

AI505
Optimization

Marco Chiarandini

Department of Mathematics & Computer Science
University of Southern Denmark

List of Lectures

- 1: Introduction
- 2: Derivatives
- 3: Bracketing
- 4: Local Descent
- 5: First-Order Methods
- 6: Second-Order Methods
- 7: Direct Methods
- 8: Stochastic Methods
- 9: Population Methods
- 10: Machine Learning Applications
- 11: Machine Learning Applications
- 12: Machine Learning Applications
- 13: onstrained Optimization
- 14: Linear Programming
- 15: Sampling Plans

1. Introduction

Who is here?

42 in total registered in ItsLearning

- **AI505 (7.5 ECTS)**
34 from Bachelor in AI
- **AI505 (7.5 ECTS) + IAAI501 (2.5 ECTS)**
8 from Master in Mathematics and Economics

Prerequisites

- Calculus
- Linear Algebra
- Programming

Outline

1. Course Organization

2. Introduction

Aims of the course

Learn about **optimization**:

- continuous multivariate optimization
- discrete optimization

Optimization is an important tool in **machine learning**, **decision making** and in analyzing physical systems.

In mathematical terms, an **optimization problem** is the problem of finding the **best solution** from the set of all **feasible solutions**.

The first step in the optimization process is constructing an appropriate **mathematical formulation**. The second is devising an **algorithm** for solving the mathematical formulation.

Contents of the Course (Pensum)

Unit	Main topic
1	Introduction, Univariate Problems
2	Multivariate Problems, Gradient-Based Methods
3	Derivative-Free Methods
4	Optimization for Machine Learning
5	Constrained Optimization, Linear Programming
6	Sampling Methods
7	Discrete Optimization and Heuristics

Practical Information

Teacher: Marco Chiarandini (imada.sdu.dk/u/marco/)

Instructor: Bonnie Lifting (H21)

Schedule, alternative views:

- mitsdu.sdu.dk, SDU Mobile
- Official course description (læserplanen)
- ItsLearning
- <https://ai-505.github.io>

Schedule (16 weeks):

- Introductory classes: 40 hours (20 classes)
- Training classes: 30 hours (15 classes)
- Scheduled: $16 \times 4 = 64$ hours
- No classes in week 8

Communication Means

- ItsLearning ⇔ External Web Page
(link <https://ai-505.github.io>)
- Announcements in ItsLearning
- Write to Marco (marco@imada.sdu.dk) or to instructor
- Collaborate with peers

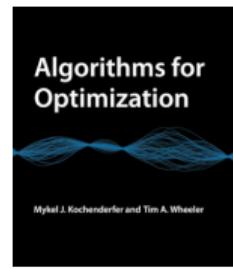
~~ It is good to ask questions!!

~~ Let me know if you think we should do things differently!

Sources — Reading Material

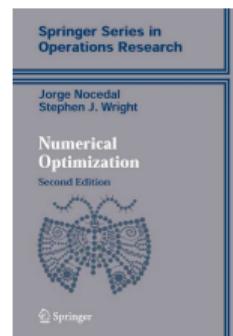
Main reference:

[KW] Mykel J. Kochenderfer, Tim A. Wheeler. Algorithms for Optimization 2019. The MIT Press.



Others

[NW] J. Nocedal and S. J. Wright, Numerical Optimization, Second Edition. Springer Series in Operations Research, 2006



Course Material

External Web Page is the main reference for list of contents
(ie¹, syllabus, pensum).

It will collect:

- slides
- list of topics and references
- exercises
- links
- tutorials for programming tasks

¹ie = id est = that is, eg = exempli gratia = for example, wrt = with respect to, et al. = et alii = and others

Assessment

Portfolio consisting of:

- mandatory assignments in groups of 2:

1. assignment
2. assignment
3. assignment
4. assignment

Time line to be announced.

- oral exam on June 11-13, 2025

The oral examination consists of questions on the basis of the handed in assignments and can be extended to other parts of the course curriculum.

Final grade: starting from the grade of the assignments overall evaluation can move up or down by at most two levels.

(language: Danish and/or English)

Exercise Sessions (1/2)

For exercises in general:

- Both theoretical, modeling and programming tasks
- No mandatory hand-ins, voluntary participation, but all lectures and exercises are relevant for the exam
- Help each other! Teaching others is the best form of learning
- Bonnie will be there to help
- Exercises should help you to learn - don't hesitate to find a work setting or extend with other material that works better for you

Exercises are group work

- Exercises best done in pairs of 2-3 people
- Try to gather in different groups every now and then

Exercise Sessions (1/2)

Notation: \dagger , $*$, \ddagger

- plus exercises are to be done before the class
- starred exercises are done in class
- unmarked exercises are for self study
- starred exercises are examples of assignment questions

Coding

Set up your local Python programming environment and use the available tutorial to:

- Brush up some knowledge on Python-IDEs and Git
- Grasp some basics on jupyter notebooks
- Recap basics of data processing and visualization

We will span several programming modes, functional, imperative, object oriented, small scripts
large code bases, use pf modules and packages. We will use git.

Outline

1. Course Organization

2. Introduction

Introduction

- Applications of Optimization
 - Physics
 - Business
 - Biology
 - Engineering
 - Machine Learning
 - Logistics and Scheduling
- Objectives to Optimize
 - Efficiency
 - Safety
 - Accuracy
- Constraints
 - Cost
 - Weight
 - Structural Integrity
- Challenges
 - High-Dimensional Search Spaces
 - Multiple Competing Objectives
 - Model Uncertainty

A Brief History

- Algebra, study of the rules for manipulating mathematical symbols
- Calculus, study of continuous change It stems from the developments of Leibniz (1646–1716) and Newton (1642–1727). Both differential and integral calculus make use of the notion of convergence of infinite series to a well-defined limit.
- Computers mid-twentieth and numerical algorithms for optimization
- Linear programming, which is an optimization problem with a linear objective function and linear constraints. Leonid Kantorovich (1912–1986) presented a formulation for linear programming and an algorithm to solve it.
- It was applied to optimal resource allocation problems during World War II. George Dantzig (1914–2005) developed the simplex algorithm, which represented a significant advance in solving linear
- Richard Bellman (1920–1984) developed the notion of dynamic programming, which is a commonly used method for optimally solving complex problems by breaking them down into simpler problems
- Artificial Intelligence (⇌ Optimization)

Real World vs Model vs Representation vs Implementation

The Real World: That messy thing we are trying to study (with computers).

Model: Mathematical object in some class M.

Representation: An object of an abstract data type R used to store the model.

Implementation: An object of a concrete type used to store the model.

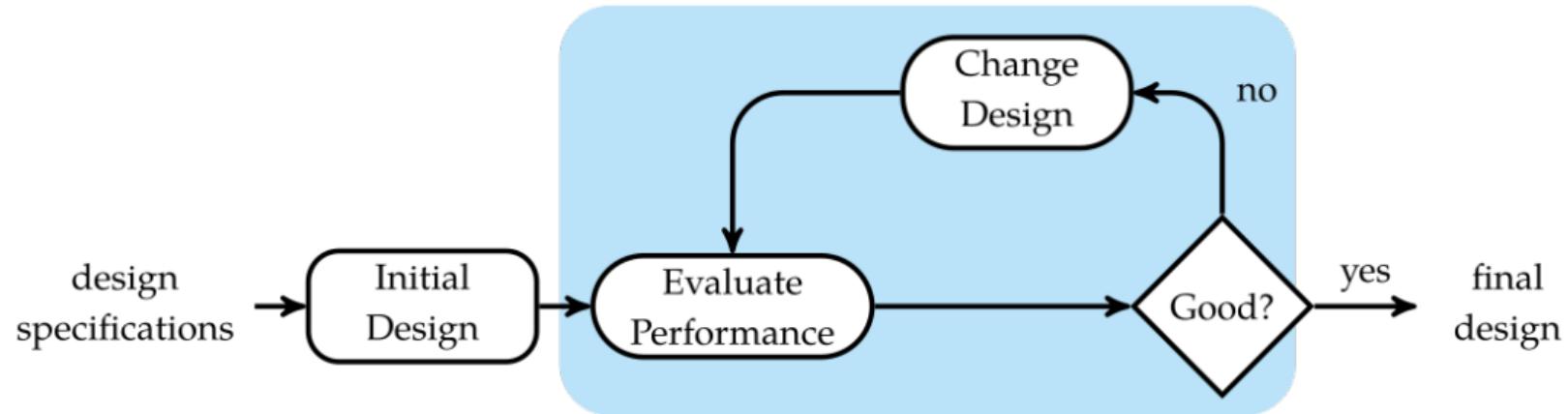
Any object from the real world might have different models.

Any model might have several representations (exact).

And representation might have different implementations (exact).

We will focus on the algorithmic aspects of optimization that arise after the problem has been properly formulated

Optimization Process



An optimization algorithm is used to incrementally improve the design until it can no longer be improved or until the budgeted time or cost has been reached.

Optimization Process

Search the space of possible designs with the aim of finding the best one.

Depending on the application, this search may involve:

- evaluating an analytical expression (white or glass box)
- running physical experiments, such as wind tunnel tests (black box)
- running computer simulations

Modern optimization techniques can be applied to problems with millions of variables and constraints.

Basic Optimization Problem

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad f(x) \\ & \text{subject to } x \in \mathcal{X} \end{aligned}$$

- Design Point
- Design Variables
- Objective Function
- Feasible Set
- Minimizer

Any value of x from among all points in the feasible set \mathcal{X} that minimizes the objective function is called a **solution** or **minimizer**. A particular solution is written x^* .

$$x^* = \operatorname{argmin}_x f(x) \text{ subject to } x \in \mathcal{X}$$

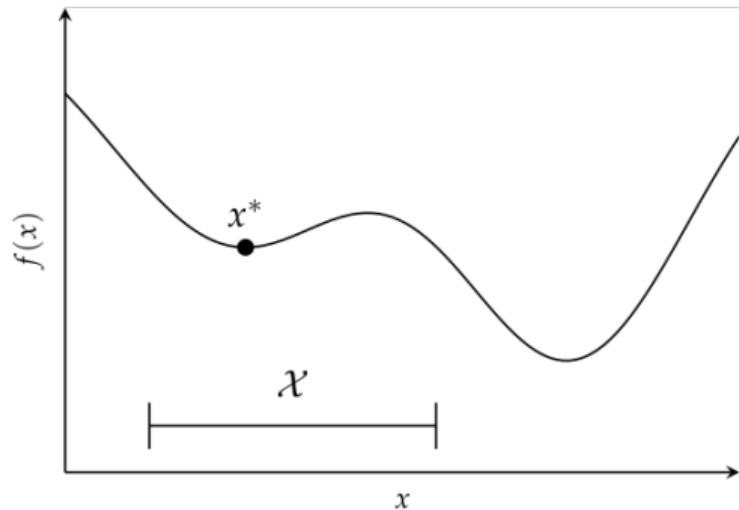
There is only one minimum but there can be many minimizers

$$\text{maximize } x \ f(x) \text{ subject to } x \in \mathcal{X} \equiv \text{minimize } x - f(x) \text{ subject to } x \in \mathcal{X}$$

Basic Optimization Problem

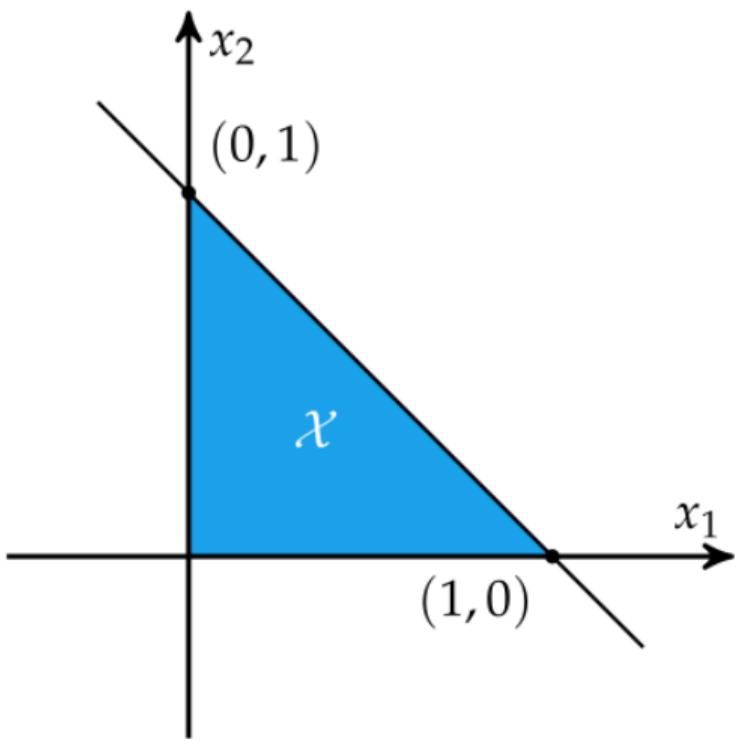
minimize
$$_x f(x)$$

subject to $x \in \mathcal{X}$



Constraints

$$\begin{array}{ll}\text{minimize}_{x_1, x_2} & f(x_1, x_2) \\ \text{subject to} & x_1 \geq 0 \\ & x_2 \geq 0 \\ & x_1 + x_2 \leq 1\end{array}$$



Constraints

$$\begin{aligned} & \underset{x}{\text{minimize}} \quad f(x) \\ & \text{subject to } x > 1 \end{aligned}$$

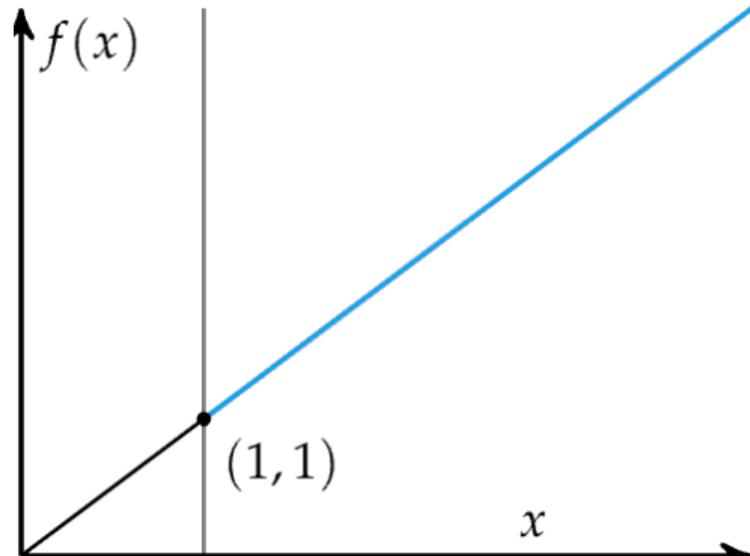
The problem has no solution.

$x = 1$ would not be feasible.

$x = 1$ would be the solutions to

$$\inf_x f(x) \text{ subject to } x > 1$$

infimum of a subset \mathcal{X} of a partially ordered set \mathcal{P} is the greatest element in \mathcal{P} that is less than or equal to each element of \mathcal{X} , if such an element exists. Aka, **greatest lower bound**.

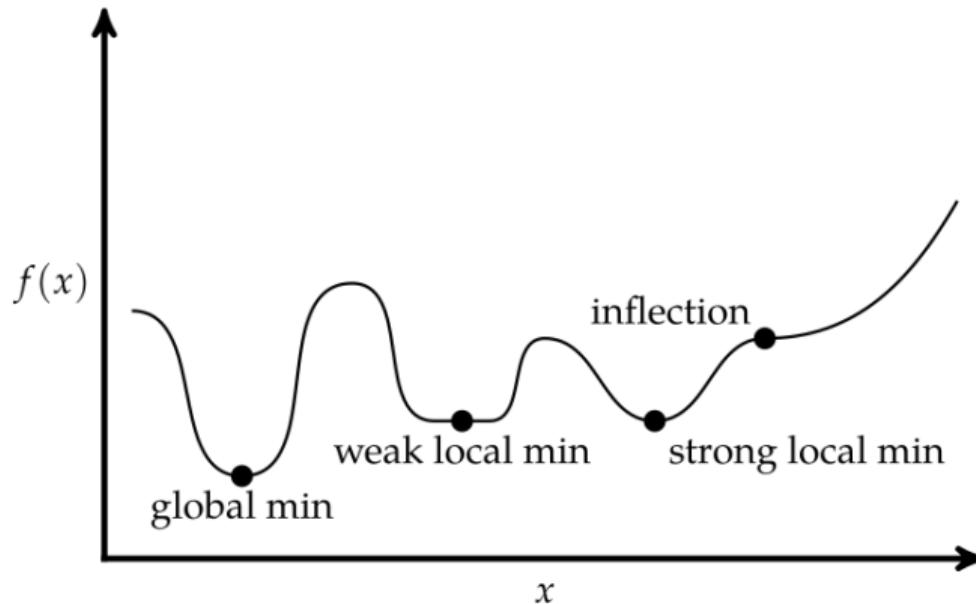


Critical Points

Univariate function

global vs local minimum

Def. A point x^* is at a **local minimum** (or is a local minimizer) if there exists a $\delta > 0$ such that $f(x^*) \leq f(x)$ for all x with $\|x - x^*\| < \delta$.

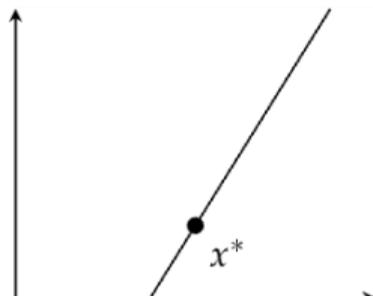


Conditions for Local Minima

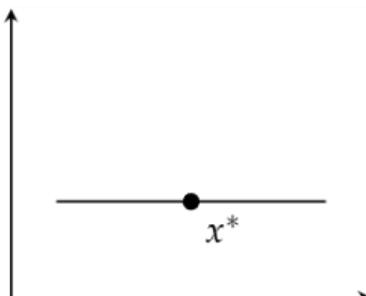
Univariate objective functions, assuming they are differentiable (Derivatives exist), without constraints

Local minimum: Necessary condition but not sufficient condition:

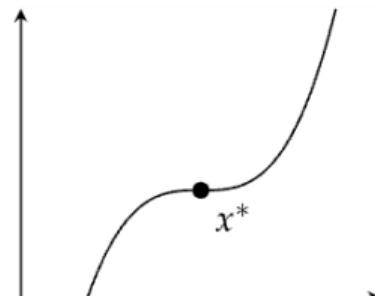
1. $f'(x^*) = 0$, the first-order necessary condition (FONC)
2. $f''(x^*) \geq 0$, the second-order necessary condition (SONC)



SONC but not FONC



FONC and SONC



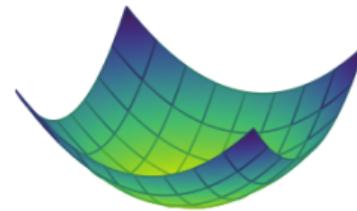
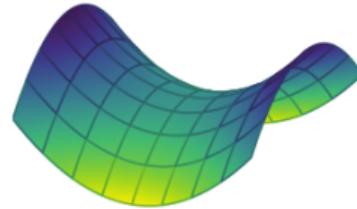
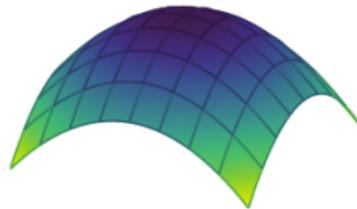
FONC and SONC

Conditions for Local Minima

Multivariate objective functions, assuming they are differentiable (gradients and Hessians exist), without constraints

Local minimum: Necessary condition but not sufficient condition:

