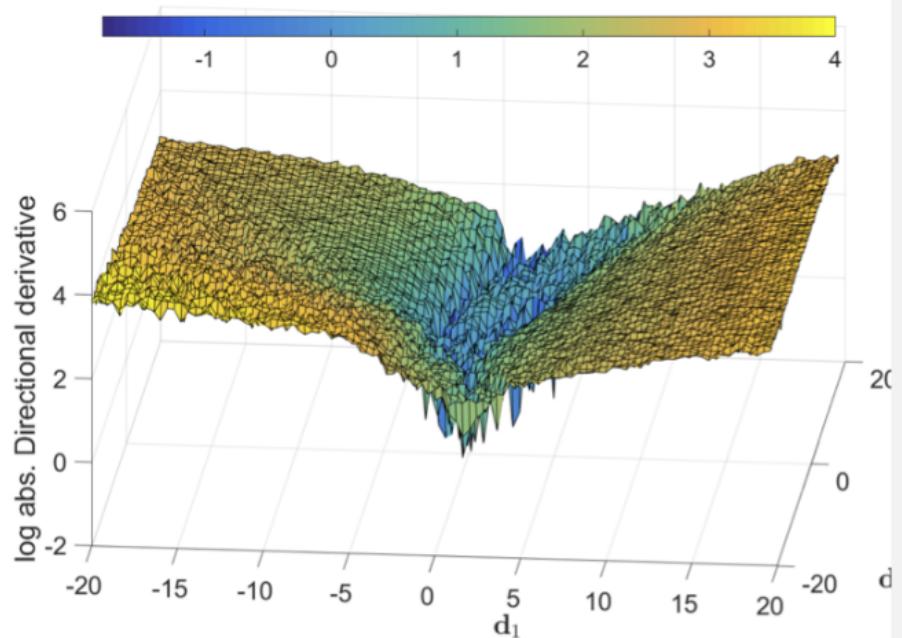
(a) Function value, $|\mathcal{M}| = 150$



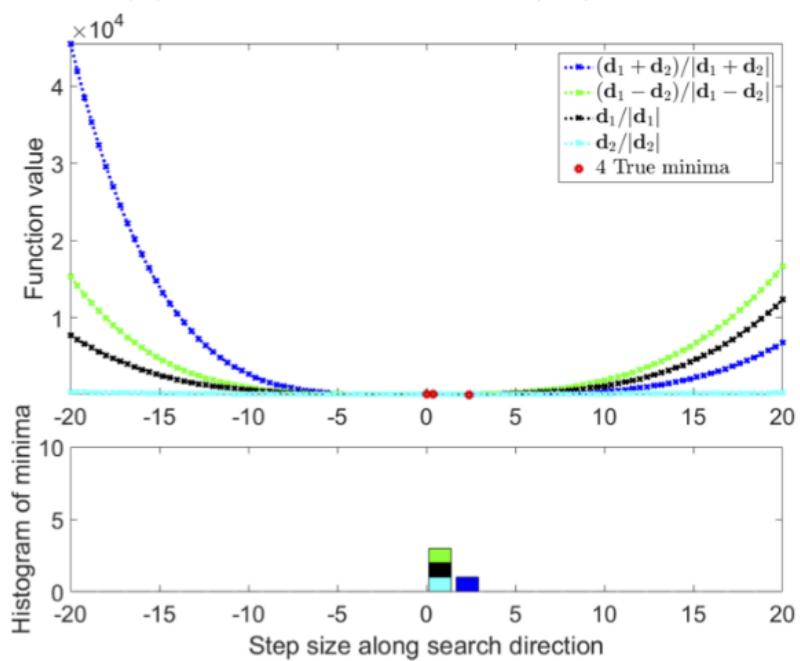
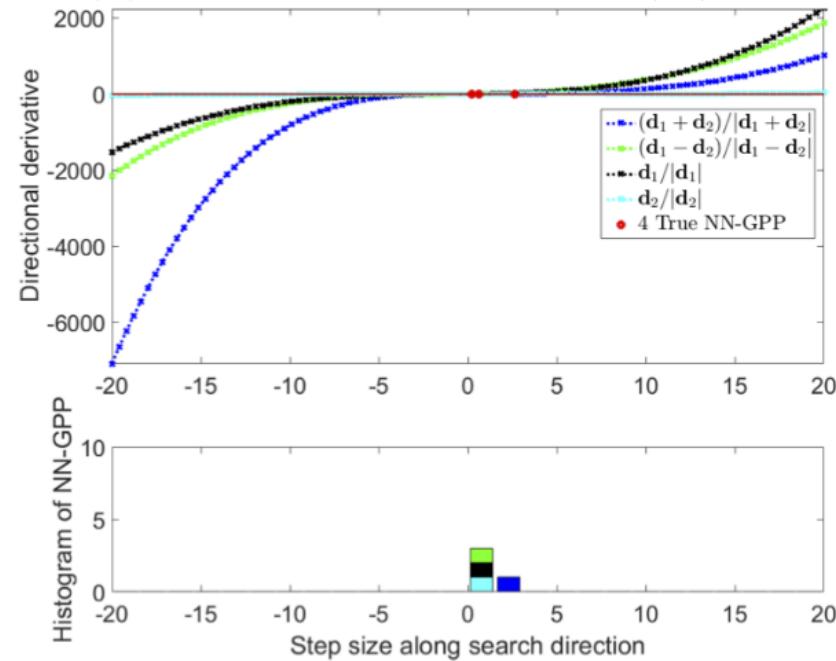
(b) Directional derivative, $|\mathcal{M}| = 150$



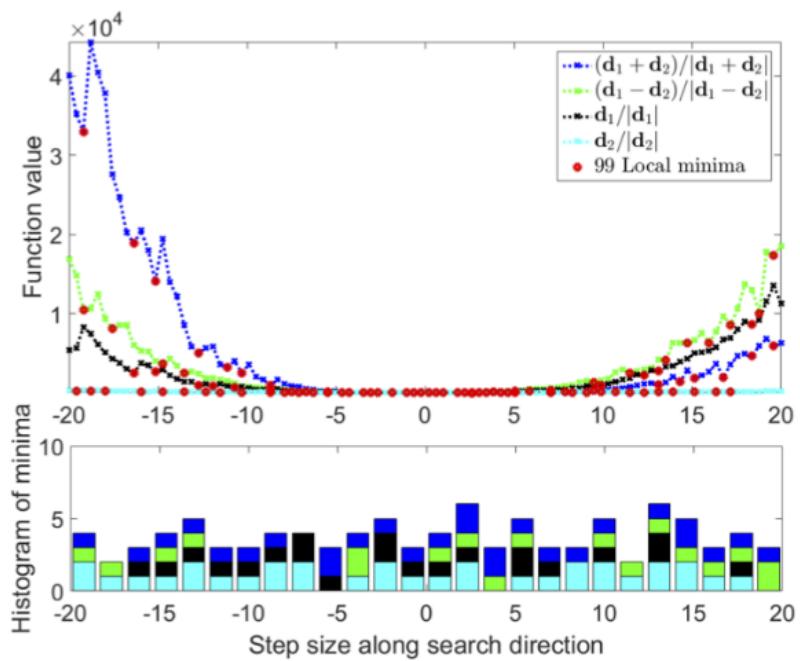
(c) Function value, $|\mathcal{B}| = 10$



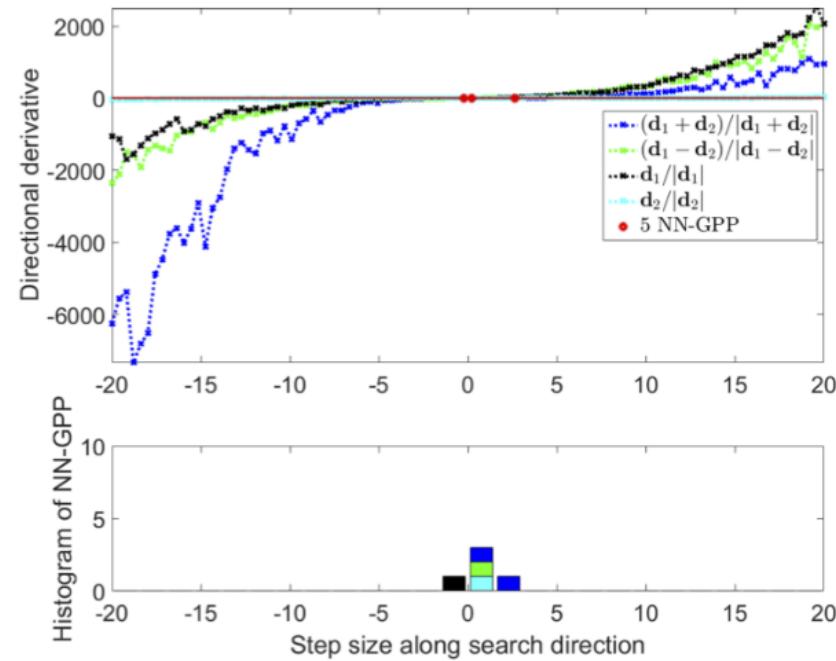
(d) Directional derivative, $|\mathcal{B}| = 10$

