

---

# Nonsmooth, Nonconvex Optimization Algorithms and Examples

Michael L. Overton  
Courant Institute of Mathematical Sciences  
New York University

Convex and Nonsmooth Optimization Class, Spring 2018, Final Lecture

Based on my research work with Jim Burke (Washington), Adrian Lewis (Cornell)  
and others



Introduction

Nonsmooth,

Nonconvex

Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

# Introduction



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

## ■ Continuous



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

## Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex
- Usually, but not always, locally Lipschitz: for all  $x$  there exists  $L_x$  s.t.  $|f(y) - f(z)| \leq L_x \|y - z\|$  for all  $y, z$  near  $x$



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex
- Usually, but not always, locally Lipschitz: for all  $x$  there exists  $L_x$  s.t.  $|f(y) - f(z)| \leq L_x \|y - z\|$  for all  $y, z$  near  $x$
- Otherwise, no structure assumed



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex
- Usually, but not always, locally Lipschitz: for all  $x$  there exists  $L_x$  s.t.  $|f(y) - f(z)| \leq L_x \|y - z\|$  for all  $y, z$  near  $x$
- Otherwise, no structure assumed

Lots of interesting applications



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex
- Usually, but not always, locally Lipschitz: for all  $x$  there exists  $L_x$  s.t.  $|f(y) - f(z)| \leq L_x \|y - z\|$  for all  $y, z$  near  $x$
- Otherwise, no structure assumed

Lots of interesting applications

Any locally Lipschitz function is differentiable almost everywhere on its domain. So, whp, can evaluate gradient at any given point.



# Nonsmooth, Nonconvex Optimization

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Problem: find  $x$  that locally minimizes  $f$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- Continuous
- Not differentiable everywhere, in particular often not differentiable at local minimizers
- Not convex
- Usually, but not always, locally Lipschitz: for all  $x$  there exists  $L_x$  s.t.  $|f(y) - f(z)| \leq L_x \|y - z\|$  for all  $y, z$  near  $x$
- Otherwise, no structure assumed

Lots of interesting applications

Any locally Lipschitz function is differentiable almost everywhere on its domain. So, whp, can evaluate gradient at any given point.

What happens if we simply use gradient descent (steepest descent) with a standard line search?



# A Simple Nonconvex Example

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

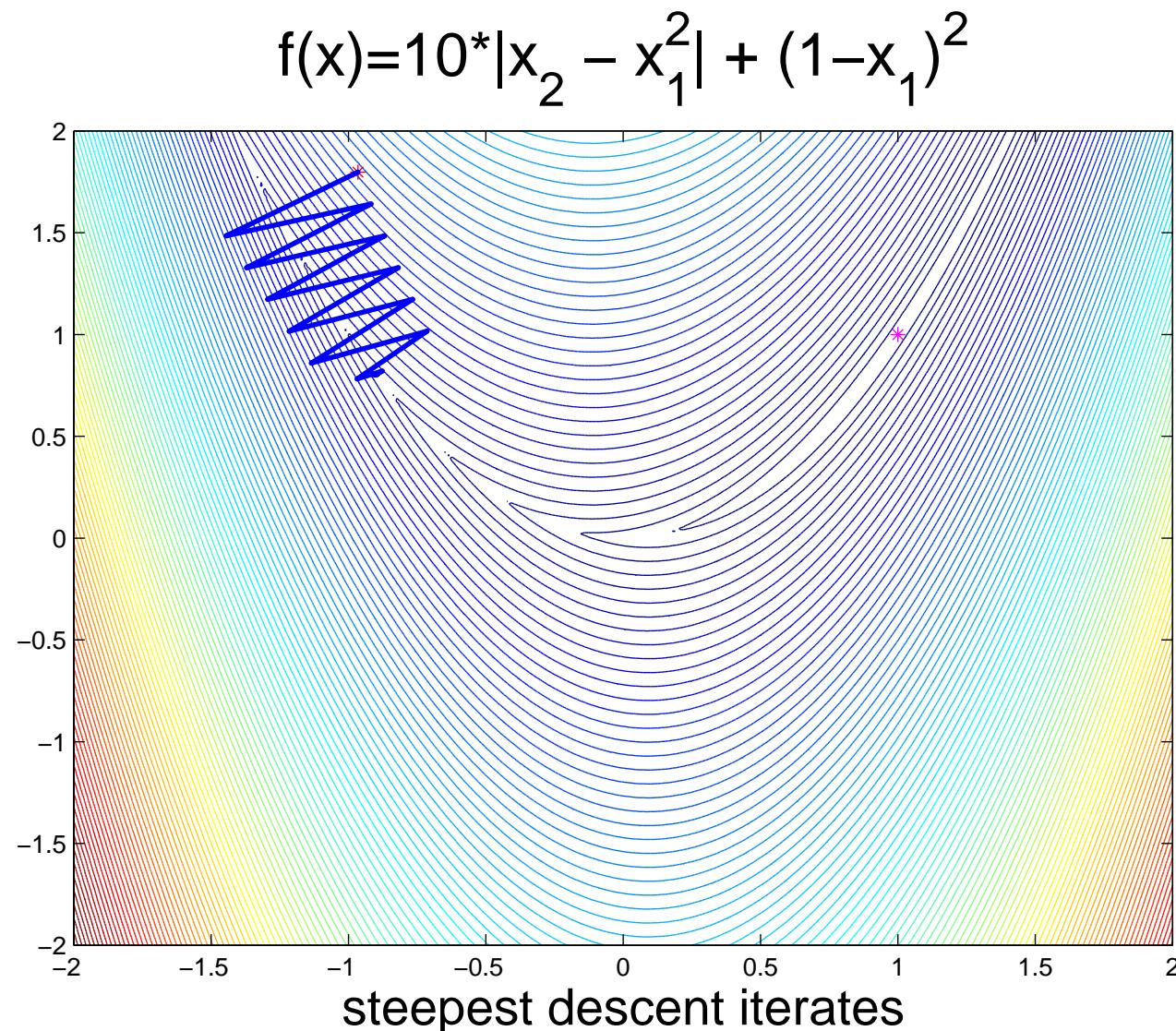
Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks





# A Simple Nonconvex Example

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

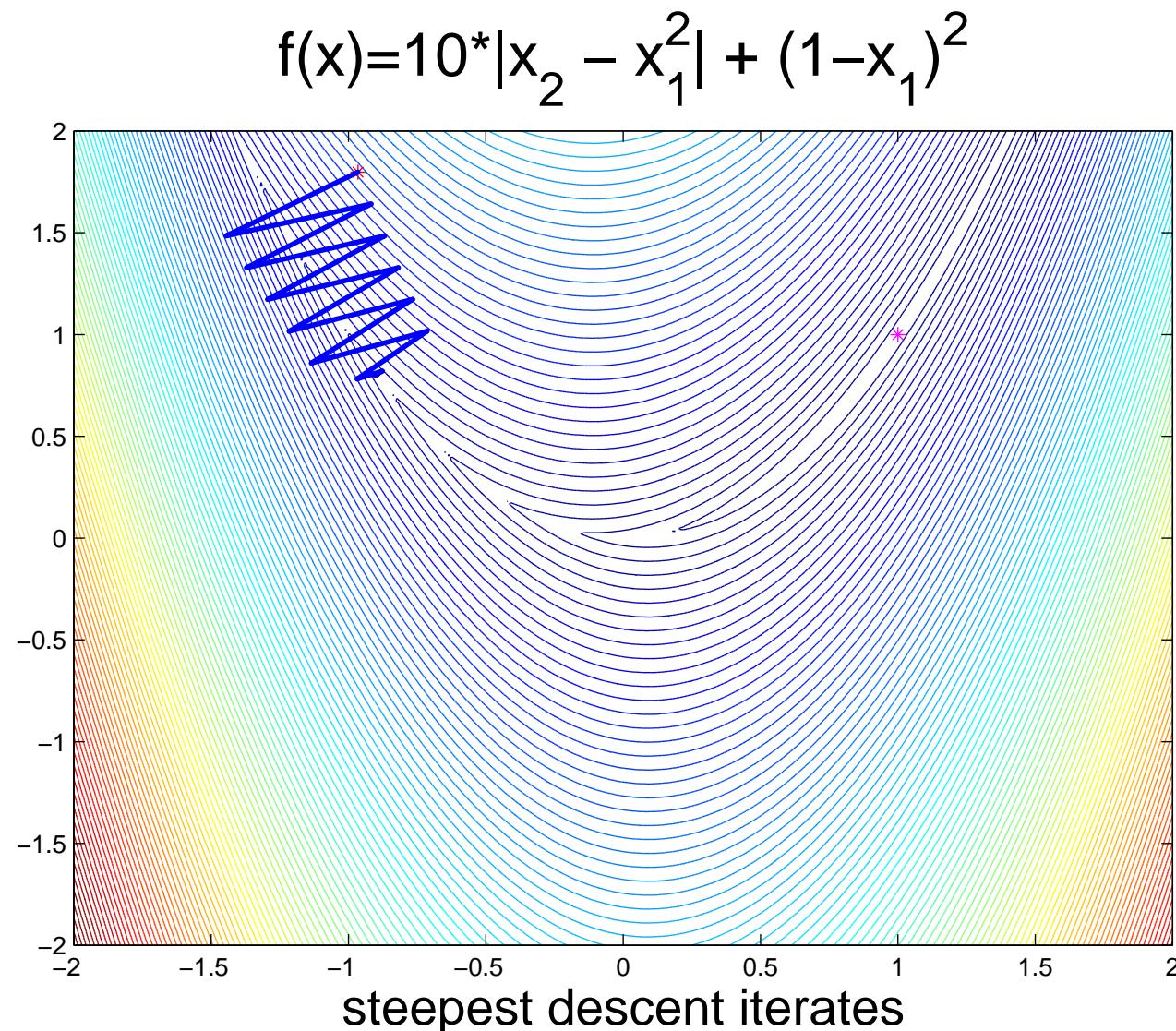
Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



On this example, iterates invariably converge to a nonstationary point



# Failure of Gradient Descent in Nonsmooth Case

Known for decades that gradient descent may converge to nonstationary points when  $f$  is nonsmooth, even if it is convex.

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



# Failure of Gradient Descent in Nonsmooth Case

Known for decades that gradient descent may converge to nonstationary points when  $f$  is nonsmooth, even if it is convex.

- V.F. Dem'janov and V.N. Malozemov, 1970
- P. Wolfe, 1975
- J.-B. Hiriart-Urruty and C. Lemaréchal, 1993

But these are all examples cooked up to defeat exact line searches from a specific starting point.

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



# Failure of Gradient Descent in Nonsmooth Case

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Known for decades that gradient descent may converge to nonstationary points when  $f$  is nonsmooth, even if it is convex.

- V.F. Dem'janov and V.N. Malozemov, 1970
- P. Wolfe, 1975
- J.-B. Hiriart-Urruty and C. Lemaréchal, 1993

But these are all examples cooked up to defeat exact line searches from a specific starting point.

Failure can be avoided by using sufficiently short steplengths (N.Z. Shor, 1970s), but this is slow.



# Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



## Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

- sufficient decrease in function value:

$$f(x + td) < f(x) + c_1 t \nabla f(x)^T d \text{ (L. Armijo, 1966)}$$

Introduction  
Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
**Armijo-Wolfe Line  
Search**  
Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



## Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

- sufficient decrease in function value:  
$$f(x + td) < f(x) + c_1 t \nabla f(x)^T d \quad (\text{L. Armijo, 1966})$$
- sufficient increase in directional derivative:  $f$  is differentiable at  $x + td$  and  $\nabla f(x + td)^T d > c_2 \nabla f(x)^T d$  ( $P.$  Wolfe, 1969)

---

Introduction  
Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

---

Gradient Sampling

---

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods

---

Concluding Remarks



# Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

- sufficient decrease in function value:  
$$f(x + td) < f(x) + c_1 t \nabla f(x)^T d \quad (\text{L. Armijo, 1966})$$
- sufficient increase in directional derivative:  $f$  is differentiable at  $x + td$  and  $\nabla f(x + td)^T d > c_2 \nabla f(x)^T d$  ( $P.$  Wolfe, 1969)

Assuming  $\inf_t f(x + td)$  is bounded below,

- the Armijo condition holds for sufficiently small  $t$  as long as  $f$  is continuous
- the Wolfe condition holds for sufficiently large  $t$  as long as  $f$  is differentiable
- the intervals where each holds overlap

so combining the two conditions leads to a convenient, convergent bracketing line search (M.J.D. Powell, 1976)

---

Introduction

Nonsmooth,

Nonconvex

Optimization

A Simple Nonconvex

Example

Failure of Gradient

Descent in

Nonsmooth Case

Armijo-Wolfe Line

Search

Failure of Gradient

Method: Simple

Convex Example

Illustration of Failure

and Success

Methods Suitable for

Nonsmooth

Functions

---

Gradient Sampling

---

Quasi-Newton

Methods

---

A Difficult

Nonconvex Problem

from Nesterov

---

Limited Memory

Methods

---

Concluding Remarks



# Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

- sufficient decrease in function value:  
$$f(x + td) < f(x) + c_1 t \nabla f(x)^T d \quad (\text{L. Armijo, 1966})$$
- sufficient increase in directional derivative:  $f$  is differentiable at  $x + td$  and  $\nabla f(x + td)^T d > c_2 \nabla f(x)^T d$  ( $P.$  Wolfe, 1969)

Assuming  $\inf_t f(x + td)$  is bounded below,

- the Armijo condition holds for sufficiently small  $t$  as long as  $f$  is continuous
- the Wolfe condition holds for sufficiently large  $t$  as long as  $f$  is differentiable
- the intervals where each holds overlap

so combining the two conditions leads to a convenient, convergent bracketing line search (M.J.D. Powell, 1976)

Extends to locally Lipschitz case (A.S. Lewis and M.L.O., 2013)

---

Introduction

Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in

Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

---

Gradient Sampling

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods

---

Concluding Remarks



## Armijo-Wolfe Line Search

Given  $x$  with  $f$  differentiable at  $x$  and  $d$  with  $\nabla f(x)^T d < 0$ , and parameters  $0 < c_1 < c_2 < 1$ , find steplength  $t$  so that

- sufficient decrease in function value:  
$$f(x + td) < f(x) + c_1 t \nabla f(x)^T d \quad (\text{L. Armijo, 1966})$$
- sufficient increase in directional derivative:  $f$  is differentiable at  $x + td$  and  $\nabla f(x + td)^T d > c_2 \nabla f(x)^T d$  ( $P.$  Wolfe, 1969)

Assuming  $\inf_t f(x + td)$  is bounded below,

- the Armijo condition holds for sufficiently small  $t$  as long as  $f$  is continuous
- the Wolfe condition holds for sufficiently large  $t$  as long as  $f$  is differentiable
- the intervals where each holds overlap

so combining the two conditions leads to a convenient, convergent bracketing line search (M.J.D. Powell, 1976)

Extends to locally Lipschitz case (A.S. Lewis and M.L.O., 2013)

Searching for “Armijo-Wolfe” on the web, we found  
Melissa Armijo-Wolfe’s LinkedIn page!

Introduction

Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example

Failure of Gradient  
Descent in  
Nonsmooth Case

Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

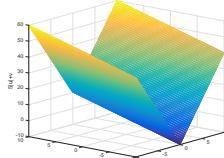
Limited Memory  
Methods

Concluding Remarks



## Failure of Gradient Method: Simple Convex Example

Let  $f(x) = a|x_1| + x_2$ , with  $a \geq 1$ . Although  $f$  is unbounded below, it is bounded below along any direction  $d = -\nabla f(x)$ .



Introduction  
Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

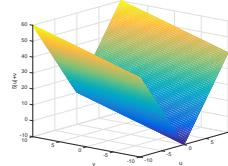
Limited Memory  
Methods

Concluding Remarks



## Failure of Gradient Method: Simple Convex Example

Let  $f(x) = a|x_1| + x_2$ , with  $a \geq 1$ . Although  $f$  is unbounded below, it is bounded below along any direction  $d = -\nabla f(x)$ .



**Theorem.** Let  $x^{(0)}$  satisfy  $x_1^{(0)} \neq 0$  and define  $x^{(k)} \in \mathbb{R}^2$  by

$$x^{(k+1)} = x^{(k)} + t_k d^{(k)} \text{ where } d^{(k)} = -\nabla f(x^{(k)})$$

and  $t_k$  is any steplength satisfying the Armijo and Wolfe conditions with Armijo parameter  $c_1$ . If

$$c_1(a^2 + 1) > 1$$

then  $x^{(k)}$  converges to  $\bar{x}$  with  $\bar{x}_1 = 0$ , even though  $f$  is unbounded below.

Azam Asl and M.L.O., 2017

Introduction  
Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success  
Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



# Illustration of Failure and Success

## Introduction

Nonsmooth,  
Nonconvex  
Optimization  
A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

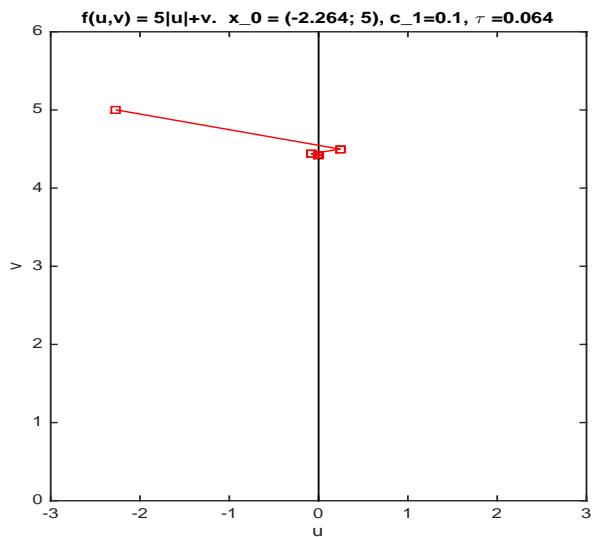
## Gradient Sampling

Quasi-Newton  
Methods

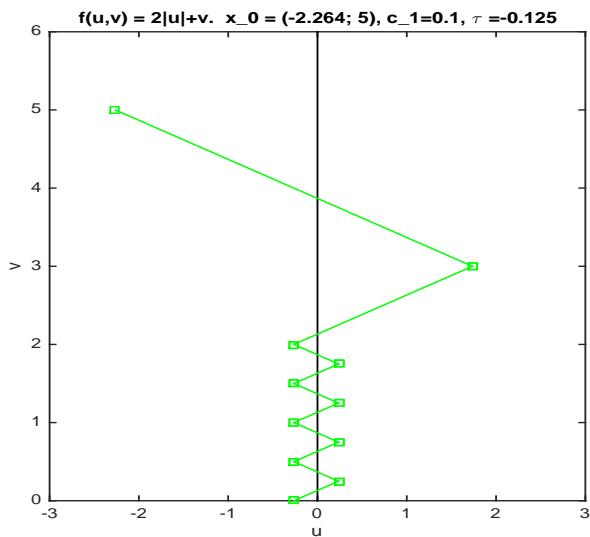
A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks



$$a = 5, c_1 = 0.1 \\ x^{(k)} \rightarrow \bar{x}$$



$$a = 2, c_1 = 0.1 \\ f(x^{(k)}) \downarrow -\infty$$



# Methods Suitable for Nonsmooth Functions

Exploit the gradient information obtained at several points, not just at one point:

[Introduction](#)  
[Nonsmooth,  
Nonconvex  
Optimization](#)  
[A Simple Nonconvex  
Example](#)  
[Failure of Gradient  
Descent in  
Nonsmooth Case](#)  
[Armijo-Wolfe Line  
Search](#)  
[Failure of Gradient  
Method: Simple  
Convex Example](#)  
[Illustration of Failure  
and Success](#)  
**Methods Suitable for  
Nonsmooth  
Functions**

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

[Limited Memory  
Methods](#)

[Concluding Remarks](#)



# Methods Suitable for Nonsmooth Functions

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example

Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Exploit the gradient information obtained at several points, not just at one point:

- Bundle methods (C. Lemaréchal, K.C. Kiwiel, 1980s –) extensive practical use and theoretical analysis, but complicated in nonconvex case



# Methods Suitable for Nonsmooth Functions

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search  
Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Exploit the gradient information obtained at several points, not just at one point:

- Bundle methods (C. Lemaréchal, K.C. Kiwiel, 1980s –) extensive practical use and theoretical analysis, but complicated in nonconvex case
- Gradient sampling: an easily stated method with nice convergence theory (J.V. Burke, A.S. Lewis, M.L.O., 2005; K.C. Kiwiel, 2007), but computationally intensive



# Methods Suitable for Nonsmooth Functions

Introduction  
Nonsmooth,  
Nonconvex  
Optimization

A Simple Nonconvex  
Example  
Failure of Gradient  
Descent in  
Nonsmooth Case  
Armijo-Wolfe Line  
Search

Failure of Gradient  
Method: Simple  
Convex Example  
Illustration of Failure  
and Success

Methods Suitable for  
Nonsmooth  
Functions

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Concluding Remarks

Exploit the gradient information obtained at several points, not just at one point:

- Bundle methods (C. Lemaréchal, K.C. Kiwiel, 1980s –) extensive practical use and theoretical analysis, but complicated in nonconvex case
- Gradient sampling: an easily stated method with nice convergence theory (J.V. Burke, A.S. Lewis, M.L.O., 2005; K.C. Kiwiel, 2007), but computationally intensive
- BFGS: traditional workhorse for smooth optimization, works amazingly well for nonsmooth optimization too, but very limited convergence theory



## Introduction

---

### Gradient Sampling

The Gradient Sampling Method  
With First Phase of Gradient Sampling  
With Second Phase of Gradient

Sampling  
The Clarke Subdifferential

Example  
Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

Extensions  
Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods

# Gradient Sampling



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

[Introduction](#)

[Gradient Sampling](#)

**The Gradient Sampling Method**

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

[Introduction](#)

[Gradient Sampling](#)

**The Gradient Sampling Method**

[With First Phase of Gradient Sampling](#)

[With Second Phase of Gradient Sampling](#)

[The Clarke Subdifferential](#)

[Example](#)

[Note that](#)

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

[Grad. Samp.: A Stabilized Steepest Descent Method](#)  
[Convergence of Gradient Sampling Method](#)

[Extensions](#)

[Some Gradient Sampling Success Stories](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

[Introduction](#)

[Gradient Sampling](#)

[The Gradient Sampling Method](#)

[With First Phase of Gradient Sampling](#)

[With Second Phase of Gradient Sampling](#)

[The Clarke Subdifferential](#)

[Example](#)

[Note that](#)

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

[Grad. Samp.: A Stabilized Steepest Descent Method](#)  
[Convergence of Gradient Sampling Method](#)

[Extensions](#)

[Some Gradient Sampling Success Stories](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

[Introduction](#)

[Gradient Sampling](#)

[The Gradient Sampling Method](#)

[With First Phase of Gradient Sampling](#)

[With Second Phase of Gradient Sampling](#)

[Sampling](#)

[The Clarke Subdifferential](#)

[Example](#)

[Note that](#)

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

[Grad. Samp.: A Stabilized Steepest Descent Method](#)  
[Convergence of Gradient Sampling Method](#)

[Extensions](#)

[Some Gradient Sampling Success Stories](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of

Gradient Sampling

Method

Extensions

Some Gradient Sampling Success Stories

---

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$
- ◆ If  $\|g\| \leq \tau$ , break out of loop.

Introduction

Gradient Sampling

The Gradient  
Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$
- ◆ If  $\|g\| \leq \tau$ , break out of loop.
- ◆ Backtracking line search: set  $d = -g$  and replace  $x$  by  $x + td$ , with  $t \in \{1, \frac{1}{2}, \frac{1}{4}, \dots\}$  and  $f(x + td) < f(x) - \beta t \|g\|^2$

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$
- ◆ If  $\|g\| \leq \tau$ , break out of loop.
- ◆ Backtracking line search: set  $d = -g$  and replace  $x$  by  $x + td$ , with  $t \in \{1, \frac{1}{2}, \frac{1}{4}, \dots\}$  and  $f(x + td) < f(x) - \beta t \|g\|^2$
- ◆ If  $f$  is not differentiable at  $x + td$ , replace  $x + td$  by a nearby point where  $f$  is differentiable.<sup>1</sup>

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that  $0 \in \partial^C f(x)$  at  $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$
- ◆ If  $\|g\| \leq \tau$ , break out of loop.
- ◆ Backtracking line search: set  $d = -g$  and replace  $x$  by  $x + td$ , with  $t \in \{1, \frac{1}{2}, \frac{1}{4}, \dots\}$  and  $f(x + td) < f(x) - \beta t \|g\|$
- ◆ If  $f$  is not differentiable at  $x + td$ , replace  $x + td$  by a nearby point where  $f$  is differentiable.<sup>1</sup>

■ New phase: set  $\epsilon = \mu\epsilon$  and  $\tau = \theta\tau$ .

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# The Gradient Sampling Method

Fix sample size  $m \geq n + 1$ , line search parameter  $\beta \in (0, 1)$ , reduction factors  $\mu \in (0, 1)$  and  $\theta \in (0, 1)$ .

Initialize sampling radius  $\epsilon > 0$ , tolerance  $\tau > 0$ , iterate  $x$ .

Repeat (outer loop)

■ Repeat (inner loop: gradient sampling with fixed  $\epsilon$ ):

- ◆ Set  $G = \{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\}$ , sampling  $u_1, \dots, u_m$  uniformly from the unit ball
- ◆ Set  $g = \arg \min\{||g|| : g \in \text{conv}(G)\}$
- ◆ If  $\|g\| \leq \tau$ , break out of loop.
- ◆ Backtracking line search: set  $d = -g$  and replace  $x$  by  $x + td$ , with  $t \in \{1, \frac{1}{2}, \frac{1}{4}, \dots\}$  and  $f(x + td) < f(x) - \beta t \|g\|^2$
- ◆ If  $f$  is not differentiable at  $x + td$ , replace  $x + td$  by a nearby point where  $f$  is differentiable.<sup>1</sup>

■ New phase: set  $\epsilon = \mu\epsilon$  and  $\tau = \theta\tau$ .

J.V. Burke, A.S. Lewis and M.L.O., SIOPT, 2005.

---

<sup>1</sup>Needed in theory, but typically not in practice.