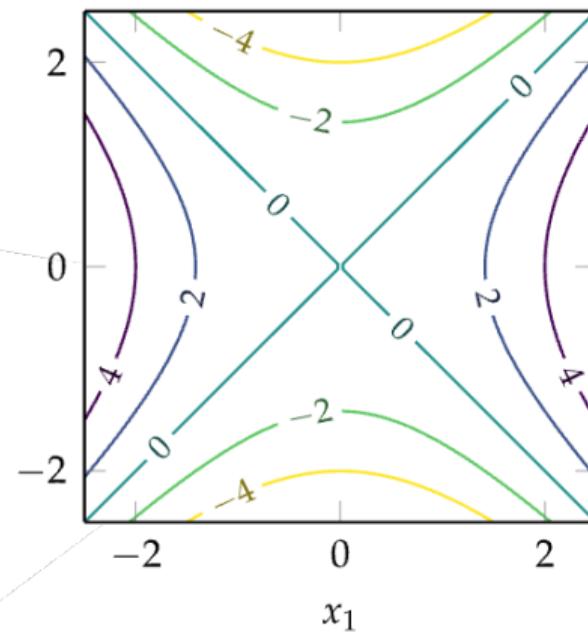
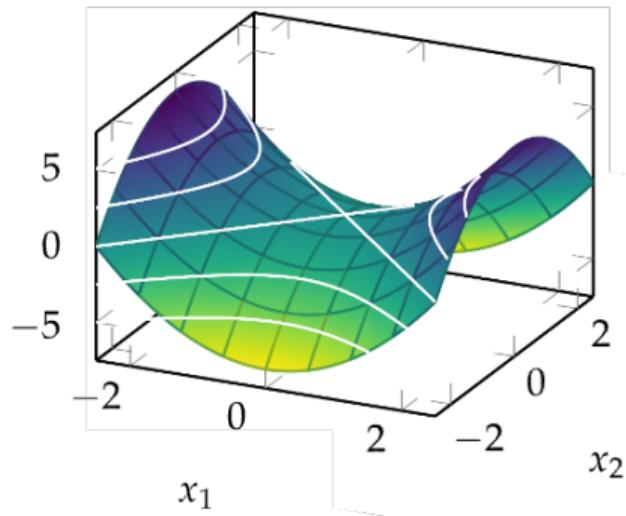
1. $\nabla f(\mathbf{x}^*) = \mathbf{0}$, the first-order necessary condition (FONC)
2. $\nabla^2 f(\mathbf{x}^*) \geq \mathbf{0}$, the second-order necessary condition (SONC)



Contour Plots

$$f(x_1, x_2) = x_1^2 - x_2^2$$

can be rendered in a 3D space but convenient to represent it also in 2D showing the lines of constant output value



Taylor Expansion

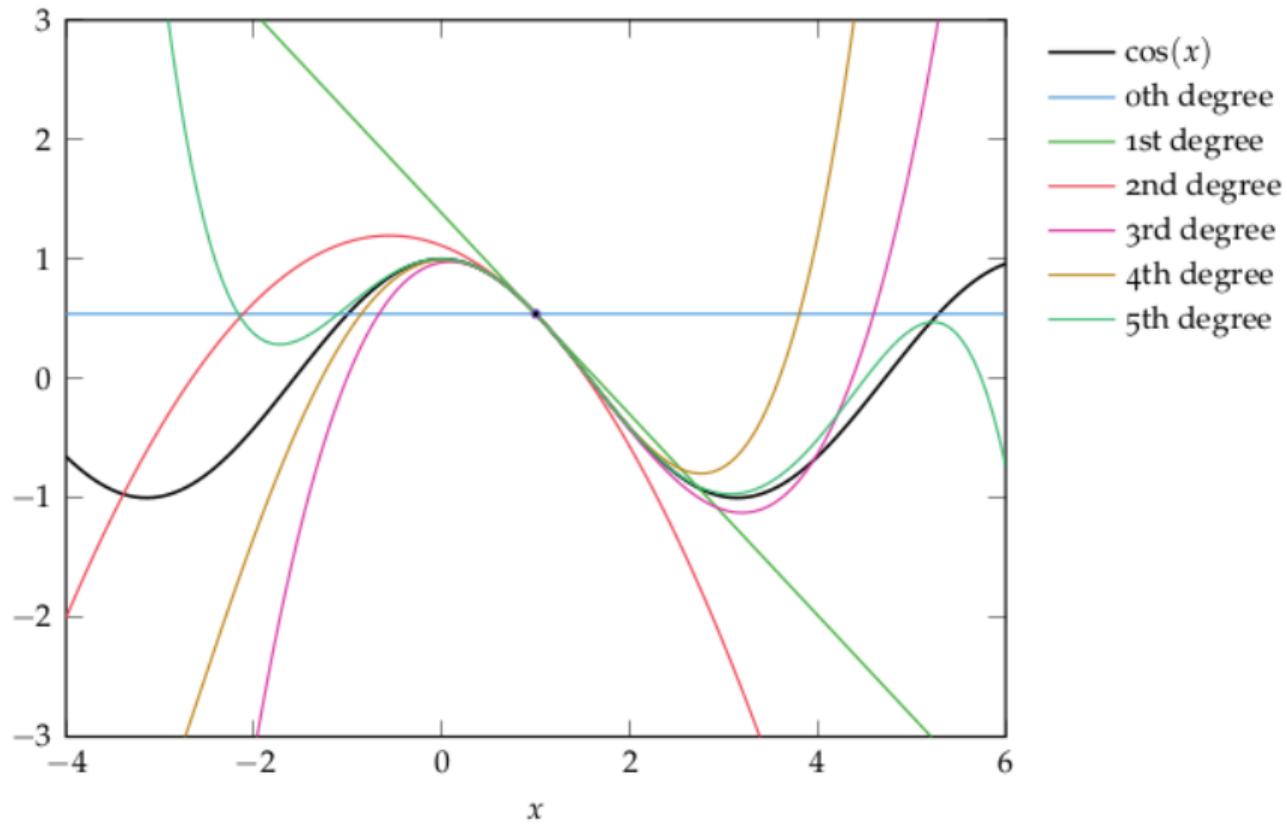
From the *first fundamental theorem of calculus*,² we know that

$$f(x+h) = f(x) + \int_0^h f'(x+a) da$$

Nesting this definition produces the Taylor expansion of f about x :

$$\begin{aligned} f(x+h) &= f(x) + \int_0^h \left(f'(x) + \int_0^a f''(x+b) db \right) da \\ &= f(x) + f'(x)h + \int_0^h \int_0^a f''(x+b) db da \\ &= f(x) + f'(x)h + \int_0^h \int_0^a \left(f''(x) + \int_0^b f'''(x+c) dc \right) db da \\ &= f(x) + f'(x)h + \frac{f''(x)}{2!}h^2 + \int_0^h \int_0^a \int_0^b f'''(x+c) dc db da \\ &\quad \vdots \\ &= f(x) + \frac{f'(x)}{1!}h + \frac{f''(x)}{2!}h^2 + \frac{f'''(x)}{3!}h^3 + \dots \\ &= \sum_{n=0}^{\infty} \frac{f^{(n)}(x)}{n!}h^n \end{aligned}$$

Taylor Expansion



Taylor Expansion Multidim

In multiple dimensions, the Taylor expansion about \mathbf{a} generalizes to

$$f(\mathbf{x}) = f(\mathbf{a}) + \nabla f(\mathbf{a})^\top (\mathbf{x} - \mathbf{a}) + \frac{1}{2} (\mathbf{x} - \mathbf{a})^\top \nabla^2 f(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

Example: Rosenbrock function

$$f(x, y) = (a - x)^2 + b(y - x^2)^2$$

It has a global minimum at $(x, y) = (a, a^2)$, where $f(x, y) = 0$. Usually, these parameters are set such that $a = 1$ and $b = 100$. Only in the trivial case where $a = 0$ the function is symmetric and the minimum is at the origin.

Multivariate generalization
sum of $N/2$ uncoupled 2D Rosenbrock problems, and defined only for even N :

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_N) = \sum_{i=1}^{N/2} [100(x_{2i-1}^2 - x_{2i})^2 + (x_{2i-1} - 1)^2].$$

This variant has predictably simple solutions.

Example: Rosenbrock function

Consider the Rosenbrock banana function,

$$f(\mathbf{x}) = (1 - x_1)^2 + 5(x_2 - x_1^2)^2$$

Does the point $(1, 1)$ satisfy the FONC and SONC?

The gradient is:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2(10x_1^3 - 10x_1x_2 + x_1 - 1) \\ 10(x_2 - x_1^2) \end{bmatrix}$$

and the Hessian is:

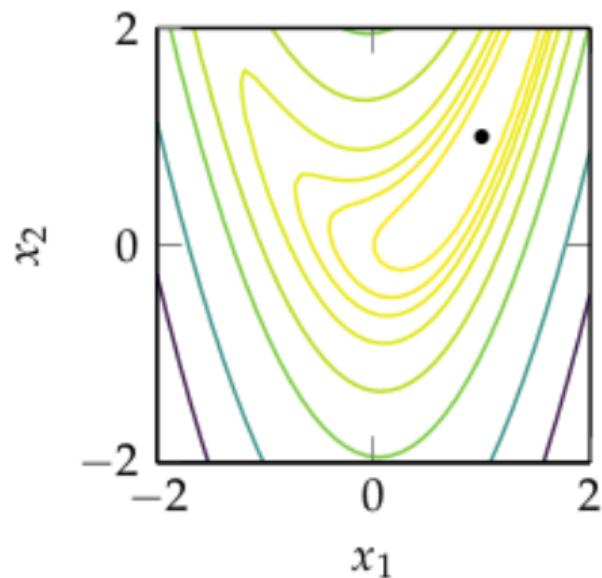
$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} -20(x_2 - x_1^2) + 40x_1^2 + 2 & -20x_1 \\ -20x_1 & 10 \end{bmatrix}$$

We compute $\nabla(f)([1, 1]) = 0$, so the FONC is satisfied. The Hessian at $[1, 1]$ is:

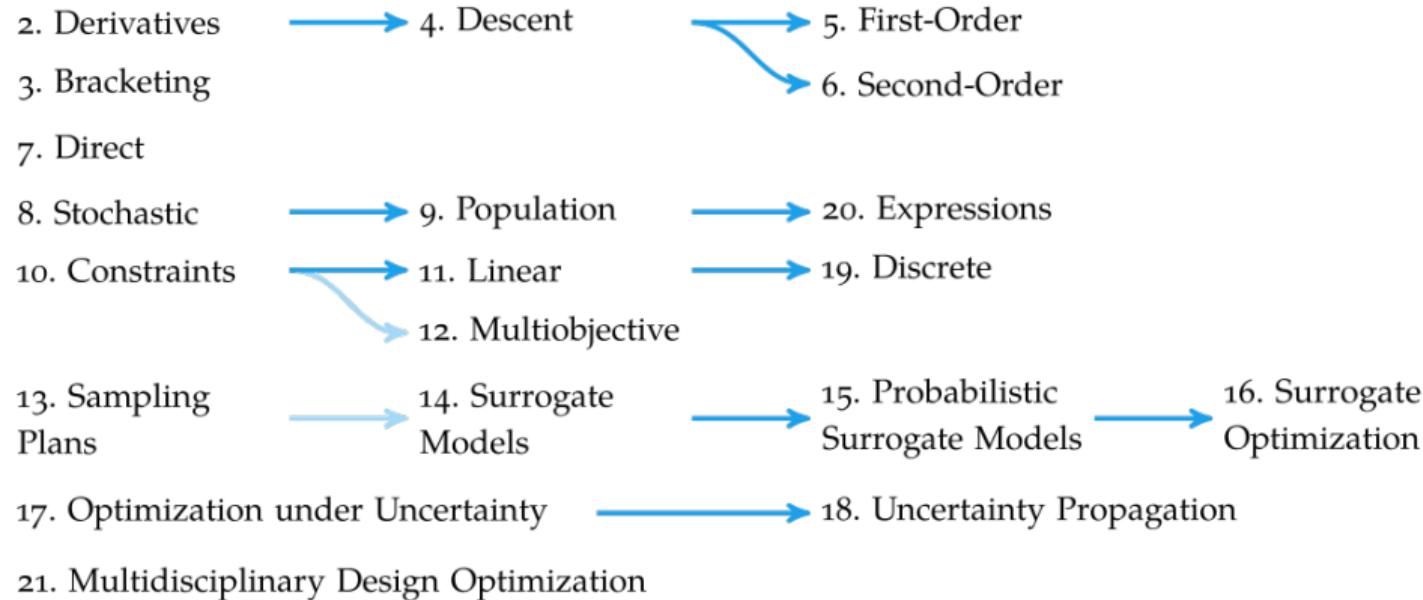
$$\begin{bmatrix} 42 & -20 \\ -20 & 10 \end{bmatrix}$$

which is positive definite, so the SONC is satisfied.

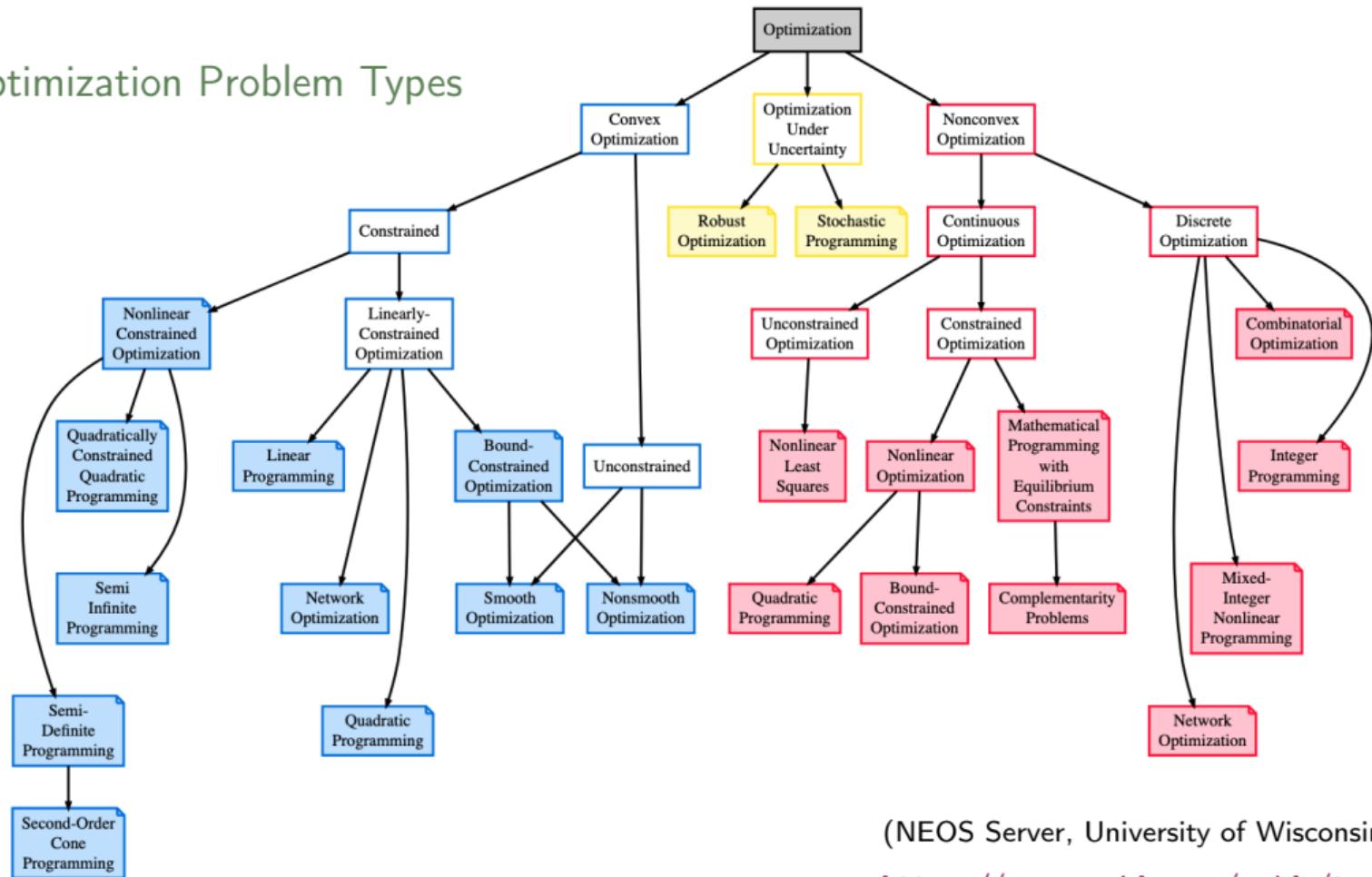
Example: Rosenbrock function



Overview



Optimization Problem Types



(NEOS Server, University of Wisconsin)

<https://neos-guide.org/guide/types/>

Problem classification

Different classification parameters:

- Univariate $f : \mathbb{R} \rightarrow \mathbb{R}$ vs Multivariate $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- Real-valued $f : \mathbb{R}^n \rightarrow \mathbb{R}$ vs vector functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Linear vs Nonlinear
- Nonlinear: Convex vs Nonconvex, unimodal vs multimodal
- Constrained vs unconstrained
- Smooth (differentiable) vs non smooth (non differentiable)
- Deterministic vs Uncertain
- Continuous vs Discrete

Application Example

In robotics, an autonomous agent (e.g., a drone, robotic arm, or self-driving car) needs to determine an **optimal path** from a starting position to a target while satisfying constraints such as avoiding obstacles, minimizing energy consumption, or optimizing smoothness.

Formulating the Optimization Problem

The goal is to find a smooth trajectory $\mathbf{x}(t)$ that minimizes a cost function:

$$J = \int_{t_0}^{t^f} C(\mathbf{x}(t), \mathbf{u}(t)) dt$$

$\mathbf{x}(t)$ is the state (e.g., position, velocity), $\mathbf{u}(t)$ is the control input (e.g., force, acceleration), $C(\mathbf{x}, \mathbf{u})$ is the cost function, which could represent energy usage, time, or distance.

Constraints:

- Dynamic Constraints: Governed by the system's physics (e.g., Newton's laws for a robot arm or quadcopter): $\mathbf{x}' = f(\mathbf{x}, \mathbf{u})$
- Obstacle Avoidance: Ensures that the trajectory does not collide with obstacles:
 $g(\mathbf{x}(t)) \geq 0, \forall t$
- Boundary Conditions: The system must start and end at given positions with certain velocities.

Summary

- Optimization is the process of finding the best system design subject to a set of constraints
- Optimization is concerned with finding global minima of a function
- Minima can occur where the gradient is zero, but zero-gradient does not imply optimality

2. Derivatives and Gradients

Definitions

- $[a, b] = \{x \in \mathbb{R} \mid a \leq x \leq b\}$ closed interval
 $(a, b) = \{x \in \mathbb{R} \mid a < x < b\}$ open interval
- column vectors and matrices
scalar product: $\mathbf{y}^T \mathbf{x} = \sum_{i=1}^n y_i x_i$
- \mathbf{Ax} column vector combination of the columns of A ;
 $\mathbf{u}^T A$ row vector combination of the rows of A

Definitions

- **linear combination**

$$\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k \in \mathbb{R}^n$$
$$\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_k]^T \in \mathbb{R}^k$$

$$\mathbf{x} = \lambda_1 \mathbf{v}_1 + \cdots + \lambda_k \mathbf{v}_k = \sum_{i=1}^k \lambda_i \mathbf{v}_i$$

moreover:

$$\boldsymbol{\lambda} \geq 0$$

$$\boldsymbol{\lambda}^T \mathbf{1} = 1$$

$$\boldsymbol{\lambda} \geq 0 \quad \text{and} \quad \boldsymbol{\lambda}^T \mathbf{1} = 1$$

conic combination

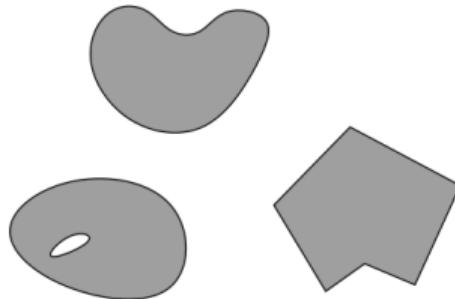
affine combination

convex combination

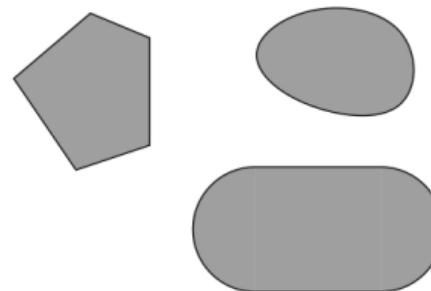
$$\left(\sum_{i=1}^k \lambda_i = 1 \right)$$

Definitions

- **convex set**: if $\mathbf{x}, \mathbf{y} \in S$ and $0 \leq \lambda \leq 1$ then $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in S$



nonconvex



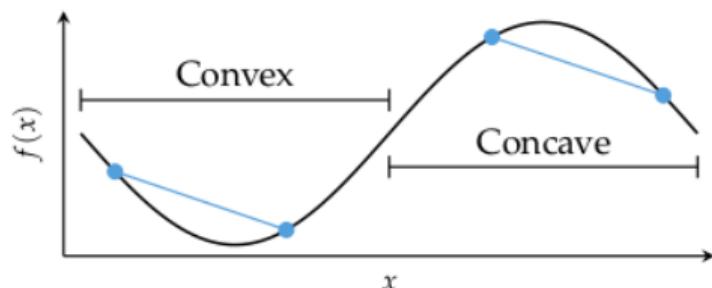
convex

- **convex function** if its epigraph

$\{(x, y) \in \mathbb{R}^2 : y \geq f(x)\}$ is a convex set or if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and

if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \alpha \in [0, 1]$ it holds that

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})$$



Definitions

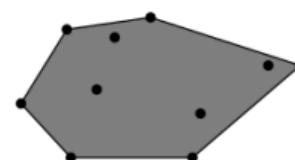
- For a set of points $S \subseteq \mathbb{R}^n$

$\text{lin}(S)$ linear hull (span)

$\text{cone}(S)$ conic hull

$\text{aff}(S)$ affine hull

$\text{conv}(S)$ convex hull



the convex hull of X

$$\text{conv}(X) = \left\{ \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n \mid x_i \in X, \lambda_1, \dots, \lambda_n \geq 0 \text{ and } \sum_i \lambda_i = 1 \right\}$$

Norms

Def. A **norm** is a function that assigns a length to a vector.