(c) Function value, $|\mathcal{B}| = 10$ (d) Directional derivative, $|\mathcal{B}| = 10$ 

(e) True minima along search directions



(f) True NN-GPP along search directions

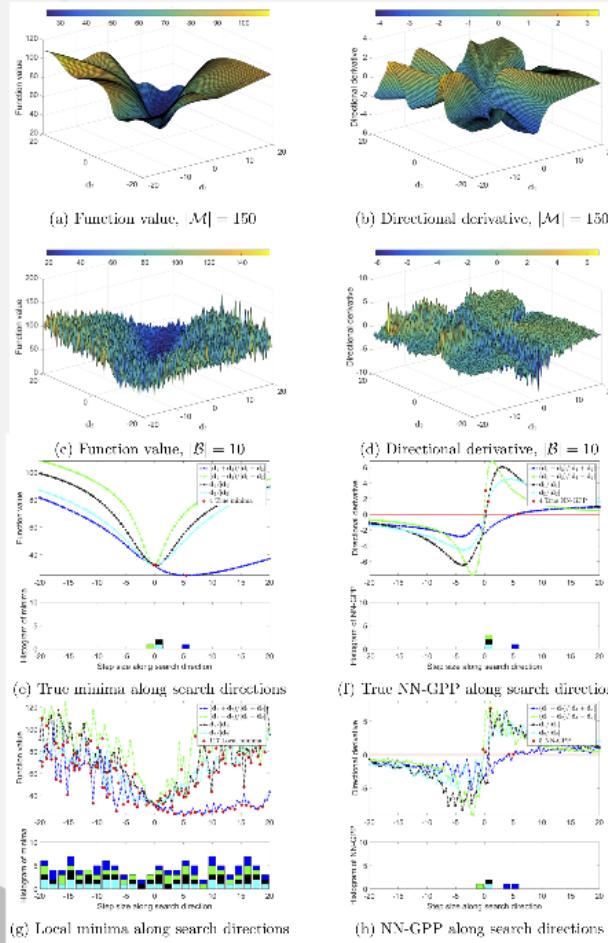


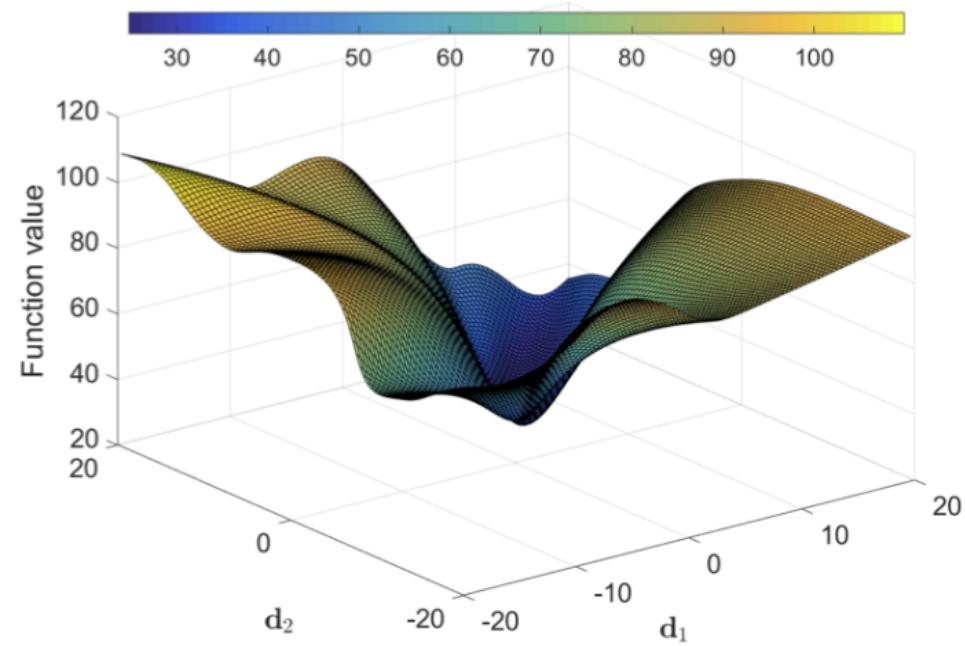
(g) Local minima along search directions