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# With First Phase of Gradient Sampling

[Introduction](#)

[Gradient Sampling](#)

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

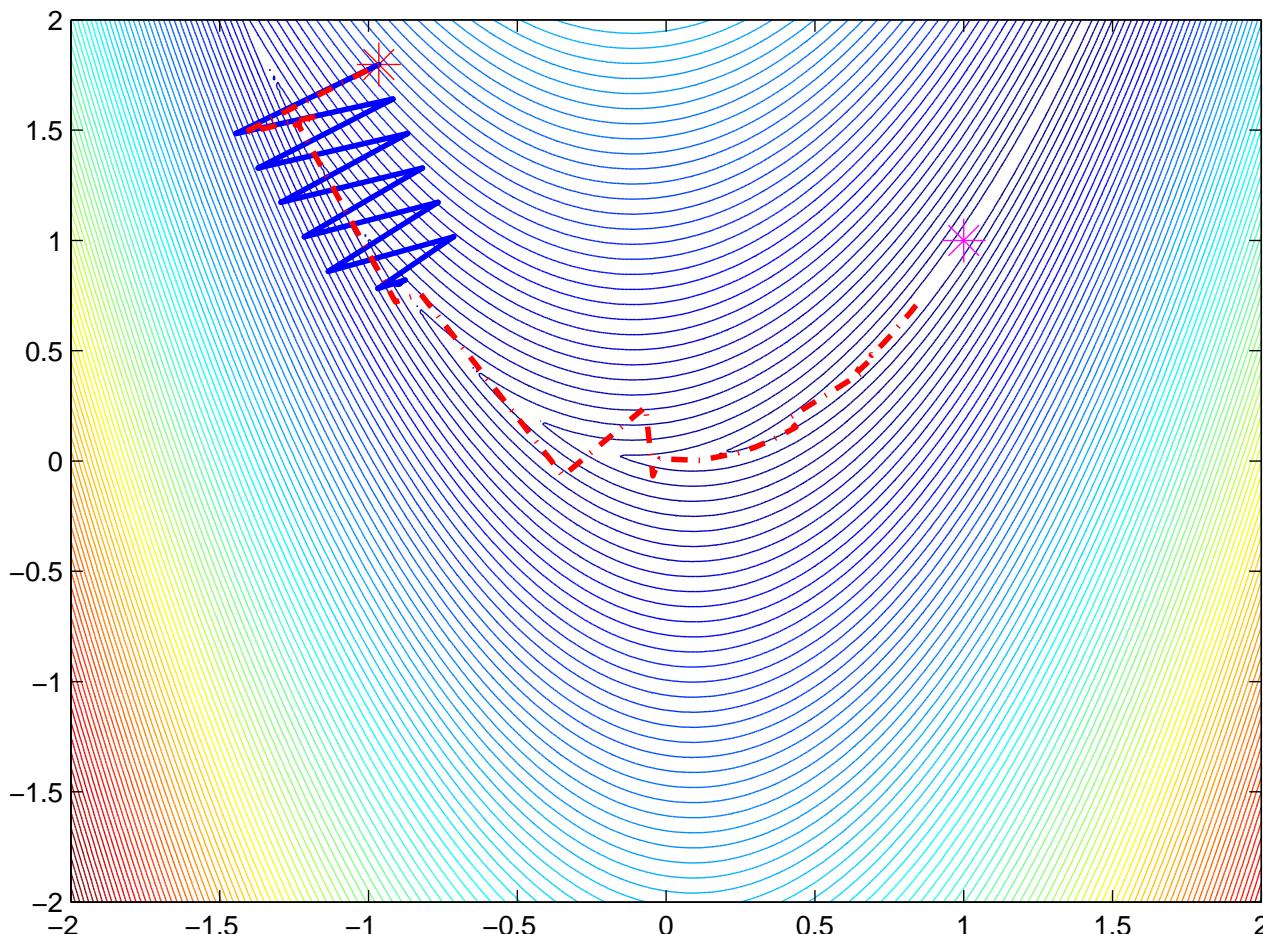
Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods

$$f(x) = 10|x_2 - x_1^2| + (1-x_1)^2$$





# With Second Phase of Gradient Sampling

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

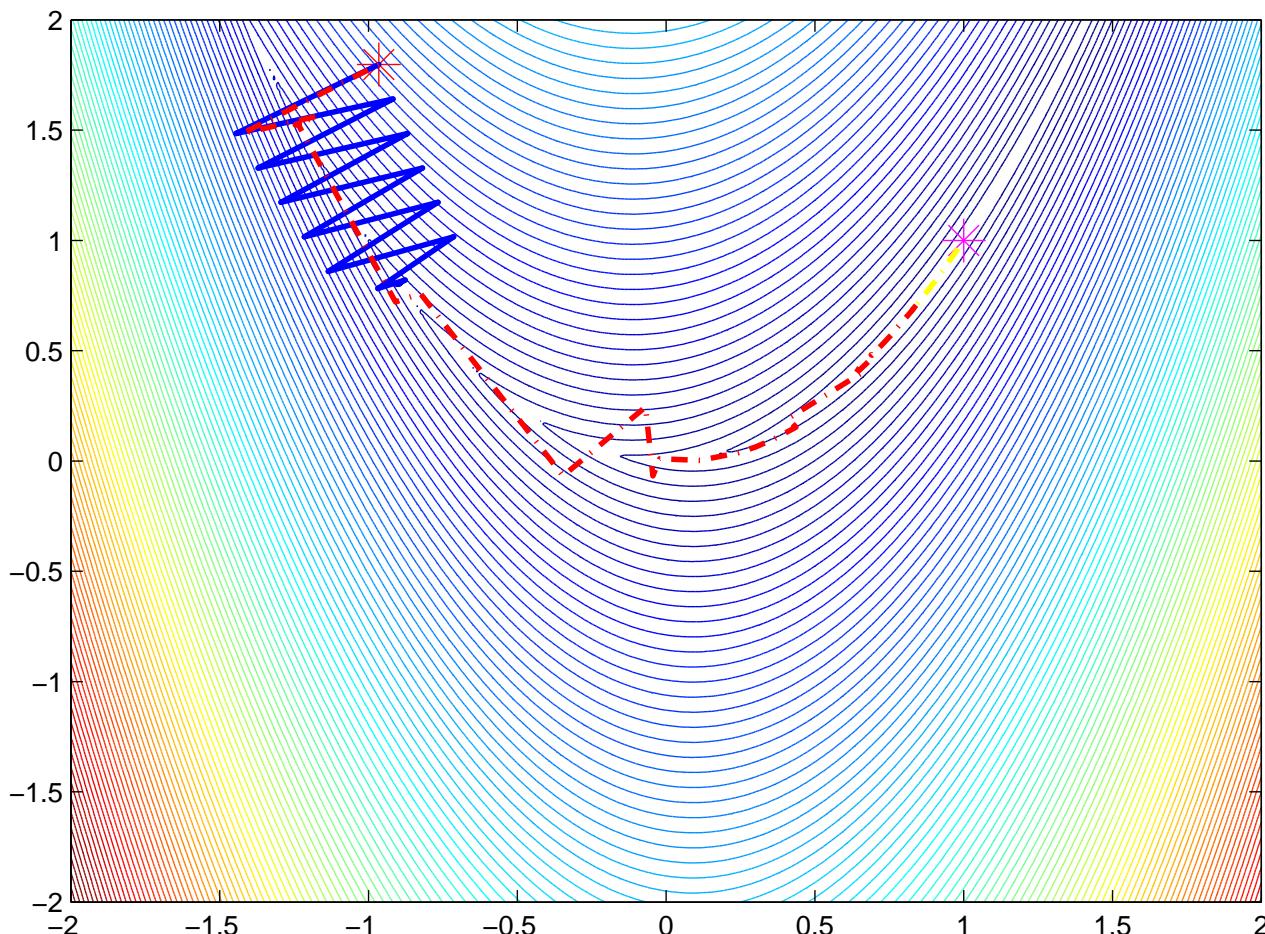
Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

$$f(x) = 10|x_2 - x_1^2| + (1-x_1)^2$$





# The Clarke Subdifferential

[Introduction](#)

---

[Gradient Sampling](#)

---

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

**The Clarke  
Subdifferential**

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

**The Clarke  
Subdifferential**

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$

F.H. Clarke, 1973 (he used the name “generalized gradient”).



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

**The Clarke  
Subdifferential**

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$

F.H. Clarke, 1973 (he used the name “generalized gradient”).

If  $f$  is continuously differentiable at  $\bar{x}$ , then  $\partial^C f(\bar{x}) = \{\nabla f(\bar{x})\}$ .



# The Clarke Subdifferential

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$

F.H. Clarke, 1973 (he used the name “generalized gradient”).

If  $f$  is continuously differentiable at  $\bar{x}$ , then  $\partial^C f(\bar{x}) = \{\nabla f(\bar{x})\}$ .

If  $f$  is convex,  $\partial^C f$  is the subdifferential of convex analysis.



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

**The Clarke  
Subdifferential**

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and  
let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$

F.H. Clarke, 1973 (he used the name “generalized gradient”).

If  $f$  is continuously differentiable at  $\bar{x}$ , then  $\partial^C f(\bar{x}) = \{\nabla f(\bar{x})\}$ .

If  $f$  is convex,  $\partial^C f$  is the subdifferential of convex analysis.

We say  $\bar{x}$  is Clarke stationary for  $f$  if  $0 \in \partial^C f(\bar{x})$ .



# The Clarke Subdifferential

[Introduction](#)

[Gradient Sampling](#)

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

**The Clarke  
Subdifferential**

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is locally Lipschitz, and let  $D = \{x \in \mathbb{R}^n : f \text{ is differentiable at } x\}$ .

Rademacher's Theorem:  $\mathbb{R}^n \setminus D$  has measure zero.

The Clarke subdifferential of  $f$  at  $\bar{x}$  is

$$\partial^C f(\bar{x}) = \text{conv} \left\{ \lim_{x \rightarrow \bar{x}, x \in D} \nabla f(x) \right\}.$$

F.H. Clarke, 1973 (he used the name “generalized gradient”).

If  $f$  is continuously differentiable at  $\bar{x}$ , then  $\partial^C f(\bar{x}) = \{\nabla f(\bar{x})\}$ .

If  $f$  is convex,  $\partial^C f$  is the subdifferential of convex analysis.

We say  $\bar{x}$  is Clarke stationary for  $f$  if  $0 \in \partial^C f(\bar{x})$ .

Key point: the convex hull of the set  $G$  generated by Gradient Sampling is a surrogate for  $\partial^C f$ .



## Example

Let

$$f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of

Gradient Sampling

With Second Phase

of Gradient

Sampling

The Clarke

Subdifferential

**Example**

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A

Stabilized Steepest

Descent Method

Convergence of

Gradient Sampling

Method

Extensions

Some Gradient

Sampling Success

Stories

Quasi-Newton

Methods

A Difficult

Nonconvex Problem

from Nesterov

Limited Memory

Methods



## Example

Let

$$f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2$$

For  $x$  with  $x_2 \neq x_1^2$ ,  $f$  is differentiable with gradient

$$\nabla f(x) = 10 \operatorname{sgn}\{x_2 - x_1^2\} \begin{bmatrix} -2x_1 \\ 1 \end{bmatrix} + \begin{bmatrix} -2(1 - x_1) \\ 0 \end{bmatrix}$$

so  $\partial^C f(x) = \{\nabla f(x)\}$ .

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

The Clarke Subdifferential

**Example**

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



## Example

Let

$$f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2$$

For  $x$  with  $x_2 \neq x_1^2$ ,  $f$  is differentiable with gradient

$$\nabla f(x) = 10 \operatorname{sgn}\{x_2 - x_1^2\} \begin{bmatrix} -2x_1 \\ 1 \end{bmatrix} + \begin{bmatrix} -2(1 - x_1) \\ 0 \end{bmatrix}$$

so  $\partial^C f(x) = \{\nabla f(x)\}$ . For  $x$  with  $x_2 = x_1^2$ , there are two limiting gradients, namely

$$\pm 10 \begin{bmatrix} -2x_1 \\ 1 \end{bmatrix} + \begin{bmatrix} -2(1 - x_1) \\ 0 \end{bmatrix}$$

so  $\partial^C f(x)$  consists of the convex hull of these two vectors.

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

The Clarke Subdifferential

**Example**

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Method

Extensions  
Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



## Example

Let

$$f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2$$

For  $x$  with  $x_2 \neq x_1^2$ ,  $f$  is differentiable with gradient

$$\nabla f(x) = 10 \operatorname{sgn}\{x_2 - x_1^2\} \begin{bmatrix} -2x_1 \\ 1 \end{bmatrix} + \begin{bmatrix} -2(1 - x_1) \\ 0 \end{bmatrix}$$

so  $\partial^C f(x) = \{\nabla f(x)\}$ . For  $x$  with  $x_2 = x_1^2$ , there are two limiting gradients, namely

$$\pm 10 \begin{bmatrix} -2x_1 \\ 1 \end{bmatrix} + \begin{bmatrix} -2(1 - x_1) \\ 0 \end{bmatrix}$$

so  $\partial^C f(x)$  consists of the convex hull of these two vectors. The unique  $x$  for which  $0 \in \partial^C f(x)$  is  $x = [1; 1]^T$ , so this is the unique Clarke stationary point of  $f$  (it follows that it is the global minimizer).

[Introduction](#)

[Gradient Sampling](#)

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

The Clarke Subdifferential

**Example**

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



Note that  $0 \in \partial^C f(x)$  at  $x = [1; 1]^T$

Introduction

Gradient Sampling

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling  
With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential  
Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

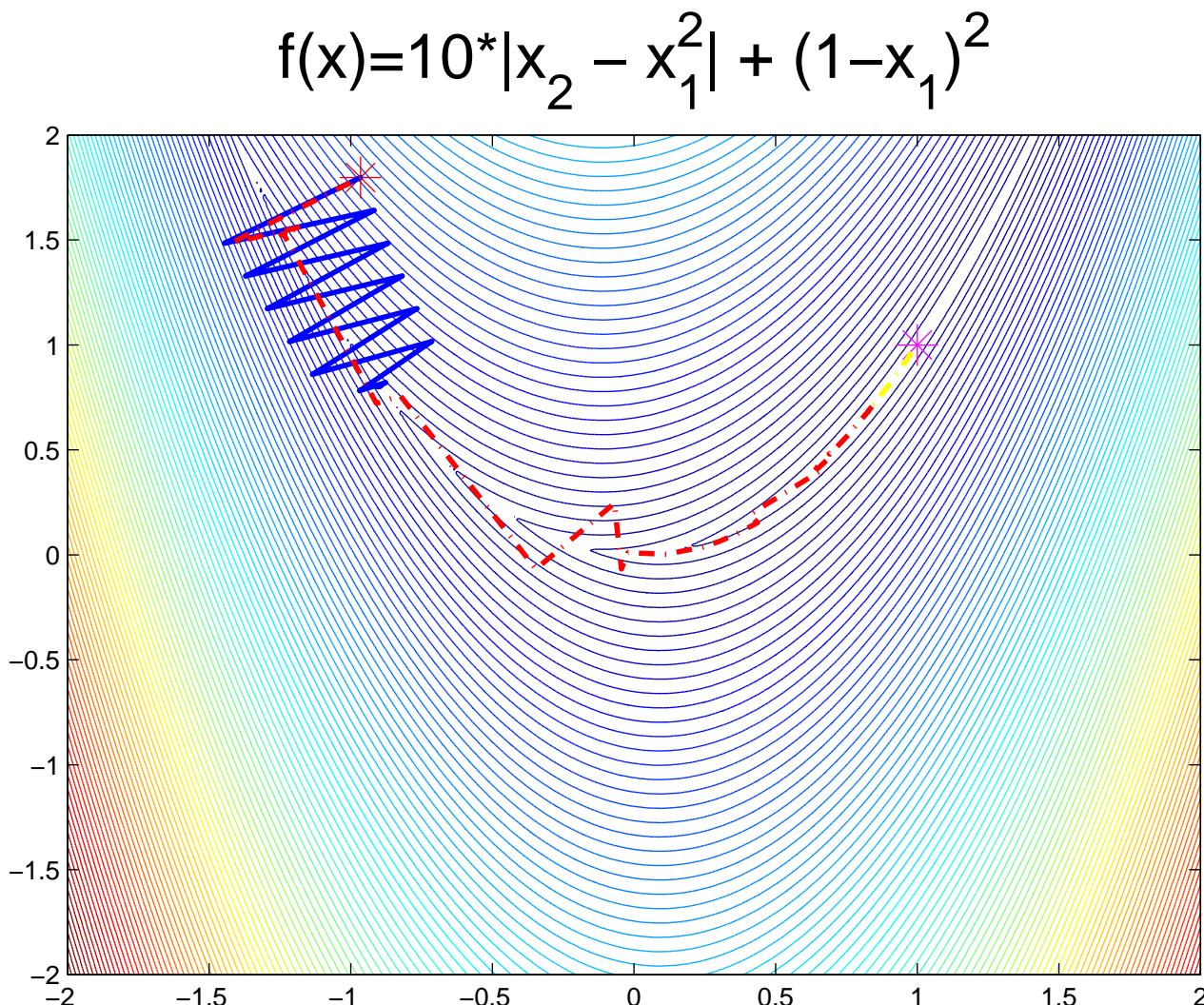
Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions  
Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods





# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling

Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

## Proof.

$$-\text{dist}(0, C) = -\min_{g \in C} \|g\|$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

## Proof.

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \end{aligned}$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling

Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

## Proof.

$$\begin{aligned}-\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\&= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\&= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d\end{aligned}$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient

Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling

Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

## Proof.

$$\begin{aligned}-\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\&= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\&= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\&= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d)\end{aligned}$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

## Proof.

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d) \\ &= \min_{\|d\| \leq 1} \max_{g \in C} g^T d. \end{aligned}$$

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

**Proof.**

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d) \\ &= \min_{\|d\| \leq 1} \max_{g \in C} g^T d. \end{aligned}$$

Note: the distance is nonnegative, and zero iff  $0 \in C$ .

---

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling

Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

**Proof.**

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d) \\ &= \min_{\|d\| \leq 1} \max_{g \in C} g^T d. \end{aligned}$$

Note: the distance is nonnegative, and zero iff  $0 \in C$ .

Otherwise, equality is attained by  $g = \Pi_C(0)$ ,  $d = -g/\|g\|$ .

---

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling

Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

**Proof.**

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d) \\ &= \min_{\|d\| \leq 1} \max_{g \in C} g^T d. \end{aligned}$$

Note: the distance is nonnegative, and zero iff  $0 \in C$ .

Otherwise, equality is attained by  $g = \Pi_C(0)$ ,  $d = -g/\|g\|$ .

Ordinary steepest descent:  $C = \{\nabla f(x)\}$ .

---

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

---

Limited Memory  
Methods



# Grad. Samp.: A Stabilized Steepest Descent Method

**Lemma.** Let  $C$  be a compact convex set and  $\|\cdot\| = \|\cdot\|_2$ . Then

$$-\text{dist}(0, C) = \min_{\|d\| \leq 1} \max_{g \in C} g^T d$$

**Proof.**

$$\begin{aligned} -\text{dist}(0, C) &= -\min_{g \in C} \|g\| \\ &= -\min_{g \in C} \max_{\|d\| \leq 1} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T d \\ &= -\max_{\|d\| \leq 1} \min_{g \in C} g^T (-d) \\ &= \min_{\|d\| \leq 1} \max_{g \in C} g^T d. \end{aligned}$$

Note: the distance is nonnegative, and zero iff  $0 \in C$ .

Otherwise, equality is attained by  $g = \Pi_C(0)$ ,  $d = -g/\|g\|$ .

Ordinary steepest descent:  $C = \{\nabla f(x)\}$ .

Gradient sampling:  $C = \text{conv}(G)$

$$= \text{conv}(\{\nabla f(x), \nabla f(x + \epsilon u_1), \dots, \nabla f(x + \epsilon u_m)\})$$

---

Introduction

Gradient Sampling

The Gradient

Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

---

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz

Introduction

Gradient Sampling

The Gradient  
Sampling Method

With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method

Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Method

---

Extensions

Some Gradient Sampling Success Stories

Stories

---

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke

Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at

$x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

- the inner loop always terminates, so the sequences of sampling radii  $\{\epsilon\}$  and tolerances  $\{\tau\}$  converge to zero, and

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke

Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

- the inner loop always terminates, so the sequences of sampling radii  $\{\epsilon\}$  and tolerances  $\{\tau\}$  converge to zero, and
- $\bar{x}$  is Clarke stationary for  $f$ , i.e.,  $0 \in \partial^C f(\bar{x})$ .

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

- the inner loop always terminates, so the sequences of sampling radii  $\{\epsilon\}$  and tolerances  $\{\tau\}$  converge to zero, and
- $\bar{x}$  is Clarke stationary for  $f$ , i.e.,  $0 \in \partial^C f(\bar{x})$ .

J.V. Burke, A.S. Lewis and M.L.O., SIOPT, 2005.

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

---

Extensions

Some Gradient Sampling Success Stories

---

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

- the inner loop always terminates, so the sequences of sampling radii  $\{\epsilon\}$  and tolerances  $\{\tau\}$  converge to zero, and
- $\bar{x}$  is Clarke stationary for  $f$ , i.e.,  $0 \in \partial^C f(\bar{x})$ .

J.V. Burke, A.S. Lewis and M.L.O., SIOPT, 2005.

Drop the assumption that  $f$  has bounded level sets. Then, wp 1, either the sequence  $\{f(x)\} \rightarrow -\infty$ , or every cluster point of the sequence of iterates  $\{x\}$  is Clarke stationary.

---

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke Subdifferential Example

Note that  $0 \in \partial^C f(x)$  at  $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions  
Some Gradient Sampling Success Stories

Quasi-Newton Methods

---

A Difficult Nonconvex Problem from Nesterov

---

Limited Memory Methods



# Convergence of Gradient Sampling Method

**Theorem.** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

- is locally Lipschitz
- is cont. differentiable on an open full-measure subset of  $\mathbb{R}^n$
- has bounded level sets

Then, with probability one,  $f$  is differentiable at the sampled points, the line search always terminates, and if the sequence of iterates  $\{x\}$  converges to some point  $\bar{x}$ , then, with probability 1

- the inner loop always terminates, so the sequences of sampling radii  $\{\epsilon\}$  and tolerances  $\{\tau\}$  converge to zero, and
- $\bar{x}$  is Clarke stationary for  $f$ , i.e.,  $0 \in \partial^C f(\bar{x})$ .

J.V. Burke, A.S. Lewis and M.L.O., SIOPT, 2005.

Drop the assumption that  $f$  has bounded level sets. Then, wp 1, either the sequence  $\{f(x)\} \rightarrow -\infty$ , or every cluster point of the sequence of iterates  $\{x\}$  is Clarke stationary.

K.C. Kiwiel, SIOPT, 2007.

Introduction

Gradient Sampling

The Gradient Sampling Method

With First Phase of Gradient Sampling

With Second Phase of Gradient

Sampling

The Clarke Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method

Convergence of Gradient Sampling Method

Extensions  
Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



## Extensions

A more efficient version: Adaptive Gradient Sampling (F.E. Curtis and X. Que, 2013).

[Introduction](#)

[Gradient Sampling](#)

The Gradient Sampling Method  
With First Phase of Gradient Sampling

With Second Phase of Gradient Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

[Extensions](#)

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



## Extensions

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling  
With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential  
Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

[Extensions](#)

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

A more efficient version: Adaptive Gradient Sampling  
(F.E. Curtis and X. Que, 2013).

## Problems with Nonsmooth Constraints

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$



## Extensions

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling  
With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential  
Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

[Extensions](#)

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

A more efficient version: Adaptive Gradient Sampling  
(F.E. Curtis and X. Que, 2013).

## Problems with Nonsmooth Constraints

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.



## Extensions

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

**Extensions**

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

A more efficient version: Adaptive Gradient Sampling  
(F.E. Curtis and X. Que, 2013).

## Problems with Nonsmooth Constraints

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

A successive quadratic programming gradient sampling method with convergence theory.



## Extensions

[Introduction](#)

[Gradient Sampling](#)

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling

With Second Phase  
of Gradient  
Sampling

The Clarke  
Subdifferential

Example

Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

[Extensions](#)

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

A more efficient version: Adaptive Gradient Sampling  
(F.E. Curtis and X. Que, 2013).

## Problems with Nonsmooth Constraints

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

A successive quadratic programming gradient sampling method with convergence theory.

F.E. Curtis and M.L.O., SIOPT, 2012.



# Some Gradient Sampling Success Stories

- Non-Lipschitz eigenvalue optimization for non-normal matrices (J.V. Burke, A.S. Lewis and M.L.O., 2002 – )

Introduction

Gradient Sampling

The Gradient Sampling Method  
With First Phase of Gradient Sampling  
With Second Phase of Gradient Sampling

Sampling  
The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# Some Gradient Sampling Success Stories

- Non-Lipschitz eigenvalue optimization for non-normal matrices (J.V. Burke, A.S. Lewis and M.L.O., 2002 – )
- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)

Introduction

Gradient Sampling

The Gradient Sampling Method  
With First Phase of Gradient Sampling  
With Second Phase of Gradient Sampling

The Clarke Subdifferential

Example

Note that

$0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A Stabilized Steepest Descent Method  
Convergence of Gradient Sampling Method

Extensions

Some Gradient Sampling Success Stories

Quasi-Newton Methods

A Difficult Nonconvex Problem from Nesterov

Limited Memory Methods



# Some Gradient Sampling Success Stories

Introduction

Gradient Sampling

The Gradient  
Sampling Method  
With First Phase of  
Gradient Sampling  
With Second Phase  
of Gradient  
Sampling  
The Clarke  
Subdifferential

Example  
Note that  
 $0 \in \partial^C f(x)$  at  
 $x = [1; 1]^T$

Grad. Samp.: A  
Stabilized Steepest  
Descent Method  
Convergence of  
Gradient Sampling  
Method

Extensions

Some Gradient  
Sampling Success  
Stories

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

- Non-Lipschitz eigenvalue optimization for non-normal matrices (J.V. Burke, A.S. Lewis and M.L.O., 2002 – )
- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)
- Design of path planning for robots: avoids “chattering” that otherwise arises from nonsmoothness (I. Mitchell et al, 2017)



[Introduction](#)

[Gradient Sampling](#)

**Quasi-Newton  
Methods**

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
(“Full” Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity  
Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example  
Again

Relationships

# Quasi-Newton Methods



## Bill Davidon

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

W. Davidon, a physicist at Argonne, had the breakthrough idea in 1959: since it's too expensive to compute and factor the Hessian  $\nabla^2 f(x)$  at every iteration, update an approximation to its inverse using information from gradient differences, and multiply this onto the negative gradient to approximate Newton's method.



## Bill Davidon

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

W. Davidon, a physicist at Argonne, had the breakthrough idea in 1959: since it's too expensive to compute and factor the Hessian  $\nabla^2 f(x)$  at every iteration, update an approximation to its inverse using information from gradient differences, and multiply this onto the negative gradient to approximate Newton's method.

Each inverse Hessian approximation differs from the previous one by a rank-two correction.



## Bill Davidon

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

**Bill Davidon**

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

W. Davidon, a physicist at Argonne, had the breakthrough idea in 1959: since it's too expensive to compute and factor the Hessian  $\nabla^2 f(x)$  at every iteration, update an approximation to its inverse using information from gradient differences, and multiply this onto the negative gradient to approximate Newton's method.

Each inverse Hessian approximation differs from the previous one by a rank-two correction.

Ahead of its time: the paper was rejected by the physics journals, but published 30 years later in the first issue of SIAM J. Optimization.



## Bill Davidon

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Positive](#)

W. Davidon, a physicist at Argonne, had the breakthrough idea in 1959: since it's too expensive to compute and factor the Hessian  $\nabla^2 f(x)$  at every iteration, update an approximation to its inverse using information from gradient differences, and multiply this onto the negative gradient to approximate Newton's method.

Each inverse Hessian approximation differs from the previous one by a rank-two correction.

Ahead of its time: the paper was rejected by the physics journals, but published 30 years later in the first issue of SIAM J. Optimization.

Davidon was a well known active anti-war protester during the Vietnam War. In December 2013, it was revealed that he was the mastermind behind the break-in at the FBI office in Media, PA, on March 8, 1971, during the Muhammad Ali - Joe Frazier world heavyweight boxing championship.



## Fletcher and Powell

In 1963, R. Fletcher and M.J.D. Powell improved Davidon's method and established convergence for convex quadratic functions.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Positive](#)



## Fletcher and Powell

In 1963, R. Fletcher and M.J.D. Powell improved Davidon's method and established convergence for convex quadratic functions.

They applied it to solve problems in 100 variables: a lot at the time.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relatives of Partial](#)



## Fletcher and Powell

In 1963, R. Fletcher and M.J.D. Powell improved Davidon's method and established convergence for convex quadratic functions.

They applied it to solve problems in 100 variables: a lot at the time.

The method became known as the DFP method.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relatives of Partial](#)



## Fletcher and Powell

In 1963, R. Fletcher and M.J.D. Powell improved Davidon's method and established convergence for convex quadratic functions.

They applied it to solve problems in 100 variables: a lot at the time.

The method became known as the DFP method.

Davidon, Fletcher and Powell all died during 2013–2016.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relatives of Partial](#)



## BFGS

In 1970, C.G. Broyden, R. Fletcher, D. Goldfarb and D. Shanno all independently proposed the BFGS method, which is a kind of dual of the DFP method. It was soon recognized that this was a remarkably effective method for smooth optimization.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

**BFGS**

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Positive Definite Matrices](#)



## BFGS

In 1970, C.G. Broyden, R. Fletcher, D. Goldfarb and D. Shanno all independently proposed the BFGS method, which is a kind of dual of the DFP method. It was soon recognized that this was a remarkably effective method for smooth optimization.

In 1973, C.G. Broyden, J.E. Dennis and J.J. Moré proved generic local superlinear convergence of BFGS and DFP and other quasi-Newton methods.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon Fletcher and Powell](#)

**BFGS**

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)



## BFGS

In 1970, C.G. Broyden, R. Fletcher, D. Goldfarb and D. Shanno all independently proposed the BFGS method, which is a kind of dual of the DFP method. It was soon recognized that this was a remarkably effective method for smooth optimization.

In 1973, C.G. Broyden, J.E. Dennis and J.J. Moré proved generic local superlinear convergence of BFGS and DFP and other quasi-Newton methods.

In 1976, M.J.D. Powell established convergence of BFGS with an inexact Armijo-Wolfe line search for a general class of smooth convex functions for BFGS. In 1987, this was extended by R.H. Byrd, J. Nocedal and Y.-X. Yuan to include the whole “Broyden” class of methods interpolating BFGS and DFP: *except for the DFP end point.*

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon Fletcher and Powell](#)

**BFGS**

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



## BFGS

In 1970, C.G. Broyden, R. Fletcher, D. Goldfarb and D. Shanno all independently proposed the BFGS method, which is a kind of dual of the DFP method. It was soon recognized that this was a remarkably effective method for smooth optimization.

In 1973, C.G. Broyden, J.E. Dennis and J.J. Moré proved generic local superlinear convergence of BFGS and DFP and other quasi-Newton methods.