A function f is a norm if:

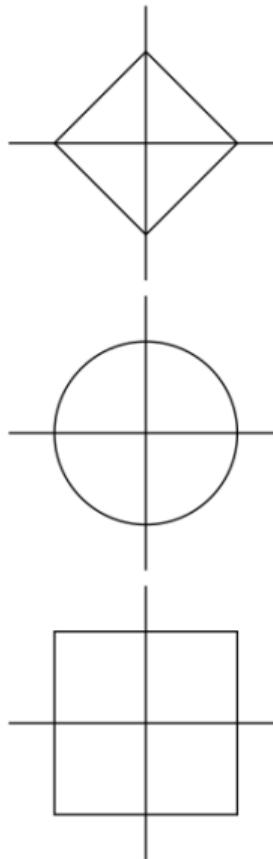
1. $f(\mathbf{x}) = 0$ if and only if \mathbf{x} is the zero vector
2. $f(a\mathbf{x}) = |a|f(\mathbf{x})$, such that lengths scale
3. $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$, also known as triangle inequality

L_p norms are commonly used set of norms parameterized by a scalar $p \geq 1$:

$$\|\mathbf{x}\|_p = \lim_{p \rightarrow p} (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}$$

L_∞ is also called the max norm, Chebyshev distance or chessboard distance.

$$L_1: \|\mathbf{x}\|_1 = |x_1| + |x_2| + \cdots + |x_n|$$



$$L_2: \|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2}$$

$$L_\infty: \|\mathbf{x}\|_\infty = \max(|x_1|, |x_2|, \dots, |x_n|)$$

Outline

3. Derivatives

4. Symbolic Differentiation

5. Numerical Differentiation

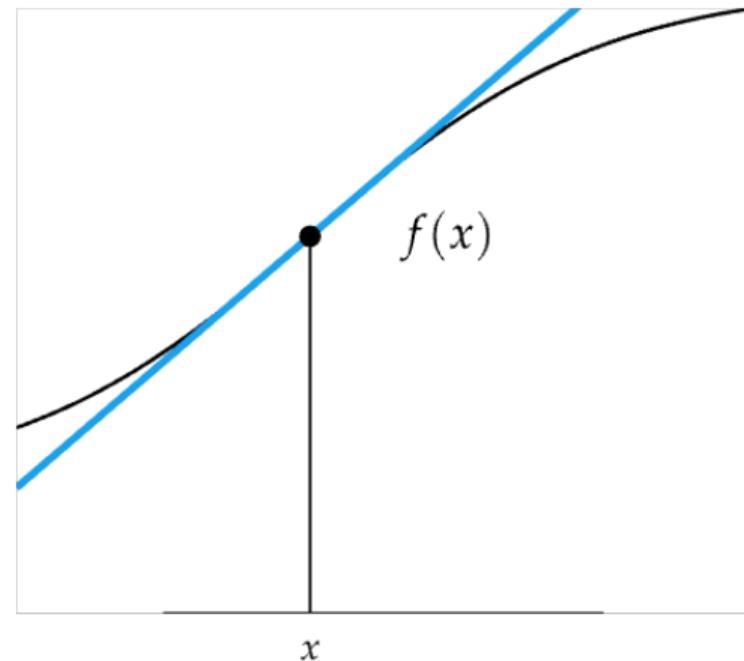
6. Automatic Differentiation

Derivatives

- Derivatives tell us which direction to search for a solution
- Slope of Tangent Line

$$f'(x) := \frac{df(x)}{dx}$$

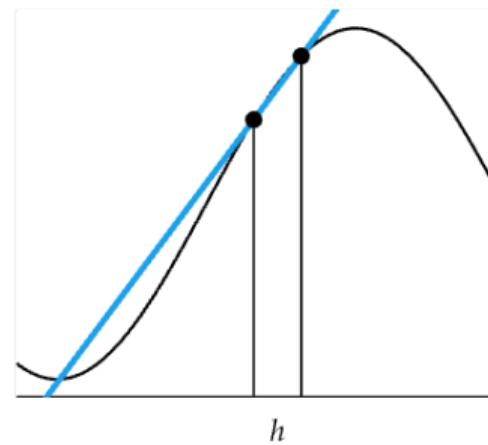
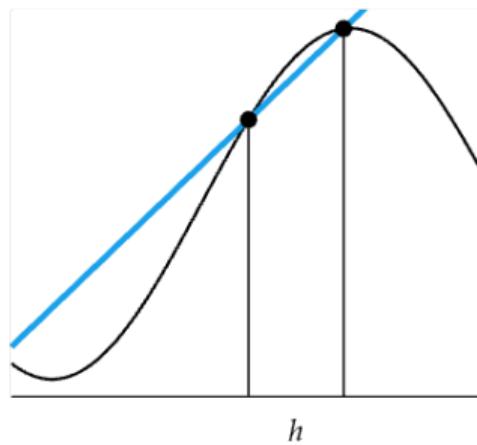
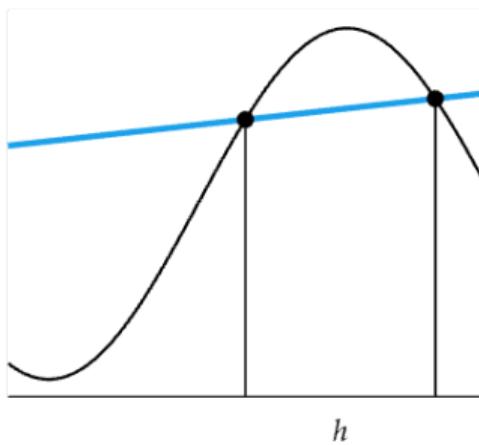
(Leibniz notation)



Derivatives

$$f(x + \Delta x) \approx f(x) + f'(x)\Delta x$$

$$f'(x) = \frac{\Delta x}{\Delta x}$$



Symbolic Differentiation

$$f'(x) \equiv \underbrace{\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}}_{\text{forward difference}} = \underbrace{\lim_{h \rightarrow 0} \frac{f(x+h/2) - f(x-h/2)}{h}}_{\text{central difference}} = \underbrace{\lim_{h \rightarrow 0} \frac{f(x) - f(x-h)}{h}}_{\text{backward difference}}$$

Symbolic Differentiation

```
import sympy as sp

# Define the variable
x = sp.symbols('x')

# Define the function
f = x**2 + x/2 - sp.sin(x)/x

# Compute the derivative
df_dx = sp.diff(f, x)

# Display the result
print("The symbolic derivative of f is:")
print(df_dx)
```

derivative.py

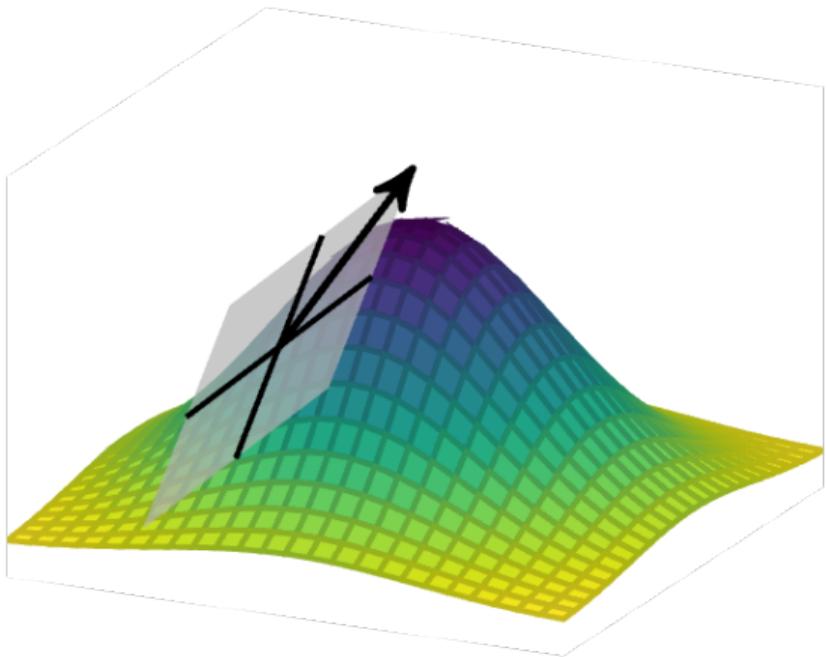
Derivatives in Multiple Dimensions

- Gradient Vector

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]$$

- Hessian Matrix

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix}$$



Directional derivative

The **directional derivative** $\nabla_s f(\mathbf{x})$ of a multivariate function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the instantaneous rate of change of $f(\mathbf{x})$ as $\mathbf{x} = [x_1, x_2, \dots, x_n]$ is moved with velocity $\mathbf{s} = [s_1, s_2, \dots, s_n]$.

$$\nabla_s f(\mathbf{x}) \equiv \underbrace{\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{s}) - f(\mathbf{x})}{h}}_{\text{forward difference}} = \underbrace{\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{s}/2) - f(\mathbf{x} - h\mathbf{s}/2)}{h}}_{\text{central difference}} = \underbrace{\lim_{h \rightarrow 0} \frac{f(\mathbf{x}) - f(\mathbf{x} - h\mathbf{s})}{h}}_{\text{backward difference}}$$

To compute $\nabla_s f(\mathbf{x})$:

- compute $\nabla_s f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial x_1} s_1 + \frac{\partial f(\mathbf{x})}{\partial x_2} s_2 + \dots + \frac{\partial f(\mathbf{x})}{\partial x_n} s_n = \nabla f(\mathbf{x})^T \mathbf{s} = \nabla f(\mathbf{x}) \cdot \mathbf{s}$
- $g(\alpha) := f(\mathbf{x} + \alpha \mathbf{s})$ and then compute $g'(0)$

We wish to compute the directional derivative of $f(\mathbf{x}) = x_1x_2$ at $\mathbf{x} = [1, 0]$ in the direction $\mathbf{s} = [-1, -1]$:

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}, \quad \frac{\partial f}{\partial x_2} \right] = [x_2, x_1]$$

$$\nabla_{\mathbf{s}} f(\mathbf{x}) = \nabla f(\mathbf{x})^\top \mathbf{s} = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix} = -1$$

We can also compute the directional derivative as follows:

$$g(\alpha) = f(\mathbf{x} + \alpha \mathbf{s}) = (1 - \alpha)(-\alpha) = \alpha^2 - \alpha$$

$$g'(\alpha) = 2\alpha - 1$$

$$g'(0) = -1$$

Matrix Calculus

Common gradient:

$$\nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} = ?$$

$$\mathbf{b}^T \mathbf{x} = [b_1 x_1 + b_2 x_2 + \dots + b_n x_n]$$

$$\frac{\partial \mathbf{b}^T \mathbf{x}}{\partial x_i} = b_i$$

$$\nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} = \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{b} = \mathbf{b}$$

Matrix Calculus

Common gradient:

$$\nabla_x x^T A x = ?$$

$$x^T A x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}^T \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}^T \begin{bmatrix} x_1 a_{11} + x_2 a_{12} + \dots + x_n a_{1n} \\ x_1 a_{21} + x_2 a_{22} + \dots + x_n a_{2n} \\ \vdots \\ x_1 a_{n1} + x_2 a_{n2} + \dots + x_n a_{nn} \end{bmatrix}$$

$$x_1^2 a_{11} + x_1 x_2 a_{12} + \dots + x_1 x_n a_{1n} +$$

$$= x_1 x_2 a_{21} + x_2^2 a_{22} + \dots + x_2 x_n a_{2n} +$$

⋮

$$x_1 x_n a_{n1} + x_2 x_n a_{n2} + \dots + x_n^2 a_{nn}$$

$$\frac{\partial}{\partial x_i} \mathbf{x}^T A \mathbf{x} = \sum_{j=1}^n x_j (a_{ij} + a_{ji})$$

$$\nabla_{\mathbf{x}} \mathbf{x}^T A \mathbf{x} = \begin{bmatrix} \sum_{j=1}^n x_j (a_{1j} + a_{j1}) \\ \sum_{j=1}^n x_j (a_{2j} + a_{j2}) \\ \vdots \\ \sum_{j=1}^n x_j (a_{nj} + a_{jn}) \end{bmatrix} = \begin{bmatrix} a_{11} + a_{11} & a_{12} + a_{21} & \dots & a_{1n} + a_{n1} \\ a_{21} + a_{12} & a_{22} + a_{22} & \dots & a_{2n} + a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} + a_{1n} & a_{n2} + a_{2n} & \dots & a_{nn} + a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = (A + A^T) \mathbf{x}$$

Smoothness

Def. The **smoothness** of a function is a property measured by the number of continuous derivatives (differentiability class) it has over its domain.

A function of **class C^k** is a function of smoothness at least k ; that is, a function of class C^k is a function that has a k th derivative that is continuous in its domain.

The term **smooth function** refers to a C^∞ -function. However, it may also mean “sufficiently differentiable” for the problem under consideration.

Smoothness

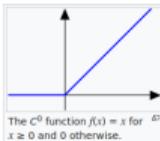
- Let U be an open set on the real line and a function f defined on U with real values. Let k be a non-negative integer.
- The function f is said to be of **differentiability class C^k** if the derivatives $f', f'', \dots, f^{(k)}$ exist and are continuous on U .
- If f is k -differentiable on U , then it is at least in the class C^{k-1} since $f', f'', \dots, f^{(k-1)}$ are continuous on U .
- The function f is said to be **infinitely differentiable, smooth**, or of **class C^∞** , if it has derivatives of all orders (continuous) on U .
- The function f is said to be of **class C^ω** , or **analytic**, if f is smooth and its Taylor series expansion around any point in its domain converges to the function in some neighborhood of the point.
- There exist functions that are smooth but not analytic; C^ω is thus strictly contained in C^∞ . Bump functions are examples of functions with this property.

Example: continuous (C^0) but not differentiable [edit]

The function

$$f(x) = \begin{cases} x & \text{if } x \geq 0, \\ 0 & \text{if } x < 0 \end{cases}$$

is continuous, but not differentiable at $x = 0$, so it is of class C^0 , but not of class C^1 .



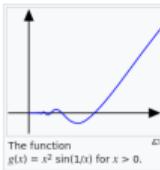
The C^0 function $f(x) = x$ for $x \geq 0$ and 0 otherwise.

Example: finitely-times differentiable (C^k) [edit]

For each even integer k , the function

$$f(x) = |x|^{k+1}$$

is continuous and k times differentiable at all x . At $x = 0$, however, f is not $(k+1)$ times differentiable, so f is of class C^k , but not of class C^l where $l > k$.



The function $f(x) = |x|^{k+1}$ for even k .

Example: differentiable but not continuously differentiable (not C^1) [edit]

The function

$$g(x) = \begin{cases} x^2 \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0, \\ 0 & \text{if } x = 0 \end{cases}$$

is differentiable, with derivative

$$g'(x) = \begin{cases} -\cos\left(\frac{1}{x}\right) + 2x \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0, \\ 0 & \text{if } x = 0. \end{cases}$$

Because $\cos(1/x)$ oscillates as $x \rightarrow 0$, $g'(x)$ is not continuous at zero.

Therefore, $g(x)$ is differentiable but not of class C^1 .



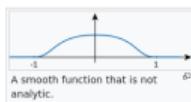
The function $f : \mathbb{R} \rightarrow \mathbb{R}$ with $f(x) = x^2 \sin\left(\frac{1}{x}\right)$ for $x \neq 0$ and $f(0) = 0$ is differentiable. However, this function is not continuously differentiable.

Example: differentiable but not Lipschitz continuous [edit]

The function

$$h(x) = \begin{cases} x^{4/3} \sin\left(\frac{1}{x}\right) & \text{if } x \neq 0, \\ 0 & \text{if } x = 0 \end{cases}$$

is differentiable but its derivative is unbounded on a compact set. Therefore, h is an example of a function that is differentiable but not locally Lipschitz continuous.



A smooth function that is not analytic.

Example: analytic (C^ω) [edit]

The exponential function e^x is analytic, and hence falls into the class C^ω (where ω is the smallest transfinite ordinal). The trigonometric functions are also analytic wherever they are defined, because they are linear combinations of complex exponential functions e^{ix} and e^{-ix} .

Example: smooth (C^∞) but not analytic (C^ω) [edit]

The bump function

$$f(x) = \begin{cases} e^{-\frac{1}{1-x^2}} & \text{if } |x| < 1, \\ 0 & \text{otherwise} \end{cases}$$

is smooth, so of class C^∞ , but it is not analytic at $x = \pm 1$, and hence is not of class C^ω . The function f is an example of a smooth function with compact support.

Positive Definiteness

Def. A symmetric matrix A is **positive definite** if $x^T Ax$ is positive for all points other than the origin: $x^T Ax > 0$ for all $x \neq 0$.

Def. A symmetric matrix A is **positive semidefinite** if $x^T Ax$ is always non-negative: $x^T Ax \geq 0$ for all x .

A matrix A is positive definite if and only all its eigenvalues are positive.

If the matrix A is positive definite in the function $f(x) = x^T Ax$, then f has a unique global minimum.

Recall that the second order Taylor approximation of a twice-differentiable function f at x_0 is

$$f(x) \approx f(x_0) + \nabla f(x_0)^T (x - x_0) + \frac{1}{2} (x - x_0)^T H_0 (x - x_0)$$

where H_0 is the Hessian evaluated at x_0 . If $(x - x_0)^T H_0 (x - x_0)$ has a unique global minimum, then the overall approximation has a unique global minimum.