(h) NN-GPP along search directions

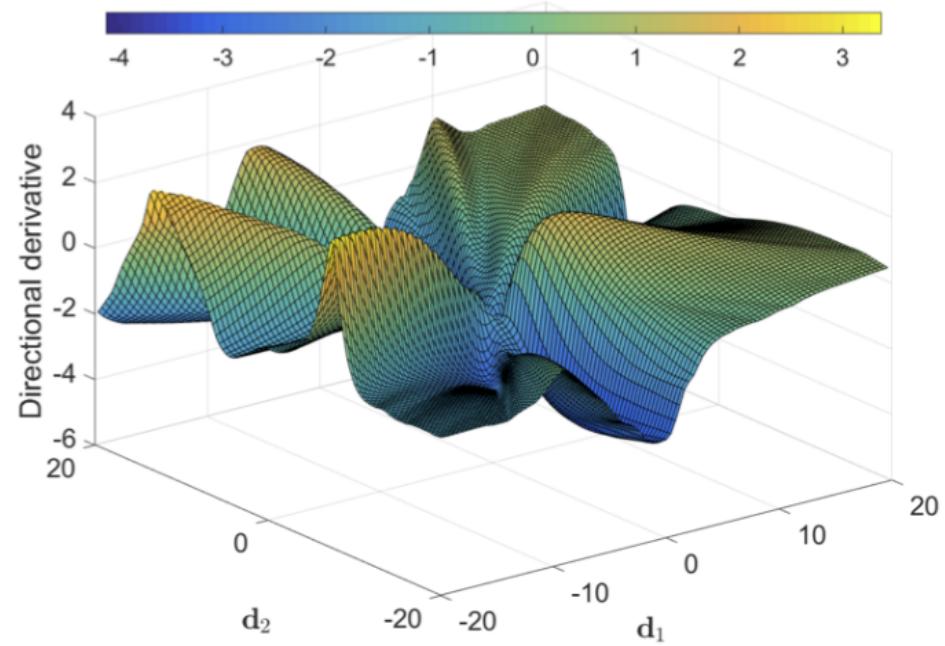
Role of Activation Functions

SOFTSIGN

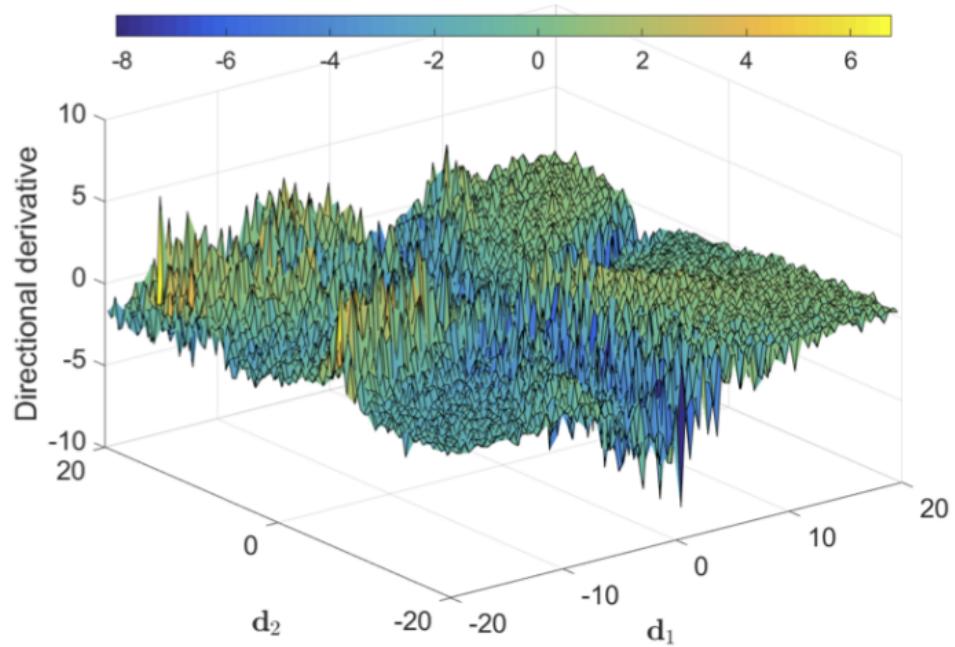
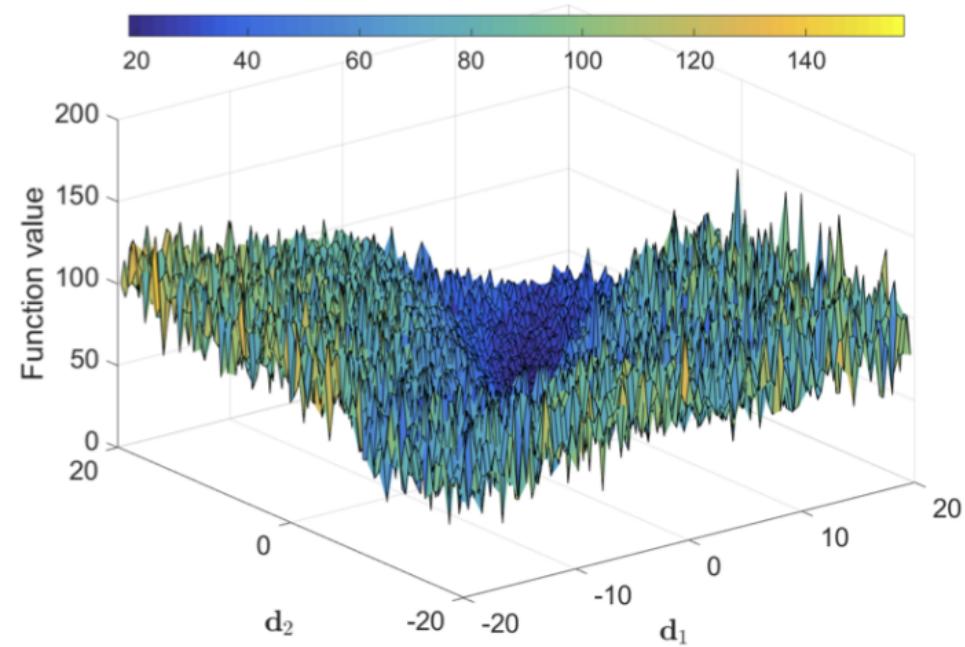




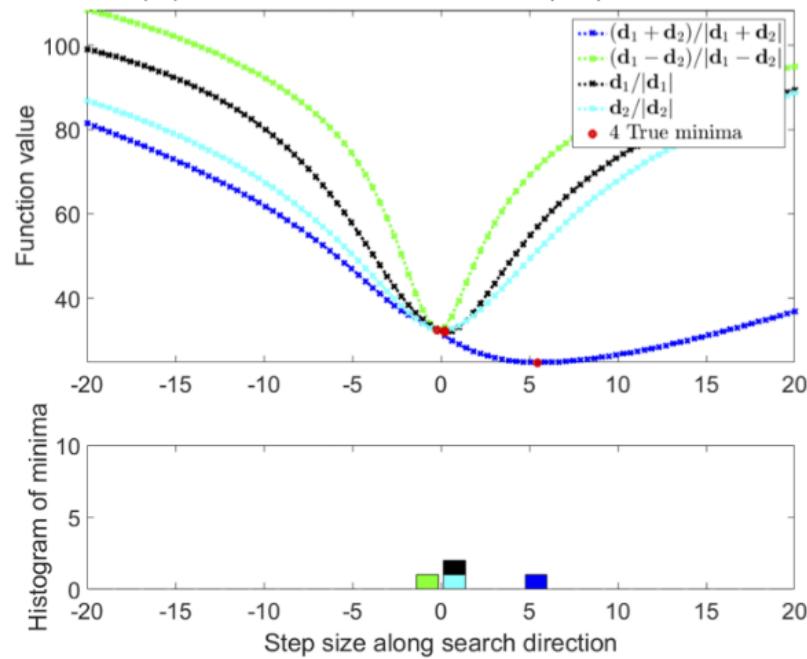
(a) Function value, $|\mathcal{M}| = 150$



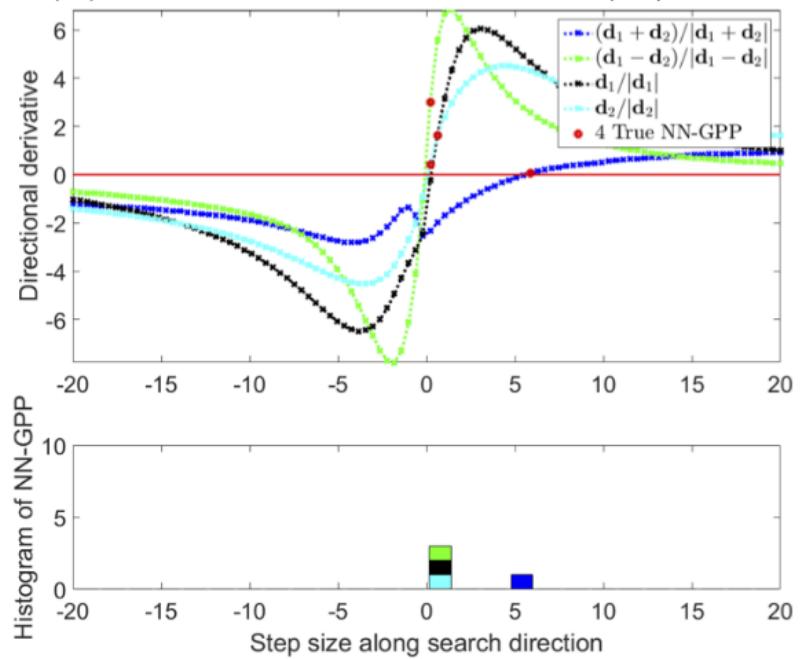
(b) Directional derivative, $|\mathcal{M}| = 150$



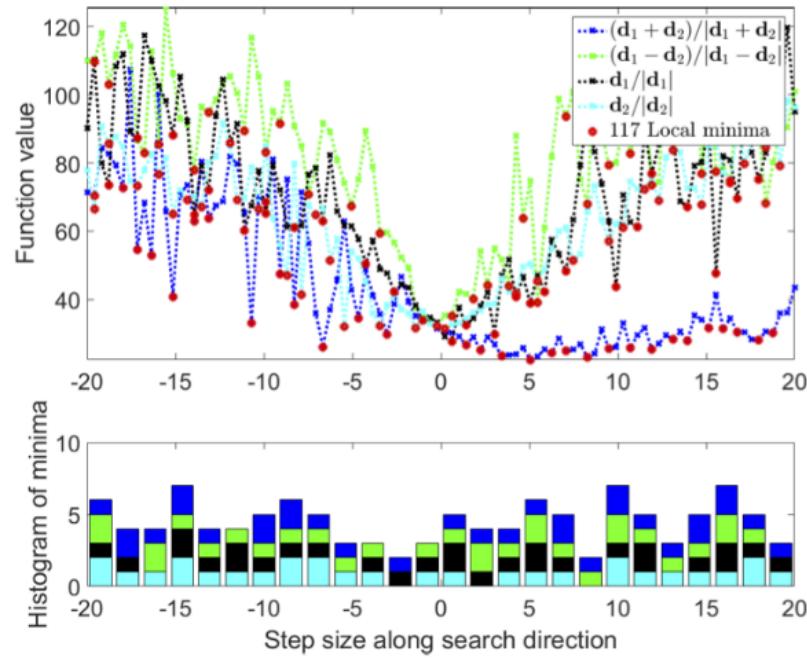
(c) Function value, $|\mathcal{B}| = 10$