In 1976, M.J.D. Powell established convergence of BFGS with an inexact Armijo-Wolfe line search for a general class of smooth convex functions for BFGS. In 1987, this was extended by R.H. Byrd, J. Nocedal and Y.-X. Yuan to include the whole “Broyden” class of methods interpolating BFGS and DFP: *except for the DFP end point.*

Pathological counterexamples to convergence in the smooth, nonconvex case are known to exist (Y.-H. Dai, 2002, 2013; W. Mascarenhas 2004), but it is widely accepted that the method works well in practice in the smooth, nonconvex case.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon Fletcher and Powell](#)

**BFGS**

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Positive Definite Matrices](#)



# The BFGS Method (“Full” Version)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$   
(which is supposed to approximate the *inverse Hessian* of  $f$ )

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for](#)

[Nonsmooth](#)

[Optimization](#)

[With BFGS](#)

[Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10](#)

[Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships of Partial](#)



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)  
[With BFGS Example:](#)  
[Minimizing a Product of Eigenvalues](#)  
[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)  
[Partly Smooth Functions, continued](#)  
[Same Example Again](#)  
[Properties of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse Hessian* of  $f$ )

Repeat



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for](#)

[Nonsmooth Optimization](#)

[With BFGS](#)

[Example:  
Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse Hessian* of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

**The BFGS Method (“Full” Version)**

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse Hessian* of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example:  
Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse Hessian* of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search
- Set  $s = td$ ,  $y = \nabla f(x + td) - \nabla f(x)$



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse Hessian* of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search
- Set  $s = td$ ,  $y = \nabla f(x + td) - \nabla f(x)$
- Replace  $x$  by  $x + td$



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

**The BFGS Method (“Full” Version)**

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse* Hessian of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search
- Set  $s = td$ ,  $y = \nabla f(x + td) - \nabla f(x)$
- Replace  $x$  by  $x + td$
- Replace  $H$  by  $VHV^T + \frac{1}{s^T y} ss^T$ , where  $V = I - \frac{1}{s^T y} sy^T$



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

**The BFGS Method (“Full” Version)**

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)  
[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse* Hessian of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search
- Set  $s = td$ ,  $y = \nabla f(x + td) - \nabla f(x)$
- Replace  $x$  by  $x + td$
- Replace  $H$  by  $VHV^T + \frac{1}{s^T y} ss^T$ , where  $V = I - \frac{1}{s^T y} sy^T$

Note that  $H$  can be computed in  $O(n^2)$  operations since  $V$  is a rank one perturbation of the identity



# The BFGS Method (“Full” Version)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

**The BFGS Method (“Full” Version)**

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of](#)

[Eigenvalues of](#)

[A  \$\circ\$  X](#)

[Evolution of Eigenvalues of H](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

Initialize iterate  $x$  and positive-definite symmetric matrix  $H$  (which is supposed to approximate the *inverse* Hessian of  $f$ )

Repeat

- Set  $d = -H\nabla f(x)$ .
- Obtain  $t$  from Armijo-Wolfe line search
- Set  $s = td$ ,  $y = \nabla f(x + td) - \nabla f(x)$
- Replace  $x$  by  $x + td$
- Replace  $H$  by  $VHV^T + \frac{1}{s^T y} ss^T$ , where  $V = I - \frac{1}{s^T y} sy^T$

Note that  $H$  can be computed in  $O(n^2)$  operations since  $V$  is a rank one perturbation of the identity

The Wolfe condition guarantees that  $s^T y > 0$  and hence that the new  $H$  is positive definite.



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10](#)

[Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partly](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS](#)

[Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued  
Same Example  
Again](#)

[Properties of Partial](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

Key point: use an Armijo-Wolfe line search. Do not insist on reducing the magnitude of the directional derivative along the line!

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

Key point: use an Armijo-Wolfe line search. Do not insist on reducing the magnitude of the directional derivative along the line!

In the nonsmooth case, BFGS builds a very ill-conditioned inverse “Hessian” approximation, with some tiny eigenvalues converging to zero, corresponding to “infinitely large” curvature in the directions defined by the associated eigenvectors.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

Key point: use an Armijo-Wolfe line search. Do not insist on reducing the magnitude of the directional derivative along the line!

In the nonsmooth case, BFGS builds a very ill-conditioned inverse “Hessian” approximation, with some tiny eigenvalues converging to zero, corresponding to “infinitely large” curvature in the directions defined by the associated eigenvectors.

Remarkably, the condition number of the inverse Hessian approximation typically reaches  $10^{16}$  before the method breaks down.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

Key point: use an Armijo-Wolfe line search. Do not insist on reducing the magnitude of the directional derivative along the line!

In the nonsmooth case, BFGS builds a very ill-conditioned inverse “Hessian” approximation, with some tiny eigenvalues converging to zero, corresponding to “infinitely large” curvature in the directions defined by the associated eigenvectors.

Remarkably, the condition number of the inverse Hessian approximation typically reaches  $10^{16}$  before the method breaks down.

We have never seen convergence to non-stationary points that cannot be explained by numerical difficulties.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



# BFGS for Nonsmooth Optimization

In 1982, C. Lemaréchal observed that quasi-Newton methods can be effective for nonsmooth optimization, but dismissed them as there was no theory behind them and no good way to terminate them.

Otherwise, there is not much in the literature on the subject until A.S. Lewis and M.L.O. (Math. Prog., 2013): we address both issues in detail, but our convergence results are limited to special cases.

Key point: use an Armijo-Wolfe line search. Do not insist on reducing the magnitude of the directional derivative along the line!

In the nonsmooth case, BFGS builds a very ill-conditioned inverse “Hessian” approximation, with some tiny eigenvalues converging to zero, corresponding to “infinitely large” curvature in the directions defined by the associated eigenvectors.

Remarkably, the condition number of the inverse Hessian approximation typically reaches  $10^{16}$  before the method breaks down.

We have never seen convergence to non-stationary points that cannot be explained by numerical difficulties.

Convergence rate of BFGS is typically linear (not superlinear) in the nonsmooth case.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \(“Full” Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)

[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

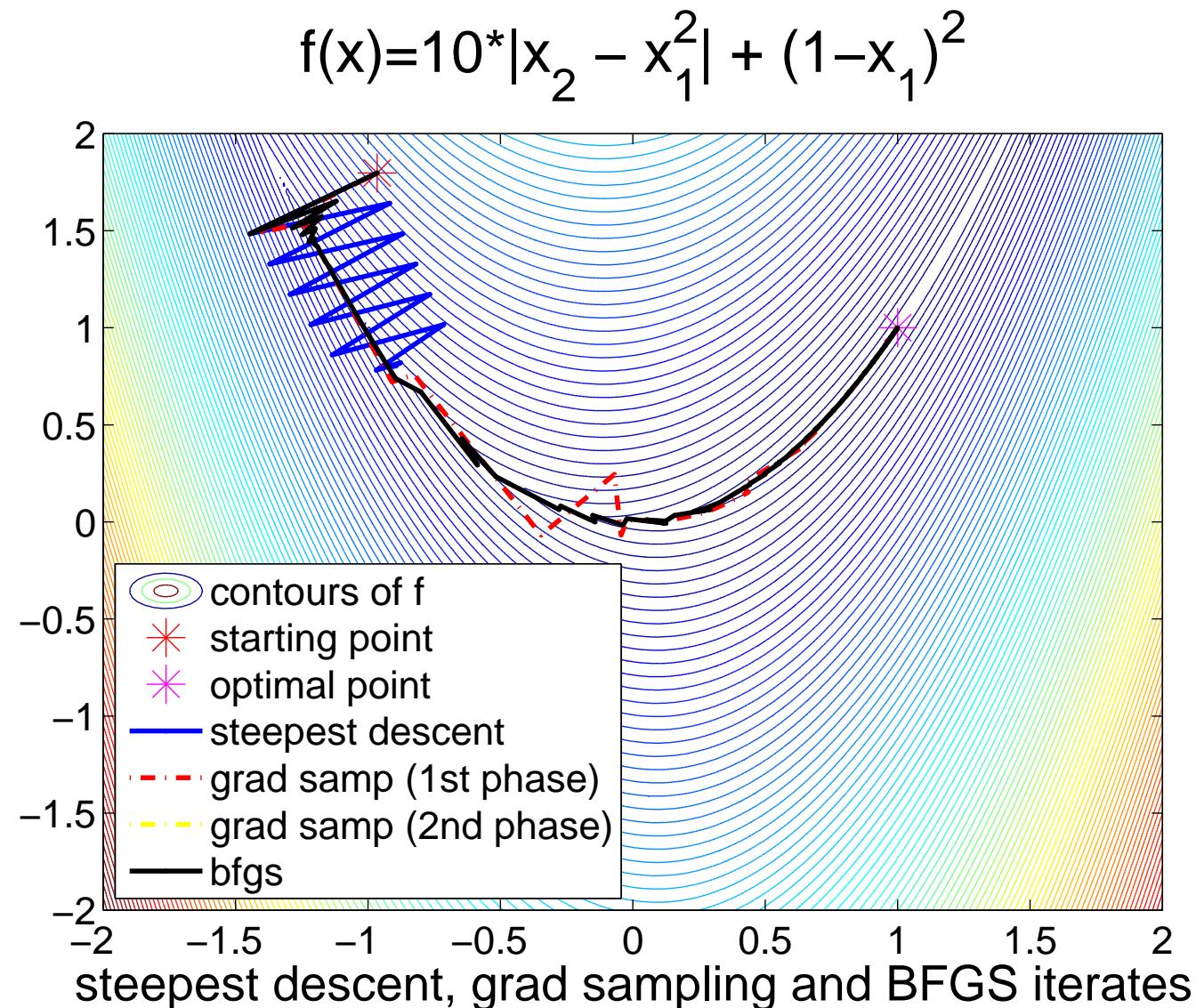
[Same Example Again](#)

[Relations of Partial](#)



# With BFGS

Introduction  
Gradient Sampling  
Quasi-Newton Methods  
Bill Davidson  
Fletcher and Powell  
BFGS  
The BFGS Method ("Full" Version)  
BFGS for Nonsmooth Optimization  
**With BFGS**  
Example:  
Minimizing a Product of Eigenvalues  
BFGS from 10 Randomly Generated Starting Points  
Evolution of Eigenvalues of  $A \circ X$   
Evolution of Eigenvalues of  $H$   
Regularity  
Partly Smooth Functions  
Partly Smooth Functions, continued  
Same Example Again  
Relations of Partial





## Example: Minimizing a Product of Eigenvalues

Let  $S^N$  denote the space of real symmetric  $N \times N$  matrices, and

$$\lambda_1(X) \geq \lambda_2(X) \geq \cdots \lambda_N(X)$$

denote the eigenvalues of  $X \in S^N$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relatives of Partial](#)



## Example: Minimizing a Product of Eigenvalues

Let  $S^N$  denote the space of real symmetric  $N \times N$  matrices, and

$$\lambda_1(X) \geq \lambda_2(X) \geq \cdots \lambda_N(X)$$

denote the eigenvalues of  $X \in S^N$ . We wish to minimize

$$f(X) = \log \prod_{i=1}^{N/2} \lambda_i(A \circ X)$$

where  $A \in S^N$  is fixed and  $\circ$  is the Hadamard (componentwise) matrix product, subject to the constraints that  $X$  is positive semidefinite and has diagonal entries equal to 1.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)

[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)  
[Properties of Partly Smooth Functions](#)



## Example: Minimizing a Product of Eigenvalues

Let  $S^N$  denote the space of real symmetric  $N \times N$  matrices, and

$$\lambda_1(X) \geq \lambda_2(X) \geq \cdots \lambda_N(X)$$

denote the eigenvalues of  $X \in S^N$ . We wish to minimize

$$f(X) = \log \prod_{i=1}^{N/2} \lambda_i(A \circ X)$$

where  $A \in S^N$  is fixed and  $\circ$  is the Hadamard (componentwise) matrix product, subject to the constraints that  $X$  is positive semidefinite and has diagonal entries equal to 1.

If we replace  $\prod$  by  $\sum$  we would have a semidefinite program.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)

[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)

[Properties of Partial](#)



## Example: Minimizing a Product of Eigenvalues

Let  $S^N$  denote the space of real symmetric  $N \times N$  matrices, and

$$\lambda_1(X) \geq \lambda_2(X) \geq \cdots \lambda_N(X)$$

denote the eigenvalues of  $X \in S^N$ . We wish to minimize

$$f(X) = \log \prod_{i=1}^{N/2} \lambda_i(A \circ X)$$

where  $A \in S^N$  is fixed and  $\circ$  is the Hadamard (componentwise) matrix product, subject to the constraints that  $X$  is positive semidefinite and has diagonal entries equal to 1.

If we replace  $\prod$  by  $\sum$  we would have a semidefinite program.

Since  $f$  is not convex, may as well replace  $X$  by  $YY^T$  where  $Y \in \mathbb{R}^{N \times N}$ : eliminates psd constraint, and then also easy to eliminate diagonal constraint.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon  
Fletcher and Powell

BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example  
Again

Relationships



## Example: Minimizing a Product of Eigenvalues

Let  $S^N$  denote the space of real symmetric  $N \times N$  matrices, and

$$\lambda_1(X) \geq \lambda_2(X) \geq \cdots \lambda_N(X)$$

denote the eigenvalues of  $X \in S^N$ . We wish to minimize

$$f(X) = \log \prod_{i=1}^{N/2} \lambda_i(A \circ X)$$

where  $A \in S^N$  is fixed and  $\circ$  is the Hadamard (componentwise) matrix product, subject to the constraints that  $X$  is positive semidefinite and has diagonal entries equal to 1.

If we replace  $\prod$  by  $\sum$  we would have a semidefinite program.

Since  $f$  is not convex, may as well replace  $X$  by  $YY^T$  where  $Y \in \mathbb{R}^{N \times N}$ : eliminates psd constraint, and then also easy to eliminate diagonal constraint.

Application: entropy minimization in an environmental application (K.M. Anstreicher and J. Lee, 2004)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)

[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)  
[Properties of Partly Smooth Functions](#)



# BFGS from 10 Randomly Generated Starting Points

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

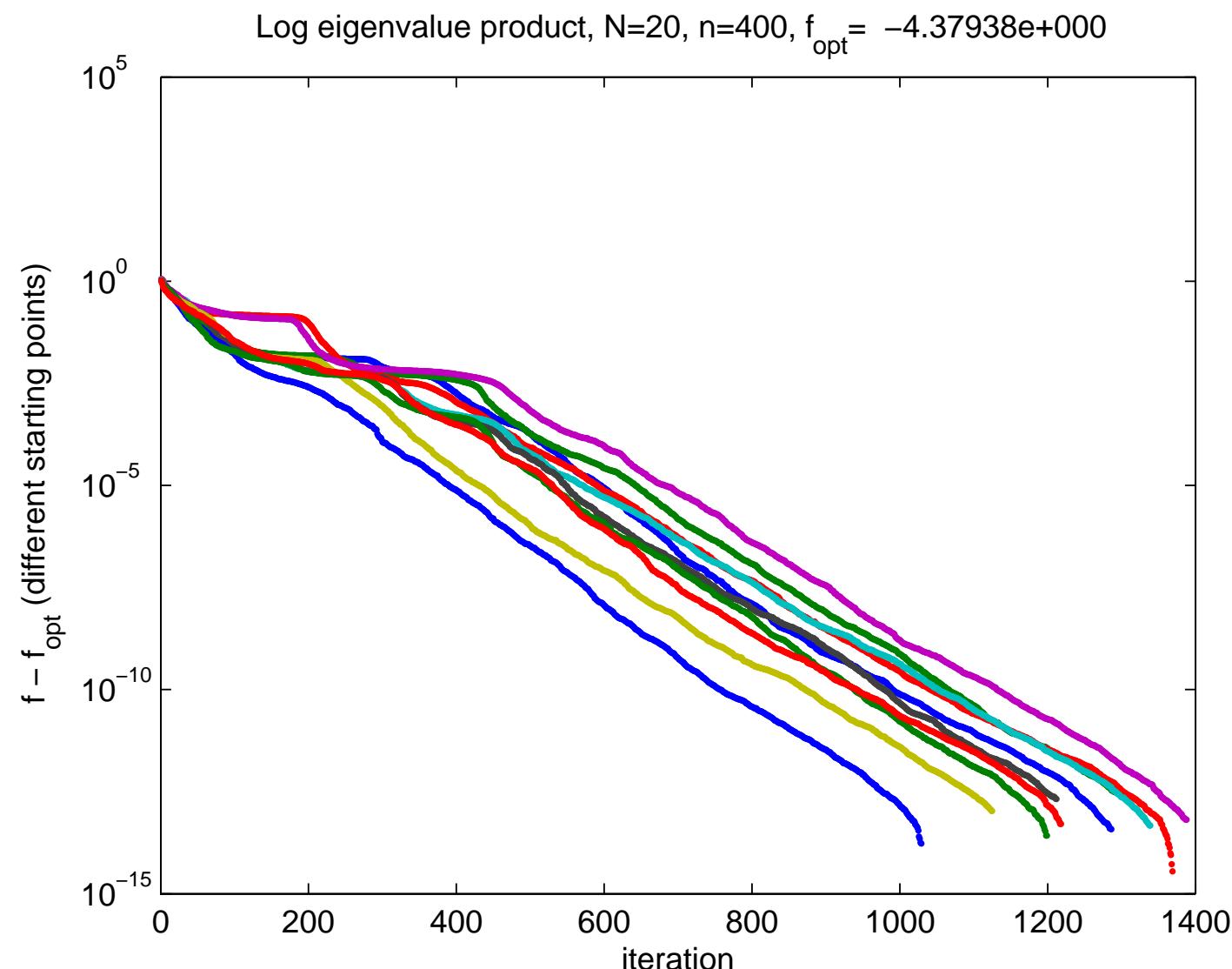
Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example

Again

Relations of Partial



$f - f_{\text{opt}}$ , where  $f_{\text{opt}}$  is least value of  $f$  found over all runs



# Evolution of Eigenvalues of $A \circ X$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

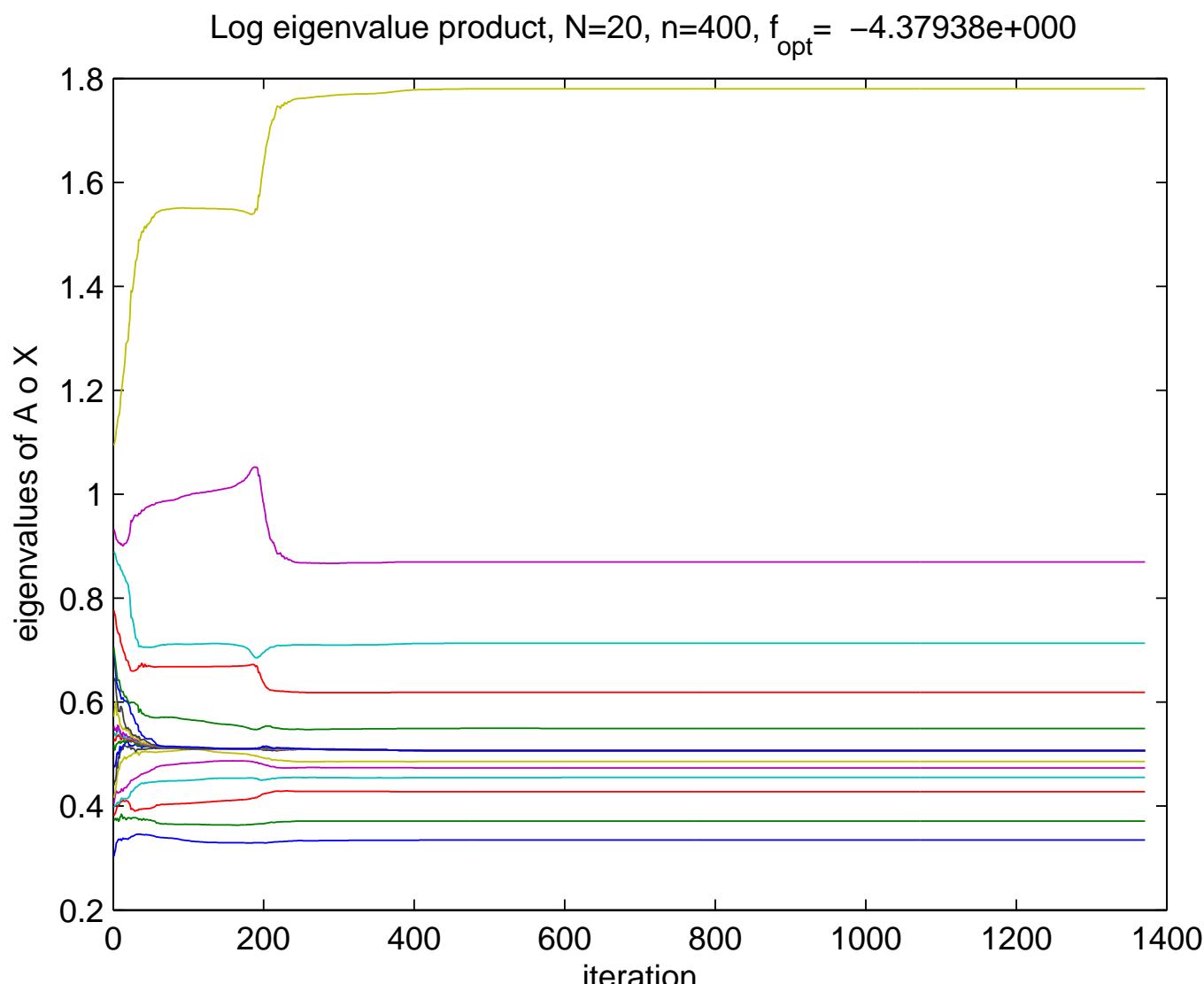
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity  
Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again





# Evolution of Eigenvalues of $A \circ X$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davison

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

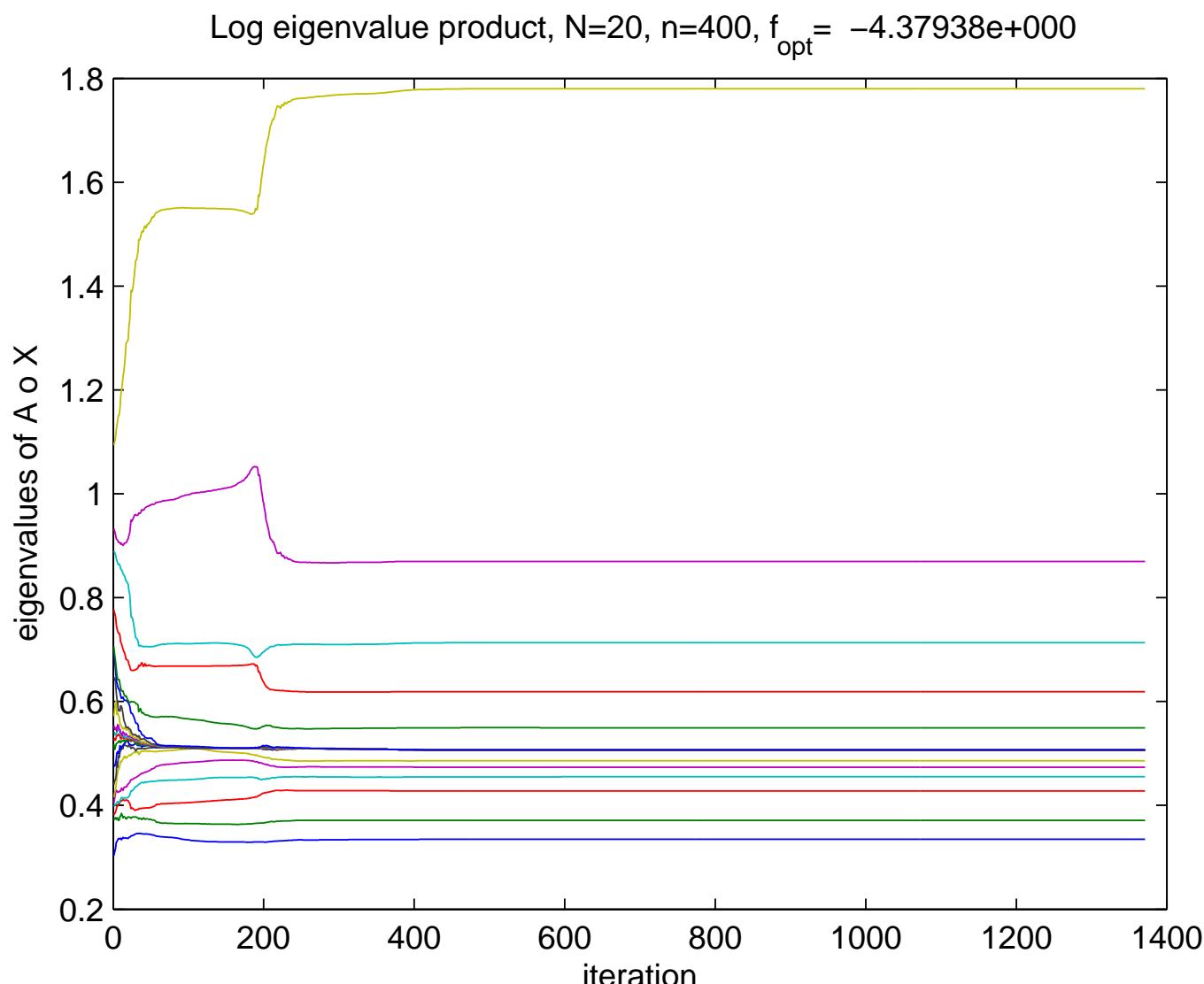
Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example

Again



Note that  $\lambda_6(X), \dots, \lambda_{14}(X)$  coalesce



# Evolution of Eigenvalues of $H$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidson  
Fletcher and Powell  
BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

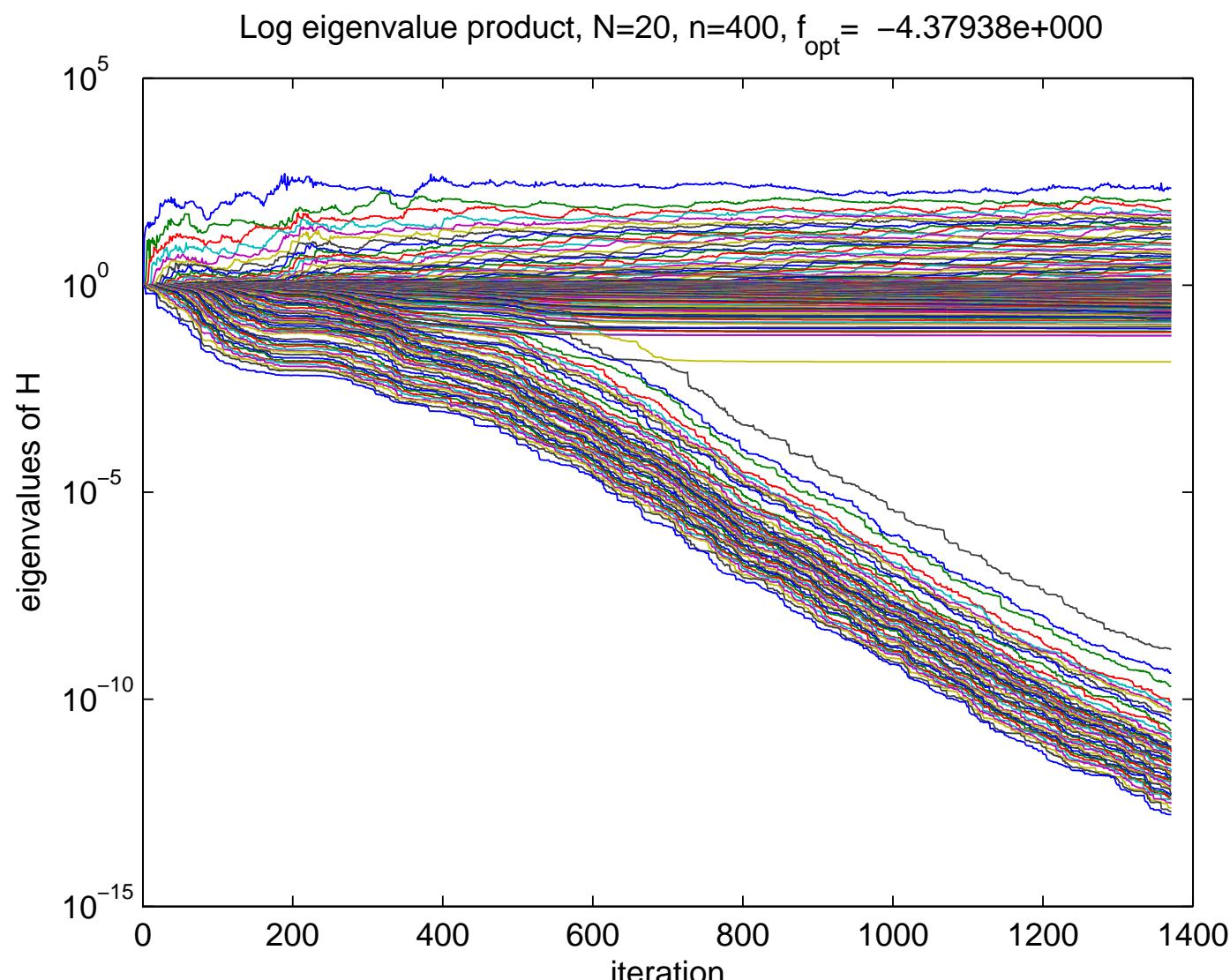
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity  
Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again  
Relations of Partial





# Evolution of Eigenvalues of $H$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidson  
Fletcher and Powell

BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

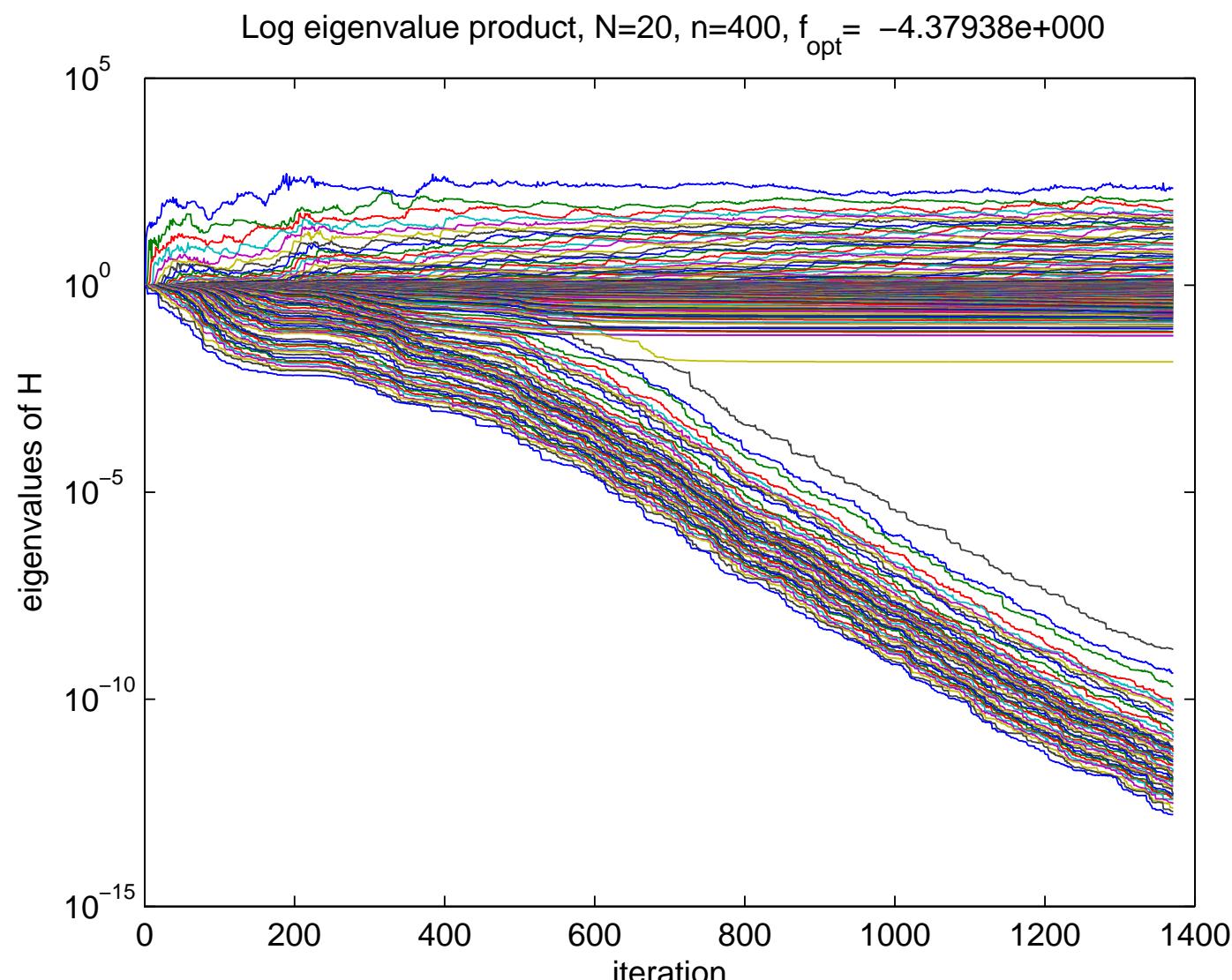
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity  
Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again



44 eigenvalues of  $H$  converge to zero...why???