Outline

3. Derivatives

4. Symbolic Differentiation

5. Numerical Differentiation

6. Automatic Differentiation

Symbolic Derivatives

- Symbolic derivatives can give valuable insight into the structure of the problem domain and, in some cases, produce analytical solutions of extrema (e.g., solving for $\frac{d}{dx} f(x) = 0$) that can eliminate the need for derivative calculation altogether.
- But they do not lend themselves to efficient runtime calculation of derivative values, as they can get exponentially larger than the expression whose derivative they represent

Outline

- 3. Derivatives
- 4. Symbolic Differentiation
- 5. Numerical Differentiation
- 6. Automatic Differentiation

Numerical Differentiation

Finite Difference Method

- Neighboring points are used to approximate the derivative

$$f'(x) \approx \underbrace{\frac{f(x+h) - f(x)}{h}}_{\text{forward difference}} \approx \underbrace{\frac{f(x+h/2) - f(x-h/2)}{h}}_{\text{central difference}} \approx \underbrace{\frac{f(x) - f(x-h)}{h}}_{\text{backward difference}}$$

- h too small causes numerical cancellation errors (square root or cube root of the machine precision for floating point values: `sys.float_info.epsilon` difference between 1 and closest representable number)

Derivation

from Taylor series expansion:

$$f(x+h) = f(x) + \frac{f'(x)}{1!}h + \frac{f''(x)}{2!}h^2 + \frac{f'''(x)}{3!}h^3 + \dots$$

We can rearrange and solve for the first derivative:

$$f'(x)h = f(x+h) - f(x) - \frac{f''(x)}{2!}h^2 - \frac{f'''(x)}{3!}h^3 - \dots$$

$$f'(x) = \frac{f(x+h) - f(x)}{h} - \frac{f''(x)}{2!}h - \frac{f'''(x)}{3!}h^2 - \dots$$

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

- forward difference has error term $O(h)$, linear error as h approaches zero
- central difference has error term is $O(h^2)$

```
import sys
import numpy as np

def diff_forward(f, x: float, h: float=np.sqrt(sys.float_info.epsilon)) -> float:
    return (f(x+h) - f(x))/h

def diff_central(f, x: float, h: float=np.cbrt(sys.float_info.epsilon)) -> float:
    return (f(x+h/2) - f(x-h/2))/h

def diff_backward(f, x: float, h: float=np.sqrt(sys.float_info.epsilon)) -> float:
    return (f(x) - f(x-h))/h

# Example usage
def func(x):
    return x**2 + np.sin(x)

x0 = 1.0
print(f"The derivative at x = {x0} is {diff_forward(func, x0)})
```

finite_diff.py

Numerical Differentiation

Complex step method

Uses one single function evaluation after taking a step in the imaginary direction.

$$f(x + ih) = f(x) + ihf'(x) - h^2 \frac{f''(x)}{2!} - ih^3 \frac{f'''(x)}{3!} + \dots$$

$$\operatorname{Im}(f(x + ih)) = hf'(x) - h^3 \frac{f'''(x)}{3!} + \dots$$

$$\begin{aligned}\Rightarrow f'(x) &= \frac{\operatorname{Im}(f(x + ih))}{h} + h^2 \frac{f'''(x)}{3!} - \dots \\ &= \frac{\operatorname{Im}(f(x + ih))}{h} + O(h^2) \text{ as } h \rightarrow 0\end{aligned}$$

$$\operatorname{Re}(f(x + ih)) = f(x) - h^2 \frac{f''(x)}{2!} + \dots$$

$$\Rightarrow f(x) = \operatorname{Re}(f(x + ih)) + h^2 \frac{f''(x)}{2!} - \dots$$

```
import numpy as np

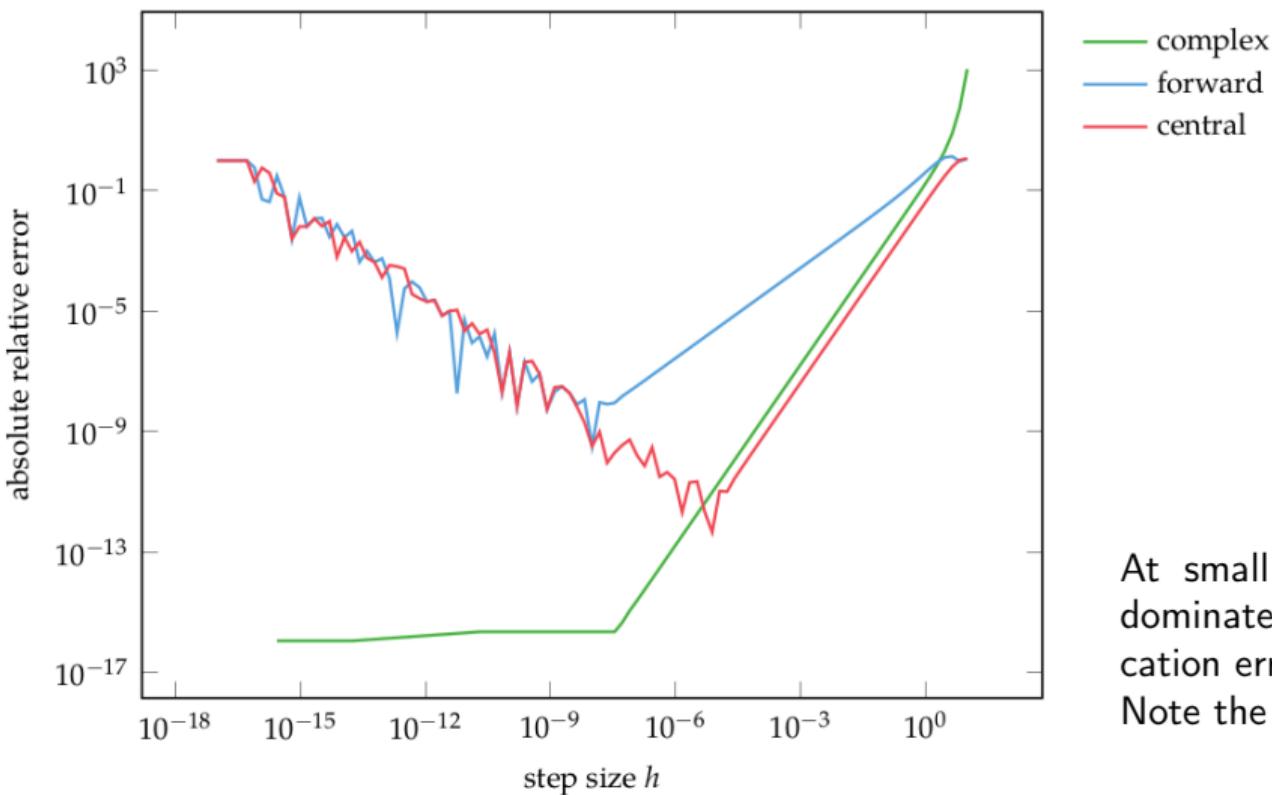
def diff_complex(f, x: float, h: float=1e-20) -> float:
    return np.imag(f(x + h * 1j)) / h

# Example usage
def func(x):
    return x**2 + np.sin(x)

x0 = 1.0
print(f"The derivative at x = {x0} is {diff_complex(func, x0)}")
```

complex_diff.py

Numerical Differentiation Error Comparison



At small h , round off errors dominate, and at large h , truncation errors dominate.

Note the log transformation.

Numerical Differentiation in ML

- Approximation errors would be tolerated in a deep learning setting thanks to the well-documented error resiliency of neural network architectures (Gupta et al., 2015).
- The $O(n)$ complexity of numerical differentiation for a gradient in n dimensions is the main obstacle to its usefulness in machine learning, where n can be as large as millions or billions in state-of-the-art deep learning models (Shazeer et al., 2017).

Outline

- 3. Derivatives
- 4. Symbolic Differentiation
- 5. Numerical Differentiation
- 6. Automatic Differentiation

Automatic Differentiation

Automatic differentiation techniques are founded on the observation that any function is evaluated by performing a sequence of simple elementary operations involving just one or two arguments at a time:

- addition
- multiplication
- division
- power operation a^b
- trigonometric functions
- exponential functions
- logarithmic
- chain rule:

$$\frac{d}{dx} f(g(x)) = \frac{d}{dx} f \circ g(x) = \frac{df}{dg} \frac{dg}{dx}$$

- Forward Accumulation is equivalent to expanding a function using the chain rule and computing the derivatives inside-out
- Requires n -passes to compute n -dimensional gradient
- Example:

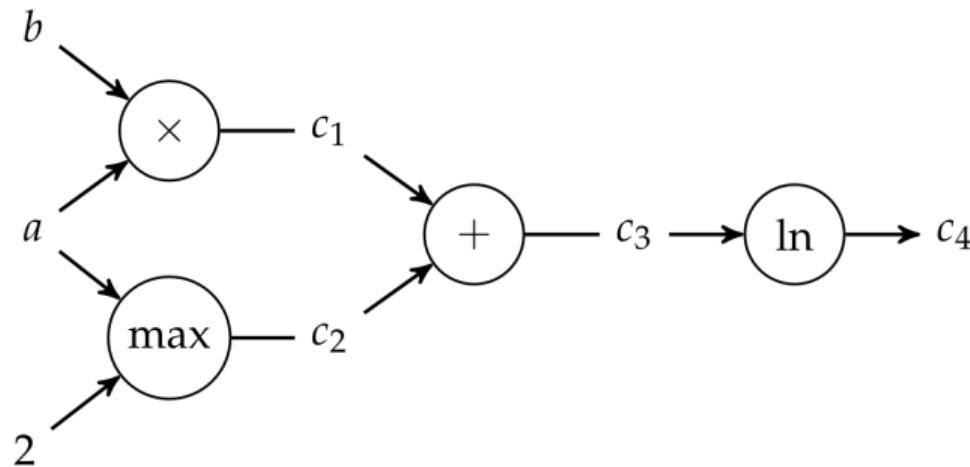
$$f(a, b) = \ln(ab + \max(a, 2))$$

$$\begin{aligned}\frac{\partial f}{\partial a} &= \frac{\partial}{\partial a} \ln(ab + \max(a, 2)) \\ &= \frac{1}{ab + \max(a, 2)} \frac{\partial}{\partial a} (ab + \max(a, 2)) \\ &= \frac{1}{ab + \max(a, 2)} \left[\frac{\partial(ab)}{\partial a} + \frac{\partial \max(a, 2)}{\partial a} \right] \\ &= \frac{1}{ab + \max(a, 2)} \left[\left(b \frac{\partial a}{\partial a} + a \frac{\partial b}{\partial a} \right) + \left((2 > a) \frac{\partial 2}{\partial a} + (2 < a) \frac{\partial a}{\partial a} \right) \right] \\ &= \frac{1}{ab + \max(a, 2)} [b + (2 < a)]\end{aligned}$$

Automatic Differentiation

Computational graph: nodes are operations and the edges are input-output relations. leaf nodes of a computational graph are input variables or constants, and terminal nodes are values output by the function

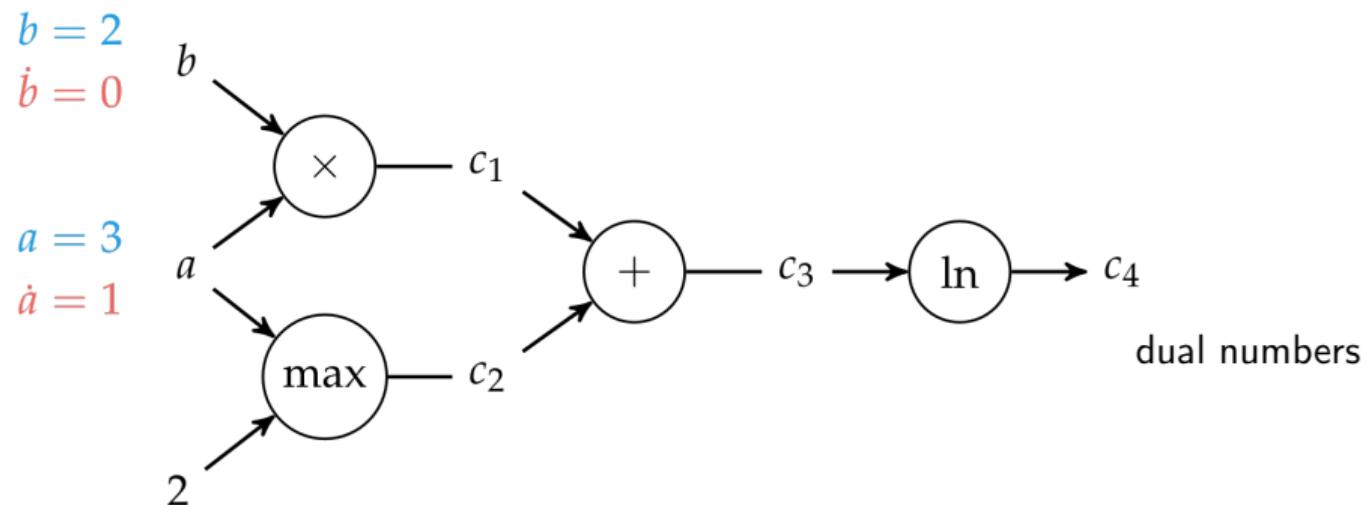
Forward accumulation for $f(a, b) = \ln(ab + \max(a, 2))$



Automatic Differentiation

Computational graph: nodes are operations and the edges are input-output relations. leaf nodes of a computational graph are input variables or constants, and terminal nodes are values output by the function

Forward accumulation for $f(a, b) = \ln(ab + \max(a, 2))$

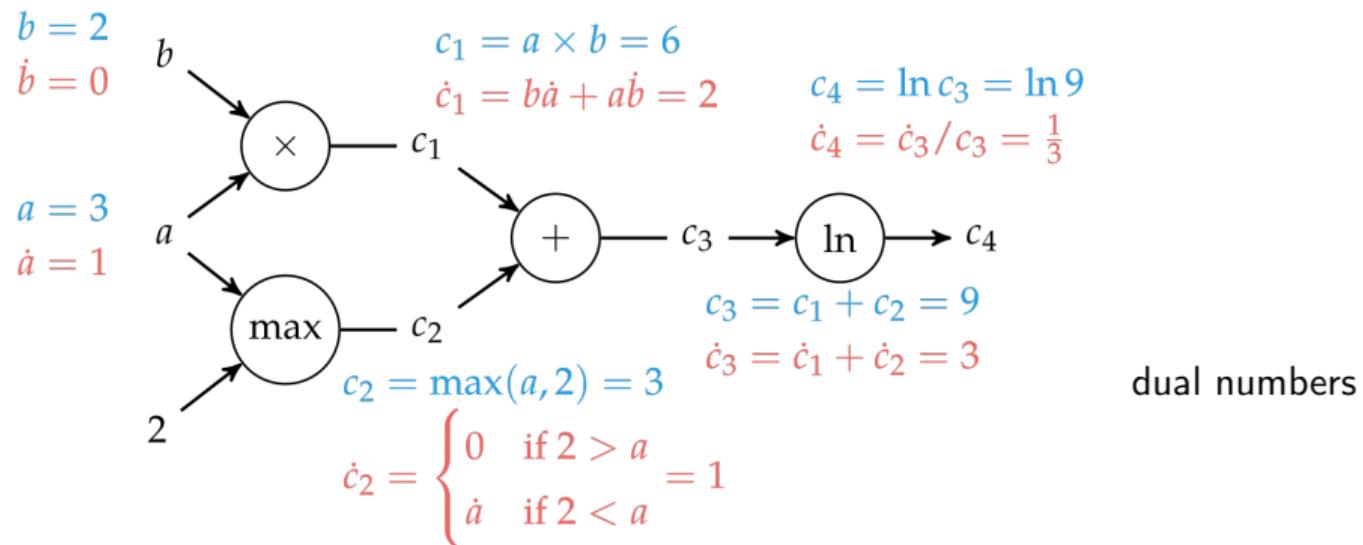


$$\frac{\partial b}{\partial a} = \dot{b} \text{ Newton notation}$$

Automatic Differentiation

Computational graph: nodes are operations and the edges are input-output relations. leaf nodes of a computational graph are input variables or constants, and terminal nodes are values output by the function

Forward accumulation for $f(a, b) = \ln(ab + \max(a, 2))$

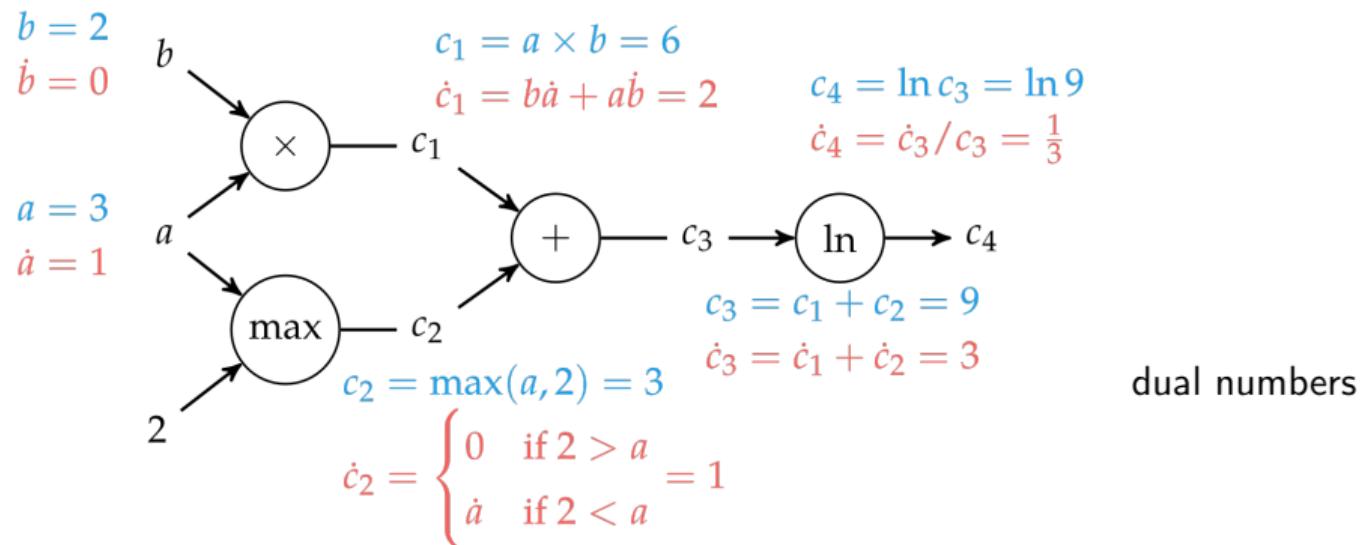


$$\frac{\partial b}{\partial a} = \dot{b} \text{ Newton notation}$$

Automatic Differentiation

Computational graph: nodes are operations and the edges are input-output relations. leaf nodes of a computational graph are input variables or constants, and terminal nodes are values output by the function

Forward accumulation for $f(a, b) = \ln(ab + \max(a, 2))$



$$\frac{\partial b}{\partial a} = \dot{b} \text{ Newton notation}$$

$$\text{for } \frac{\partial f}{\partial b} \text{ set } \dot{a} = 0, \dot{b} = 1$$

Dual numbers

- Dual numbers can be expressed mathematically by including the abstract quantity ϵ , where ϵ^2 is defined to be 0.
- Like a complex number, a dual number is written $a + b\epsilon$ where a and b are both real values.
- $(a + b\epsilon) + (c + d\epsilon) = (a + c) + (b + d)\epsilon$
 $(a + b\epsilon) \times (c + d\epsilon) = (ac) + (ad + bc)\epsilon$
- by passing a dual number into any smooth function f , we get the evaluation and its derivative. We can show this using the Taylor series:

$$\begin{aligned}f(x) &= \sum_{k=0}^{\infty} \frac{f^{(k)}(a)}{k!} (x - a)^k &= f(a) + bf'(a)\epsilon + \epsilon^2 \sum_{k=2}^{\infty} \frac{f^{(k)}(a)b^k}{k!} \epsilon^{(k-2)} \\f(a + b\epsilon) &= \sum_{k=0}^{\infty} \frac{f^{(k)}(a)}{k!} (a + b\epsilon - a)^k &= f(a) + bf'(a)\epsilon \\&= \sum_{k=0}^{\infty} \frac{f^{(k)}(a)b^k \epsilon^k}{k!} &\end{aligned}$$

Note that

$$\begin{aligned}(v + \dot{v}\epsilon) + (u + \dot{u}\epsilon) &= (v + u) + (\dot{v} + \dot{u})\epsilon \\(v + \dot{v}\epsilon)(u + \dot{u}\epsilon) &= (vu) + (v\dot{u} + \dot{v}u)\epsilon,\end{aligned}$$

satisfies the rules of differentiation

Setting:

$$f(v + \dot{v}\epsilon) = f(v) + f'(v)\dot{v}\epsilon$$

The chain rule follows:

$$\begin{aligned}f(g(v + \dot{v}\epsilon)) &= f(g(v) + g'(v)\dot{v}\epsilon) \\&= f(g(v)) + f'(g(v))g'(v)\dot{v}\epsilon.\end{aligned}$$

Automatic Differentiation

- **Reverse accumulation** is performed in a single run using two passes $O(m \cdot \text{ops}(f))$ (forward and back) for $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Note: this is central to the backpropagation algorithm used to train neural networks because it needs only one pass for the n -dimensional function to find the gradient.
- implemented through two different operation overloading functions (for forward and backward)
- Many open-source software implementations are available: eg, Tensorflow

Forward implements:

$$\frac{df}{dx} = \frac{df}{dc_4} \frac{dc_4}{dx} = \frac{df}{dc_4} \left(\frac{dc_4}{dc_3} \frac{dc_3}{dx} \right) = \frac{df}{dc_4} \left(\frac{dc_4}{dc_3} \left(\frac{dc_3}{dc_2} \frac{dc_2}{dx} + \frac{dc_3}{dc_1} \frac{dc_1}{dx} \right) \right)$$

Backward implements:

$$\frac{df}{dx} = \frac{df}{dc_4} \frac{dc_4}{dx} = \left(\frac{df}{dc_3} \frac{dc_3}{dc_4} \right) \frac{dc_4}{dx} = \left(\left(\frac{df}{dc_2} \frac{dc_2}{dc_3} + \frac{df}{dc_1} \frac{dc_1}{dc_3} \right) \frac{dc_3}{dc_4} \right) \frac{dc_4}{dx}$$

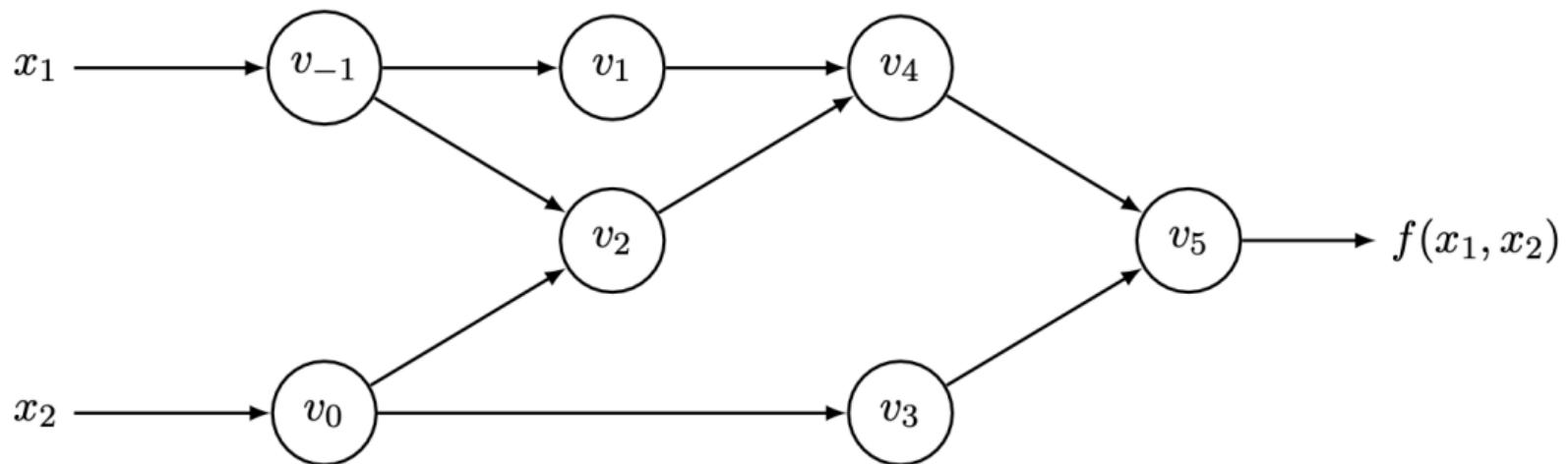
Complementing each intermediate variable v_i with an **adjoint**

$$\bar{v}_i = \frac{\partial y_j}{\partial v_i}$$

which represents the sensitivity of a considered output y_j with respect to changes in v_i .