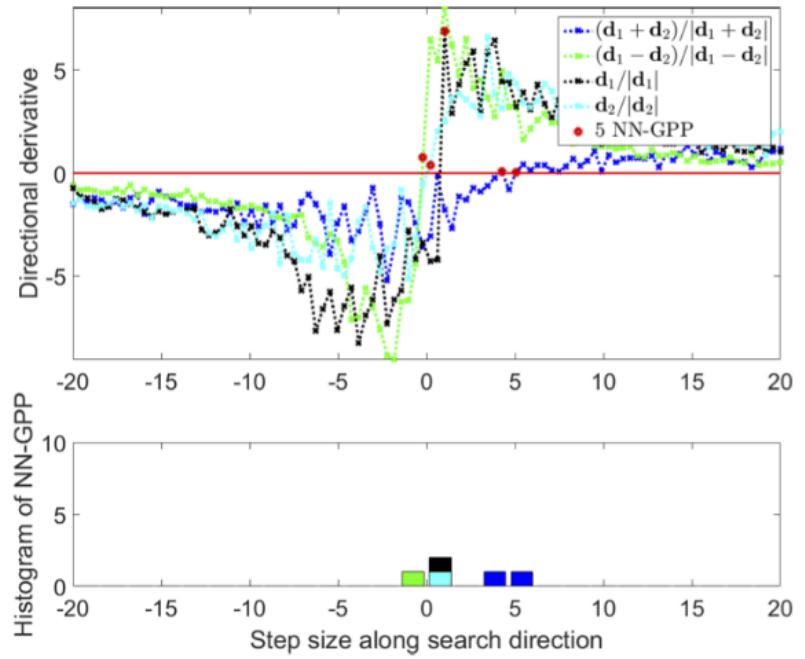
(d) Directional derivative, $|\mathcal{B}| = 10$