# Regularity

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Properties of Partial](#)

A locally Lipschitz, directionally differentiable function  $f$  is *regular* (Clarke 1970s) near a point  $x$  when its directional derivative  $f'(\cdot; d)$  is upper semicontinuous there for every fixed direction  $d$ .



# Regularity

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:  
Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partial](#)

A locally Lipschitz, directionally differentiable function  $f$  is *regular* (Clarke 1970s) near a point  $x$  when its directional derivative  $f'(\cdot; d)$  is upper semicontinuous there for every fixed direction  $d$ .

In this case  $0 \in \partial^C f(x)$  is equivalent to the first-order optimality condition  $f'(x, d) \geq 0$  for all directions  $d$ .



# Regularity

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon Fletcher and Powell](#)

[BFGS  
The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:  
Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

A locally Lipschitz, directionally differentiable function  $f$  is *regular* (Clarke 1970s) near a point  $x$  when its directional derivative  $f'(\cdot; d)$  is upper semicontinuous there for every fixed direction  $d$ .

In this case  $0 \in \partial^C f(x)$  is equivalent to the first-order optimality condition  $f'(x, d) \geq 0$  for all directions  $d$ .

■ All convex functions are regular



# Regularity

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:  
Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partial](#)

A locally Lipschitz, directionally differentiable function  $f$  is *regular* (Clarke 1970s) near a point  $x$  when its directional derivative  $f'(\cdot; d)$  is upper semicontinuous there for every fixed direction  $d$ .

In this case  $0 \in \partial^C f(x)$  is equivalent to the first-order optimality condition  $f'(x, d) \geq 0$  for all directions  $d$ .

- All convex functions are regular
- All smooth functions are regular



# Regularity

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:  
Minimizing a Product of Eigenvalues  
BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$   
Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)

A locally Lipschitz, directionally differentiable function  $f$  is *regular* (Clarke 1970s) near a point  $x$  when its directional derivative  $f'(\cdot; d)$  is upper semicontinuous there for every fixed direction  $d$ .

In this case  $0 \in \partial^C f(x)$  is equivalent to the first-order optimality condition  $f'(x, d) \geq 0$  for all directions  $d$ .

- All convex functions are regular
- All smooth functions are regular
- Nonsmooth concave functions are not regular

Example:  $f(x) = -|x|$

Note: this is a somewhat simpler definition of regularity than the one in Lecture 12, but it is less precise: it defines regularity in a neighborhood, not at a point.



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partly Smooth Functions](#)



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$
- the Clarke subdifferential  $\partial f$  is continuous on  $\mathcal{M}$  near  $x$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Relations of Partial](#)



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$
- the Clarke subdifferential  $\partial f$  is continuous on  $\mathcal{M}$  near  $x$
- $\text{par } \partial f(x)$ , the subspace parallel to the affine hull of the subdifferential of  $f$  at  $x$ , is exactly the subspace normal to  $\mathcal{M}$  at  $x$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)

[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)  
[Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relations of Partial](#)



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$
- the Clarke subdifferential  $\partial f$  is continuous on  $\mathcal{M}$  near  $x$
- $\text{par } \partial f(x)$ , the subspace parallel to the affine hull of the subdifferential of  $f$  at  $x$ , is exactly the subspace normal to  $\mathcal{M}$  at  $x$ .

We refer to  $\text{par } \partial f(x)$  as the *V-space* for  $f$  at  $x$  (with respect to  $\mathcal{M}$ ), and to its orthogonal complement, the subspace tangent to  $\mathcal{M}$  at  $x$ , as the *U-space* for  $f$  at  $x$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)  
[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)  
[With BFGS](#)  
[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)

[Relationships](#)



## Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$
- the Clarke subdifferential  $\partial f$  is continuous on  $\mathcal{M}$  near  $x$
- $\text{par } \partial f(x)$ , the subspace parallel to the affine hull of the subdifferential of  $f$  at  $x$ , is exactly the subspace normal to  $\mathcal{M}$  at  $x$ .

We refer to  $\text{par } \partial f(x)$  as the *V-space* for  $f$  at  $x$  (with respect to  $\mathcal{M}$ ), and to its orthogonal complement, the subspace tangent to  $\mathcal{M}$  at  $x$ , as the *U-space* for  $f$  at  $x$ .

When we refer to the *V-space* and *U-space* without reference to a point  $x$ , we mean at a minimizer.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)

[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)  
[With BFGS](#)  
[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

**Partly Smooth Functions**

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



# Partly Smooth Functions

A regular function  $f$  is *partly smooth* at  $x$  relative to a manifold  $\mathcal{M}$  containing  $x$  (A.S. Lewis 2003) if

- its restriction to  $\mathcal{M}$  is twice continuously differentiable near  $x$
- the Clarke subdifferential  $\partial f$  is continuous on  $\mathcal{M}$  near  $x$
- $\text{par } \partial f(x)$ , the subspace parallel to the affine hull of the subdifferential of  $f$  at  $x$ , is exactly the subspace normal to  $\mathcal{M}$  at  $x$ .

We refer to  $\text{par } \partial f(x)$  as the *V-space* for  $f$  at  $x$  (with respect to  $\mathcal{M}$ ), and to its orthogonal complement, the subspace tangent to  $\mathcal{M}$  at  $x$ , as the *U-space* for  $f$  at  $x$ .

When we refer to the *V-space* and *U-space* without reference to a point  $x$ , we mean at a minimizer.

For nonzero  $y$  in the *V-space*, the mapping  $t \mapsto f(x + ty)$  is necessarily nonsmooth at  $t = 0$ , while for nonzero  $y$  in the *U-space*,  $t \mapsto f(x + ty)$  is differentiable at  $t = 0$  as long as  $f$  is locally Lipschitz.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)  
[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)  
[With BFGS](#)  
[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)



## Partly Smooth Functions, continued

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example

Again

Properties of Partly

Example:  $f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2$ .

**Question:** What is  $\mathcal{M}$  and what are the  $U$  and  $V$  spaces at the minimizer?



## Partly Smooth Functions, continued

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partly Smooth Functions](#)

**Example:**  $f(x) = 10|x_2 - x_1^2| + (1 - x_1)^2.$

**Question:** What is  $\mathcal{M}$  and what are the  $U$  and  $V$  spaces at the minimizer?

**Example:**  $f(x) = \|x\|_2.$

**Question:** What is  $\mathcal{M}$  and what are the  $U$  and  $V$  spaces at the minimizer?



## Same Example Again

[Introduction](#)

[Gradient Sampling](#)

Quasi-Newton  
Methods

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

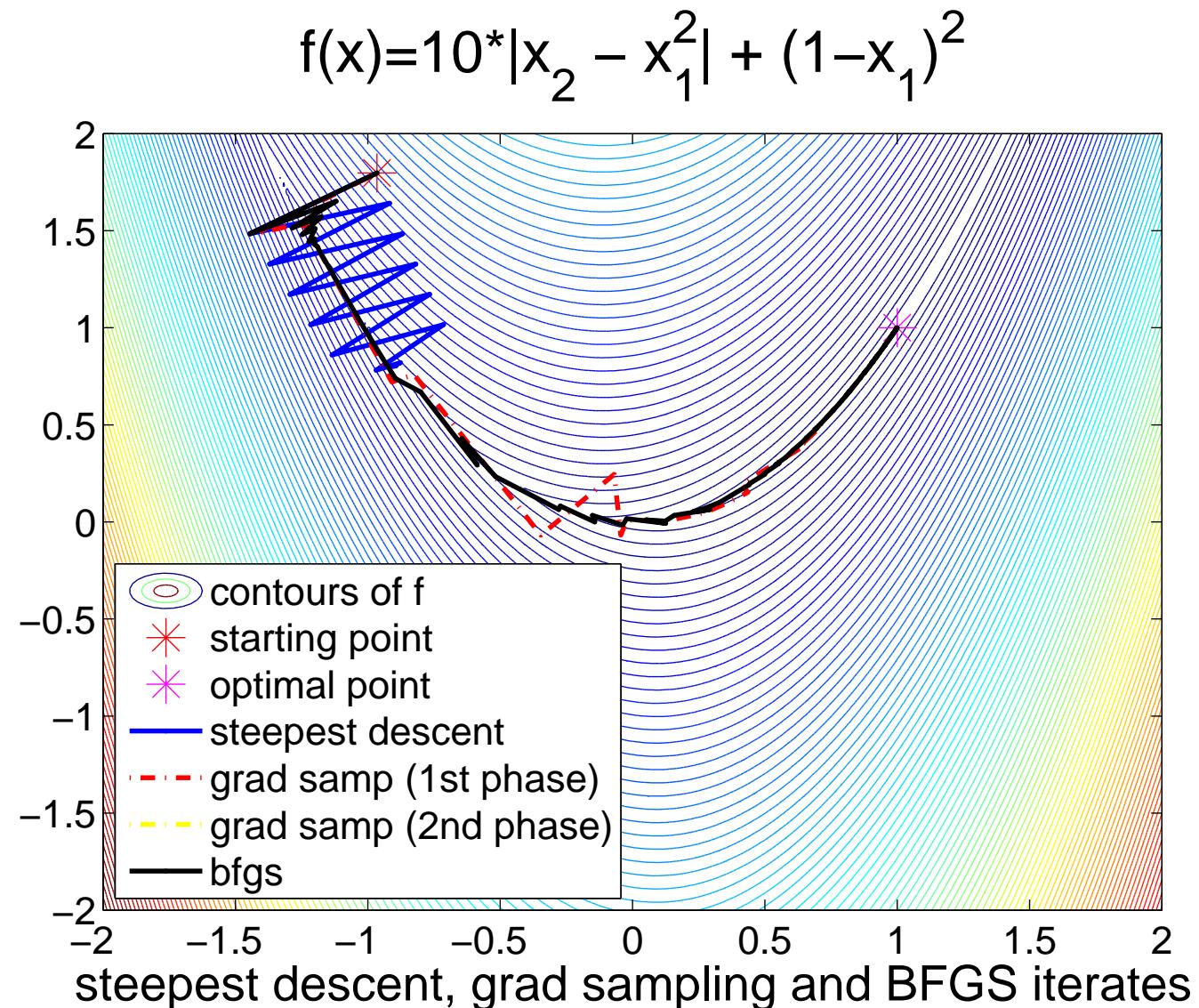
Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example  
Again





## Relation of Partial Smoothness to Earlier Work

Partial smoothness is closely related to earlier work of J.V. Burke and J.J. Moré (1990,1994) and S.J. Wright (1993) on identification of constraint structure by algorithms.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



## Relation of Partial Smoothness to Earlier Work

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

Partial smoothness is closely related to earlier work of J.V. Burke and J.J. Moré (1990,1994) and S.J. Wright (1993) on identification of constraint structure by algorithms.

When  $f$  is convex, the partly smooth nomenclature is consistent with the usage of  $V$ -space and  $U$ -space by C. Lemaréchal, F. Oustry and C. Sagastizábal (2000), but partial smoothness does not imply convexity and convexity does not imply partial smoothness.



# Why Did 44 Eigenvalues of $H$ Converge to Zero?

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

The eigenvalue product is *regular* and also *partly smooth* (in the sense of A.S. Lewis, 2003) with respect to the manifold of matrices with an eigenvalue with given multiplicity. This implies that *tangent* to this manifold (preserving the multiplicity to first-order) the function is *smooth* ("U-shaped") and *normal* to it, the function is *nonsmooth* ("V-shaped").



# Why Did 44 Eigenvalues of $H$ Converge to Zero?

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for](#)

[Nonsmooth Optimization](#)

[With BFGS](#)

[Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10](#)

[Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Properties of Positive](#)

The eigenvalue product is *regular* and also *partly smooth* (in the sense of A.S. Lewis, 2003) with respect to the manifold of matrices with an eigenvalue with given multiplicity. This implies that *tangent* to this manifold (preserving the multiplicity to first-order) the function is *smooth* ("U-shaped") and *normal* to it, the function is *nonsmooth* ("V-shaped").

Recall that at the computed minimizer,

$$\lambda_6(A \circ X) \approx \dots \approx \lambda_{14}(A \circ X).$$

Matrix theory says that imposing multiplicity  $m$  on an eigenvalue a matrix  $\in S^N$  is  $\frac{m(m+1)}{2} - 1$  conditions, or 44 when  $m = 9$ , so the dimension of the  $V$ -space at this minimizer is 44.



# Why Did 44 Eigenvalues of $H$ Converge to Zero?

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10](#)

[Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)

The eigenvalue product is *regular* and also *partly smooth* (in the sense of A.S. Lewis, 2003) with respect to the manifold of matrices with an eigenvalue with given multiplicity. This implies that *tangent* to this manifold (preserving the multiplicity to first-order) the function is *smooth* ("U-shaped") and *normal* to it, the function is *nonsmooth* ("V-shaped").

Recall that at the computed minimizer,

$$\lambda_6(A \circ X) \approx \dots \approx \lambda_{14}(A \circ X).$$

Matrix theory says that imposing multiplicity  $m$  on an eigenvalue a matrix  $\in S^N$  is  $\frac{m(m+1)}{2} - 1$  conditions, or 44 when  $m = 9$ , so the dimension of the  $V$ -space at this minimizer is 44.

Tiny eigenvalues of  $H$  correspond to huge curvature, which corresponds to  $V$ -space directions.



# Why Did 44 Eigenvalues of $H$ Converge to Zero?

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)  
[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)  
Example:  
Minimizing a Product of Eigenvalues

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)   
Regularity

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)  
[Same Example Again](#)  
[Properties of Partly Smooth Functions](#)

The eigenvalue product is *regular* and also *partly smooth* (in the sense of A.S. Lewis, 2003) with respect to the manifold of matrices with an eigenvalue with given multiplicity. This implies that *tangent* to this manifold (preserving the multiplicity to first-order) the function is *smooth* ("U-shaped") and *normal* to it, the function is *nonsmooth* ("V-shaped").

Recall that at the computed minimizer,

$$\lambda_6(A \circ X) \approx \dots \approx \lambda_{14}(A \circ X).$$

Matrix theory says that imposing multiplicity  $m$  on an eigenvalue a matrix  $\in S^N$  is  $\frac{m(m+1)}{2} - 1$  conditions, or 44 when  $m = 9$ , so the dimension of the  $V$ -space at this minimizer is 44.

Tiny eigenvalues of  $H$  correspond to huge curvature, which corresponds to  $V$ -space directions.

Thus BFGS *automatically* detected the  $U$  and  $V$  space partitioning without knowing anything about the mathematical structure of  $f$ !



# Variation of $f$ from Minimizer, along EigVecs of $H$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidson

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

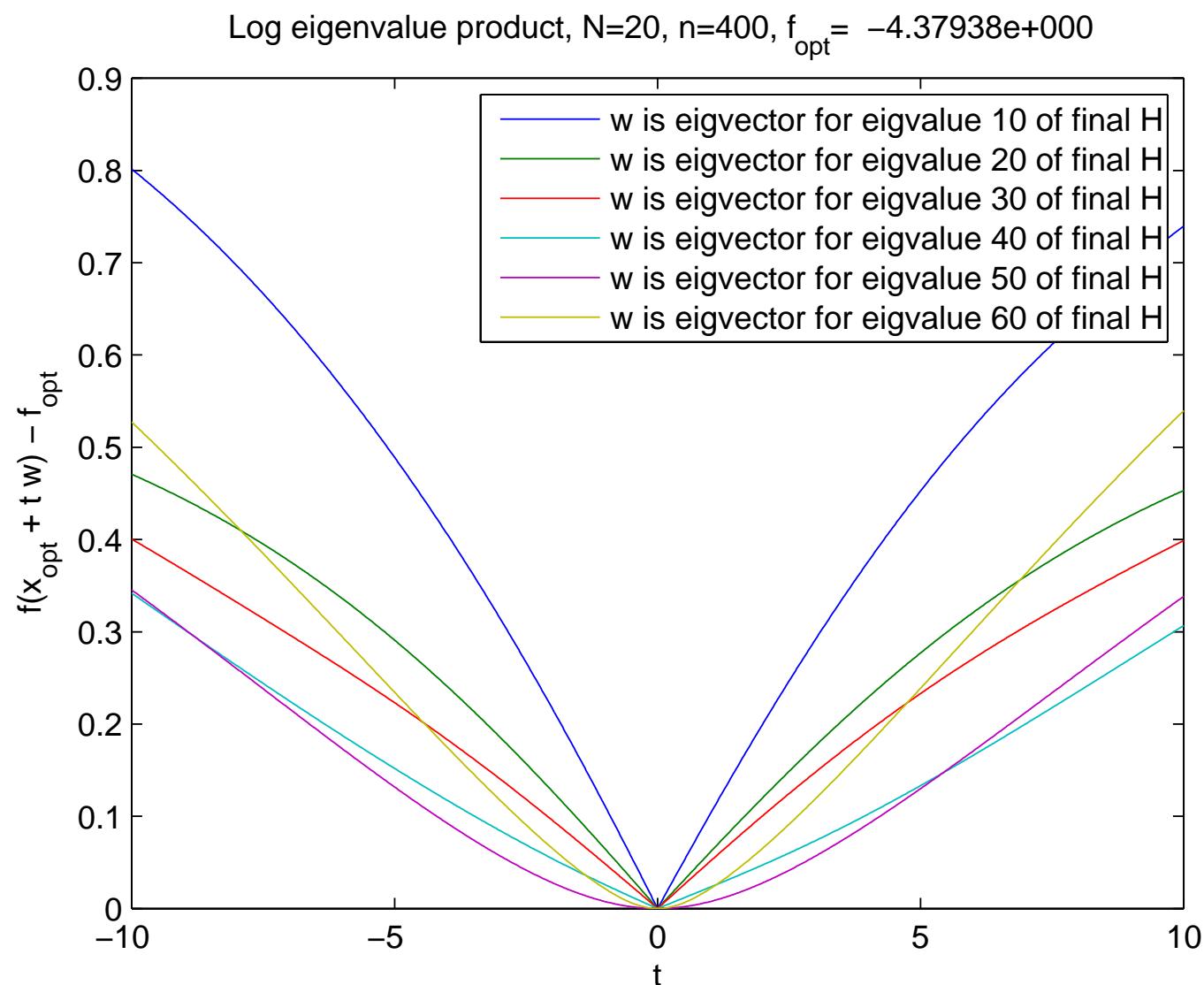
Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued  
Same Example  
Again



Eigenvalues of  $H$  numbered smallest to largest



# BFGS Theory for Special Nonsmooth Functions

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued  
Same Example  
Again](#)

Convergence results for BFGS with Armijo-Wolfe line search when  $f$  is nonsmooth are limited to very special cases.

- $f(x) = |x|$  (one variable!): sequence generated converging to 0 is related to a certain binary expansion of the starting point (A.S. Lewis and M.L.O., 2013)



# BFGS Theory for Special Nonsmooth Functions

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Properties of Partial](#)

Convergence results for BFGS with Armijo-Wolfe line search when  $f$  is nonsmooth are limited to very special cases.

- $f(x) = |x|$  (one variable!): sequence generated converging to 0 is related to a certain binary expansion of the starting point (A.S. Lewis and M.L.O., 2013)
- $f(x) = |x_1| + x_2$ :  $f(x) \downarrow -\infty$  (A.S. Lewis and Shanshan Zhang, 2015)



# BFGS Theory for Special Nonsmooth Functions

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS](#)

[Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Properties of Partial](#)

Convergence results for BFGS with Armijo-Wolfe line search when  $f$  is nonsmooth are limited to very special cases.

- $f(x) = |x|$  (one variable!): sequence generated converging to 0 is related to a certain binary expansion of the starting point (A.S. Lewis and M.L.O., 2013)
- $f(x) = |x_1| + x_2$ :  $f(x) \downarrow -\infty$  (A.S. Lewis and Shanshan Zhang, 2015)
- $f(x) = |x_1| + \sum_{i=2}^n x_i$ : eventually a direction is identified on which  $f$  is unbounded below (Yuchen Xie and A. Waechter, 2017)



# BFGS Theory for Special Nonsmooth Functions

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example](#)

[Again](#)

[Properties of Partial](#)

Convergence results for BFGS with Armijo-Wolfe line search when  $f$  is nonsmooth are limited to very special cases.

- $f(x) = |x|$  (one variable!): sequence generated converging to 0 is related to a certain binary expansion of the starting point (A.S. Lewis and M.L.O., 2013)
- $f(x) = |x_1| + x_2$ :  $f(x) \downarrow -\infty$  (A.S. Lewis and Shanshan Zhang, 2015)
- $f(x) = |x_1| + \sum_{i=2}^n x_i$ : eventually a direction is identified on which  $f$  is unbounded below (Yuchen Xie and A. Waechter, 2017)
- $f(x) = \sqrt{\sum_{i=1}^n x_i^2}$ : iterates converge to  $[0, \dots, 0]$  (Jiayi Guo and A.S. Lewis, 2017) (proof based on Powell (1976))



# BFGS Theory for Special Nonsmooth Functions

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell  
BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$   
Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued  
Same Example  
Again](#)

[Properties of Partly](#)

Convergence results for BFGS with Armijo-Wolfe line search when  $f$  is nonsmooth are limited to very special cases.

- $f(x) = |x|$  (one variable!): sequence generated converging to 0 is related to a certain binary expansion of the starting point (A.S. Lewis and M.L.O., 2013)
- $f(x) = |x_1| + x_2$ :  $f(x) \downarrow -\infty$  (A.S. Lewis and Shanshan Zhang, 2015)
- $f(x) = |x_1| + \sum_{i=2}^n x_i$ : eventually a direction is identified on which  $f$  is unbounded below (Yuchen Xie and A. Waechter, 2017)
- $f(x) = \sqrt{\sum_{i=1}^n x_i^2}$ : iterates converge to  $[0, \dots, 0]$  (Jiayi Guo and A.S. Lewis, 2017) (proof based on Powell (1976))
- $f(x) = |x_1| + x_2^2$ : remains open!



# Challenge: General Nonsmooth Case

[Introduction](#)

---

[Gradient Sampling](#)

---

[Quasi-Newton  
Methods](#)

---

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method  
\(“Full” Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example](#)

[Again](#)

[Relatives of Partial](#)



## Challenge: General Nonsmooth Case

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relatives of Partial](#)



## Challenge: General Nonsmooth Case

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Relationships](#)



## Challenge: General Nonsmooth Case

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)  
[Fletcher and Powell](#)  
[BFGS](#)  
[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example:](#)  
[Minimizing a Product of Eigenvalues](#)  
[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)  
[Partly Smooth Functions](#)  
[Partly Smooth Functions, continued](#)  
[Same Example Again](#)  
[Properties of Partial](#)

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

Prove or disprove that the following hold with probability one:



## Challenge: General Nonsmooth Case

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS](#)

[Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example](#)

[Again](#)

[Properties of Partial](#)

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

Prove or disprove that the following hold with probability one:

1. BFGS generates an infinite sequence  $\{x\}$  with  $f$  differentiable at all iterates



## Challenge: General Nonsmooth Case

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization  
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example](#)

[Again](#)

[Properties of Partial](#)

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

Prove or disprove that the following hold with probability one:

1. BFGS generates an infinite sequence  $\{x\}$  with  $f$  differentiable at all iterates
2. Any cluster point  $\bar{x}$  is Clarke stationary



## Challenge: General Nonsmooth Case

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization  
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again  
Properties of Partial](#)

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

Prove or disprove that the following hold with probability one:

1. BFGS generates an infinite sequence  $\{x\}$  with  $f$  differentiable at all iterates
2. Any cluster point  $\bar{x}$  is Clarke stationary
3. The sequence of function values generated (including all of the line search iterates) converges to  $f(\bar{x})$  R-linearly



## Challenge: General Nonsmooth Case

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization  
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$   
Evolution of  
Eigenvalues of  \$H\$   
Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued  
Same Example  
Again  
Properties of Partial](#)

Assume  $f$  is locally Lipschitz with bounded level sets and is semi-algebraic (its graph is a finite union of sets each defined by a finite list of polynomial inequalities)

Assume the initial  $x$  and  $H$  are generated randomly (e.g. from normal and Wishart distributions)

Prove or disprove that the following hold with probability one:

1. BFGS generates an infinite sequence  $\{x\}$  with  $f$  differentiable at all iterates
2. Any cluster point  $\bar{x}$  is Clarke stationary
3. The sequence of function values generated (including all of the line search iterates) converges to  $f(\bar{x})$  R-linearly
4. If  $\{x\}$  converges to  $\bar{x}$  where  $f$  is "partly smooth" w.r.t. a manifold  $\mathcal{M}$  then the subspace defined by the eigenvectors corresponding to eigenvalues of  $H$  converging to zero converges to the "V-space" of  $f$  w.r.t.  $\mathcal{M}$  at  $\bar{x}$



# Some BFGS Nonsmooth Success Stories

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method  
\(“Full” Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued](#)

[Same Example  
Again](#)

[Relationships](#)



## Some BFGS Nonsmooth Success Stories

- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth  
Optimization](#)

[With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues](#)

[BFGS from 10  
Randomly Generated  
Starting Points](#)

[Evolution of  
Eigenvalues of  
 \$A \circ X\$](#)

[Evolution of  
Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth  
Functions](#)

[Partly Smooth  
Functions, continued  
Same Example  
Again](#)

[Properties of Partial](#)



# Some BFGS Nonsmooth Success Stories

- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)
- Shape optimization for spectral functions of Dirichlet-Laplacian operators (B. Osting, 2010)

[Introduction](#)

[Gradient Sampling](#)

Quasi-Newton  
Methods

Bill Davidon  
Fletcher and Powell

BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity  
Partly Smooth  
Functions

Partly Smooth  
Functions, continued  
Same Example  
Again

Properties of Partial



## Some BFGS Nonsmooth Success Stories

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth Optimization](#)

[With BFGS  
Example:  
Minimizing a Product of Eigenvalues  
BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$   
Evolution of Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth Functions  
Partly Smooth Functions, continued  
Same Example Again  
Properties of Partial](#)

- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)
- Shape optimization for spectral functions of Dirichlet-Laplacian operators (B. Osting, 2010)
- Condition metric optimization (P. Boito and J. Dedieu, 2010)



# Some BFGS Nonsmooth Success Stories

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon  
Fletcher and Powell](#)

[BFGS  
The BFGS Method  
\("Full" Version\)](#)

[BFGS for  
Nonsmooth Optimization](#)

[With BFGS  
Example:  
Minimizing a Product of Eigenvalues  
BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$   
Evolution of Eigenvalues of  \$H\$](#)

[Regularity  
Partly Smooth Functions  
Partly Smooth Functions, continued  
Same Example Again  
Properties of Partial](#)

- Design of fixed-order controllers for linear dynamical systems with input and output (D. Henrion and M.L.O., 2006, and many subsequent users of our HIFOO (H-Infinity Fixed Order Optimization) toolbox)
- Shape optimization for spectral functions of Dirichlet-Laplacian operators (B. Osting, 2010)
- Condition metric optimization (P. Boito and J. Dedieu, 2010)

Software is available: HANSO



# Extensions of BFGS for Nonsmooth Optimization

A combined BFGS-Gradient Sampling method with convergence theory (F.E. Curtis and X. Que, 2015)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[Bill Davidon](#)

[Fletcher and Powell](#)

[BFGS](#)

[The BFGS Method \("Full" Version\)](#)

[BFGS for Nonsmooth Optimization](#)

[With BFGS Example: Minimizing a Product of Eigenvalues](#)

[BFGS from 10 Randomly Generated Starting Points](#)

[Evolution of Eigenvalues of  \$A \circ X\$](#)

[Evolution of Eigenvalues of  \$H\$](#)

[Regularity](#)

[Partly Smooth Functions](#)

[Partly Smooth Functions, continued](#)

[Same Example Again](#)

[Relationships](#)



# Extensions of BFGS for Nonsmooth Optimization

A combined BFGS-Gradient Sampling method with convergence theory (F.E. Curtis and X. Que, 2015)

## Constrained Problems

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS

Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example

Again

Properties of Partial



# Extensions of BFGS for Nonsmooth Optimization

A combined BFGS-Gradient Sampling method with convergence theory (F.E. Curtis and X. Que, 2015)

## Constrained Problems

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

A successive quadratic programming (SQP) BFGS method applied to challenging problems in static-output-feedback control design (F.E. Curtis, T. Mitchell and M.L.O., 2015).