Example

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$



Example: Forward Accumulation

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Forward Primal Trace

$v_{-1} = x_1$	= 2
$v_0 = x_2$	= 5
<hr/>	
$v_1 = \ln v_{-1}$	= $\ln 2$
$v_2 = v_{-1} \times v_0$	= 2×5
$v_3 = \sin v_0$	= $\sin 5$
$v_4 = v_1 + v_2$	= $0.693 + 10$
$v_5 = v_4 - v_3$	= $10.693 + 0.959$
<hr/>	
$y = v_5$	= 11.652

Forward Tangent (Derivative) Trace

$\dot{v}_{-1} = \dot{x}_1$	= 1
$\dot{v}_0 = \dot{x}_2$	= 0
<hr/>	
$\dot{v}_1 = \dot{v}_{-1}/v_{-1}$	= $1/2$
$\dot{v}_2 = \dot{v}_{-1} \times v_0 + \dot{v}_0 \times v_{-1}$	= $1 \times 5 + 0 \times 2$
$\dot{v}_3 = \dot{v}_0 \times \cos v_0$	= $0 \times \cos 5$
$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$	= $0.5 + 5$
$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$	= $5.5 - 0$
<hr/>	
$\dot{y} = \dot{v}_5$	= 5.5

$O(n \cdot \text{ops}(f))$

Example: Reverse Accumulation

Forward Primal Trace

$$v_{-1} = x_1 = 2$$

$$v_0 = x_2 = 5$$

$$v_1 = \ln v_{-1} = \ln 2$$

$$v_2 = v_{-1} \times v_0 = 2 \times 5$$

$$v_3 = \sin v_0 = \sin 5$$

$$v_4 = v_1 + v_2 = 0.693 + 10$$

$$v_5 = v_4 - v_3 = 10.693 + 0.959$$

$$y = v_5 = 11.652$$

Reverse Adjoint (Derivative) Trace

$$\bar{x}_1 = \bar{v}_{-1} = 5.5$$

$$\bar{x}_2 = \bar{v}_0 = 1.716$$

$$\bar{v}_{-1} = \bar{v}_{-1} + \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} = \bar{v}_{-1} + \bar{v}_1 / v_{-1} = 5.5$$

$$\bar{v}_0 = \bar{v}_0 + \bar{v}_2 \frac{\partial v_2}{\partial v_0} = \bar{v}_0 + \bar{v}_2 \times v_{-1} = 1.716$$

$$\bar{v}_{-1} = \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_2 \times v_0 = 5$$

$$\bar{v}_0 = \bar{v}_3 \frac{\partial v_3}{\partial v_0} = \bar{v}_3 \times \cos v_0 = -0.284$$

$$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$$

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$$

$$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$$

$$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$$

$$\bar{v}_5 = \bar{y} = 1$$

$O(m \cdot \text{ops}(f))$

Summary

- Derivatives are useful in optimization because they provide information about how to change a given point in order to improve the objective function
- For multivariate functions, various derivative-based concepts are useful for directing the search for an optimum, including the gradient, the Hessian, and the directional derivative
- computation of derivatives in computer programs can be classified into four categories:
 1. manually working out derivatives and coding them (error prone and time consuming)
 2. numerical differentiation using finite difference approximations
Complex step method can eliminate the effect of subtractive cancellation error when taking small steps
 3. symbolic differentiation using expression manipulation in computer algebra systems
 4. automatic differentiation, (aka algorithmic differentiation)
forward and reverse accumulation on computational graphs

3. Bracketing

Bracketing

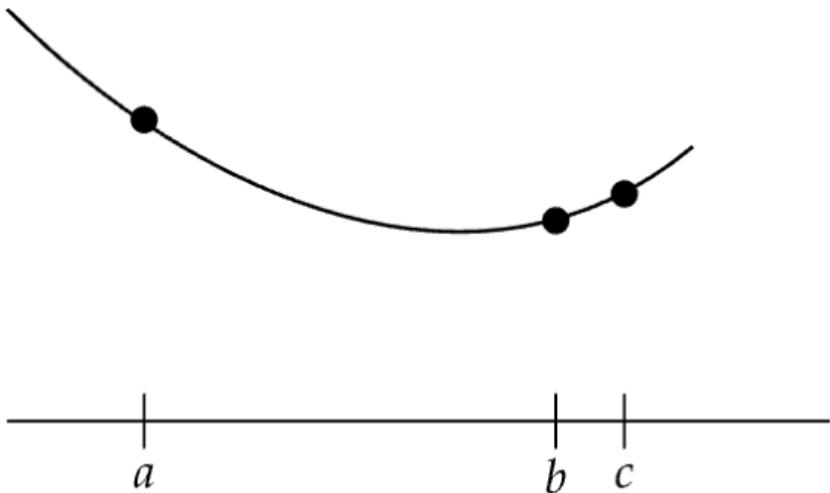
A derivative-free method to identify an interval containing a local minimum and then successively shrinking that interval

Unimodality

There exists a unique optimizer x^* such that f is monotonically decreasing for $x \leq x^*$ and monotonically increasing for $x \geq x^*$

Finding an Initial Bracket

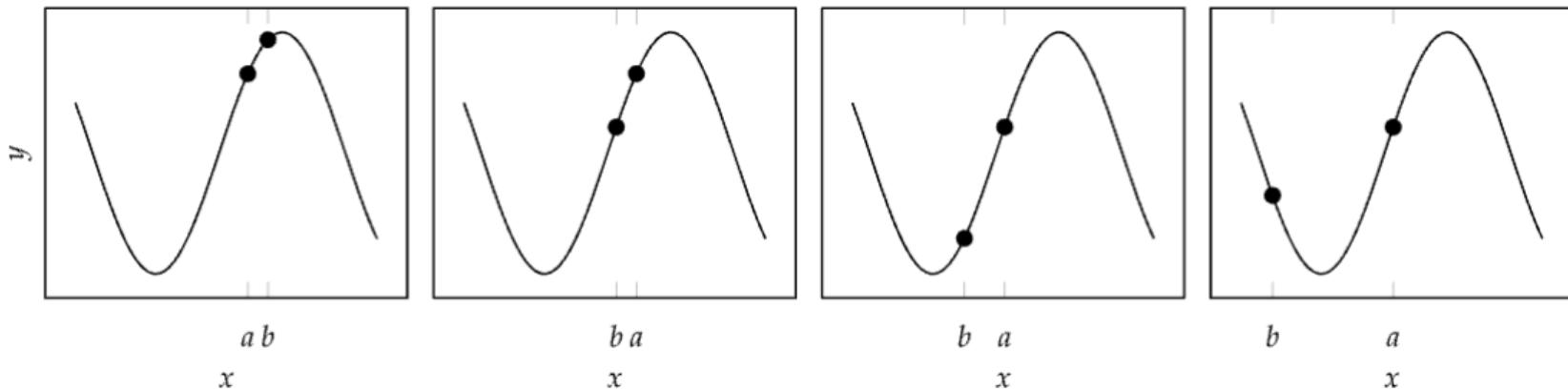
Given a unimodal function, the global minimum is guaranteed to be inside the interval $[a, c]$ if $f(a) > f(b) < f(c)$



```
function bracket_minimum(f, x=0; s=1e-2, k=2.0)
    a, ya = x, f(x)
    b, yb = a + s, f(a + s)
    if yb > ya
        a, b = b, a
        ya, yb = yb, ya
        s = -s
    end
    while true
        c, yc = b + s, f(b + s)
        if yc > yb
            return a < c ? (a, c) : (c, a)
        end
        a, ya, b, yb = b, yb, c, yc
        s *= k
    end
end
```

Finding an Initial Bracket

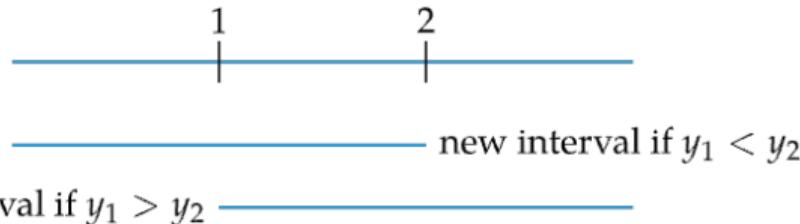
Example of `bracket_minimum` on a function



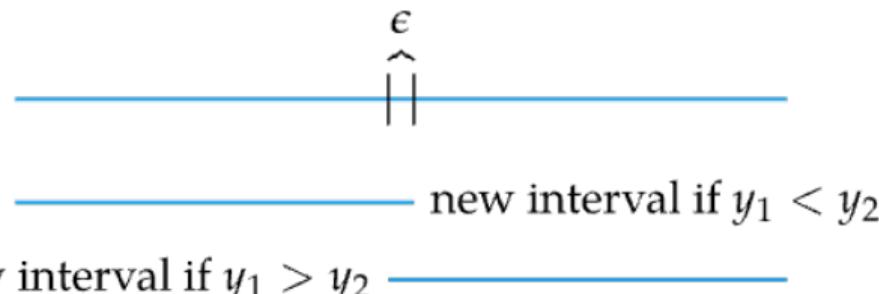
reverses direction between the first and second iteration and expands until a minimum is bracketed in the fourth iteration.

For unimodal functions, when function evaluations are limited, what is the maximal shrinkage we can achieve?

When restricted to only 2 function evaluations (queries) the most we can guarantee to shrink our interval is by just under a factor of 2.

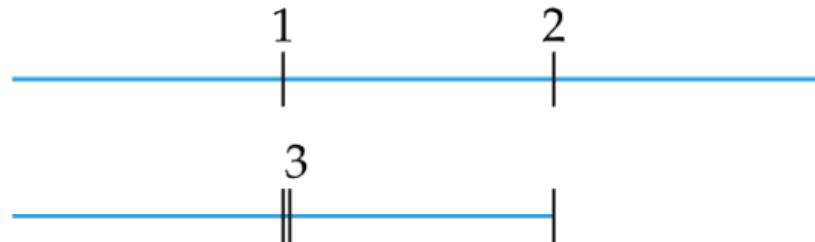


yields a factor of 3.



for $\epsilon \rightarrow 0$ yields a factor of just less than 2

When restricted to only 3 function evaluations (queries) the most we can guarantee to shrink our interval is by a factor of 3.



Fibonacci Search

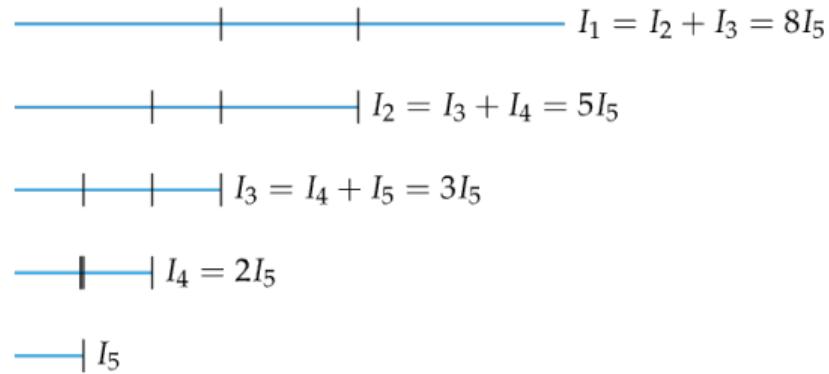
When restricted to n function evaluations following the previous strategy, we are guaranteed to shrink our interval by a factor of F_{n+1} .

Fibonacci numbers: sum of previous two,
1, 1, 2, 3, 5, 8, 13, ...

$$F_n = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1, 2 \\ F_{n-1} + F_{n-2} & \text{otherwise} \end{cases}$$

The length of every interval constructed can be expressed in terms of the final interval times a Fibonacci number.

- final, smallest interval has length I_n ,
- second smallest interval has length $I_{n-1} = F_3 I_n$
- third smallest interval has length $I_{n-2} = F_4 I_n$,
and so forth.



Fibonacci Search Algorithm

For a unimodal function f in the interval $[a, b]$, we want to shrink the interval within n iterations.
(At each iteration we want to shrink by a factor ϕ).

$$b_{k+1} - a_{k+1} = \frac{F_{n-k+1}}{F_{n-k+2}}(b_k - a_k)$$

Closed-form expression (Binet's formula):

$$F_n = \frac{\phi^n - (1-\phi)^n}{\sqrt{5}},$$

Therefore:

$$b_n - a_n = \frac{F_2}{F_3}(b_{n-1} - a_{n-1})$$

$\phi = (1 + \sqrt{5})/2 \approx 1.61803$ is the golden ratio.

$$= \frac{F_2}{F_3} \frac{F_3}{F_4} \cdots \frac{F_n}{F_{n+1}}(b_1 - a_1)$$

$$\frac{F_{n+1}}{F_n} = \phi \frac{1 - s^{n+1}}{1 - s^n}, \quad s = (1 - \sqrt{5})(1 + \sqrt{5}) \approx -0.3827$$

$$= \frac{1}{F_{n+1}}(b_1 - a_1)$$

Suppose we have an unimodal function f in the interval $[a, b]$ and a tolerance $\epsilon = 0.01$. Let $k = 1$.

1. $d_k = a_k + \frac{F_{n-k+1}}{F_{n-k+2}}(b_k - a_k)$

$$\frac{F_n}{F_{n+1}} = \rho_n = \frac{1 - s^n}{\phi(1 - s^{n+1})} \approx 0.6$$

2. if $k \neq n - 1$:

$$c_k = a_k + \left(1 - \frac{F_{n-k+1}}{F_{n-k+2}}\right)(b_k - a_k)$$

otherwise: $c_k = d_k + \epsilon(a_k - d_k)$

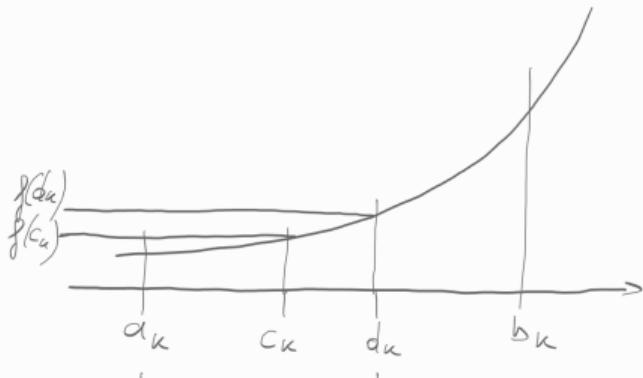
3. if $f(c_k) < f(d_k)$: $b_{k+1} = d_k, d_{k+1} = c_k, a_{k+1} = a_k$

otherwise: $a_{k+1} = b_k, b_{k+1} = c_k, d_{k+1} = d_k$

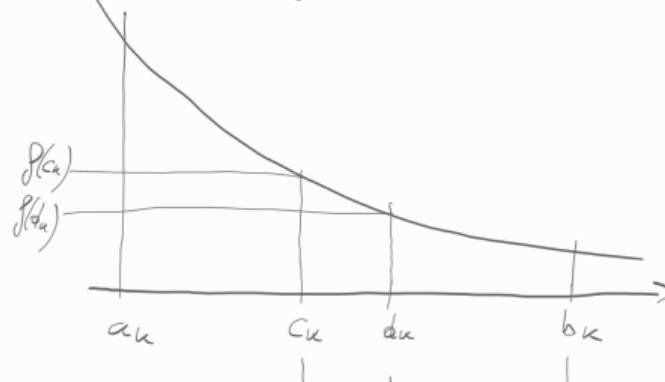
4. $k = k + 1$, if $k = n$ go to step 5, else go to step 2

5. return (a_k, b_k) if $(a_k < b_k)$ else (b_k, a_k)

$$f(c_k) < f(d_k)$$



$$f(c_k) > f(d_k)$$



$$k \neq m-1$$

$$a_{k+1} + \underbrace{(1-\varrho)}_{\approx 0,4} \underbrace{(a_{k+1} - a_{k+1})}_{< 0}$$

$$d_{k+1} + \varepsilon (d_{k+1} - d_{k+1})$$

$$k = m-1$$

$$b_{k+1} - \underbrace{\varrho}_{\approx 0,4} \underbrace{(b_{k+1} - a_{k+1})}_{< 0}$$

$$d_{k+1} + \varepsilon (a_{k+1} - d_{k+1})$$

Golden Section Search

$$\lim_{n \rightarrow \infty} \frac{F_{n+1}}{F_n} = \lim_{n \rightarrow \infty} \frac{1}{\rho_n} = \lim_{n \rightarrow \infty} \phi \frac{1 - s^{n+1}}{1 - s^n} = \phi \approx 1.61803 \quad \frac{1}{\phi} \approx 0.618$$



A horizontal blue line segment representing the search interval I_2 , which is $I_1 \varphi^{-1}$. It is divided into three segments by two vertical tick marks. The first and third segments are of equal length, while the middle segment is half their length.

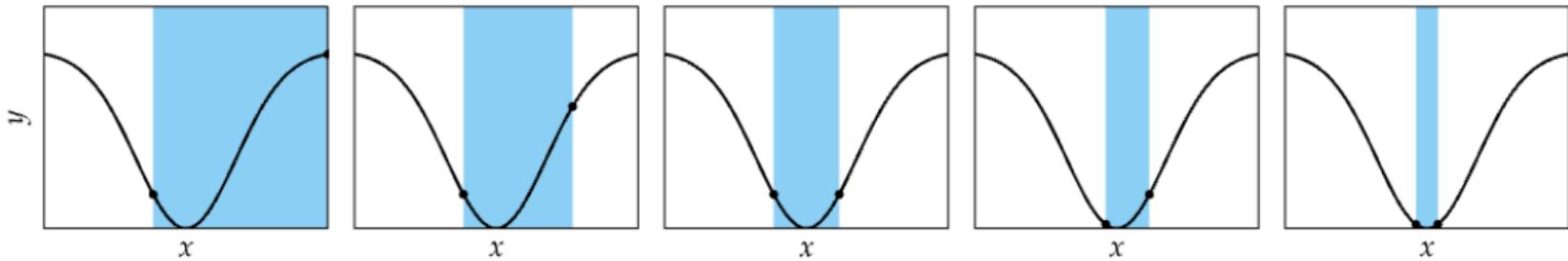
A horizontal blue line segment representing the search interval I_3 , which is $I_1 \varphi^{-2}$. It is divided into five segments by three vertical tick marks. The first and fifth segments are of equal length, while the middle three segments are each one-third their length.