(e) True minima along search directions



(f) True NN-GPP along search directions

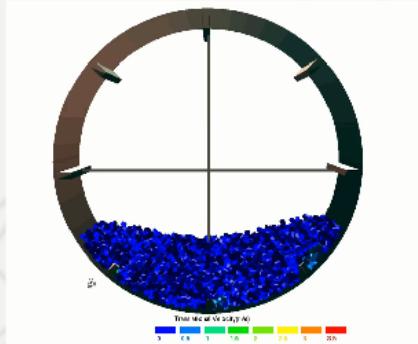


(g) Local minima along search directions

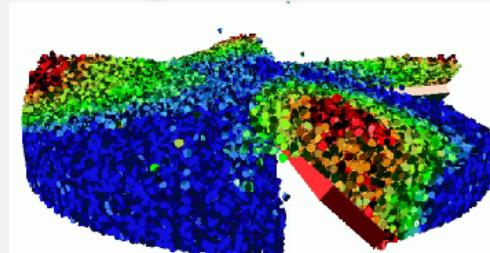
(h) NN-GPP along search directions

Actual Engineering Problems

Comminution

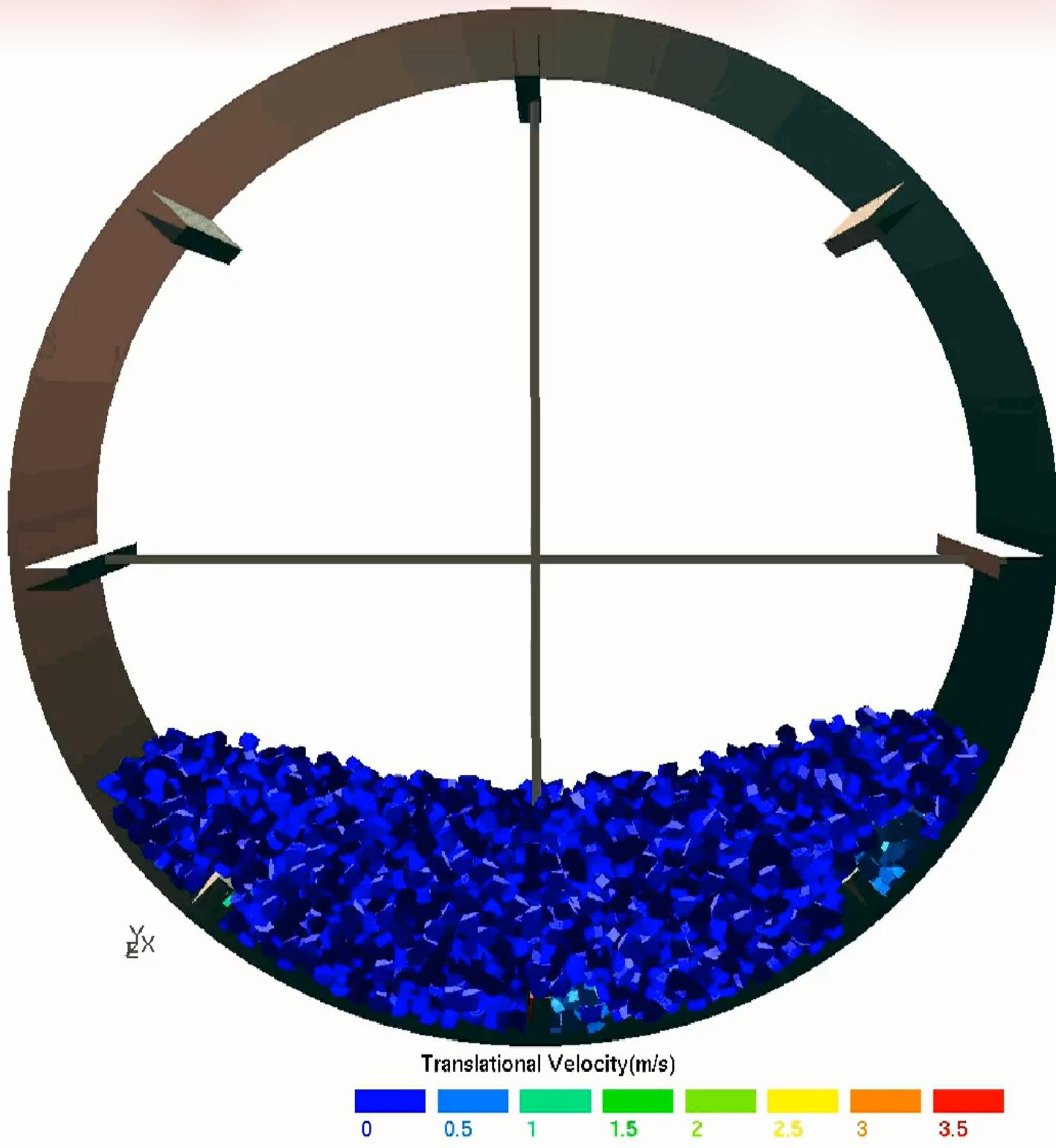


Mixing



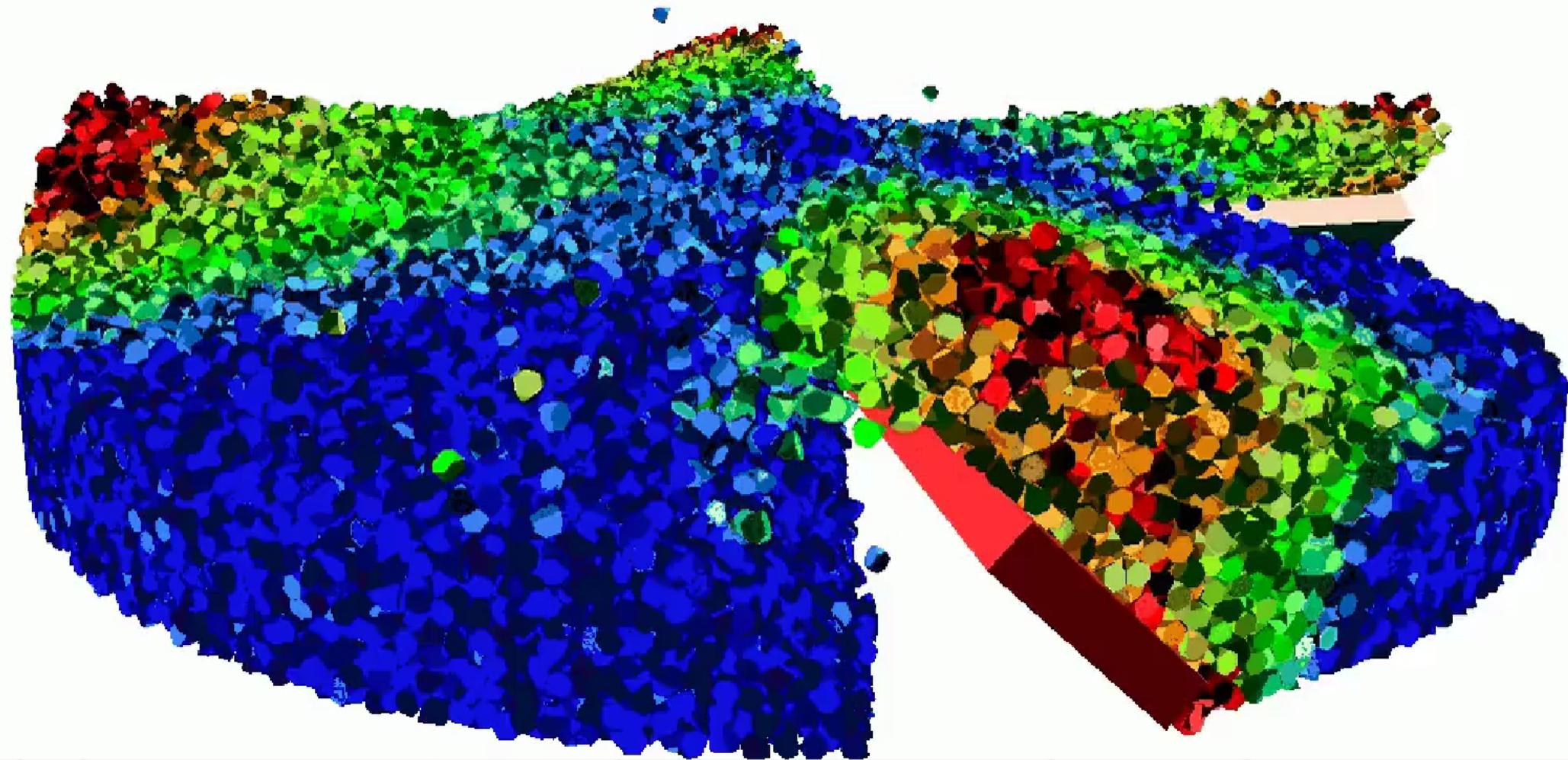
Settling

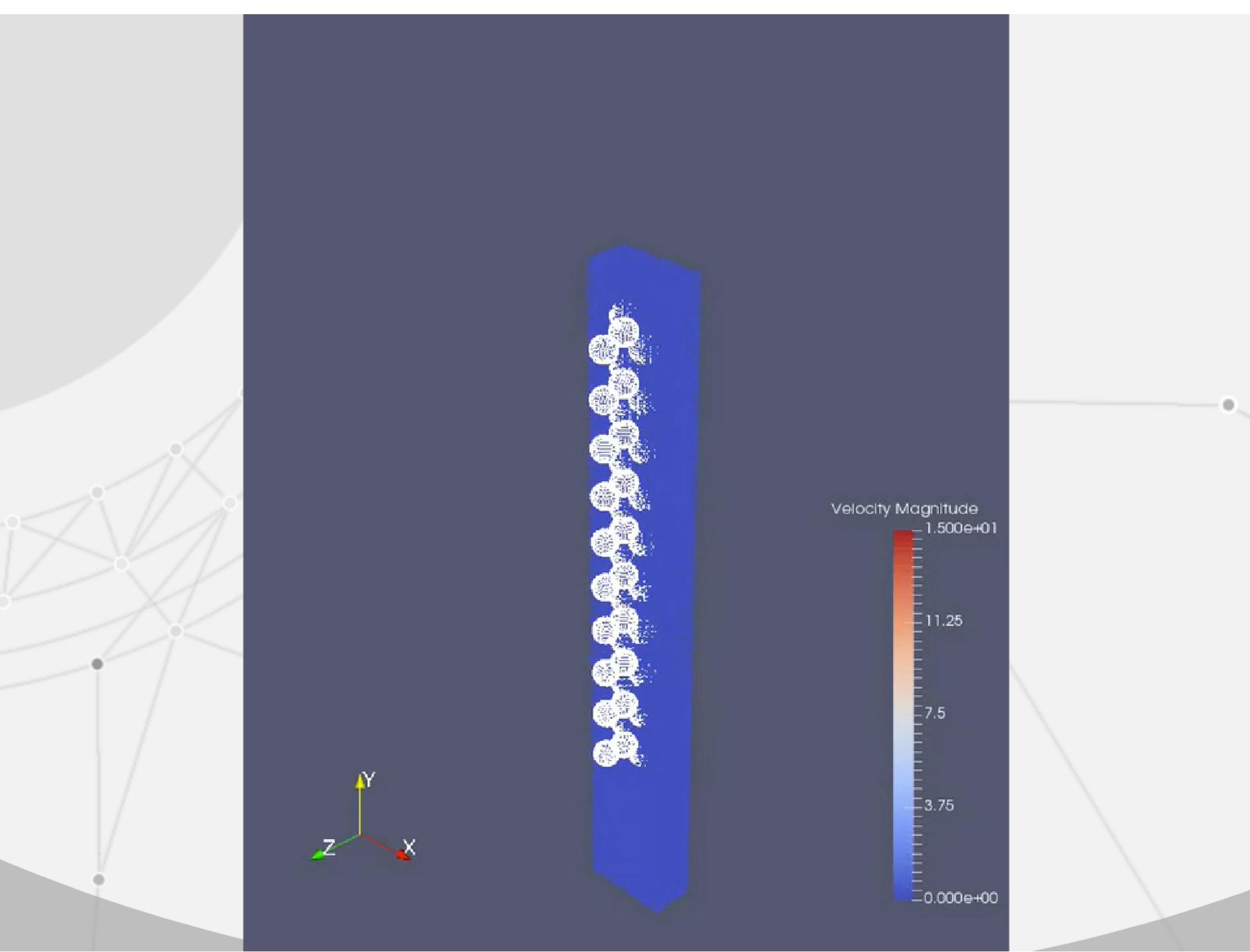




Initi

ty





Simulation BlazeDEM-GPU

Developed in conjunction with the  and UP



GPU Based Particle Solver
Discrete Element Method (Solid Particles)
Smoothed Particle Hydrodynamics (Fluid)

Problem Data Files GB - TB
What do we learn?



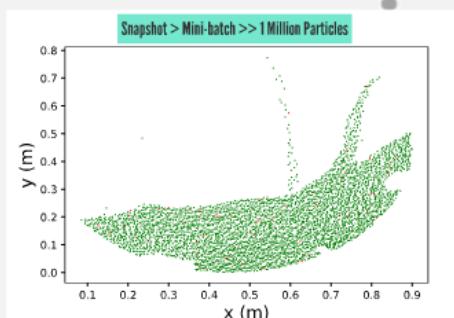
Simulation can take Days/Weeks/Months/Years What do we learn?

Mill breaks solids into pieces by grinding,
crushing and cutting.

Milling harsh environments - cannot look inside.