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon

Fletcher and Powell

BFGS

The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization

With BFGS

Example:  
Minimizing a  
Product of  
Eigenvalues

BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$

Evolution of  
Eigenvalues of  $H$

Regularity

Partly Smooth  
Functions

Partly Smooth  
Functions, continued

Same Example

Again

Properties of Partial



# Extensions of BFGS for Nonsmooth Optimization

A combined BFGS-Gradient Sampling method with convergence theory (F.E. Curtis and X. Que, 2015)

## Constrained Problems

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

A successive quadratic programming (SQP) BFGS method applied to challenging problems in static-output-feedback control design (F.E. Curtis, T. Mitchell and M.L.O., 2015).

Although there are no theoretical results, it is much more efficient and effective than the SQP Gradient Sampling method which does have convergence results.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon  
Fletcher and Powell  
BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization  
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$   
Evolution of  
Eigenvalues of  $H$   
Regularity

Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again  
Properties of Partial



# Extensions of BFGS for Nonsmooth Optimization

A combined BFGS-Gradient Sampling method with convergence theory (F.E. Curtis and X. Que, 2015)

## Constrained Problems

$$\begin{aligned} & \min f(x) \\ & \text{subject to } c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where  $f$  and  $c_1, \dots, c_p$  are locally Lipschitz but may not be differentiable at local minimizers.

A successive quadratic programming (SQP) BFGS method applied to challenging problems in static-output-feedback control design (F.E. Curtis, T. Mitchell and M.L.O., 2015).

Although there are no theoretical results, it is much more efficient and effective than the SQP Gradient Sampling method which does have convergence results.

Software is available: GRANSO

Introduction

Gradient Sampling

Quasi-Newton  
Methods

Bill Davidon  
Fletcher and Powell  
BFGS  
The BFGS Method  
("Full" Version)

BFGS for  
Nonsmooth  
Optimization  
With BFGS  
Example:  
Minimizing a  
Product of  
Eigenvalues  
BFGS from 10  
Randomly Generated  
Starting Points

Evolution of  
Eigenvalues of  
 $A \circ X$   
Evolution of  
Eigenvalues of  $H$   
Regularity

Partly Smooth  
Functions  
Partly Smooth  
Functions, continued  
Same Example  
Again  
Properties of Partial



[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

**A Difficult  
Nonconvex Problem  
from Nesterov**

An Aside:  
Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials  
Nesterov's  
Chebyshev-  
Rosenbrock  
Functions  
Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$   
First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$   
Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$   
Behavior of BFGS  
on the Second  
Nonsmooth Variant

# A Difficult Nonconvex Problem from Nesterov



## An Aside: Chebyshev Polynomials

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

**An Aside:  
Chebyshev  
Polynomials**

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

Plots of Chebyshev Polynomials

Nesterov's Chebyshev-Rosenbrock Functions

Why BFGS Takes So Many Iterations to Minimize  $N_2$

First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function  
Contour Plots of the Nonsmooth Variants for  $n = 2$

Properties of the Second Nonsmooth Variant  $\hat{N}_1$

Behavior of BFGS on the Second Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

$$\text{So } T_2(x) = 2x^2 - 1,$$



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

Plots of Chebyshev Polynomials

Nesterov's Chebyshev-Rosenbrock Functions

Why BFGS Takes So Many Iterations to Minimize  $N_2$

First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function  
Contour Plots of the Nonsmooth Variants for  $n = 2$

Properties of the Second Nonsmooth Variant  $\hat{N}_1$

Behavior of BFGS on the Second Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

So  $T_2(x) = 2x^2 - 1$ ,  $T_3(x) = 4x^3 - 3$ , etc.



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

[Plots of Chebyshev Polynomials](#)

Nesterov's Chebyshev-Rosenbrock Functions

Why BFGS Takes So Many Iterations to Minimize  $N_2$

First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function

Contour Plots of the Nonsmooth Variants for  $n = 2$

Properties of the Second Nonsmooth Variant  $\hat{N}_1$

Behavior of BFGS on the Second Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

So  $T_2(x) = 2x^2 - 1$ ,  $T_3(x) = 4x^3 - 3$ , etc.

Important properties that can be proved easily include



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

Plots of Chebyshev  
Polynomials

Nesterov's  
Chebyshev-  
Rosenbrock  
Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$

First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$

Behavior of BFGS  
on the Second  
Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

So  $T_2(x) = 2x^2 - 1$ ,  $T_3(x) = 4x^3 - 3$ , etc.

Important properties that can be proved easily include

- $T_n(x) = \cos(n \cos^{-1}(x))$



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

Plots of Chebyshev  
Polynomials

Nesterov's  
Chebyshev-  
Rosenbrock

Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$

First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$   
Behavior of BFGS  
on the Second  
Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

So  $T_2(x) = 2x^2 - 1$ ,  $T_3(x) = 4x^3 - 3$ , etc.

Important properties that can be proved easily include

- $T_n(x) = \cos(n \cos^{-1}(x))$
- $T_m(T_n(x)) = T_{mn}(x)$



## An Aside: Chebyshev Polynomials

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

A Difficult Nonconvex Problem from Nesterov

An Aside:  
Chebyshev  
Polynomials

Plots of Chebyshev  
Polynomials

Nesterov's  
Chebyshev-  
Rosenbrock

Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$

First Nonsmooth  
Variant of  
Nesterov's Function

Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$

Behavior of BFGS  
on the Second  
Nonsmooth Variant

A sequence of orthogonal polynomials defined on  $[-1, 1]$  by

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$

So  $T_2(x) = 2x^2 - 1$ ,  $T_3(x) = 4x^3 - 3$ , etc.

Important properties that can be proved easily include

- $T_n(x) = \cos(n \cos^{-1}(x))$
- $T_m(T_n(x)) = T_{mn}(x)$
- $\int_{-1}^1 \frac{1}{\sqrt{1-x^2}} T_i(x) T_j(x) dx = 0$  if  $i \neq j$



# Plots of Chebyshev Polynomials

Introduction

Gradient Sam-

Quasi-Newto-

Methods

A Difficult  
Nonconvex I-

from Nesterov's  
An Aside:  
Chebyshev  
Polynomials

Plots of Che-

Nesterov's  
Chebyshev-  
Rosenbrock

Functions  
Why BFGS

Many Iterati-

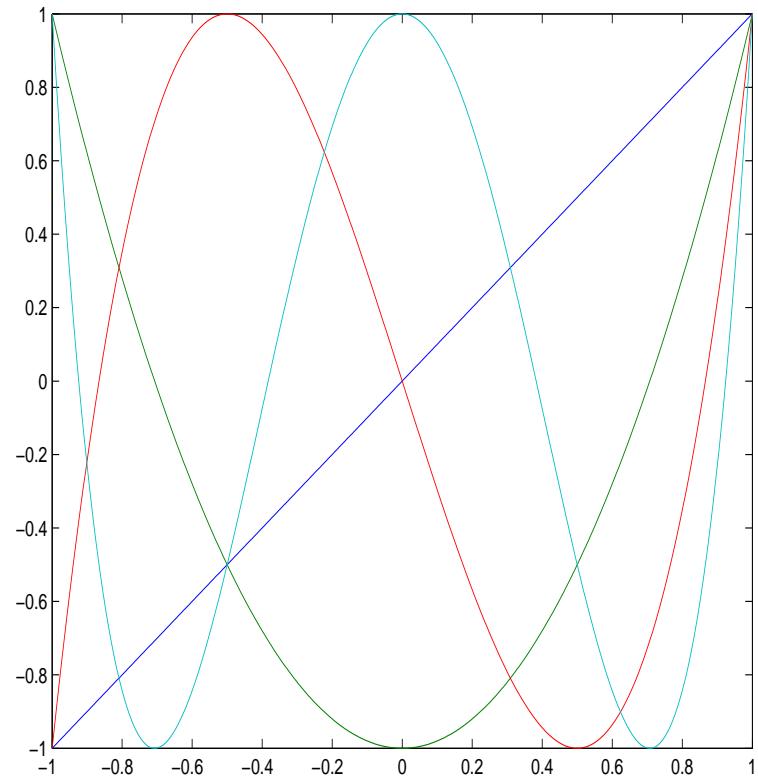
Minimize  $N$   
First Nonsm-

Variant of  
Nesterov's F

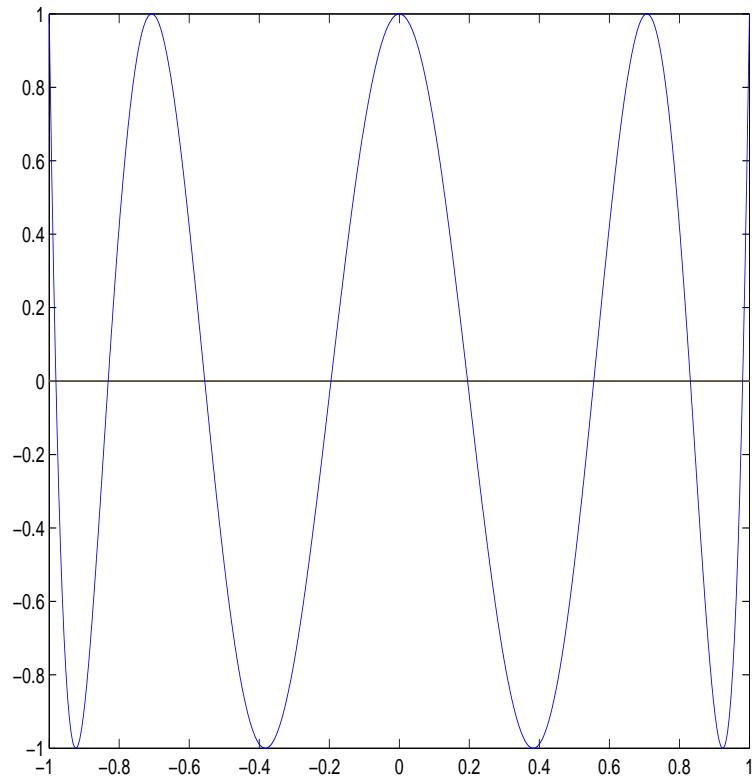
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$

Behavior of BFGS  
on the Second  
Nonsmooth Variant



Left: Plots of  $T_0(x), \dots, T_4(x)$



Right: Plot of  $T_8(x)$ .



# Plots of Chebyshev Polynomials

Introduction

Gradient Sam-

Quasi-Newto-

Methods

A Difficult  
Nonconvex I-

from Nesterov's  
An Aside:  
Chebyshev  
Polynomials

Plots of Che-

by Nesterov's  
Chebyshev-  
Rosenbrock

Functions

Why BFGS

Many Iterati-

Minimize  $N$   
First Nonsm-

Variant of  
Nesterov's F

Second Nonsmooth

Variant of  
Nesterov's Function

Contour Plots of the

Nonsmooth Variants

for  $n = 2$

Properties of the

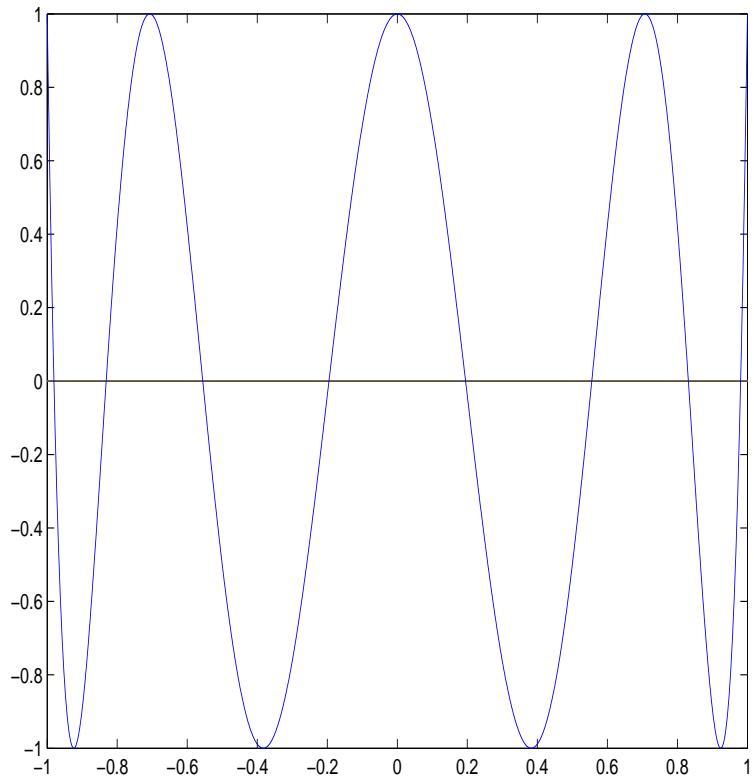
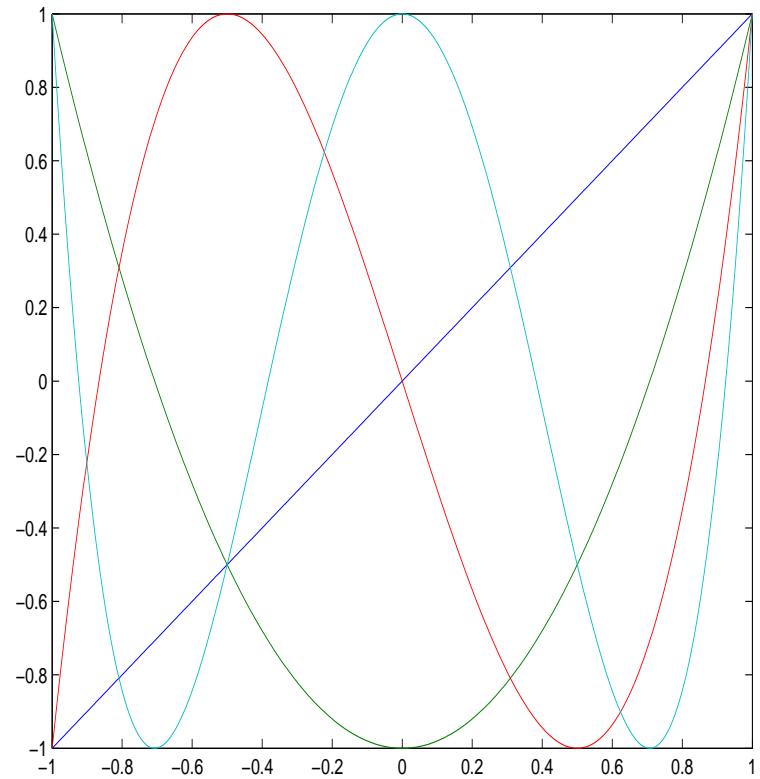
Second Nonsmooth

Variant  $\hat{N}_1$

Behavior of BFGS

on the Second

Nonsmooth Variant



Left: Plots of  $T_0(x), \dots, T_4(x)$

Right: Plot of  $T_8(x)$ .

Question: How many extrema does  $T_n(x)$  have in  $[-1, 1]$ ?



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n - 1\}$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n-1\}$$

For  $x \in \mathcal{M}_N$ , e.g.  $x = x^*$  or  $x = \hat{x}$ , the 2nd term of  $N_p$  is zero. Starting at  $\hat{x}$ , BFGS needs to approximately follow  $\mathcal{M}_N$  to reach  $x^*$  (unless it “gets lucky”).

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n - 1\}$$

For  $x \in \mathcal{M}_N$ , e.g.  $x = x^*$  or  $x = \hat{x}$ , the 2nd term of  $N_p$  is zero. Starting at  $\hat{x}$ , BFGS needs to approximately follow  $\mathcal{M}_N$  to reach  $x^*$  (unless it “gets lucky”).

When  $p = 2$ :  $N_2$  is smooth but not convex. Starting at  $\hat{x}$ :

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n-1\}$$

For  $x \in \mathcal{M}_N$ , e.g.  $x = x^*$  or  $x = \hat{x}$ , the 2nd term of  $N_p$  is zero. Starting at  $\hat{x}$ , BFGS needs to approximately follow  $\mathcal{M}_N$  to reach  $x^*$  (unless it “gets lucky”).

When  $p = 2$ :  $N_2$  is smooth but not convex. Starting at  $\hat{x}$ :

- $n = 5$ : BFGS needs 370 iterations to reduce  $N_2$  below  $10^{-15}$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n-1\}$$

For  $x \in \mathcal{M}_N$ , e.g.  $x = x^*$  or  $x = \hat{x}$ , the 2nd term of  $N_p$  is zero. Starting at  $\hat{x}$ , BFGS needs to approximately follow  $\mathcal{M}_N$  to reach  $x^*$  (unless it “gets lucky”).

When  $p = 2$ :  $N_2$  is smooth but not convex. Starting at  $\hat{x}$ :

- $n = 5$ : BFGS needs 370 iterations to reduce  $N_2$  below  $10^{-15}$
- $n = 10$ : needs  $\sim 50,000$  iterations to reduce  $N_2$  below  $10^{-15}$

even though  $N_2$  is *smooth*!

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

An Aside:  
Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials

Nesterov's  
Chebyshev-  
Rosenbrock  
Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$   
First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$   
Behavior of BFGS  
on the Second  
Nonsmooth Variant



# Nesterov's Chebyshev-Rosenbrock Functions

Consider the function

$$N_p(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|^p, \quad \text{where } p \geq 1$$

The unique minimizer is  $x^* = [1, 1, \dots, 1]^T$  with  $N_p(x^*) = 0$ .

Define  $\hat{x} = [-1, 1, 1, \dots, 1]^T$  with  $N_p(\hat{x}) = 1$  and the manifold

$$\mathcal{M}_N = \{x : x_{i+1} = 2x_i^2 - 1, \quad i = 1, \dots, n-1\}$$

For  $x \in \mathcal{M}_N$ , e.g.  $x = x^*$  or  $x = \hat{x}$ , the 2nd term of  $N_p$  is zero. Starting at  $\hat{x}$ , BFGS needs to approximately follow  $\mathcal{M}_N$  to reach  $x^*$  (unless it “gets lucky”).

When  $p = 2$ :  $N_2$  is smooth but not convex. Starting at  $\hat{x}$ :

- $n = 5$ : BFGS needs 370 iterations to reduce  $N_2$  below  $10^{-15}$
- $n = 10$ : needs  $\sim 50,000$  iterations to reduce  $N_2$  below  $10^{-15}$

even though  $N_2$  is *smooth!* . . . In the last few iterations, we observe superlinear convergence!

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials  
Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$   
First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function  
Contour Plots of the Nonsmooth Variants for  \$n = 2\$   
Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$x_{i+1} = 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1}))$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants](#)

[for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants](#)

[for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants](#)

[for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$
- $x_2 = 2x_1^2 - 1$  to trace the graph of  $T_2(x_1)$  on  $[-1, 1]$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

**Why BFGS Takes So Many Iterations to Minimize  $N_2$**

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$
- $x_2 = 2x_1^2 - 1$  to trace the graph of  $T_2(x_1)$  on  $[-1, 1]$
- $x_3 = T_2(T_2(x_1))$  to trace the graph of  $T_4(x_1)$  on  $[-1, 1]$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)  
[Nesterov's Chebyshev-Rosenbrock Functions](#)

**Why BFGS Takes So Many Iterations to Minimize  $N_2$**

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)   
[Properties of the Second Nonsmooth Variant  \$\widehat{N}\_1\$](#)   
[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$
- $x_2 = 2x_1^2 - 1$  to trace the graph of  $T_2(x_1)$  on  $[-1, 1]$
- $x_3 = T_2(T_2(x_1))$  to trace the graph of  $T_4(x_1)$  on  $[-1, 1]$
- $x_n = T_{2^{n-1}}(x_1)$  to trace the graph of  $T_{2^{n-1}}(x_1)$  on  $[-1, 1]$

which has  $2^{n-1} - 1$  extrema in  $(-1, 1)$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)  
[Nesterov's Chebyshev-Rosenbrock Functions](#)

**Why BFGS Takes So Many Iterations to Minimize  $N_2$**

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\widehat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$
- $x_2 = 2x_1^2 - 1$  to trace the graph of  $T_2(x_1)$  on  $[-1, 1]$
- $x_3 = T_2(T_2(x_1))$  to trace the graph of  $T_4(x_1)$  on  $[-1, 1]$
- $x_n = T_{2^{n-1}}(x_1)$  to trace the graph of  $T_{2^{n-1}}(x_1)$  on  $[-1, 1]$

which has  $2^{n-1} - 1$  extrema in  $(-1, 1)$ .

Even though BFGS will *not* track the manifold  $\mathcal{M}_N$  exactly, it will follow it approximately. So, since the manifold is highly oscillatory, BFGS must take relatively short steps to obtain reduction in  $N_2$  in the line search, and hence *many* iterations!

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

An Aside:  
Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials  
Nesterov's  
Chebyshev-  
Rosenbrock  
Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$

First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$

Behavior of BFGS  
on the Second  
Nonsmooth Variant



# Why BFGS Takes So Many Iterations to Minimize $N_2$

Let  $T_i(x)$  denote the  $i$ th Chebyshev polynomial. For  $x \in \mathcal{M}_N$ ,

$$\begin{aligned}x_{i+1} &= 2x_i^2 - 1 = T_2(x_i) = T_2(T_2(x_{i-1})) \\&= T_2(T_2(\dots T_2(x_1)\dots)) = T_{2^i}(x_1).\end{aligned}$$

To move from  $\hat{x}$  to  $x^*$  along the manifold  $\mathcal{M}_N$  exactly requires

- $x_1$  to change from  $-1$  to  $1$
- $x_2 = 2x_1^2 - 1$  to trace the graph of  $T_2(x_1)$  on  $[-1, 1]$
- $x_3 = T_2(T_2(x_1))$  to trace the graph of  $T_4(x_1)$  on  $[-1, 1]$
- $x_n = T_{2^{n-1}}(x_1)$  to trace the graph of  $T_{2^{n-1}}(x_1)$  on  $[-1, 1]$

which has  $2^{n-1} - 1$  extrema in  $(-1, 1)$ .

Even though BFGS will *not* track the manifold  $\mathcal{M}_N$  exactly, it will follow it approximately. So, since the manifold is highly oscillatory, BFGS must take relatively short steps to obtain reduction in  $N_2$  in the line search, and hence *many* iterations!

Newton's method is not much faster, although it converges quadratically at the end.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

An Aside:  
Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials  
Nesterov's  
Chebyshev-  
Rosenbrock  
Functions

Why BFGS Takes So  
Many Iterations to  
Minimize  $N_2$

First Nonsmooth  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$

Properties of the  
Second Nonsmooth  
Variant  $\widehat{N}_1$

Behavior of BFGS  
on the Second  
Nonsmooth Variant



# First Nonsmooth Variant of Nesterov's Function

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

[An Aside:  
Chebyshev  
Polynomials](#)

[Plots of Chebyshev  
Polynomials](#)

[Nesterov's  
Chebyshev-  
Rosenbrock  
Functions](#)

[Why BFGS Takes So  
Many Iterations to  
Minimize  \$N\_2\$](#)

[First Nonsmooth  
Variant of  
Nesterov's Function](#)

[Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  \$n = 2\$](#)

[Properties of the  
Second Nonsmooth  
Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS  
on the Second  
Nonsmooth Variant](#)



# First Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)  
[Nesterov's Chebyshev-Rosenbrock Functions](#)  
[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .



# First Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

An Aside:  
[Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .

However,  $N_1$  is regular at  $x \in \mathcal{M}_N$  and partly smooth at  $x$  w.r.t.  $\mathcal{M}_N$ .



# First Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .

However,  $N_1$  is regular at  $x \in \mathcal{M}_N$  and partly smooth at  $x$  w.r.t.  $\mathcal{M}_N$ .

We cannot initialize BFGS at  $\hat{x}$ , so starting at normally distributed random points:



# First Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

An Aside:  
Chebyshev Polynomials

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

**First Nonsmooth Variant of Nesterov's Function**

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .

However,  $N_1$  is regular at  $x \in \mathcal{M}_N$  and partly smooth at  $x$  w.r.t.  $\mathcal{M}_N$ .

We cannot initialize BFGS at  $\hat{x}$ , so starting at normally distributed random points:

- $n = 5$ : BFGS reduces  $N_1$  only to about  $10^{-2}$  in 10,000 iterations



# First Nonsmooth Variant of Nesterov's Function

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .

However,  $N_1$  is regular at  $x \in \mathcal{M}_N$  and partly smooth at  $x$  w.r.t.  $\mathcal{M}_N$ .

We cannot initialize BFGS at  $\hat{x}$ , so starting at normally distributed random points:

- $n = 5$ : BFGS reduces  $N_1$  only to about  $10^{-2}$  in 10,000 iterations
- $n = 10$ : BFGS reduces  $N_1$  only to about  $5 \times 10^{-2}$  in 10,000 iterations

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# First Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

**First Nonsmooth Variant of Nesterov's Function**

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$N_1(x) = \frac{1}{4}(x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|$$

$N_1$  is nonsmooth (though locally Lipschitz) as well as nonconvex. The second term is still zero on the manifold  $\mathcal{M}_N$ , but  $N_1$  is not differentiable on  $\mathcal{M}_N$ .

However,  $N_1$  is regular at  $x \in \mathcal{M}_N$  and partly smooth at  $x$  w.r.t.  $\mathcal{M}_N$ .

We cannot initialize BFGS at  $\hat{x}$ , so starting at normally distributed random points:

- $n = 5$ : BFGS reduces  $N_1$  only to about  $10^{-2}$  in 10,000 iterations
- $n = 10$ : BFGS reduces  $N_1$  only to about  $5 \times 10^{-2}$  in 10,000 iterations

The method appears to be converging, very slowly, but may be having numerical difficulties.



## Second Nonsmooth Variant of Nesterov's Function

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)  
[Nesterov's Chebyshev-Rosenbrock Functions](#)  
[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)   
[Properties of the Second Nonsmooth Variant  \$\widehat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

$$\widehat{N}_1(x) = \frac{1}{4}|x_1 - 1| + \sum_{i=1}^{n-1} |x_{i+1} - 2|x_i|| + 1.$$

Again, the unique global minimizer is  $x^*$ . The second term is zero on the set

$$S = \{x : x_{i+1} = 2|x_i| - 1, \quad i = 1, \dots, n-1\}$$

but  $S$  is not a manifold: it has “corners”.