A horizontal blue line segment representing the search interval I_4 , which is $I_1 \varphi^{-3}$. It is divided into seven segments by five vertical tick marks. The first and seventh segments are of equal length, while the middle five segments are each one-fifth their length.

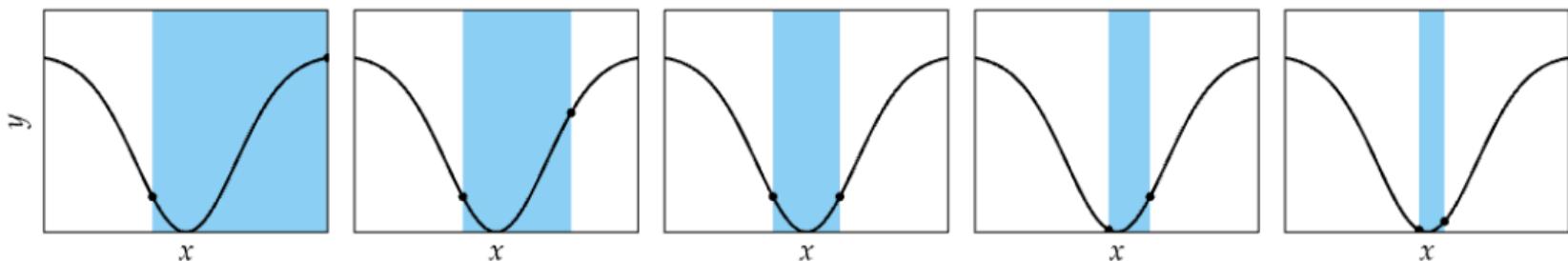
A horizontal blue line segment representing the search interval I_5 , which is $I_1 \varphi^{-4}$. It is divided into nine segments by seven vertical tick marks. The first and ninth segments are of equal length, while the middle seven segments are each one-seventh their length.

Comparison

Fibonacci Search

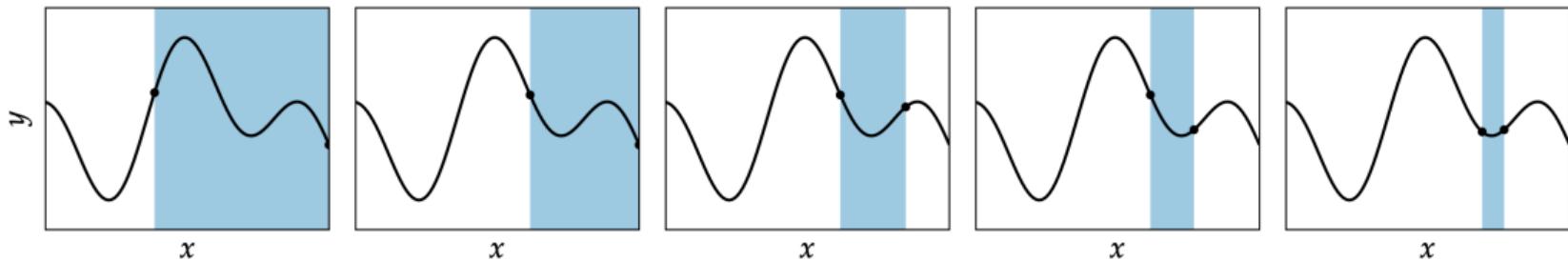


Golden Section Search

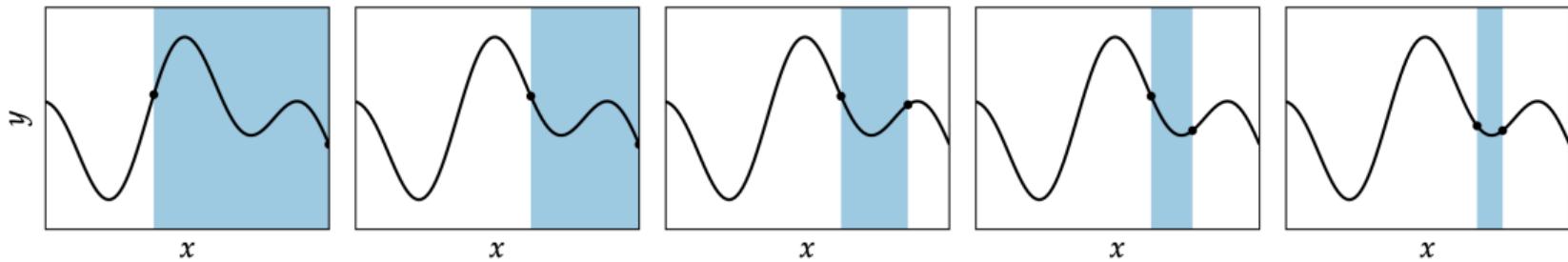


Comparison

Fibonacci Search

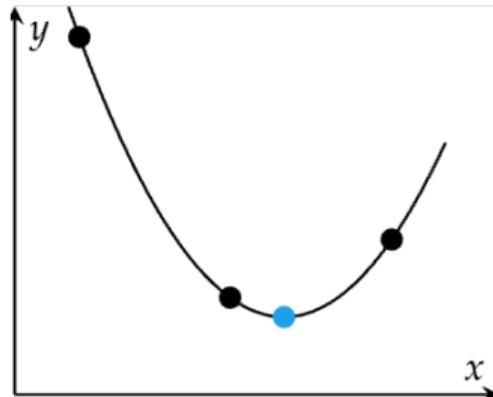


Golden Section Search



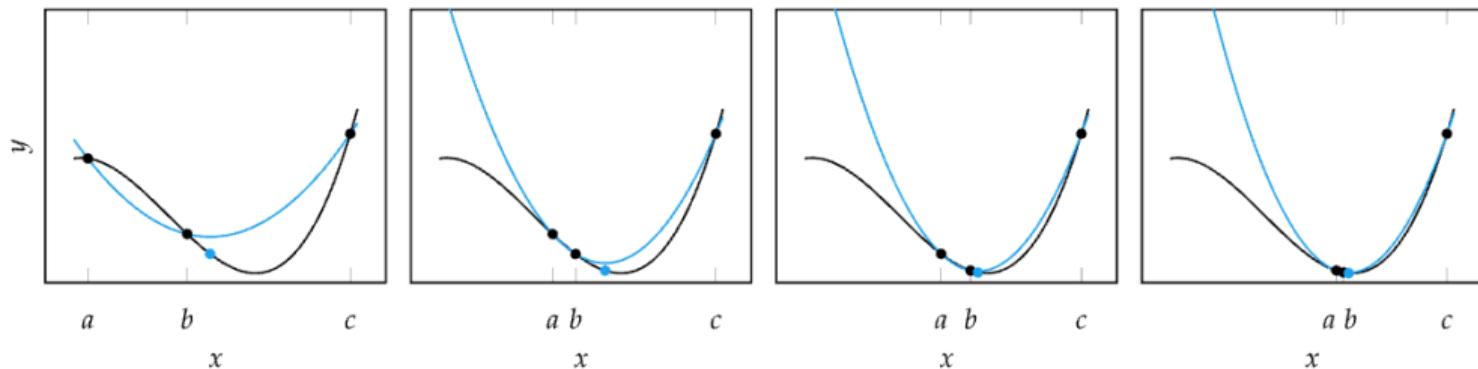
Quadratic Fit Search

- Leverages ability to analytically minimize quadratic functions
- Iteratively fits quadratic function to three bracketing points



Quadratic Fit Search

- If a function is locally nearly quadratic, the minimum can be found after several steps



Using Linear Algebra

- We assume that the variable y is related to $x \in \mathbb{R}^n$ quadratically, so for some constants b_0, b_1, b_2 :

$$y = b_0 + b_1 x + b_2 x^2$$

- Given the set of m points $(y_1, x_1), \dots, (y_3, x_3)$ in the ideal case, we have that $y_i = b_0 + b_1 x_i + b_2 x_i^2$, for all $i = 1, 2, 3$. In matrix form:

$$\begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

This can be written as $Az = y$ to emphasize that z are our unknowns and A and y are given.

In Python

In polynomial regression, the $m \times (n + 1)$ matrix A is called a **Vandermonde matrix** (a matrix with entries $a_{ij} = x_i^{n+1-j}$, $j = 1..n + 1$).

NumPy's `np.vander()` is a convenient tool for quickly constructing a Vandermonde matrix, given the values x_i , $i = 1..m$, and the number of desired columns ($n + 1$).

```
>>> print(np.vander([2, 3, 5], 2))
[[2 1]                                # [[2**1, 2**0]
 [3 1]                                # [3**1, 3**0]
 [5 1]]                               # [5**1, 5**0]]
```



```
>>> print(np.vander([2, 3, 5, 4], 3))
[[ 4  2  1]                            # [[2**2, 2**1, 2**0]
 [ 9  3  1]                            # [3**2, 3**1, 3**0]
 [25  5  1]                            # [5**2, 5**1, 5**0]
 [16  4  1]]                           # [4**2, 4**1, 4**0]]
```

In Python

```
A = np.vander(x,4)

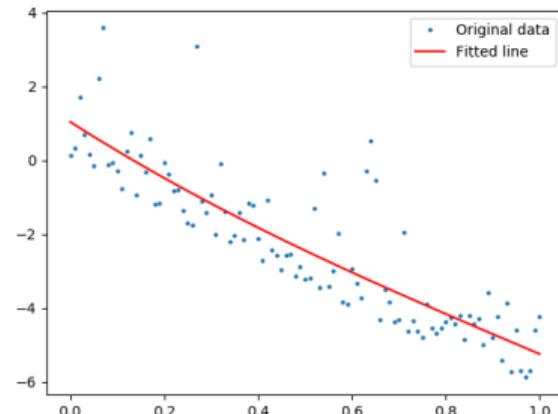
coeff = np.linalg.solve(A,y) ## Error!! Why?

B = A.T @ A
z = np.linalg.inv(B) @ A.T @ y

coeff = np.linalg.lstsq(A, y)[0]
np.allclose(z,coeff)

f=np.poly1d(coeff)
plt.plot(x, y, 'o', label='Original data', ↪
          ↪markersize=2)
plt.plot(x, f(x), 'r', label='Fitted line')
plt.legend()
plt.show()
```

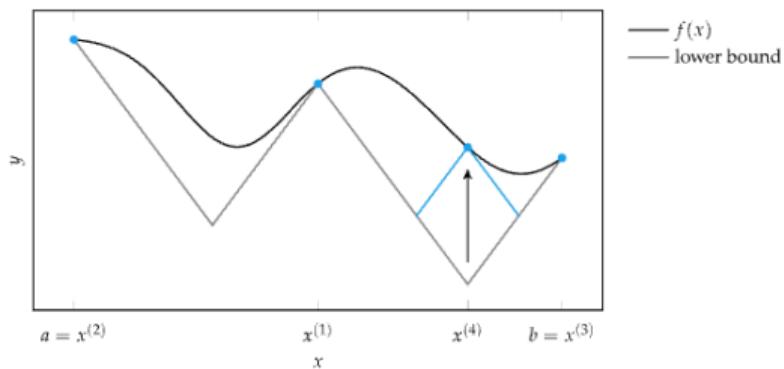
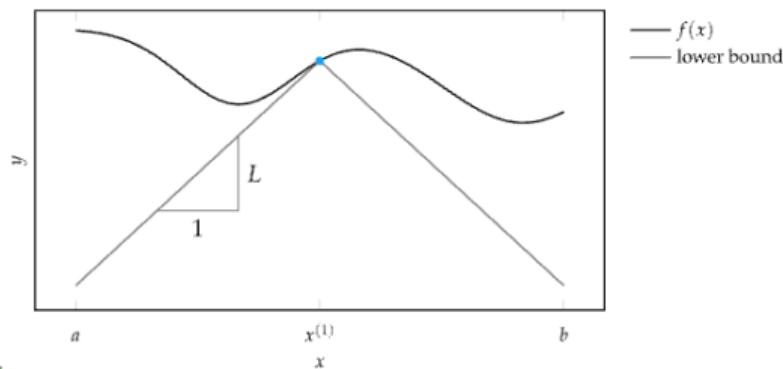
ex2.py

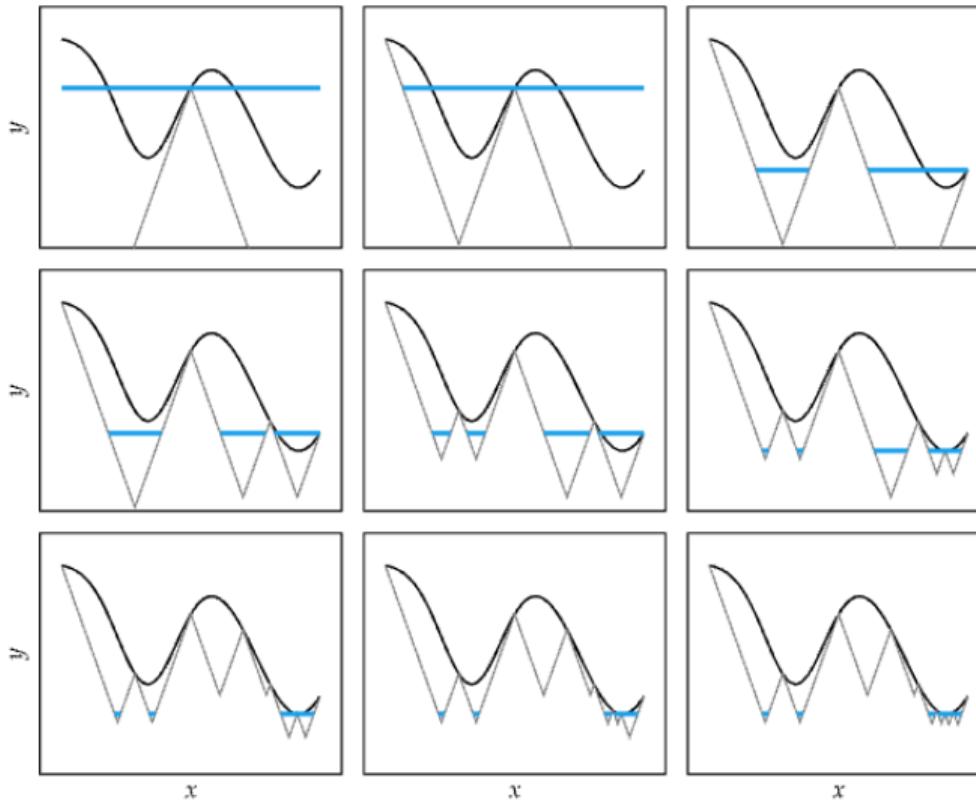


Shubert-Piyavskii Method

- The Shubert-Piyavskii method is guaranteed to find the global minimum of any bounded function
- but requires that the function be Lipschitz continuous
- A function is **Lipschitz continuous** if there is an upper bound on the magnitude of its derivative. A function f is Lipschitz continuous on $[a, b]$ if there exists an $\ell > 0$ such that:

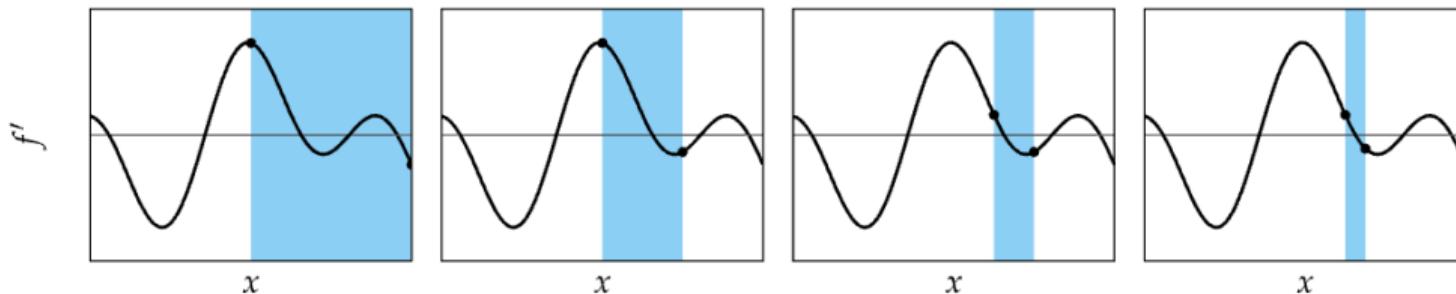
$$|f(x) - f(y)| \leq \ell|x - y|, \quad \forall x, y \in [a, b]$$





Bisection Method

- **Intermediate value theorem:** If f is continuous on $[a, b]$, and there is some $y \in [f(a), f(b)]$, then there exists at least one $x \in [a, b]$, such that $f(x) = y$.
- Used in root-finding methods
- When applied to $f'(x)$, can be used to find minimum of f
- if $\text{sign}(f'(a)) \neq \text{sign}(f'(b))$, or equivalently, $f'(a)f'(b) \leq 0$ then $[a, b]$ is guaranteed to contain a zero.



Bisection method

- Cut the bracketed region $[a, b]$ in half with every iteration
- Evaluate the midpoint $(a + b)/2$
- form a new bracket from the midpoint and whichever side that continues to bracket a zero.
- Terminate after a fixed number of iterations.
- Guaranteed to converge within ϵ of x^* within $\lg_2(|b - a|/\epsilon)$

Summary

- Many optimization methods shrink a bracketing interval, including Fibonacci search, golden section search, and quadratic fit search
- The Shubert-Piyavskii method outputs a set of bracketed intervals containing the global minima, given the Lipschitz constant
- Root-finding methods like the bisection method can be used to find where the derivative of a function is zero

4. Local Descent

Preface

For multivariate functions, we have argued that:

- derivatives can have exponential growth in the resulting analytical expression
- calculating zeros might be challenging

Hence, minimizing by solving $\nabla f(\mathbf{x}) = \mathbf{0}$ may be computationally demanding.