Sentinel Copper Mine, Zambia



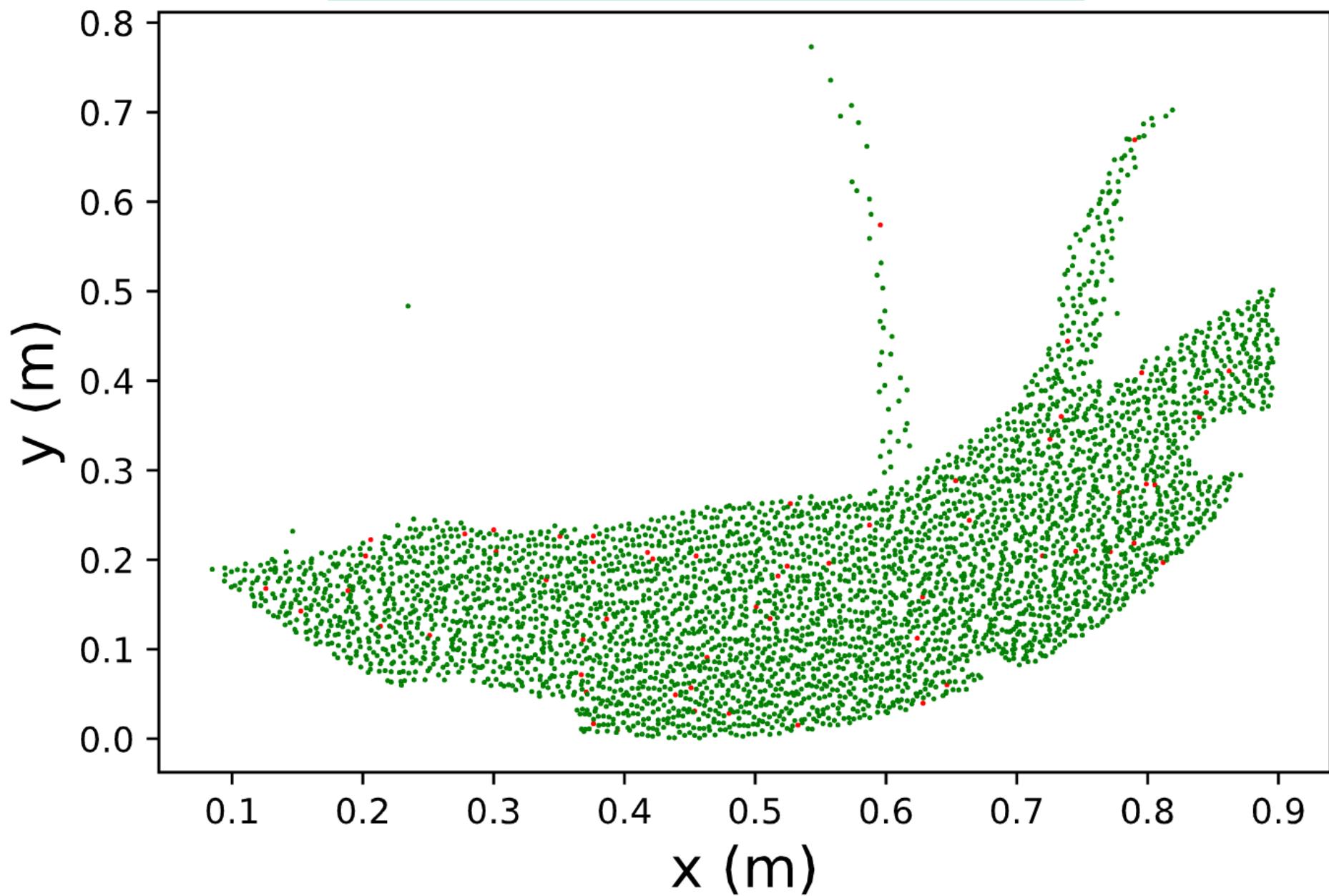


Siemens Gearless Mill Drives

High availability, lower operational costs, best performance.

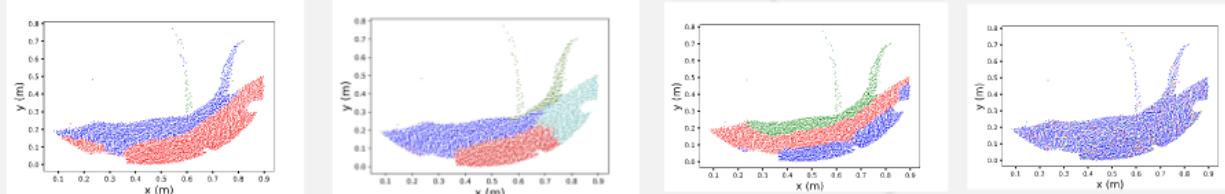
Sentinel Copper Mine, Zambia

Snapshot > Mini-batch >> 1 Million Particles

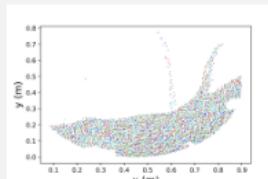


Typical Problem Weeks of Simulations What do we learn?

Use Statistical learning, Machine Learning and Deep Learning to help engineers understand

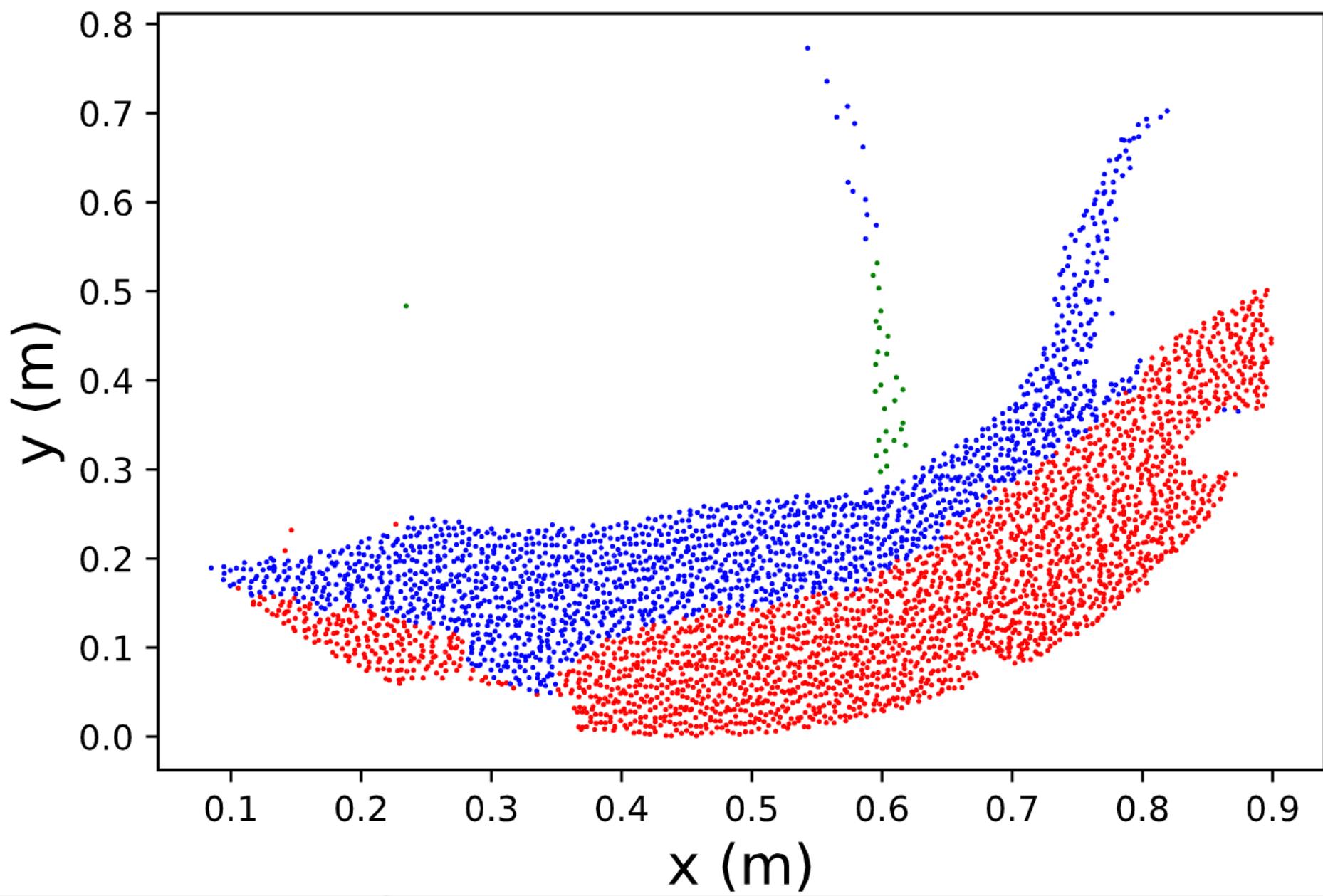


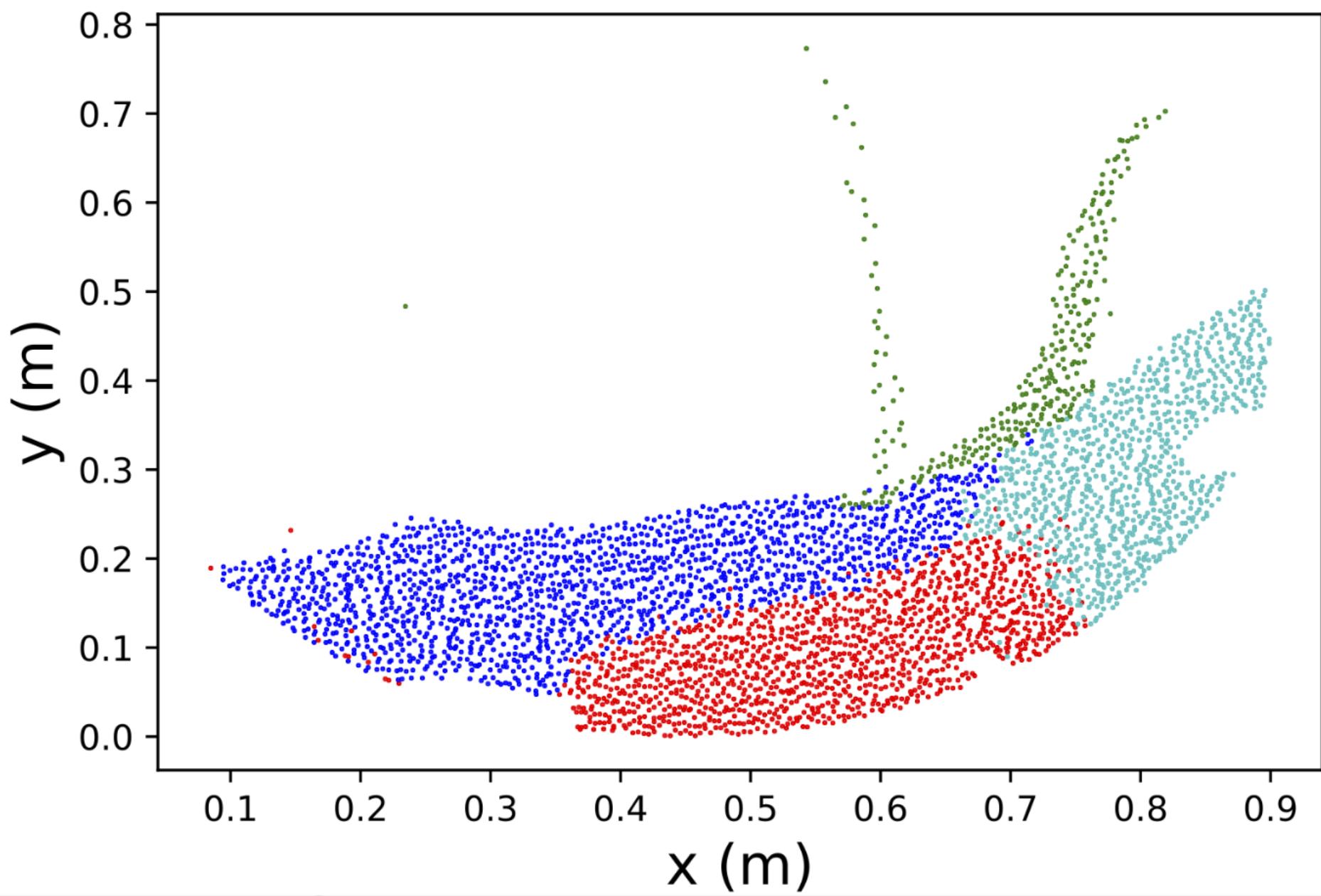
as opposed to the engineer figuring it out