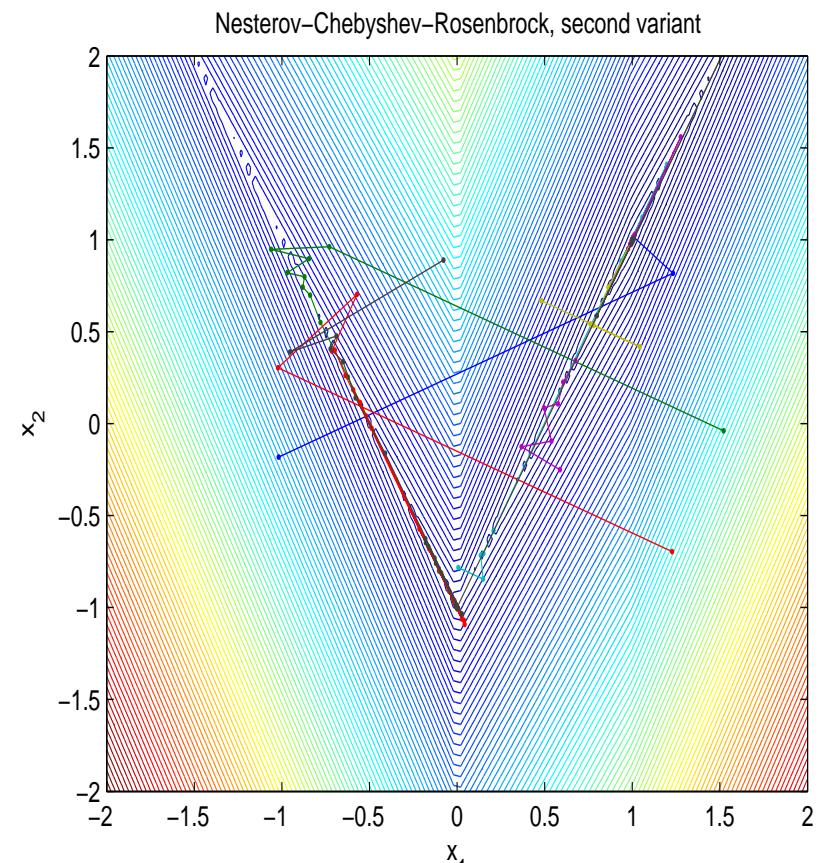
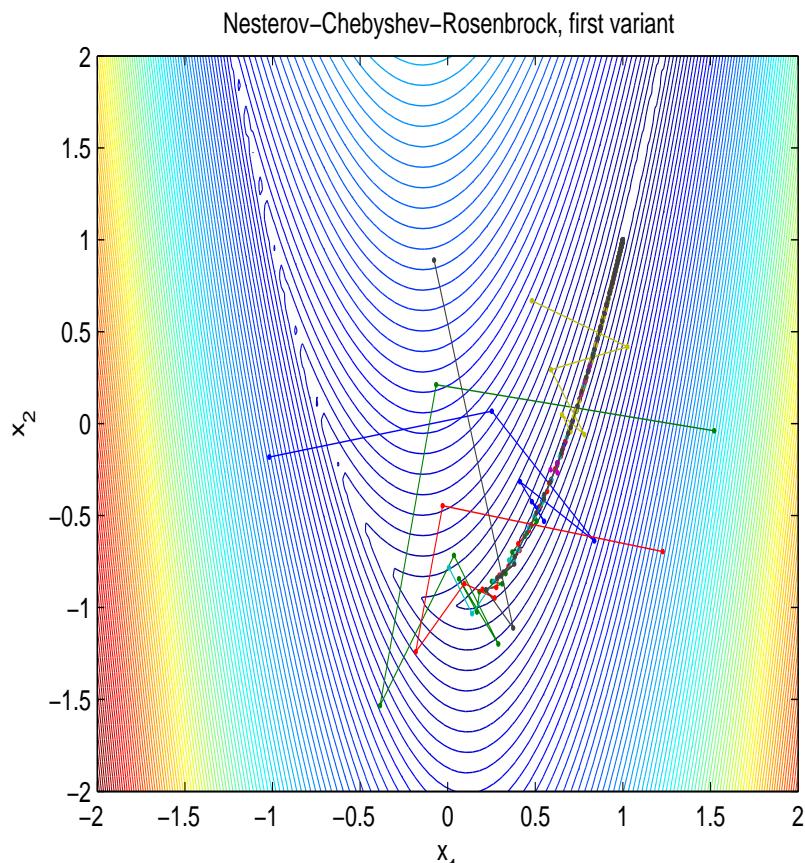


# Contour Plots of the Nonsmooth Variants for $n = 2$

Introduction  
Gradient Sam...  
Quasi-Newto...  
Methods  
A Difficult  
Nonconvex I...  
from Nesterov  
An Aside: Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials  
Nesterov's Chebyshev-Rosenbrock Functions  
Why BFGS  
Many Iterations  
Minimize  $N_1$

First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function  
Contour Plots of the Nonsmooth Variants for  $n = 2$

Properties of the Second Nonsmooth Variant  $\hat{N}_1$   
Behavior of BFGS on the Second Nonsmooth Variant



Contour plots of nonsmooth Chebyshev–Rosenbrock functions  $N_1$  (left) and  $\hat{N}_1$  (right), with  $n = 2$ , with iterates generated by BFGS initialized at 7 different randomly generated points.



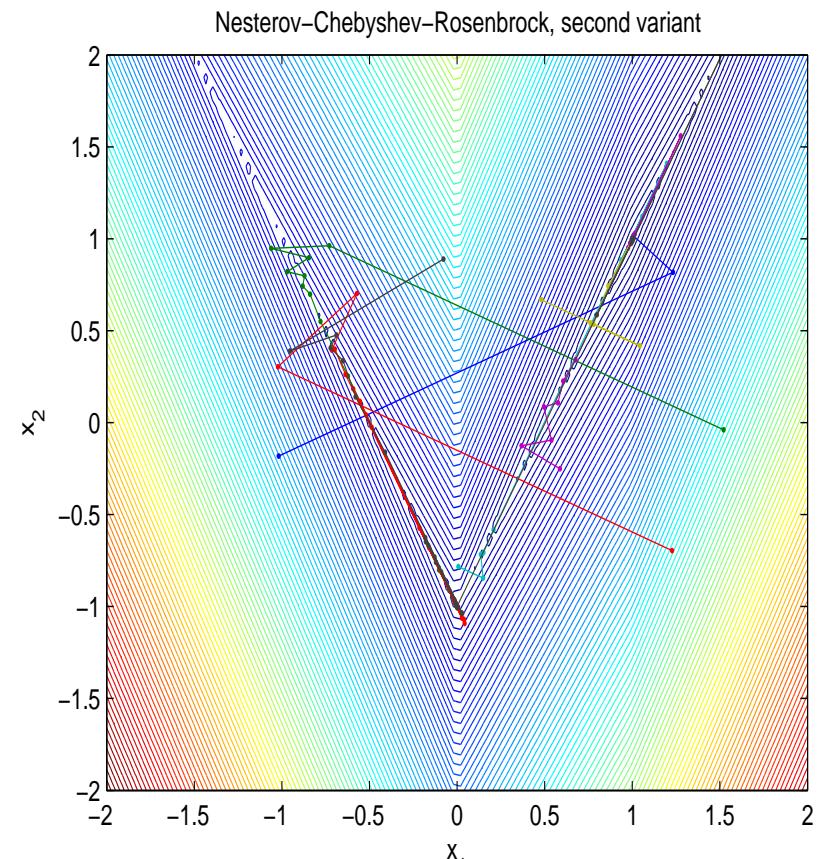
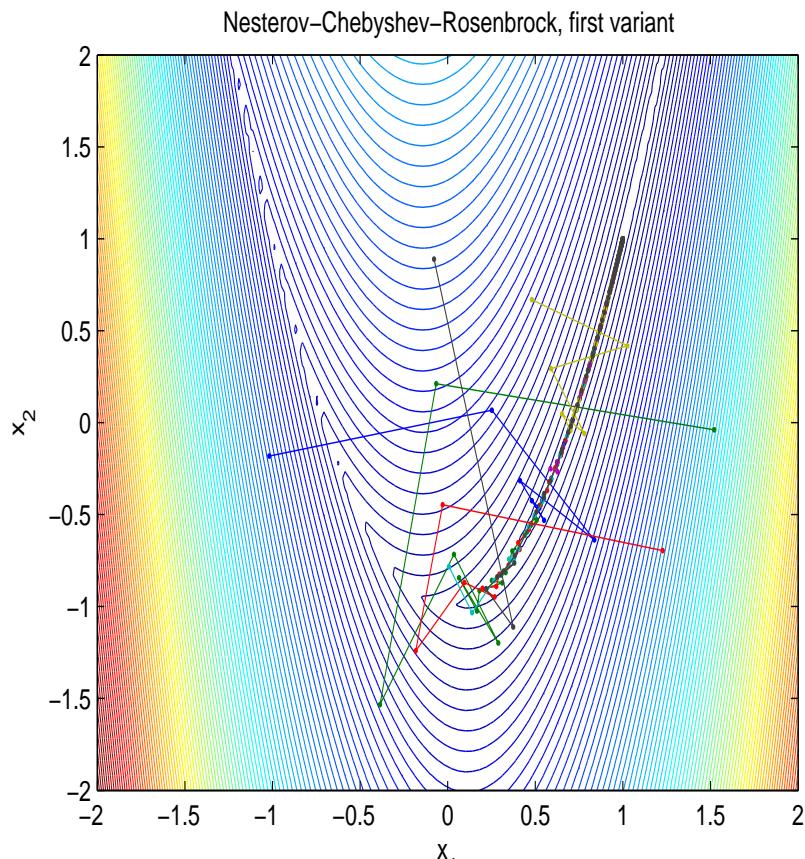
# Contour Plots of the Nonsmooth Variants for $n = 2$

Introduction  
Gradient Sam...  
Quasi-Newto...  
Methods  
A Difficult  
Nonconvex I...  
from Nesterov  
An Aside: Chebyshev  
Polynomials  
Plots of Chebyshev  
Polynomials  
Nesterov's Chebyshev-Rosenbrock Functions  
Why BFGS  
Many Iterations  
Minimize  $N_1$

First Nonsmooth Variant of Nesterov's Function  
Second Nonsmooth Variant of Nesterov's Function

Contour Plots of the Nonsmooth Variants for  $n = 2$

Properties of the Second Nonsmooth Variant  $\hat{N}_1$   
Behavior of BFGS on the Second Nonsmooth Variant



Contour plots of nonsmooth Chebyshev–Rosenbrock functions  $N_1$  (left) and  $\hat{N}_1$  (right), with  $n = 2$ , with iterates generated by BFGS initialized at 7 different randomly generated points. On the left, always get convergence to  $x^* = [1, 1]^T$ . On the right, most runs converge to  $[1, 1]$  but some go to  $x = [0, -1]^T$ .



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

These two properties mean that  $\hat{N}_1$  is *not regular* at  $[0, -1]^T$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

These two properties mean that  $\hat{N}_1$  is *not regular* at  $[0, -1]^T$ .

In fact, for  $n \geq 2$ :

- $\hat{N}_1$  has  $2^{n-1}$  Clarke stationary points

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

**Properties of the Second Nonsmooth Variant  $\hat{N}_1$**

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

These two properties mean that  $\hat{N}_1$  is *not regular* at  $[0, -1]^T$ .

In fact, for  $n \geq 2$ :

- $\hat{N}_1$  has  $2^{n-1}$  Clarke stationary points
- the only local minimizer is the global minimizer  $x^*$

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

**Properties of the Second Nonsmooth Variant  $\hat{N}_1$**

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

These two properties mean that  $\hat{N}_1$  is *not regular* at  $[0, -1]^T$ .

In fact, for  $n \geq 2$ :

- $\hat{N}_1$  has  $2^{n-1}$  Clarke stationary points
- the only local minimizer is the global minimizer  $x^*$
- $x^*$  is the only stationary point in the sense of Mordukhovich  
(i.e., with  $0 \in \partial N_1(x)$  where we defined  $\partial$  in Lecture 12)  
(see also Rockafellar and Wets, *Variational Analysis*, 1998).

(M. Gürbüzbalaban and M.L.O., 2012)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

**Properties of the Second Nonsmooth Variant  $\hat{N}_1$**

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Properties of the Second Nonsmooth Variant $\hat{N}_1$

When  $n = 2$ , the point  $x = [0, -1]^T$  is Clarke stationary for the second nonsmooth variant  $\hat{N}_1$ . We can see this because zero is in the convex hull of the gradient limits for  $\hat{N}_1$  at the point  $x$ .

However,  $x = [0, -1]^T$  is not a local minimizer, because  $d = [1, 2]^T$  is a direction of linear descent:  $\hat{N}'_1(x, d) < 0$ .

These two properties mean that  $\hat{N}_1$  is *not regular* at  $[0, -1]^T$ .

In fact, for  $n \geq 2$ :

- $\hat{N}_1$  has  $2^{n-1}$  Clarke stationary points
- the only local minimizer is the global minimizer  $x^*$
- $x^*$  is the only stationary point in the sense of Mordukhovich  
(i.e., with  $0 \in \partial N_1(x)$  where we defined  $\partial$  in Lecture 12)  
(see also Rockafellar and Wets, *Variational Analysis*, 1998).

(M. Gürbüzbalaban and M.L.O., 2012)

Furthermore, starting from enough randomly generated starting points, BFGS finds all  $2^{n-1}$  Clarke stationary points!

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

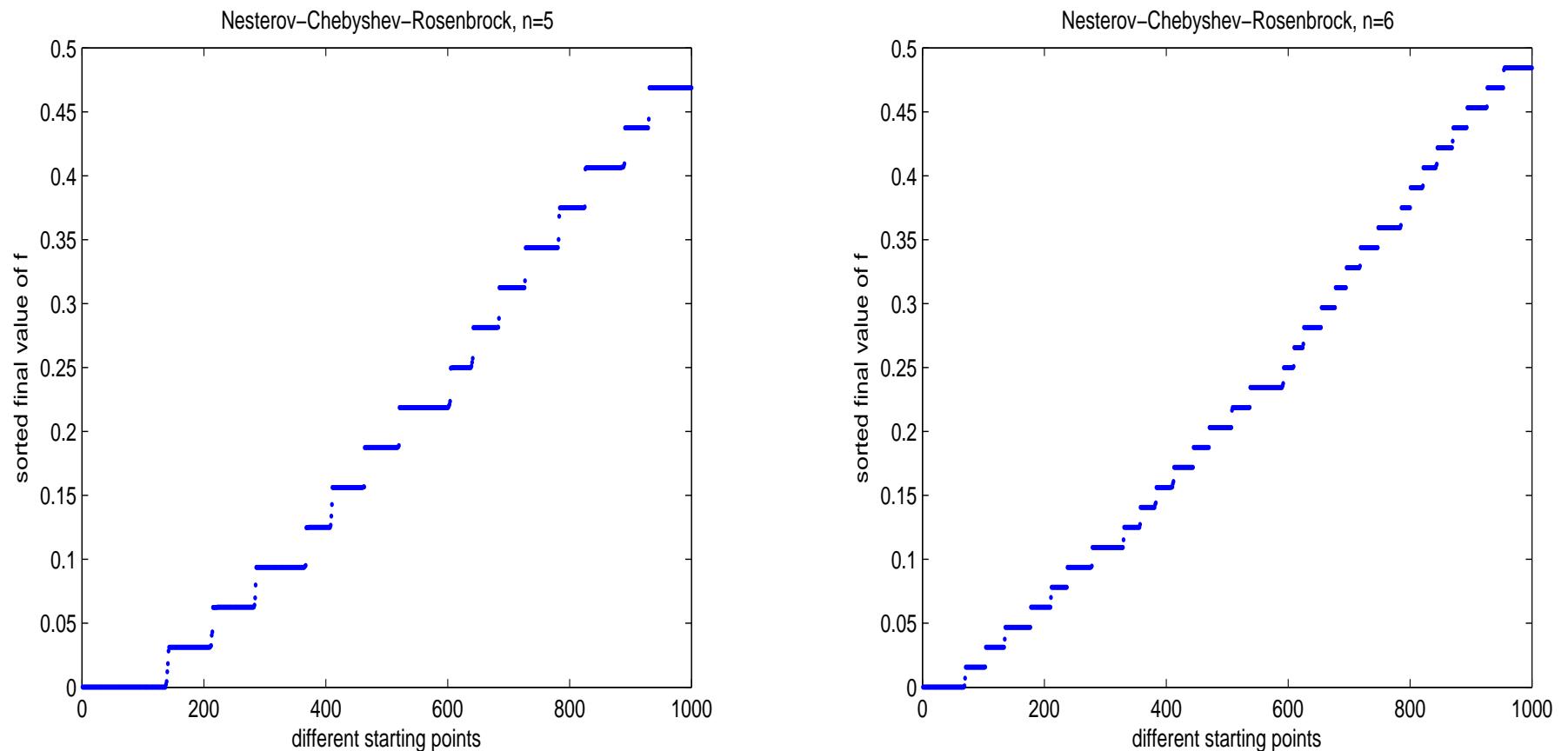
[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

# Behavior of BFGS on the Second Nonsmooth Variant

Introduction  
Gradient S...  
Quasi-Newto...  
Methods  
A Difficult  
Nonconvex I...  
from Nesterov  
An Aside:  
Chebyshev  
Polynomials  
Plots of Che...  
Polynomials  
Nesterov's  
Chebyshev-  
Rosenbrock  
Functions  
Why BFGS  
Many Iterati...  
Minimize  $N_1$   
First Nonsm...  
Variant of  
Nesterov's Function  
Second Nonsmooth  
Variant of  
Nesterov's Function  
Contour Plots of the  
Nonsmooth Variants  
for  $n = 2$   
Properties of the  
Second Nonsmooth  
Variant  $\hat{N}_1$



Left: *sorted* final values of  $\hat{N}_1$  for 1000 randomly generated starting points, when  $n = 5$ : BFGS finds all 16 Clarke stationary points. Right: same with  $n = 6$ : BFGS finds all 32 Clarke stationary points.



# Convergence to Non-Locally-Minimizing Points

When  $f$  is *smooth*, convergence of methods such as BFGS to non-locally-minimizing stationary points or local maxima is *possible* but not likely, because of the line search, and such convergence will not be stable under perturbation.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)

[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)

[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Convergence to Non-Locally-Minimizing Points

When  $f$  is *smooth*, convergence of methods such as BFGS to non-locally-minimizing stationary points or local maxima is *possible* but not likely, because of the line search, and such convergence will not be stable under perturbation.

However, this kind of convergence is what we are seeing for the non-regular, non-smooth Nesterov Chebyshev-Rosenbrock example, and it *is* stable under perturbation. The same behavior occurs for gradient sampling or bundle methods.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)

[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



## Convergence to Non-Locally-Minimizing Points

When  $f$  is *smooth*, convergence of methods such as BFGS to non-locally-minimizing stationary points or local maxima is *possible* but not likely, because of the line search, and such convergence will not be stable under perturbation.

However, this kind of convergence is what we are seeing for the non-regular, non-smooth Nesterov Chebyshev-Rosenbrock example, and it *is* stable under perturbation. The same behavior occurs for gradient sampling or bundle methods.

Kiwiel (private communication): the Nesterov example is the first he had seen which causes his bundle code to have this behavior.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)



# Convergence to Non-Locally-Minimizing Points

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[An Aside: Chebyshev Polynomials](#)  
[Plots of Chebyshev Polynomials](#)

[Nesterov's Chebyshev-Rosenbrock Functions](#)

[Why BFGS Takes So Many Iterations to Minimize  \$N\_2\$](#)

[First Nonsmooth Variant of Nesterov's Function](#)  
[Second Nonsmooth Variant of Nesterov's Function](#)  
[Contour Plots of the Nonsmooth Variants for  \$n = 2\$](#)

[Properties of the Second Nonsmooth Variant  \$\hat{N}\_1\$](#)

[Behavior of BFGS on the Second Nonsmooth Variant](#)

When  $f$  is *smooth*, convergence of methods such as BFGS to non-locally-minimizing stationary points or local maxima is *possible* but not likely, because of the line search, and such convergence will not be stable under perturbation.

However, this kind of convergence is what we are seeing for the non-regular, non-smooth Nesterov Chebyshev-Rosenbrock example, and it *is* stable under perturbation. The same behavior occurs for gradient sampling or bundle methods.

Kiwiel (private communication): the Nesterov example is the first he had seen which causes his bundle code to have this behavior.

Nonetheless, we don't know whether, in exact arithmetic, the methods would actually generate sequences converging to the nonminimizing Clarke stationary points. Experiments by Kaku (2011) suggest that the higher the precision used, the more likely BFGS is to eventually move away from such a point.



[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton  
Methods](#)

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

[Limited Memory  
Methods](#)

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.

Gradient Descent  
Convergence of the  
L-BFGS-1 Search

# Limited Memory Methods



# Limited Memory BFGS

“Full” BFGS requires storing an  $n \times n$  matrix and doing matrix-vector multiplies, which is not possible when  $n$  is large.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

**Limited Memory BFGS**

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)

$n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth, Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function,](#)

[Unbounded Below L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)



# Limited Memory BFGS

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

**Limited Memory BFGS**

[Limited Memory BFGS on the Eigenvalue Product](#)  
[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

[Nonsmooth, Nonconvex:  
 \$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

[A Nonsmooth Convex Function,](#)

[Unbounded Below](#)

[L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)

“Full” BFGS requires storing an  $n \times n$  matrix and doing matrix-vector multiplies, which is not possible when  $n$  is large.

In the 1980s, J. Nocedal and others developed a “limited memory” version of BFGS, with  $O(n)$  space and time requirements, which is very widely used for minimizing smooth functions in many variables. At the  $k$ th iteration, it applies only the most recent  $m$  rank-two updates, defined by

$$(s_j, y_j), \quad j = k - m, \dots, k - 1$$

to an initial inverse Hessian approximation  $H_0^{(k)}$ .



## Limited Memory BFGS

“Full” BFGS requires storing an  $n \times n$  matrix and doing matrix-vector multiplies, which is not possible when  $n$  is large.

In the 1980s, J. Nocedal and others developed a “limited memory” version of BFGS, with  $O(n)$  space and time requirements, which is very widely used for minimizing smooth functions in many variables. At the  $k$ th iteration, it applies only the most recent  $m$  rank-two updates, defined by

$$(s_j, y_j), \quad j = k - m, \dots, k - 1$$

to an initial inverse Hessian approximation  $H_0^{(k)}$ .

There are two variants: with “scaling” ( $H_0^{(k)} = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}} I$ ) and without scaling ( $H_0^{(k)} = I$ ).

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

**Limited Memory BFGS**

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

[Nonsmooth, Nonconvex:  
 \$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

[A Nonsmooth Convex Function, Unbounded Below  
L-BFGS-1 vs. Gradient Descent  
Convergence of the L-BFGS-1 Search](#)



## Limited Memory BFGS

“Full” BFGS requires storing an  $n \times n$  matrix and doing matrix-vector multiplies, which is not possible when  $n$  is large.

In the 1980s, J. Nocedal and others developed a “limited memory” version of BFGS, with  $O(n)$  space and time requirements, which is very widely used for minimizing smooth functions in many variables. At the  $k$ th iteration, it applies only the most recent  $m$  rank-two updates, defined by

$$(s_j, y_j), \quad j = k - m, \dots, k - 1$$

to an initial inverse Hessian approximation  $H_0^{(k)}$ .

There are two variants: with “scaling” ( $H_0^{(k)} = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}} I$ ) and without scaling ( $H_0^{(k)} = I$ ).

The convergence rate of limited memory BFGS is linear, not superlinear, on smooth problems.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product  
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.

Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## Limited Memory BFGS

“Full” BFGS requires storing an  $n \times n$  matrix and doing matrix-vector multiplies, which is not possible when  $n$  is large.

In the 1980s, J. Nocedal and others developed a “limited memory” version of BFGS, with  $O(n)$  space and time requirements, which is very widely used for minimizing smooth functions in many variables. At the  $k$ th iteration, it applies only the most recent  $m$  rank-two updates, defined by

$$(s_j, y_j), \quad j = k - m, \dots, k - 1$$

to an initial inverse Hessian approximation  $H_0^{(k)}$ .

There are two variants: with “scaling” ( $H_0^{(k)} = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}} I$ ) and without scaling ( $H_0^{(k)} = I$ ).

The convergence rate of limited memory BFGS is linear, not superlinear, on smooth problems.

Question: how effective is it on nonsmooth problems?



# Limited Memory BFGS on the Eigenvalue Product

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

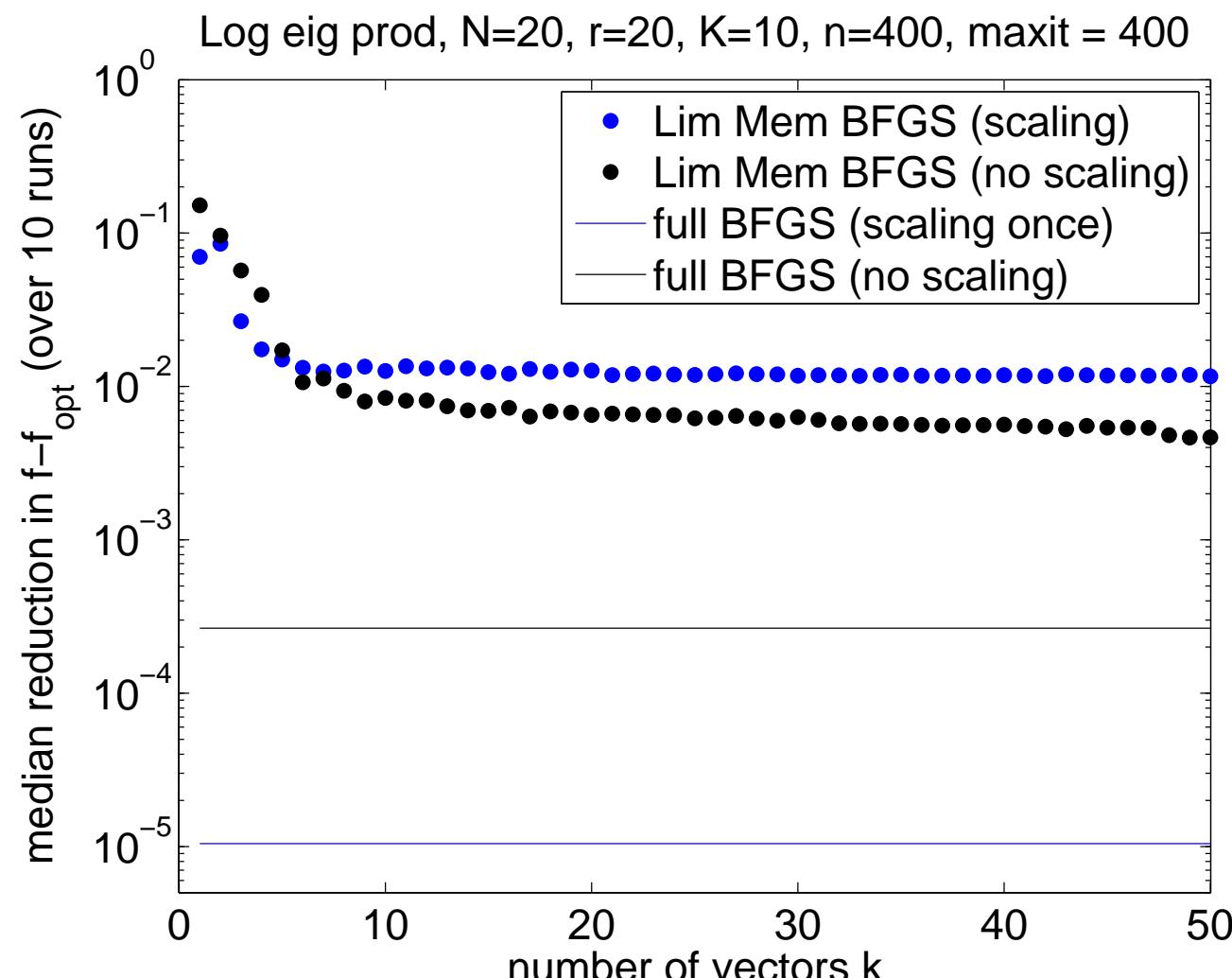
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search