Descent Direction Iteration

Descent Direction Methods use a local model to incrementally improve design point until some convergence criteria is met

1. Check termination conditions at \mathbf{x}_k ; if not met, continue.
2. Decide **descent direction** \mathbf{d}_k using local information
3. Decide **step size** (= magnitude of the overall step = α_k , since commonly $\|\mathbf{d}_k\|_2 = 1$)
4. Compute next design point \mathbf{x}_{k+1}

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

Line Search for Step Size

Assuming we have the search direction:

- Used to compute α
- Using the techniques discussed from previous classes, solve:

$$\text{minimize}_{\alpha} f(\mathbf{x} + \alpha \mathbf{d})$$

- Often this is computed approximately to reduce cost

Line Search: Alternatives

Step size:

- Fixed α called **learning rate** (commonly $\|\mathbf{d}_k\|_2 = 1$ not imposed)
- **Decaying step factor**

$$\alpha_k = \alpha_1 \gamma^{k-1} \quad \text{for } \gamma \in [0, 1]$$

Decaying step factor is often required in convergence proofs

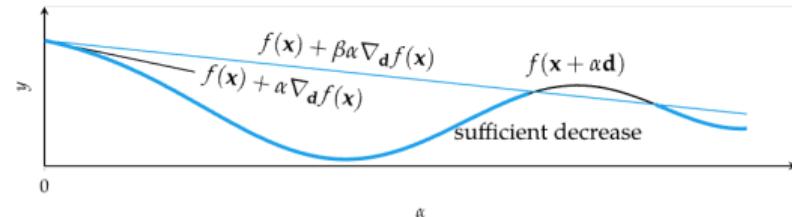
Approximate Line Search

- If function calls are expensive, rather than finding the minimum along a search direction, find a point of sufficient decrease

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \beta\alpha \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$

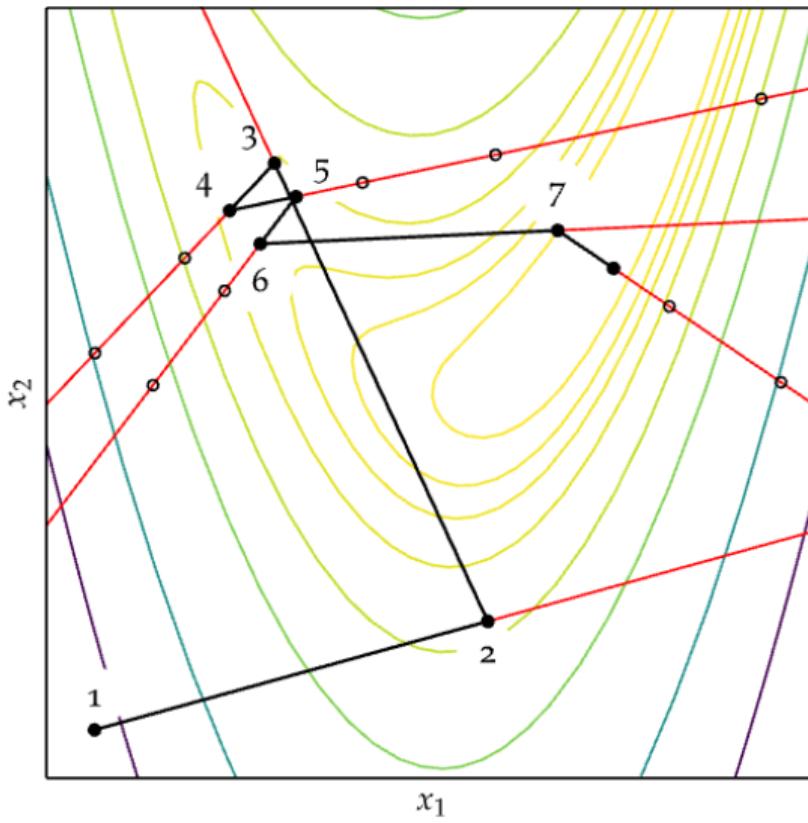
- $\beta \in [0, 1]$, usually $\beta = 1 \times 10^{-4}$

- Backtracking line search starts with a large step and then backs off



```
def backtracking_line_search(f, grad, x, d, alpha_0=1, p=0.5, beta=1e-4):
    y, g, alpha = f(x), grad(x), alpha_0
    while ( f(x + alpha * d) > y + beta * alpha * np.dot(g, d) ) :
        alpha *= p
    return alpha
```

Approximate Line Search: Example



Approximate Line Search

Building on backtracking line search are the Wolfe Conditions each sufficient to guarantee convergence to a local minimum.

1. First Wolfe Condition: Sufficient Decrease

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \beta \alpha \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$

2. Second Wolfe Condition: Curvature Condition

$$\nabla_{\mathbf{d}_k} f(\mathbf{x}_{k+1}) \geq \sigma \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$

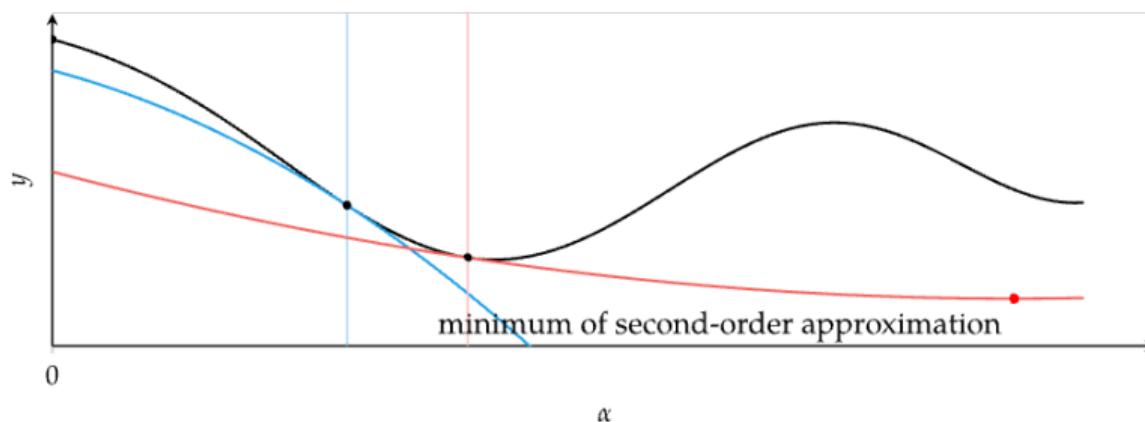
$\beta < \sigma < 1$ with

- $\sigma = 0.1$ with conjugate gradient method
- $\sigma = 0.9$ with Newton method

Approximate Line Search

The curvature condition ensures the second-order function approximations have positive curvature

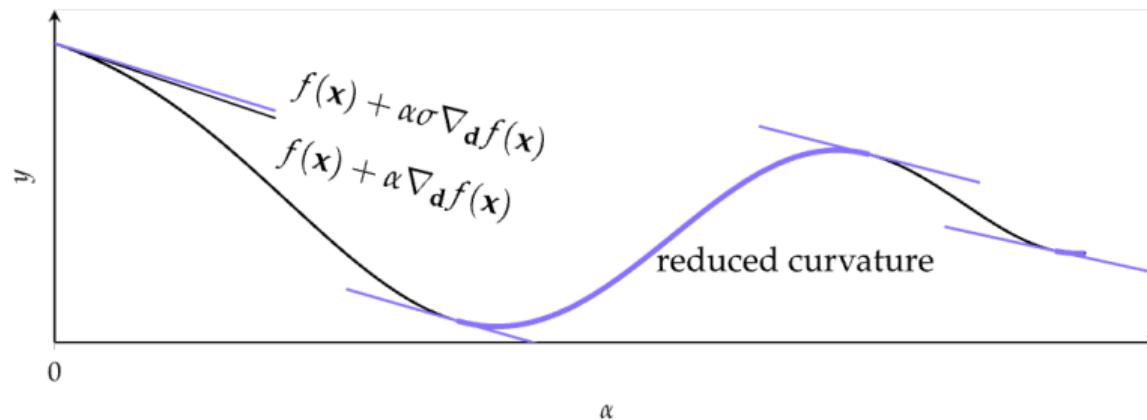
$$\nabla_{\mathbf{d}_k} f(\mathbf{x}_{k+1}) \geq \sigma \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$



Approximate Line Search

Regions satisfying the curvature condition

$$\nabla_{\mathbf{d}_k} f(\mathbf{x}_{k+1}) \geq \sigma \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$



Approximate Line Search: Example

Consider approximate line search on $f(x_1, x_2) = x_1^2 + x_1 x_2 + x_2^2$ from $\mathbf{x} = [1, 2]$ in the direction $\mathbf{d} = [-1, -1]$, gradient at \mathbf{x} is $\mathbf{g} = [4, 5]$ using a maximum step size of 10, a reduction factor of 0.5, first Wolfe condition parameter $\beta = 1 \times 10^{-4}$, second Wolfe condition parameter $\sigma = 0.9$.

first Wolfe condition ($f(\mathbf{x} + \alpha \mathbf{d}) \leq f(\mathbf{x}) + \beta \alpha (\mathbf{g}^T \cdot \mathbf{d})$):

$$\alpha = 10 : f([1, 2] + 10 \cdot [-1, -1]) \leq 7 + 1 \times 10^{-4} 10 [4, 5]^T [-1, -1] \implies 217 \not\leq 6.991$$

$$\alpha = 10 \cdot 0.5 = 5 : f([1, 2] + 5 \cdot [-1, -1]) \leq 7 + 1 \times 10^{-4} 5 [4, 5]^T [-1, -1] \implies 37 \not\leq 6.996$$

$$\alpha = 2.5 : f([1, 2] + 2.5 \cdot [-1, -1]) \leq 7 + 1 \times 10^{-4} 2.5 [4, 5]^T [-1, -1] \implies 3.25 \leq 6.998$$

The candidate design point $\mathbf{x}' = \mathbf{x} + \alpha \mathbf{d} = [-1.5, -0.5]$ is checked against the second Wolfe condition $\nabla_{\mathbf{d}} f(\mathbf{x}') \geq \sigma \nabla_{\mathbf{d}} f(\mathbf{x})$:

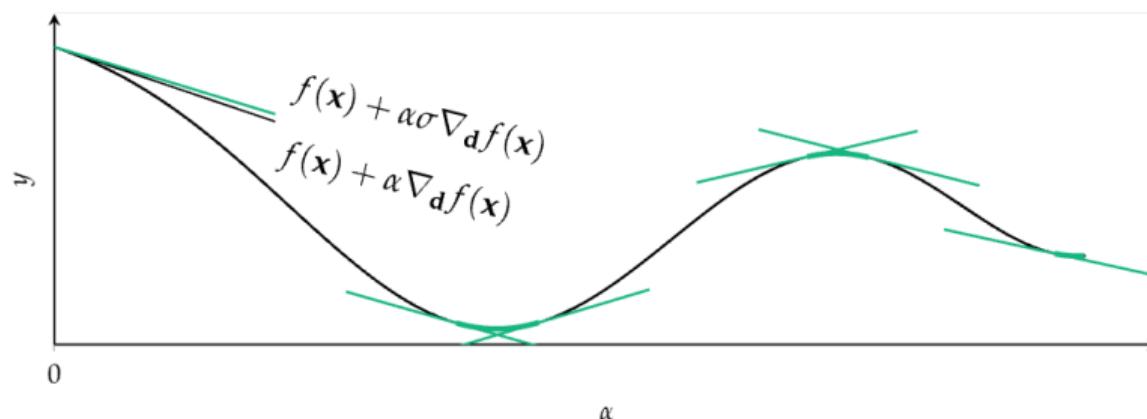
$$[-3.5, -2.5] \cdot [-1, -1] \geq 0.9 [4, 5] \cdot [-1, -1] \implies 6 \geq -8.1$$

Approximate line search terminates with $\mathbf{x} = [-1.5, -0.5]$.

Approximate Line Search

Regions where the strong curvature condition is satisfied

$$|\nabla_{\mathbf{d}_k} f(\mathbf{x}_{k+1})| \leq -\sigma \nabla_{\mathbf{d}_k} f(\mathbf{x}_k)$$



Approximate Line Search

- The sufficient decrease condition with the strong curvature condition form the strong Wolfe conditions.
- Satisfying the strong Wolfe conditions requires a more complicated algorithm

Strong backtracking line search:

1. Bracketing Phase: tests successively larger step sizes to bracket an interval $[\alpha_{k-1}, \alpha_k]$ guaranteed to contain step lengths satisfying the Wolfe conditions.
2. Zoom Phase: shrink the interval using bisection to find point satisfying the **strong** Wolfe conditions

Approximate Line Search

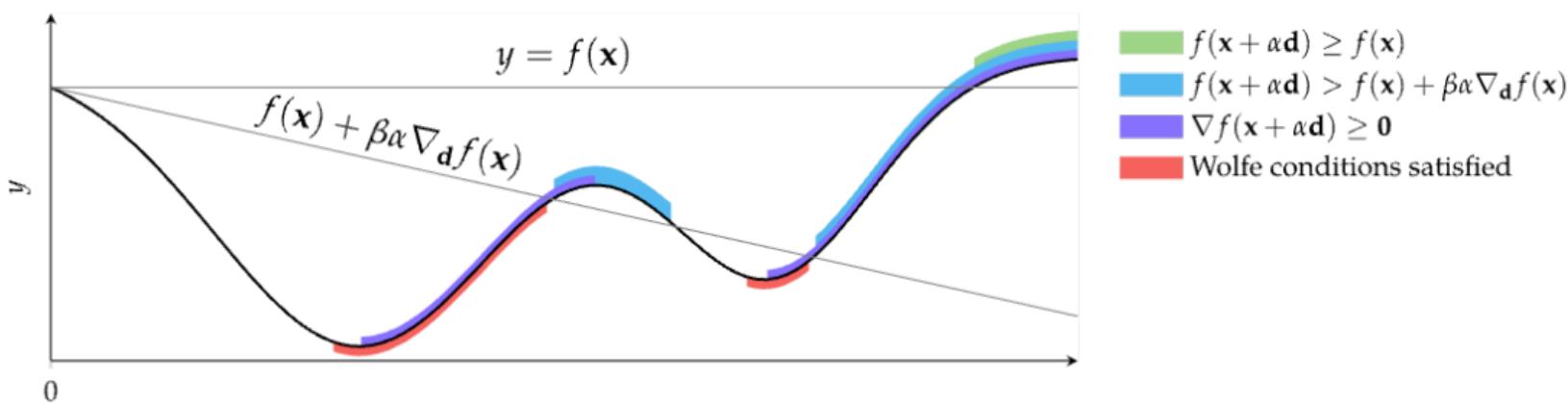
1. Bracketing Phase

An interval guaranteed to contain step lengths satisfying the Wolfe conditions is found when one of the following conditions hold:

$$f(\mathbf{x} + \alpha \mathbf{d}) \geq f(\mathbf{x})$$

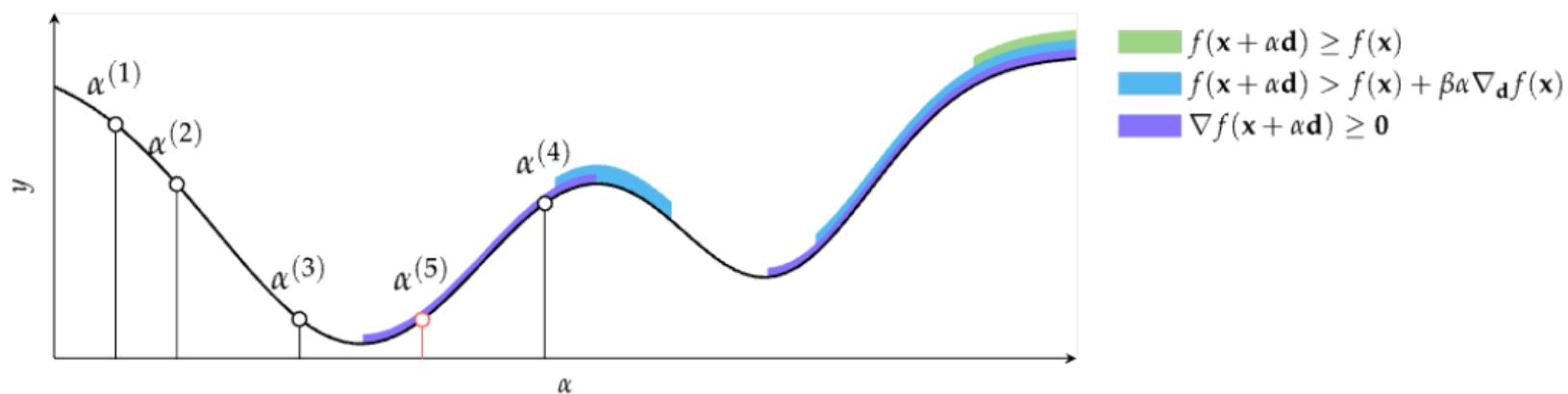
$$f(\mathbf{x} + \alpha \mathbf{d}) > f(\mathbf{x}) + \beta \alpha \nabla \mathbf{d} f(\mathbf{x})$$

$$\nabla f(\mathbf{x} + \alpha \mathbf{d}) \geq 0$$



Approximate Line Search

1. Braketing Phase + zoom phase (α_5)



Trust Region Methods

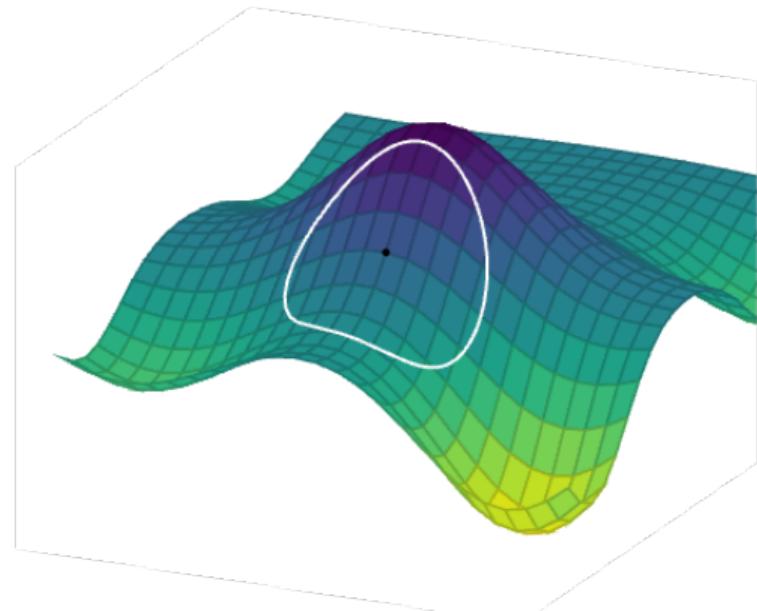
- Descent methods can place too much trust in their first and second order information
- A **trust region** is the local area of the design space where the local model is believed to be reliable.
- Trust region methods, or restricted step methods, limit the step size to ensure local approximation error is minimized
- If the improvement matches the predicted value, the trust region is expanded; otherwise it is contracted

Trust Region Methods

- \mathbf{x}' is new design point
- $\hat{f}(\mathbf{x}')$ is local function approximation, eg, second-order Taylor approximation
- δ is trust region radius

$$\begin{aligned} & \text{minimize}_{\mathbf{x}'} \hat{f}(\mathbf{x}') \\ & \text{subject to } \|\mathbf{x} - \mathbf{x}'\| \leq \delta \end{aligned}$$

Constrained optimization problem.
It can be solved efficiently if \hat{f} quadratic



Trust Region Methods

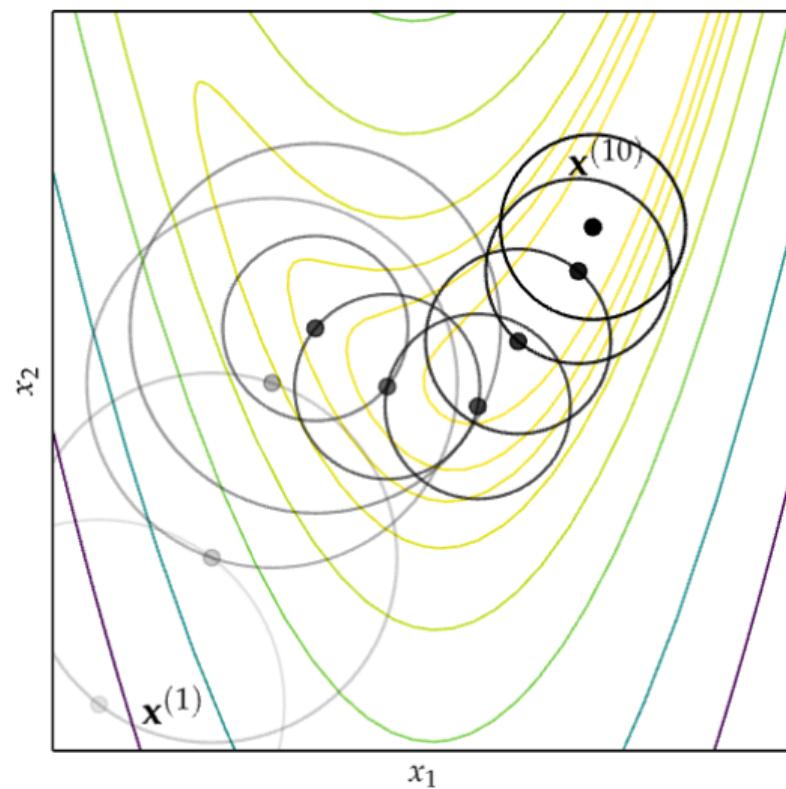
δ can be expanded or contracted based on performance

$$\eta = \frac{\text{actual improvement}}{\text{predicted improvement}} = \frac{f(\mathbf{x}) - f(\mathbf{x}')}{f(\mathbf{x}) - \hat{f}(\mathbf{x}')}$$

If $\eta < \eta_1$ contract

if $\eta > \eta_2$ expand

Trust Region Methods: Example



Trust regions can be also non circular.

Trust Region Methods

Termination Conditions (commonly used together):

- Maximum Iterations: $k > k_{\max}$
- Absolute Improvement: $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) < \epsilon_a$
- Relative Improvement: $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) < \epsilon_r |f(\mathbf{x}_k)|$
- Gradient Magnitude: $\|\nabla f(\mathbf{x}_{k+1})\| < \epsilon_g$

Then random restart.

Summary

- Descent direction methods incrementally descend toward a local optimum.
- Univariate optimization can be applied during line search.
- Approximate line search can be used to identify appropriate descent step sizes.
- Trust region methods constrain the step to lie within a local region that expands or contracts based on predictive accuracy.
- Termination conditions for descent methods can be based on criteria such as the change in the objective function value or magnitude of the gradient.