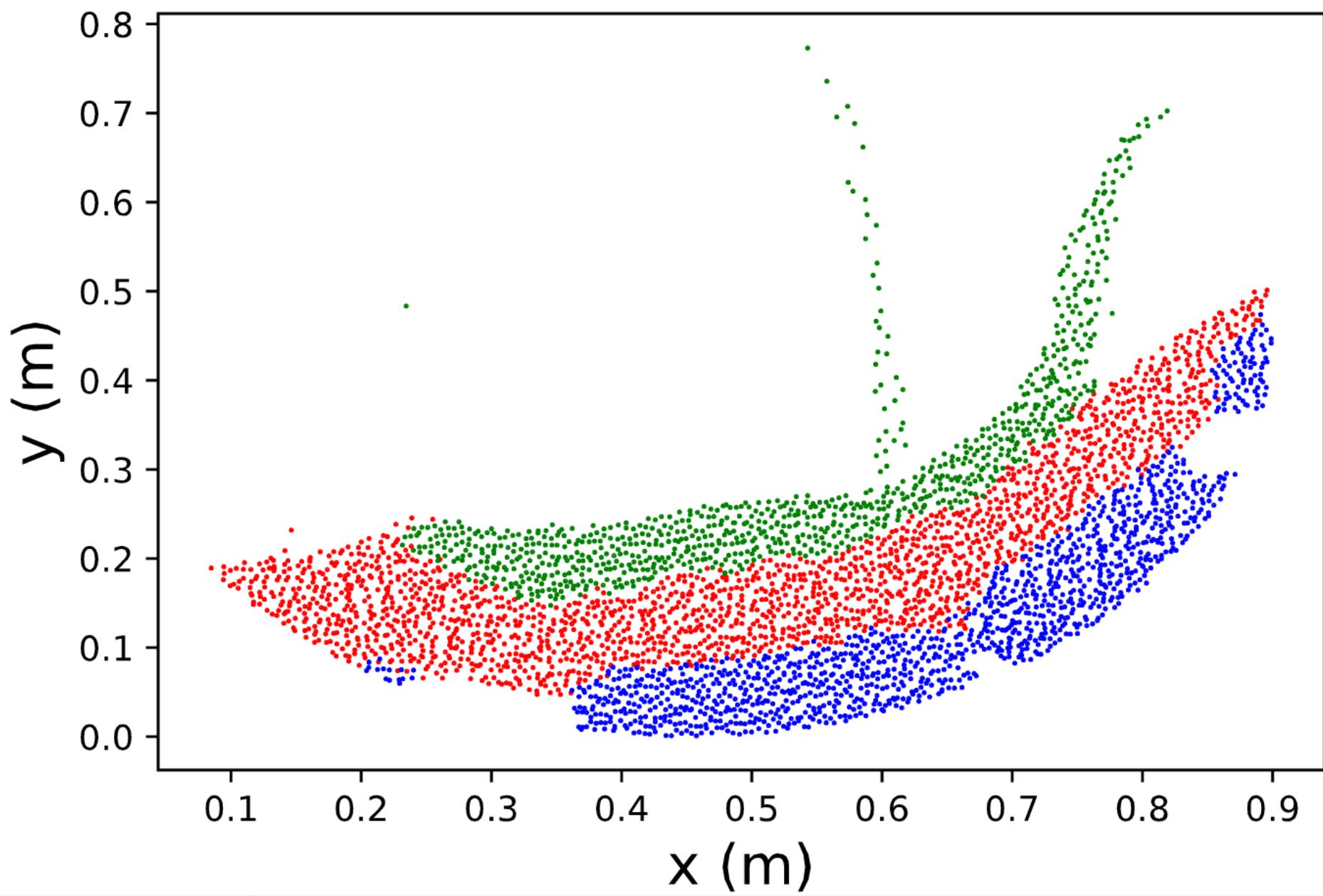


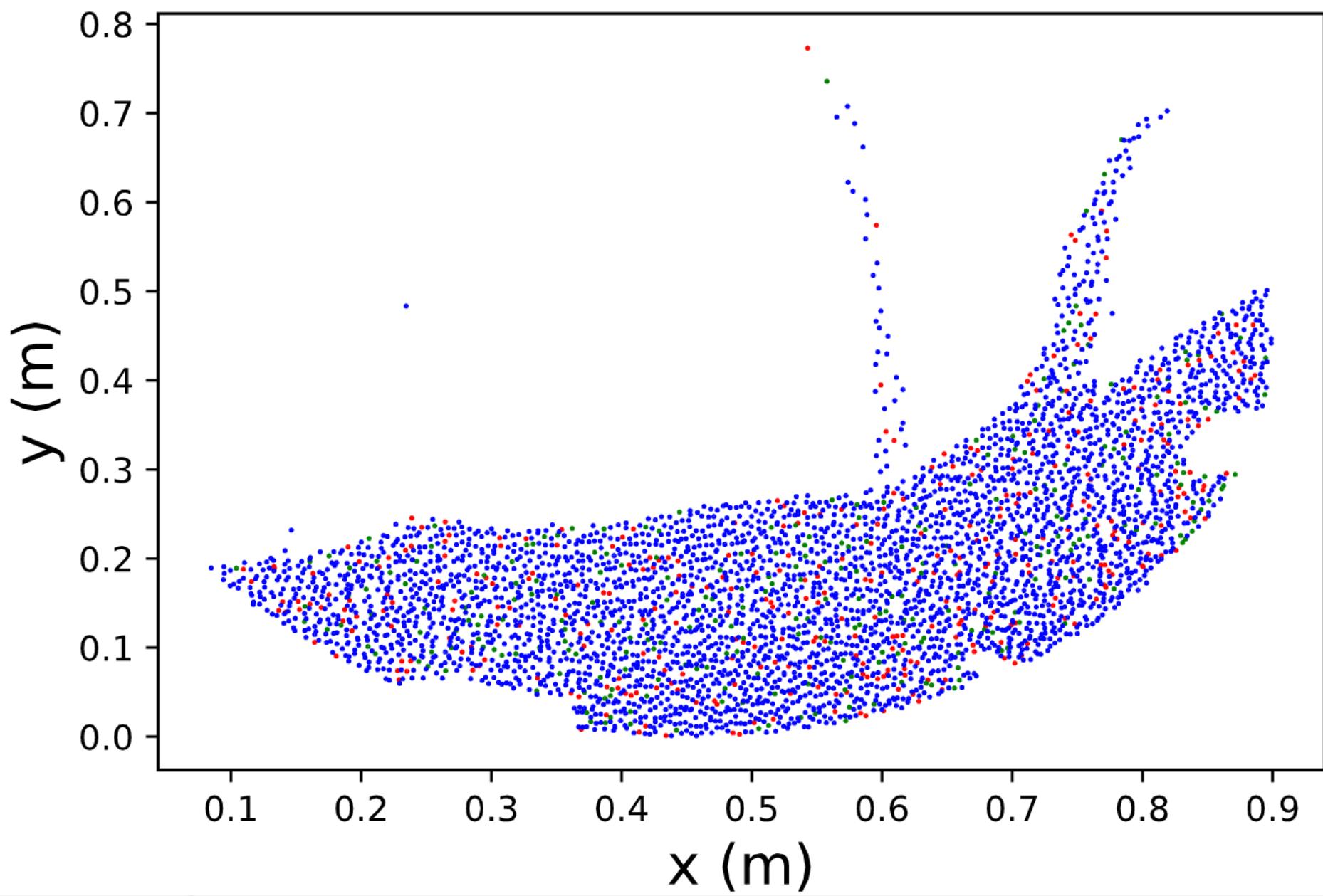
Dataset will be released
with paper

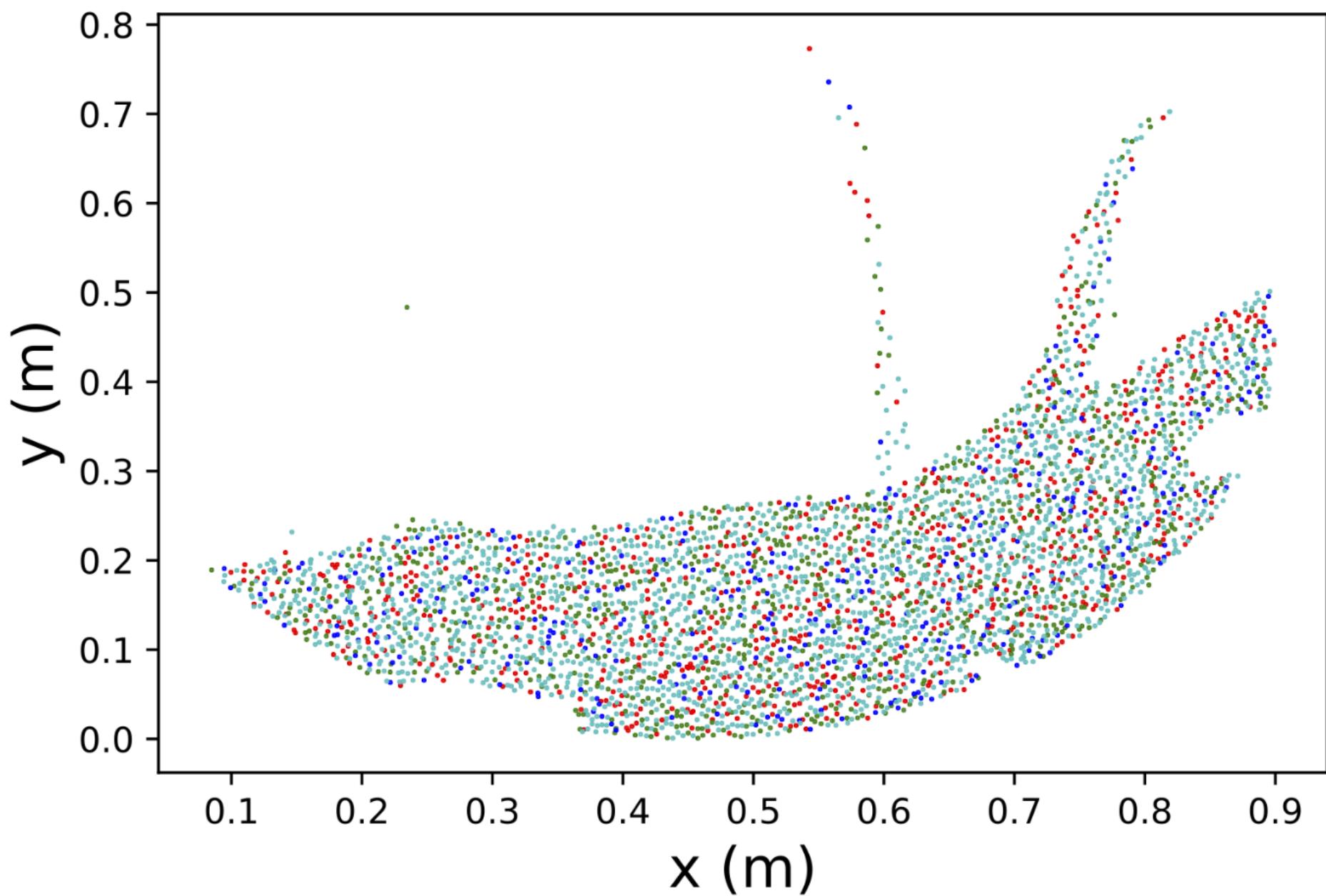












Automating Learning Rate

One step closer towards automating the processing of simulation data for engineers

Research compare idea with optimal learning rate against another idea with optimal learning rate

Take home:

Find minima in the derivative world
NN-GPP

Papers, Pytorch,TF Code Soon
gorgthelab.github.io



GORG THE LAB

gorgthelab.github.io



Dominic
Kafka (PhD)



Younghwan
Chae (PhD)



GORG THE LAB

gorgthelab.github.io



Demo

BLAZEDEM-GPU

THE TEAM

Primary Developer



Prof. Raj Rajamani, Utah University, USA

Prof. Ugur Tuzun, Cantab, UK

Dr. Patrick Pizette, IMT Lille Douai, France

Prof. Nor-Edine Abriak, IMT Lille Douai, France

Prof. Charley Wu, University of Surrey, UK

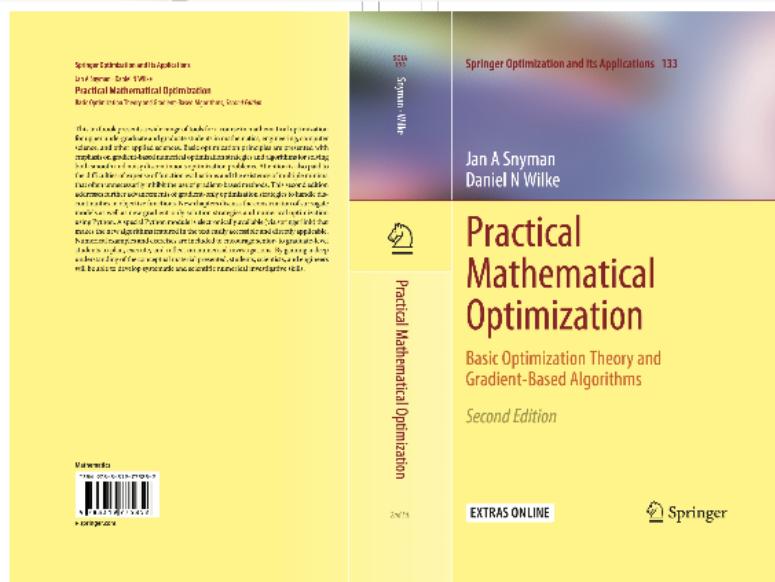
Prof. Wenjie Xu, University of Tsinghua,
China



Line search
Function minim



Do you want to learn more about
Gradient-only Optimization
Non-Negative Projection Points
Line Searches for Discontinuous Functions



Only textbook (2018)
that covers
this material



Springer Optimization and Its Applications

Jan A Snyman · Daniel N Wilke

Practical Mathematical Optimization

Basic Optimization Theory and Gradient-Based Algorithms, *Second Edition*