# Limited Memory BFGS on the Eigenvalue Product

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

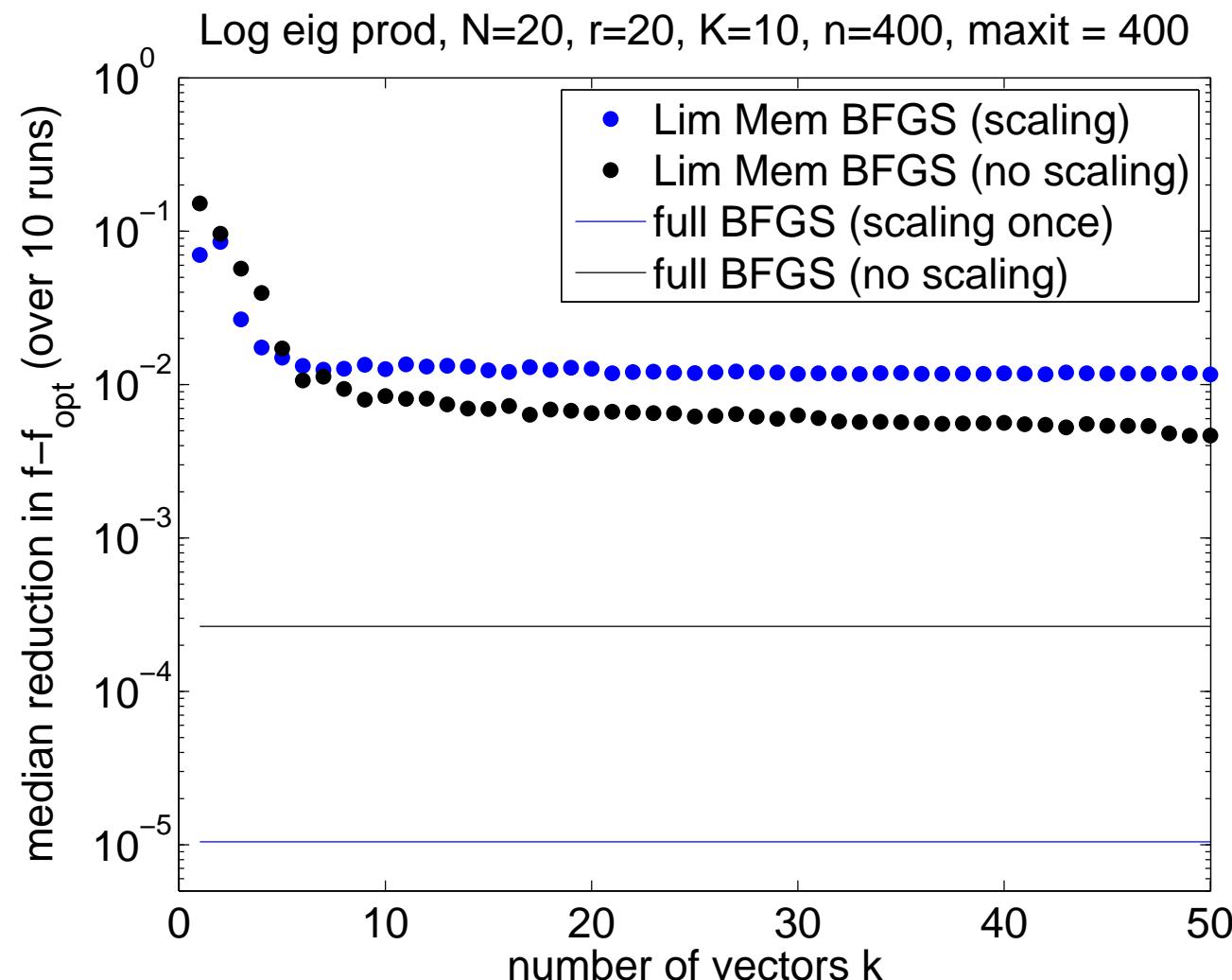
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



Limited Memory is not nearly as good as full BFGS



# Limited Memory BFGS on the Eigenvalue Product

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

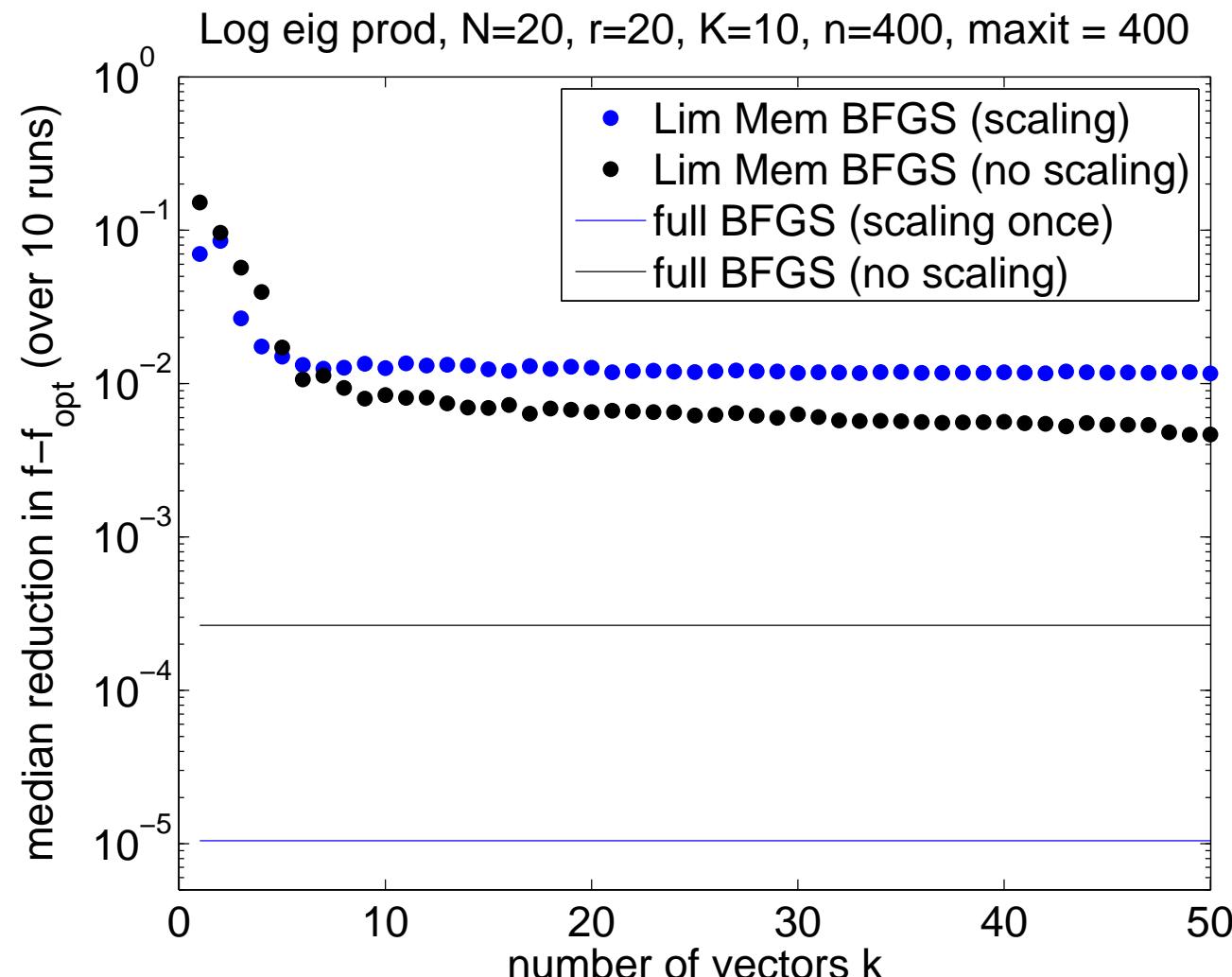
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



Limited Memory is not nearly as good as full BFGS

No significant improvement when  $k$  reaches 44



## A More Basic Example

Let  $x = [y; z; w] \in \mathbb{R}^{n_A+n_B+n_R}$  and consider the test function

$$f(x) = (y - e)^T A(y - e) + \{(z - e)^T B(z - e)\}^{1/2} + R_1(w)$$

where  $A = A^T \succ 0$ ,  $B = B^T \succ 0$ ,  $e = [1; 1; \dots; 1]$ .

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

Smooth, Convex:

$n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:

$n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth, Nonconvex:

$n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth Convex Function,

Unbounded Below  
L-BFGS-1 vs.

Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## A More Basic Example

Let  $x = [y; z; w] \in \mathbb{R}^{n_A+n_B+n_R}$  and consider the test function

$$f(x) = (y - e)^T A(y - e) + \{(z - e)^T B(z - e)\}^{1/2} + R_1(w)$$

where  $A = A^T \succ 0$ ,  $B = B^T \succ 0$ ,  $e = [1; 1; \dots; 1]$ .

The first term is quadratic, the second is nonsmooth but convex, and the third is a nonsmooth, nonconvex Rosenbrock function.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

Smooth, Convex:

$n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:

$n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth, Nonconvex:

$n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth Convex Function,

Unbounded Below L-BFGS-1 vs.

Gradient Descent Convergence of the L-BFGS-1 Search



## A More Basic Example

Let  $x = [y; z; w] \in \mathbb{R}^{n_A+n_B+n_R}$  and consider the test function

$$f(x) = (y - e)^T A(y - e) + \{(z - e)^T B(z - e)\}^{1/2} + R_1(w)$$

where  $A = A^T \succ 0$ ,  $B = B^T \succ 0$ ,  $e = [1; 1; \dots; 1]$ .

The first term is quadratic, the second is nonsmooth but convex, and the third is a nonsmooth, nonconvex Rosenbrock function.

The optimal value is 0, with  $x = e$ . The function  $f$  is partly smooth and the dimension of the V-space is  $n_B + n_R - 1$ .

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## A More Basic Example

Let  $x = [y; z; w] \in \mathbb{R}^{n_A+n_B+n_R}$  and consider the test function

$$f(x) = (y - e)^T A(y - e) + \{(z - e)^T B(z - e)\}^{1/2} + R_1(w)$$

where  $A = A^T \succ 0$ ,  $B = B^T \succ 0$ ,  $e = [1; 1; \dots; 1]$ .

The first term is quadratic, the second is nonsmooth but convex, and the third is a nonsmooth, nonconvex Rosenbrock function.

The optimal value is 0, with  $x = e$ . The function  $f$  is partly smooth and the dimension of the V-space is  $n_B + n_R - 1$ .

Set  $A = XX^T$  where  $x_{ij}$  are normally distributed, with condition number about  $10^6$  when  $n_A = 200$ . Similarly  $B$  with  $n_B < n_A$ .

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## A More Basic Example

Let  $x = [y; z; w] \in \mathbb{R}^{n_A+n_B+n_R}$  and consider the test function

$$f(x) = (y - e)^T A(y - e) + \{(z - e)^T B(z - e)\}^{1/2} + R_1(w)$$

where  $A = A^T \succ 0$ ,  $B = B^T \succ 0$ ,  $e = [1; 1; \dots; 1]$ .

The first term is quadratic, the second is nonsmooth but convex, and the third is a nonsmooth, nonconvex Rosenbrock function.

The optimal value is 0, with  $x = e$ . The function  $f$  is partly smooth and the dimension of the V-space is  $n_B + n_R - 1$ .

Set  $A = XX^T$  where  $x_{ij}$  are normally distributed, with condition number about  $10^6$  when  $n_A = 200$ . Similarly  $B$  with  $n_B < n_A$ .

Besides limited memory BFGS and full BFGS, we also compare limited memory Gradient Sampling, where we sample  $k \ll n$  gradients per iteration.

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## Smooth, Convex: $n_A = 200, n_B = 0, n_R = 1$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

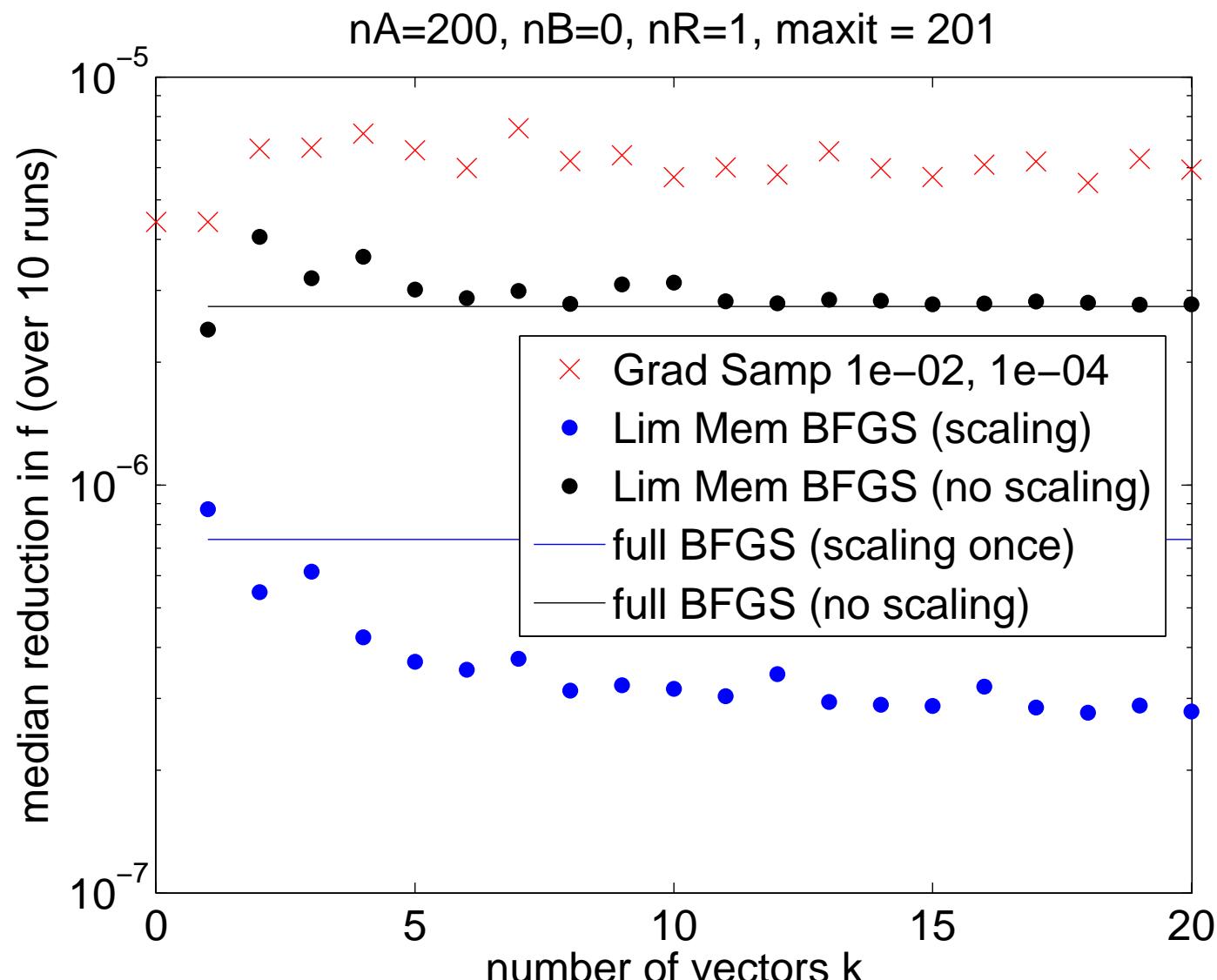
Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,

Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search





## Smooth, Convex: $n_A = 200, n_B = 0, n_R = 1$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

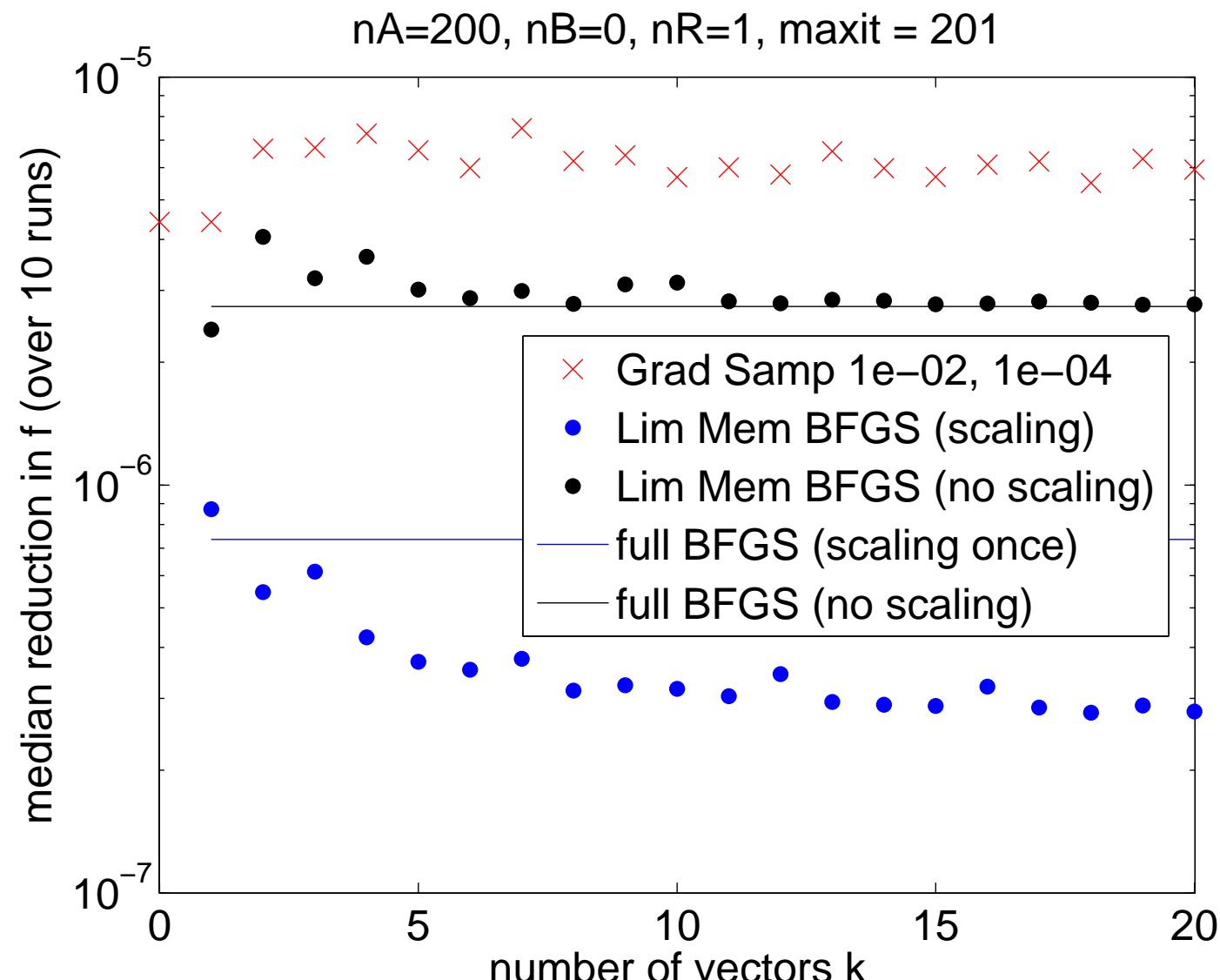
Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,

Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



LM-BFGS with scaling even better than full BFGS



## Nonsmooth, Convex: $n_A = 200, n_B = 10, n_R = 1$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

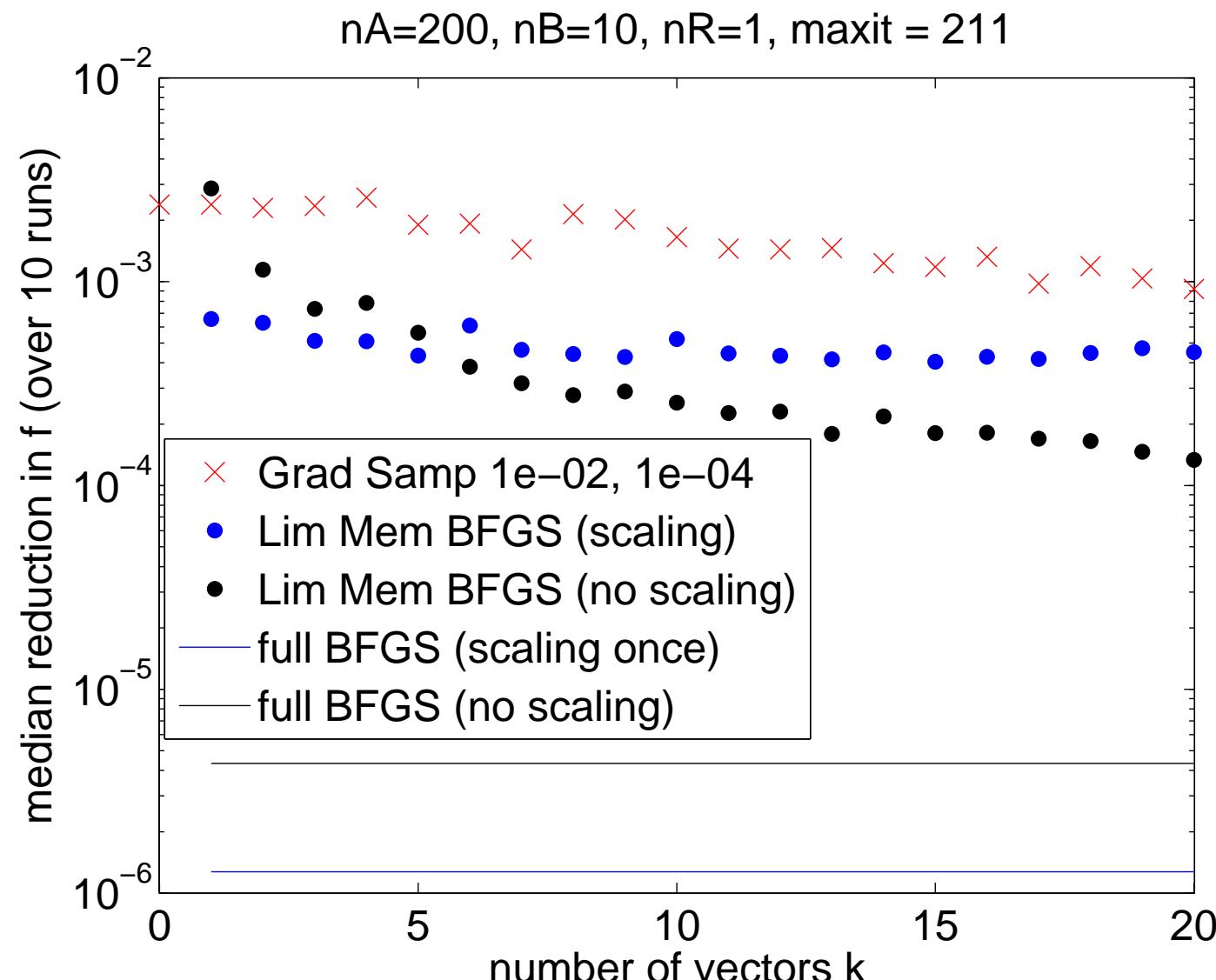
Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,

Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search





## Nonsmooth, Convex: $n_A = 200, n_B = 10, n_R = 1$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

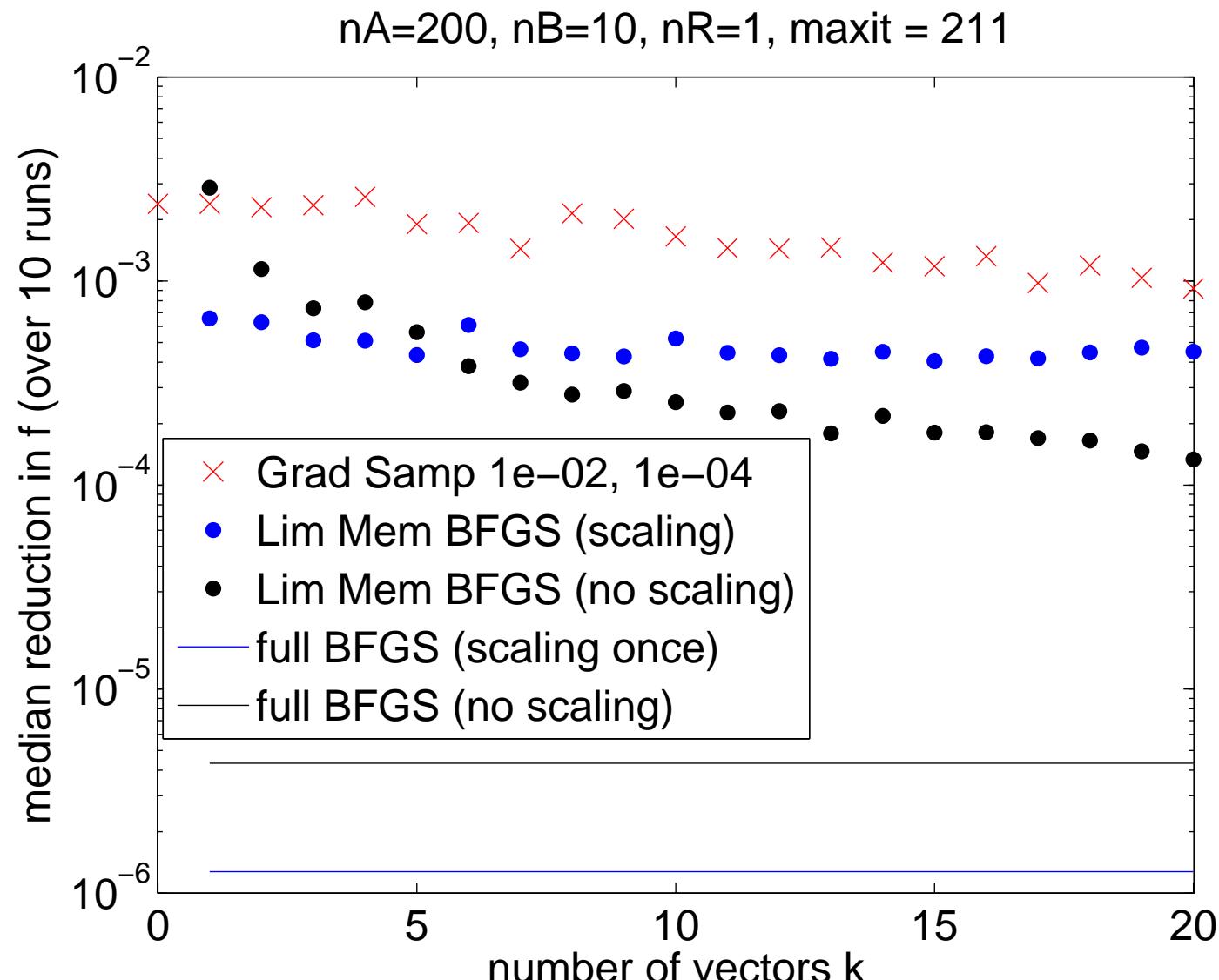
Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,

Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



LM-BFGS much worse than full BFGS



## Nonsmooth, Nonconvex: $n_A = 200, n_B = 10, n_R = 5$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

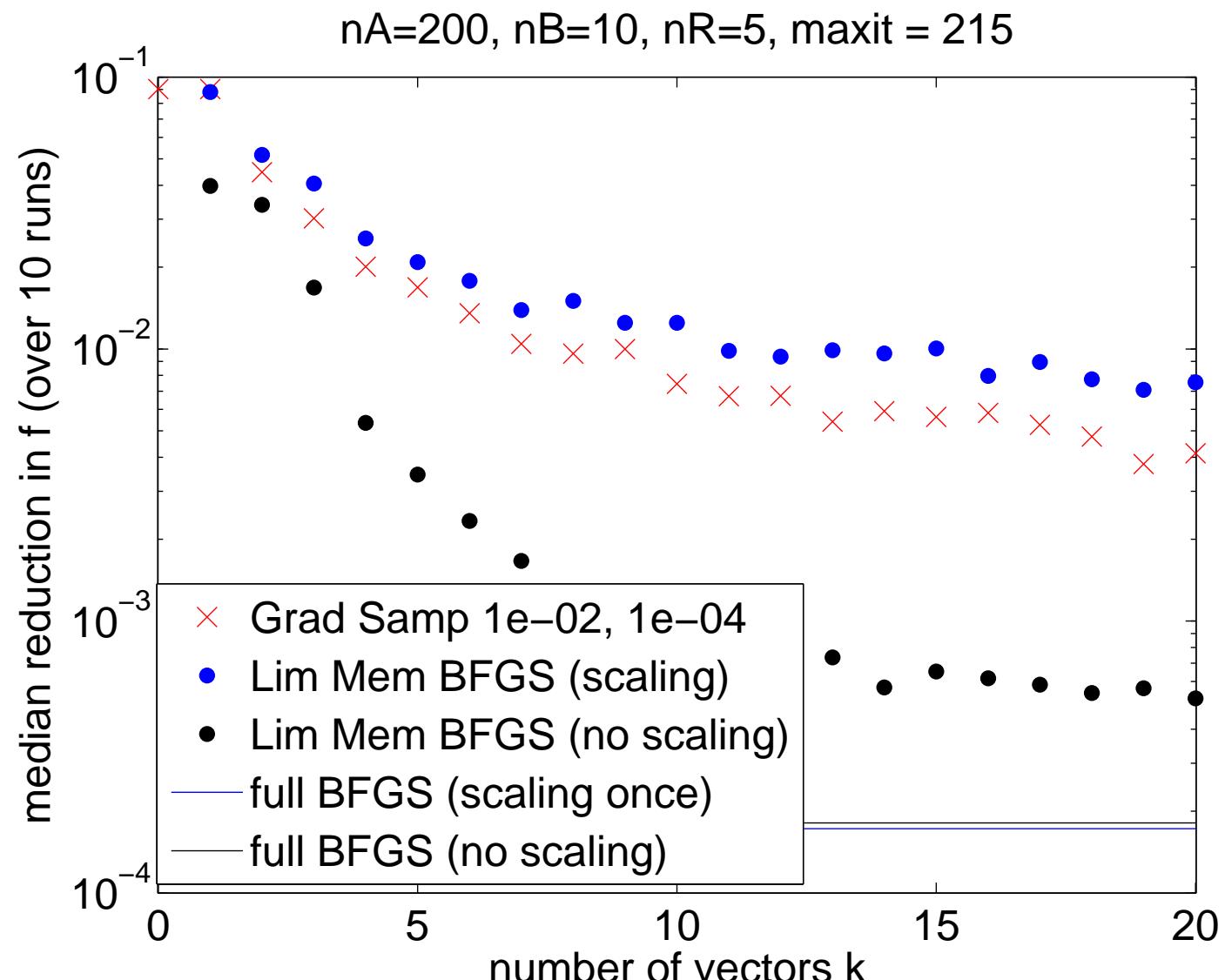
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search





## Nonsmooth, Nonconvex: $n_A = 200, n_B = 10, n_R = 5$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

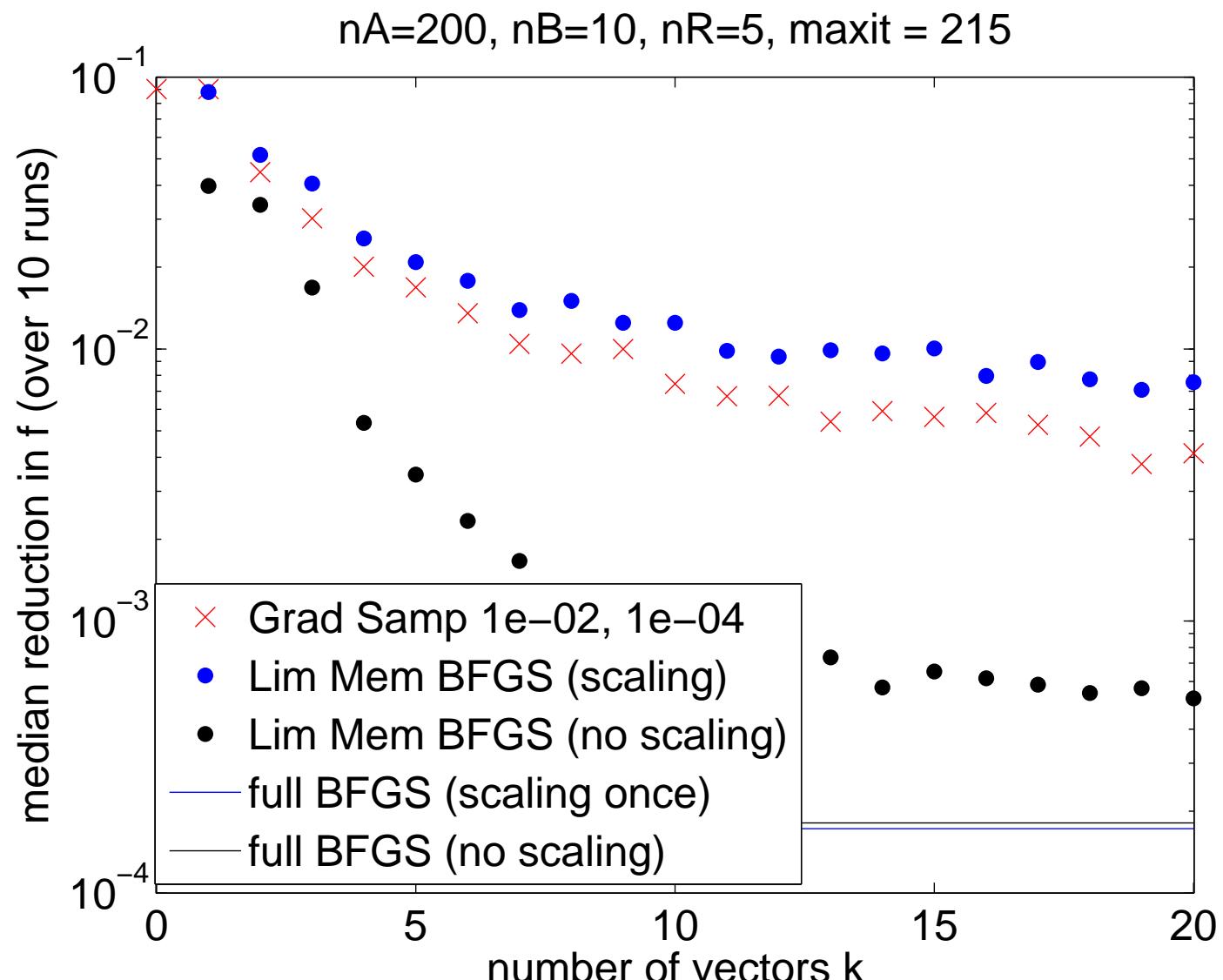
A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



LM-BFGS with scaling even worse than LM-Grad-Samp 60 / 71



# A Nonsmooth Convex Function, Unbounded Below

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

[Nonsmooth, Nonconvex:](#)

[\$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

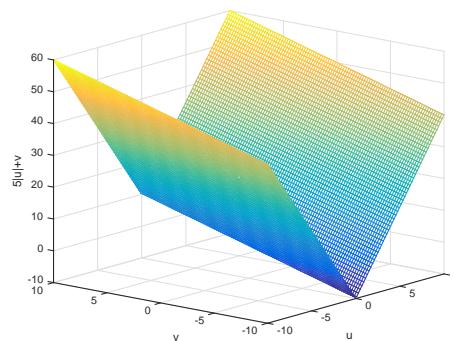
[A Nonsmooth Convex Function, Unbounded Below](#)

[L-BFGS-1 vs. Gradient Descent Convergence of the L-BFGS-1 Search](#)

Let's reconsider

$$f(x) = a|x_1| + x_2$$

with  $a \geq 1$ .



Turns out that L-BFGS-1 (saving just one update) with scaling fails for *smaller* values of  $a$  than the critical value beyond which Gradient Descent fails!



# L-BFGS-1 vs. Gradient Descent

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

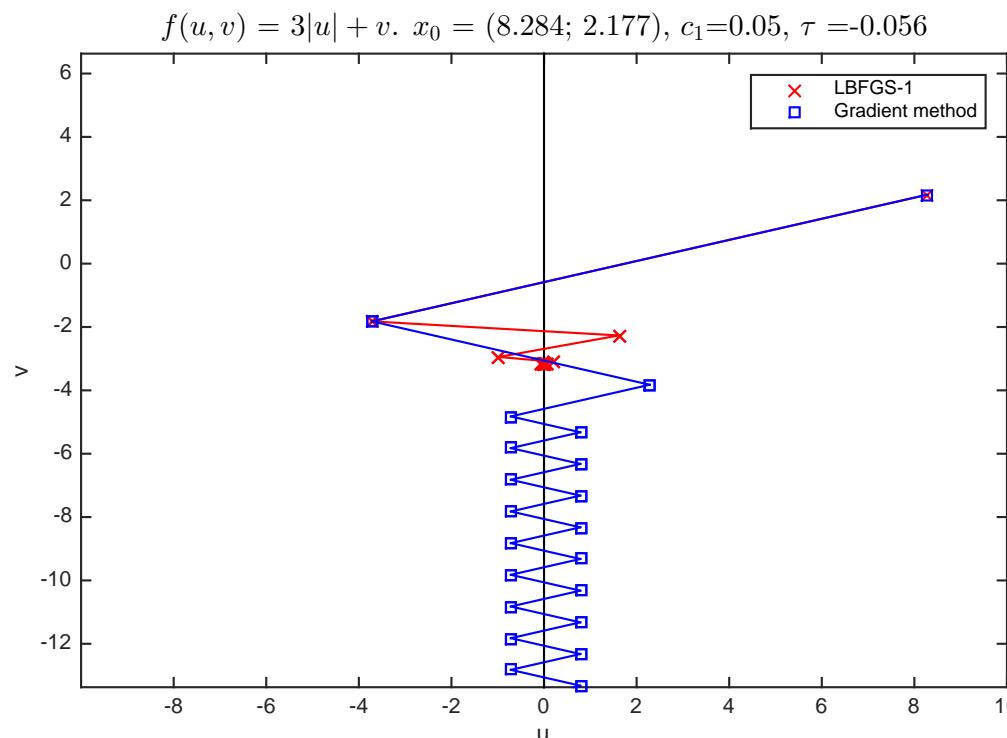
[Nonsmooth, Nonconvex:  
 \$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

[A Nonsmooth Convex Function, Unbounded Below](#)

[L-BFGS-1 vs. Gradient Descent Convergence of the L-BFGS-1 Search](#)

Red: path of L-BFGS-1 with scaling, converges to non-stationary point.

Blue: path of the gradient method with same Armijo-Wolfe line search, generates  $f(x) \downarrow -\infty$ .





# Convergence of the L-BFGS-1 Search Direction

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search

**Theorem.** Let  $d^{(k)}$  be the search direction generated by L-BFGS-1 with scaling applied to  $f(x) = a|x_1| + \sum_{i=2}^n x_i$  using an Armijo-Wolfe line search. If  $\sqrt{4(n-1)} \leq a$ , then  $\frac{|d^{(k)}|}{\|d^{(k)}\|}$  converges to some constant direction  $d$ . Furthermore, if

$$a(a + \sqrt{a^2 - 3(n-1)}) > \left(\frac{1}{c_1} - 1\right)(n-1),$$

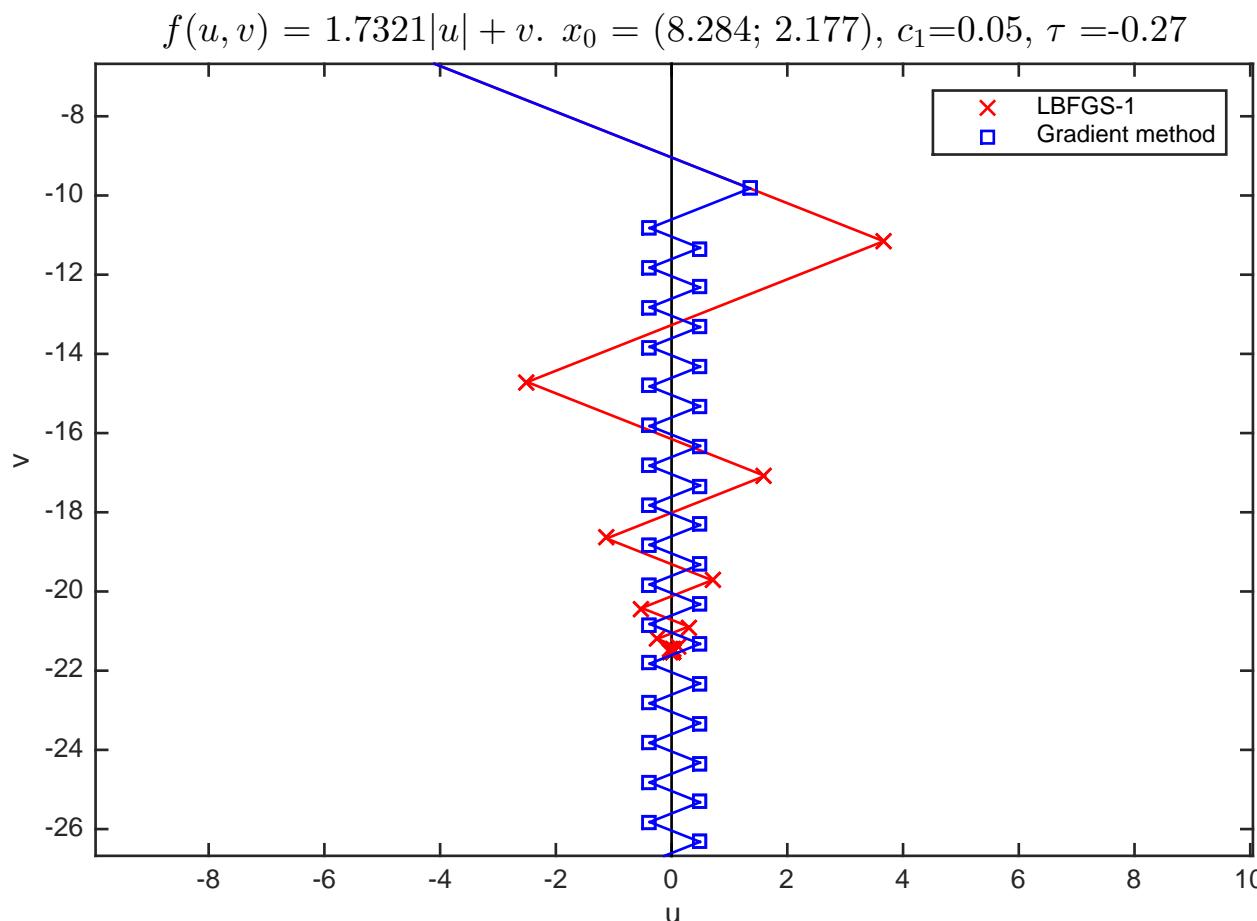
where  $c_1$  is the Armijo parameter, then the iterates  $x^{(k)}$  converge to a non-stationary point.

Azam Asl, 2018.



## Experiment, with $n = 2$ and $a = \sqrt{3}$

In practice we observe that  $\sqrt{3(n - 1)} \leq a$  suffices for the method to fail, which is a weaker condition than the previous one. Below with  $n = 2$  and  $a = \sqrt{3}$  the method fails:



Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200$ ,  $n_B = 0$ ,  $n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200$ ,  $n_B = 10$ ,  $n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search



## Experiment: slightly smaller $a$

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

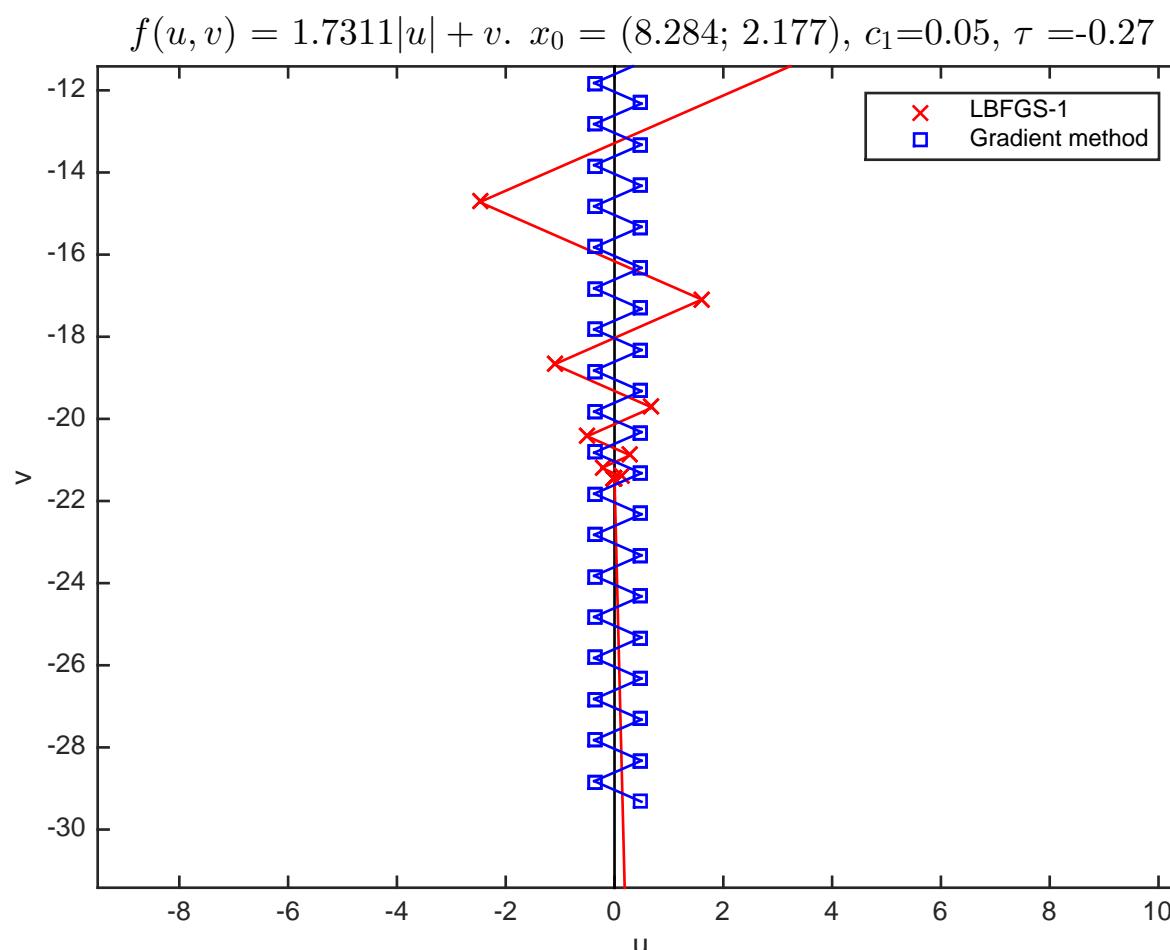
Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search

But if we set  $a = \sqrt{3} - 0.001$ , it succeeds “at the last minute”.





# Experiments: Top, scaling on; Bottom, scaling off

Introduction

Gradient Sampling

Quasi-Newton  
Methods

A Difficult  
Nonconvex Problem  
from Nesterov

Limited Memory  
Methods

Limited Memory  
BFGS

Limited Memory  
BFGS on the  
Eigenvalue Product

A More Basic  
Example

Smooth, Convex:  
 $n_A = 200, n_B = 0, n_R = 1$

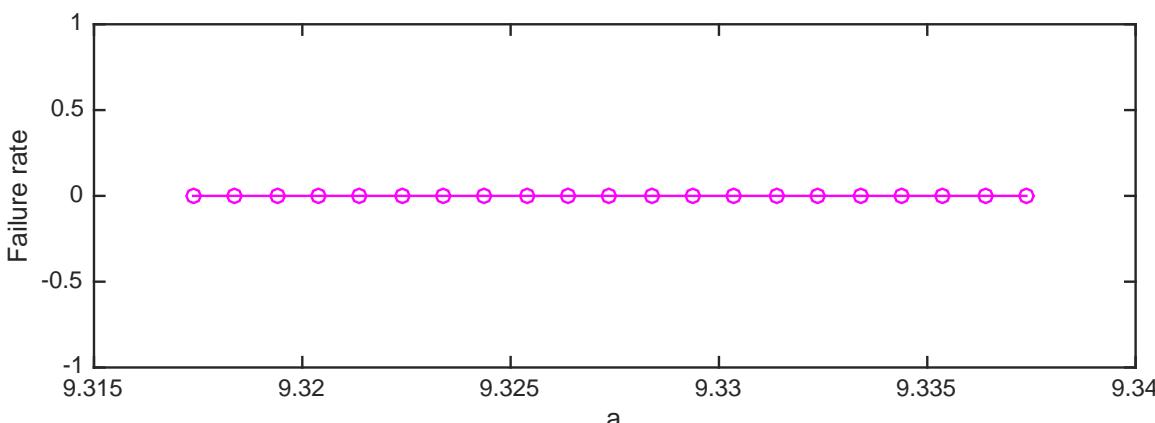
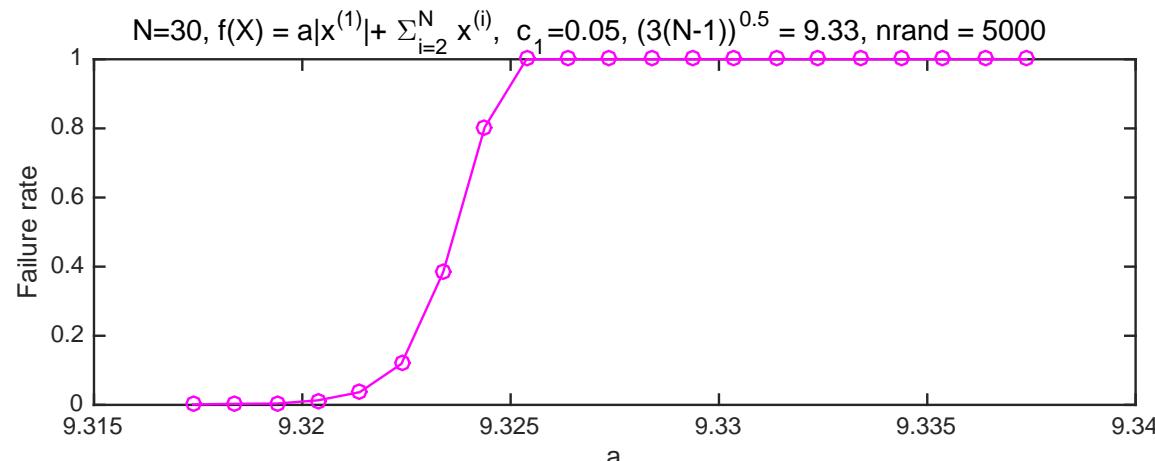
Nonsmooth, Convex:  
 $n_A = 200, n_B = 10, n_R = 1$

Nonsmooth,  
Nonconvex:  
 $n_A = 200, n_B = 10, n_R = 5$

A Nonsmooth  
Convex Function,  
Unbounded Below

L-BFGS-1 vs.  
Gradient Descent  
Convergence of the  
L-BFGS-1 Search

$$n = 30, \sqrt{3(n - 1)} = 9.327$$





# Limited Effectiveness of Limited Memory BFGS

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)

$n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth,](#)

[Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function,](#)

[Unbounded Below L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)

We have observed that that addition of nonsmoothness to a problem, convex or nonconvex, creates great difficulties for Limited Memory BFGS, with and without scaling, even when the dimension of the  $V$ -space is less than the size of the memory.



# Limited Effectiveness of Limited Memory BFGS

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)  
 $n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth, Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function, Unbounded Below](#)

[L-BFGS-1 vs. Gradient Descent Convergence of the L-BFGS-1 Search](#)

We have observed that that addition of nonsmoothness to a problem, convex or nonconvex, creates great difficulties for Limited Memory BFGS, with and without scaling, even when the dimension of the  $V$ -space is less than the size of the memory.

Azam Asl's result establishes failure of L-BFGS-1 for a specific  $f$  when scaling is on; no such result is proved yet when scaling is off.



# Limited Effectiveness of Limited Memory BFGS

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

[Nonsmooth, Nonconvex:  
 \$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

[A Nonsmooth Convex Function, Unbounded Below](#)

[L-BFGS-1 vs. Gradient Descent Convergence of the L-BFGS-1 Search](#)

We have observed that that addition of nonsmoothness to a problem, convex or nonconvex, creates great difficulties for Limited Memory BFGS, with and without scaling, even when the dimension of the  $V$ -space is less than the size of the memory.

Azam Asl's result establishes failure of L-BFGS-1 for a specific  $f$  when scaling is on; no such result is proved yet when scaling is off.

We have also investigated Limited Memory Gradient Sampling which does not work well either.



# Other Ideas for Large Scale Nonsmooth Optimization

[Introduction](#)

---

[Gradient Sampling](#)

---

[Quasi-Newton  
Methods](#)

---

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

---

[Limited Memory  
Methods](#)

---

[Limited Memory  
BFGS](#)

[Limited Memory  
BFGS on the  
Eigenvalue Product](#)

[A More Basic  
Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)

$n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth,  
Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth  
Convex Function,](#)

[Unbounded Below  
L-BFGS-1 vs.](#)

[Gradient Descent  
Convergence of the  
L-BFGS-1 Search](#)



# Other Ideas for Large Scale Nonsmooth Optimization

- Exploit structure! Lots of work on this has been done, e.g. using proximal point methods or ADMM (Alternating Direction Method of Multipliers)

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)  
 $n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth, Nonconvex:](#)  
 $n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function, Unbounded Below](#)  
L-BFGS-1 vs.  
Gradient Descent  
Convergence of the L-BFGS-1 Search



# Other Ideas for Large Scale Nonsmooth Optimization

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)

$n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth, Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function,](#)

[Unbounded Below L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)

- Exploit structure! Lots of work on this has been done, e.g. using proximal point methods or ADMM (Alternating Direction Method of Multipliers)
- Smoothing! Lots of work on this has been done too, most notably by Yu. Nesterov



# Other Ideas for Large Scale Nonsmooth Optimization

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:](#)

$n_A = 200, n_B = 0, n_R = 1$

[Nonsmooth, Convex:](#)

$n_A = 200, n_B = 10, n_R = 1$

[Nonsmooth, Nonconvex:](#)

$n_A = 200, n_B = 10, n_R = 5$

[A Nonsmooth Convex Function,](#)

[Unbounded Below L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)

- Exploit structure! Lots of work on this has been done, e.g. using proximal point methods or ADMM (Alternating Direction Method of Multipliers)
- Smoothing! Lots of work on this has been done too, most notably by Yu. Nesterov
- Automatic Differentiation (AD): (A. Griewank et. al.)



# Other Ideas for Large Scale Nonsmooth Optimization

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Limited Memory BFGS](#)

[Limited Memory BFGS on the Eigenvalue Product](#)

[A More Basic Example](#)

[Smooth, Convex:  
 \$n\_A = 200, n\_B = 0, n\_R = 1\$](#)

[Nonsmooth, Convex:  
 \$n\_A = 200, n\_B = 10, n\_R = 1\$](#)

[Nonsmooth, Nonconvex:  
 \$n\_A = 200, n\_B = 10, n\_R = 5\$](#)

[A Nonsmooth Convex Function,](#)

[Unbounded Below](#)

[L-BFGS-1 vs.](#)

[Gradient Descent Convergence of the L-BFGS-1 Search](#)

- Exploit structure! Lots of work on this has been done, e.g. using proximal point methods or ADMM (Alternating Direction Method of Multipliers)
- Smoothing! Lots of work on this has been done too, most notably by Yu. Nesterov
- Automatic Differentiation (AD): (A. Griewank et. al.)
- Stochastic Subgradient Method (D. Davis and D. Drusvyatskiy, 2018)



[Introduction](#)

---

[Gradient Sampling](#)

---

[Quasi-Newton  
Methods](#)

---

[A Difficult  
Nonconvex Problem  
from Nesterov](#)

---

[Limited Memory  
Methods](#)

---

[Concluding Remarks](#)

[Summary](#)

[Our Papers](#)

## Concluding Remarks



## Summary

Gradient descent frequently fails on nonsmooth problems.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)

**Summary**

[Our Papers](#)



## Summary

Gradient descent frequently fails on nonsmooth problems.

Gradient Sampling is a simple method for nonsmooth, nonconvex optimization for which a convergence theory is known, but it is too expensive to use in most applications.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)

**Summary**

[Our Papers](#)



## Summary

Gradient descent frequently fails on nonsmooth problems.

Gradient Sampling is a simple method for nonsmooth, nonconvex optimization for which a convergence theory is known, but it is too expensive to use in most applications.

BFGS — the full version — is remarkably effective on nonsmooth problems, but little theory is known.

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)

**Summary**

[Our Papers](#)



## Summary

[Introduction](#)

---

[Gradient Sampling](#)

---

[Quasi-Newton Methods](#)

---

[A Difficult Nonconvex Problem from Nesterov](#)

---

[Limited Memory Methods](#)

---

[Concluding Remarks](#)

**Summary**

[Our Papers](#)

Gradient descent frequently fails on nonsmooth problems.

Gradient Sampling is a simple method for nonsmooth, nonconvex optimization for which a convergence theory is known, but it is too expensive to use in most applications.

BFGS — the full version — is remarkably effective on nonsmooth problems, but little theory is known.

Limited Memory BFGS is not so effective on nonsmooth problems.



## Summary

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)

**Summary**

[Our Papers](#)

Gradient descent frequently fails on nonsmooth problems.

Gradient Sampling is a simple method for nonsmooth, nonconvex optimization for which a convergence theory is known, but it is too expensive to use in most applications.

BFGS — the full version — is remarkably effective on nonsmooth problems, but little theory is known.

Limited Memory BFGS is not so effective on nonsmooth problems.

Diabolical nonconvex problems such as Nesterov's Chebyshev-Rosenbrock problems can be very difficult, especially in the nonsmooth case.



## Summary

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)

**Summary**

[Our Papers](#)

Gradient descent frequently fails on nonsmooth problems.

Gradient Sampling is a simple method for nonsmooth, nonconvex optimization for which a convergence theory is known, but it is too expensive to use in most applications.

BFGS — the full version — is remarkably effective on nonsmooth problems, but little theory is known.

Limited Memory BFGS is not so effective on nonsmooth problems.

Diabolical nonconvex problems such as Nesterov's Chebyshev-Rosenbrock problems can be very difficult, especially in the nonsmooth case.

Our software, HANSO and GRANSO, is available (unconstrained and constrained) along with HIFOO (H-infinity fixed order optimization) for controller design, which has been used successfully in many applications.



# Our Papers

[Introduction](#)

[Gradient Sampling](#)

[Quasi-Newton Methods](#)

[A Difficult Nonconvex Problem from Nesterov](#)

[Limited Memory Methods](#)

[Concluding Remarks](#)  
[Summary](#)  
[Our Papers](#)

J. V. Burke, A. S. Lewis and M. L. Overton, A Robust Gradient Sampling Method for Nonsmooth, Nonconvex Optimization, *SIAM J. Optimization*, 2005

M. Gürbüzbalaban and M. L. Overton, On Nesterov's Nonsmooth Chebyshev-Rosenbrock Functions, *SIAM J. Optimization*, 2012

A. S. Lewis and M. L. Overton, Nonsmooth Optimization via Quasi-Newton Methods, *Math. Programming*, 2013

F. E. Curtis, T. Mitchell and M.L. Overton, A BFGS-SQP Method for Nonsmooth, Nonconvex, Constrained Optimization and its Evaluation using Relative Minimization Profiles, *Optimization Methods and Software*, 2016

A. Asl and M. L. Overton, Analysis of the Gradient Method with an Armijo-Wolfe Line Search on a Class of Nonsmooth Convex Functions, *arXiv*, 2017

J. V. Burke, F. E. Curtis, A. S. Lewis, M. L. Overton and L. Simões, Gradient Sampling Methods for Nonsmooth Optimization, *arXiv*, 2018

Papers, software are available at [www.cs.nyu.edu/overton](http://www.cs.nyu.edu/overton).