This textbook presents a wide range of tools for a course in mathematical optimization for upper undergraduate and graduate students in mathematics, engineering, computer science, and other applied sciences. Basic optimization principles are presented with emphasis on gradient-based numerical optimization strategies and algorithms for solving both smooth and noisy discontinuous optimization problems. Attention is also paid to the difficulties of expense of function evaluations and the existence of multiple minima that often unnecessarily inhibit the use of gradient-based methods. This second edition addresses further advancements of gradient-only optimization strategies to handle discontinuities in objective functions. New chapters discuss the construction of surrogate models as well as new gradient-only solution strategies and numerical optimization using Python. A special Python module is electronically available (via springerlink) that makes the new algorithms featured in the text easily accessible and directly applicable. Numerical examples and exercises are included to encourage senior- to graduate-level students to plan, execute, and reflect on numerical investigations. By gaining a deep understanding of the conceptual material presented, students, scientists, and engineers will be able to develop systematic and scientific numerical investigative skills.

SOIA
133

Snyman · Wilke

Springer Optimization and Its Applications 133

**Jan A Snyman
Daniel N Wilke**



Practical Mathematical Optimization

2nd Ed.

Practical Mathematical Optimization

**Basic Optimization Theory and
Gradient-Based Algorithms**

Second Edition

EXTRAS ONLINE

Springer

Mathematics

ISBN 978-3-319-77585-2



► springer.com