5. First-Order Methods

Descent Direction Methods

How to select the descent direction?

- first-order methods that rely on gradient
- second-order methods that rely on Hessian information

Advantages of first order methods:

- cheap iterations: good for small and large scale optimization embedded optimization
- helpful because easy to warm restart

Limitations of first order methods:

- not hard to find challenging instances for them.
- can converge slowly.

Outline

7. Gradient Descent

8. Conjugate Descent

9. Accelerated Descents

Gradient Descent

The **steepest descent** direction at \mathbf{x}_k , at k th iteration of a **local descent iterative method**, is the one opposite to the gradient (**gradient descent**):

$$\mathbf{d}_k = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|}$$

Guaranteed to lead to improvement if:

- f is smooth
- step size is sufficiently small
- \mathbf{x}_k is not a stationary point (ie, $\nabla f(\mathbf{x}_k) = 0$)

Gradient Descent: Example

- Suppose we have

$$f(\mathbf{x}) = x_1 x_2^2$$

- The gradient is $\nabla f = [x_2^2, 2x_1 x_2]$
- $\mathbf{x}_k = [1, 2]$

$$\mathbf{d}_{k+1} = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|} = \frac{[-4, -4]}{\sqrt{16 + 16}} = \left[-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right]$$

Implementation

```
class DescentMethod:  
    alpha: float  
  
class GradientDescent(DescentMethod):  
    def __init__(self, f, grad, x, alpha):  
        self.alpha = alpha  
  
    def step(self, f, grad, x):  
        alpha, g = self.alpha, grad(x)  
        return x - alpha * g
```

Gradient Descent

Theorem: The next direction is orthogonal to the current direction.

Proof:

$$\alpha_k^* = \underset{\alpha}{\operatorname{argmin}} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

$$\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k) = \nabla_{\mathbf{d}_k} f(\mathbf{x}_k) = 0 \quad \text{because } \alpha_k^* \text{ is minimum}$$

$$\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k)^T \mathbf{d}_k = 0 \quad \text{because directional derivative: } \nabla_{\mathbf{s}} f(\mathbf{x}) = \nabla f(\mathbf{x})^T \mathbf{s}$$

$$\mathbf{d}_{k+1} = -\frac{\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k)}{\|\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k)\|} \quad \text{gradient descent}$$

$$\mathbf{d}_{k+1} \cdot \mathbf{d}_k = -\frac{\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k)}{\|\nabla f(\mathbf{x}_k + \alpha_k^* \mathbf{d}_k)\|} \cdot \mathbf{d}_k = 0 \quad \mathbf{d}_{k+1}^T \mathbf{d}_k = 0 \implies \mathbf{d}_{k+1} \perp \mathbf{d}_k$$

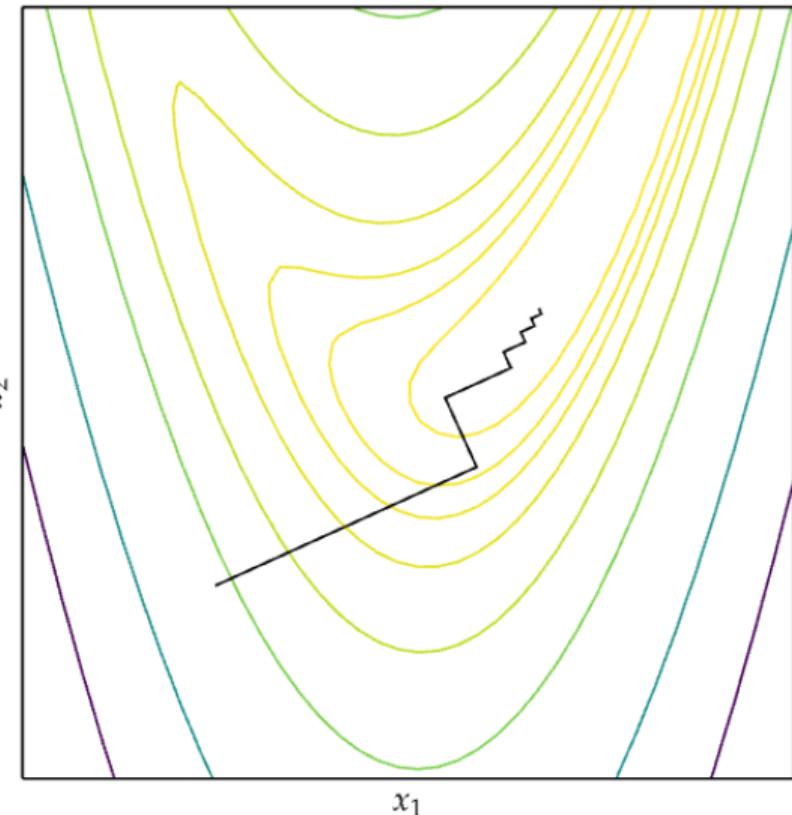
□

Gradient Descent: Example

2D Rosenbrock function

$$f(x, y) = (a - x)^2 + b(y - x^2)^2$$

Narrow valleys not aligned with gradient can be a problem



Outline

7. Gradient Descent

8. Conjugate Descent

9. Accelerated Descents

Conjugate Gradient

[Hestenes and Stiefel, 1950s]

For A symmetric positive definite:

$$Ax = b \iff \underset{x}{\text{minimize}} f(x) \stackrel{\text{def}}{=} \frac{1}{2}x^T Ax - b^T x$$

$$\nabla f(x) = Ax - b \stackrel{\text{def}}{=} r(x)$$

Conjugate Direction

Def.: A set of nonzero vectors $\{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_\ell\}$ is said to be **conjugate** with respect to the symmetric positive definite matrix A if

$$\mathbf{d}_i^T A \mathbf{d}_j = 0, \quad \text{for all } i \neq j$$

(the vectors are linearly independent. Generally, not orthogonal.)

Theorem: Given an arbitrary $\mathbf{x}_0 \in \mathbb{R}^n$ and a set of conjugate vectors $\{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1}\}$ the sequence $\{\mathbf{x}_k\}$ generated by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

where α_k is the analytical solution of $\min_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ given by:

$$\alpha_k = -\frac{\mathbf{r}_k^T \mathbf{d}_k}{\mathbf{d}_k^T A \mathbf{d}_k}$$

(aka, **conjugate direction algorithm**) converges to the solution \mathbf{x}^* of the linear system and minimization problem in at most n steps.

Proof:

$$\min_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

We can compute the derivative with respect to α :

$$\begin{aligned}\frac{\partial}{\partial \alpha} f(\mathbf{x} + \alpha \mathbf{d}) &= \frac{\partial}{\partial \alpha} (\mathbf{x} + \alpha \mathbf{d})^T A (\mathbf{x} + \alpha \mathbf{d}) - \mathbf{b}^T (\mathbf{x} + \alpha \mathbf{d}) (+c) \\ &= \mathbf{d}^T A (\mathbf{x} + \alpha \mathbf{d}) - \mathbf{d}^T \mathbf{b} \\ &= \mathbf{d}^T (A \mathbf{x} - \mathbf{b}) + \alpha \mathbf{d}^T A \mathbf{d}\end{aligned}$$

Setting $\frac{\partial f(\mathbf{x} + \alpha \mathbf{d})}{\partial \alpha} = 0$ results in:

$$\alpha_k = -\frac{\mathbf{d}_k^T (A \mathbf{x}_k - \mathbf{b})}{\mathbf{d}_k^T A \mathbf{d}_k} = -\frac{\mathbf{d}_k^T r(\mathbf{x}_k)}{\mathbf{d}_k^T A \mathbf{d}_k} \quad (1)$$

- Since the directions $\{\mathbf{d}_k\}$ are linearly independent, they must span the whole space \mathbb{R}^n . Hence, there is a set of scalars σ_k such that:

$$\mathbf{x}^* - \mathbf{x}_0 = \sigma_0 \mathbf{d}_0 + \sigma_1 \mathbf{d}_1 + \dots + \sigma_{n-1} \mathbf{d}_{n-1}$$

- By premultiplying this expression by $\mathbf{d}_k^T A$ and using the conjugacy property, we obtain:

$$\sigma_k = \frac{\mathbf{d}_k^T A(\mathbf{x}^* - \mathbf{x}_0)}{\mathbf{d}_k^T A \mathbf{d}_k} \quad (2)$$

- If \mathbf{x}_k is generated by conjugate direction algorithm, then we have

$$\mathbf{x}_k = \mathbf{x}_0 + \alpha_0 \mathbf{d}_0 + \alpha_1 \mathbf{d}_1 + \dots + \alpha_k \mathbf{d}_{k-1}$$

- By premultiplying this expression by $\mathbf{d}_k^T A$ and using the conjugacy property, we have that

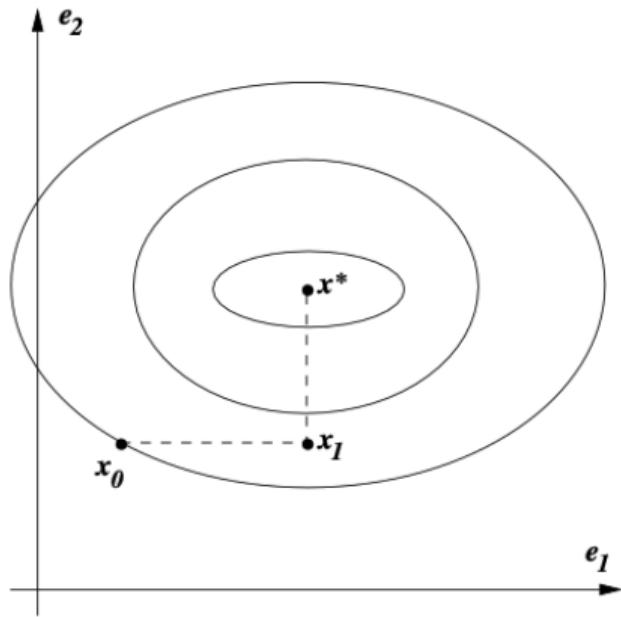
$$\mathbf{d}_k^T A(\mathbf{x}_k - \mathbf{x}_0) = 0$$

- and therefore

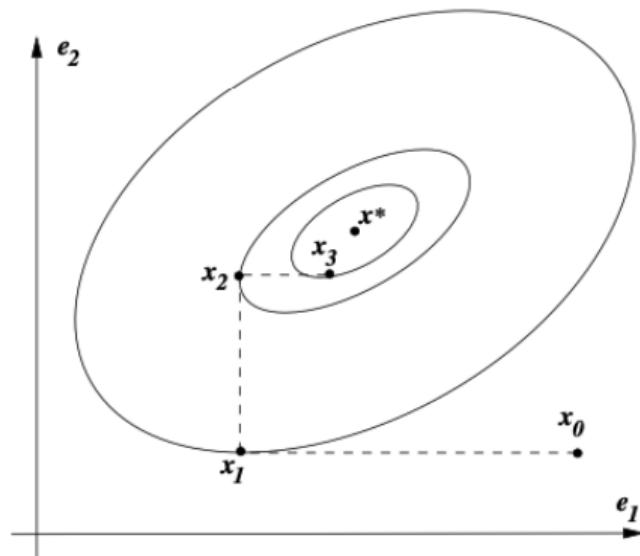
$$\begin{aligned} \mathbf{d}_k^T A(\mathbf{x}^* - \mathbf{x}_0) &= \mathbf{d}_k^T A(\mathbf{x}^* - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}_0) = \mathbf{d}_k^T A(\mathbf{x}^* - \mathbf{x}_k) + \mathbf{d}_k^T A(\mathbf{x}_k - \mathbf{x}_0) = \\ &= \mathbf{d}_k^T A(\mathbf{x}^* - \mathbf{x}_k) = \mathbf{d}_k^T (\mathbf{b} - A\mathbf{x}_k) = -\mathbf{d}_k^T \mathbf{r}_k. \end{aligned}$$

- Using this result in (2) and comparing with (1) we conclude $\alpha_k = \sigma_k$.

□



If the matrix A is diagonal, the contours of the function $f(\cdot)$ are ellipses whose axes are aligned with the coordinate directions



If A is not diagonal, its contours are elliptical, but they are usually not aligned with the coordinate directions.
Transform the problem to make A diagonal and minimize along the coordinate directions.

Conjugate Gradient Method

- The **conjugate gradient method** is a **conjugate direction method** with the property: In generating its set of conjugate vectors, it can compute a new vector \mathbf{d}_k by using only the previous vector \mathbf{d}_{k-1} . Hence, little storage and computation requirements.

$$\mathbf{d}_k = -\mathbf{r}_k + \beta_k \mathbf{d}_{k-1}$$

where β_k is to be determined such that \mathbf{d}_{k-1} and \mathbf{d}_k must be conjugate with respect to A . By premultiplying by $\mathbf{d}_{k-1}^T A$ and imposing that $\mathbf{d}_{k-1}^T A \mathbf{d}_k = 0$ we find that

$$\beta_k = \frac{\mathbf{r}_k^T A \mathbf{d}_{k-1}}{\mathbf{d}_{k-1}^T A \mathbf{d}_{k-1}}$$

- Larger values of β indicate that the previous descent direction contributes more strongly.
- \mathbf{d}_0 is commonly chosen to be the steepest descent direction at \mathbf{x}_0
- Advantage with respect to steepest descent: implicitly reuses previous information about the function and thus better convergence.

Algorithm CG

Basic version:

Input: f, x_0

Output: x^*

Set $\mathbf{r}_0 \leftarrow A\mathbf{x}_0 - b, \mathbf{d}_0 \leftarrow \mathbf{r}_0, k \leftarrow 0;$

while $\mathbf{r}_k \neq 0$ **do**

$$\alpha_k \leftarrow -\frac{\mathbf{d}_k^T \mathbf{r}(\mathbf{x}_k)}{\mathbf{d}_k^T A \mathbf{d}_k};$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k;$$

$$\mathbf{r}_{k+1} \leftarrow A\mathbf{x}_{k+1} - b;$$

$$\beta_{k+1} \leftarrow \frac{\mathbf{r}_{k+1}^T A \mathbf{d}_k}{\mathbf{d}_k^T A \mathbf{d}_k};$$

$$\mathbf{d}_{k+1} \leftarrow -\mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_k;$$

$$k \leftarrow k + 1;$$

Computationally improved version:

Input: f, x_0

Output: x^*

Set $\mathbf{r}_0 \leftarrow A\mathbf{x}_0 - b, \mathbf{d}_0 \leftarrow \mathbf{r}_0, k \leftarrow 0;$

while $\mathbf{r}_k \neq 0$ **do**

$$\alpha_k \leftarrow -\frac{\mathbf{r}(\mathbf{x}_k)^T \mathbf{r}(\mathbf{x}_k)}{\mathbf{d}_k^T A \mathbf{d}_k};$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k;$$

$$\mathbf{r}_{k+1} \leftarrow \mathbf{r}_k + \alpha_k A \mathbf{d}_k;$$

$$\beta_{k+1} \leftarrow \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k};$$

$$\mathbf{d}_{k+1} \leftarrow -\mathbf{r}_{k+1} + \beta_{k+1} \mathbf{d}_k;$$

$$k \leftarrow k + 1;$$

- we never need to know the vectors \mathbf{x}, \mathbf{r} , and \mathbf{d} for more than the last two iterations.
- major computational tasks: the matrix–vector product $A\mathbf{d}_k$, inner products $\mathbf{d}_k^T A \mathbf{d}_k$ and $\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}$, and three vector sums

NonLinear Conjugate Gradient Methods

- The conjugate gradient method can be applied to nonquadratic functions as well.
- Smooth, continuous functions behave like quadratic functions close to a local minimum
- but! we do not know the value of A that best approximates f around x_k . Instead, several choices for β_k tend to work well:
- Two changes:
 - α_k is computed by solving an approximate line search
 - the residual r , (it was simply the gradient of f), must be replaced by the gradient of the nonlinear objective f .

NonLinear Conjugate Gradient Methods

Fletcher-Reeves Method:

Input: f, \mathbf{x}_0

Output: \mathbf{x}^*

Evaluate $f_0 = f(\mathbf{x}_0), \nabla f_0 = \nabla f(\mathbf{x}_0)$;

Set $\mathbf{d}_0 \leftarrow -\nabla f_0, k \leftarrow 0$;

while $\nabla f_k \neq 0$ **do**

 Compute α_k by line search and set

$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$;

 Evaluate ∇f_{k+1} ;

$\beta_{k+1}^{FR} \leftarrow \frac{\nabla f_{k+1}^T \nabla f_{k+1}}{\nabla f_k^T \nabla f_k}$;

$\mathbf{d}_{k+1} \leftarrow -\nabla f_{k+1} + \beta_{k+1}^{FR} \mathbf{d}_k$;

$k \leftarrow k + 1$;

Polak-Ribière:

Input: f, \mathbf{x}_0

Output: \mathbf{x}^*

Evaluate $f_0 = f(\mathbf{x}_0), \nabla f_0 = \nabla f(\mathbf{x}_0)$;

Set $\mathbf{d}_0 \leftarrow -\nabla f_0, k \leftarrow 0$;

while $\nabla f_k \neq 0$ **do**

 Compute α_k by line search and set

$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$;

 Evaluate ∇f_{k+1} ;

$\beta_{k+1}^{PR} \leftarrow \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\nabla f_k^T \nabla f_k}$;

$\mathbf{d}_{k+1} \leftarrow -\nabla f_{k+1} + \beta_{k+1}^{PR} \mathbf{d}_k$;

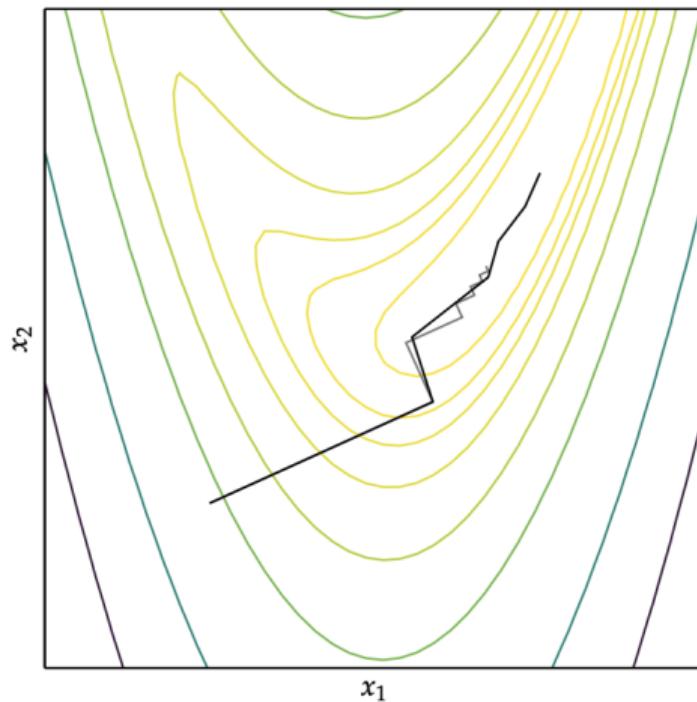
$k \leftarrow k + 1$;

PR with:

$$\beta_{k+1}^+ = \max\{\beta_{k+1}^{PR}, 0\}$$

becomes PR⁺ and guaranteed to converge (satisfies first Wolfe conditions).

The conjugate gradient method with the Polak-Ribi re update. Gradient descent is shown in gray.



Outline

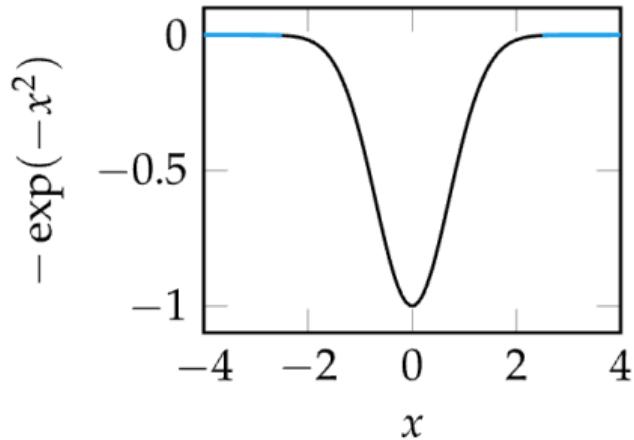
7. Gradient Descent

8. Conjugate Descent

9. Accelerated Descents

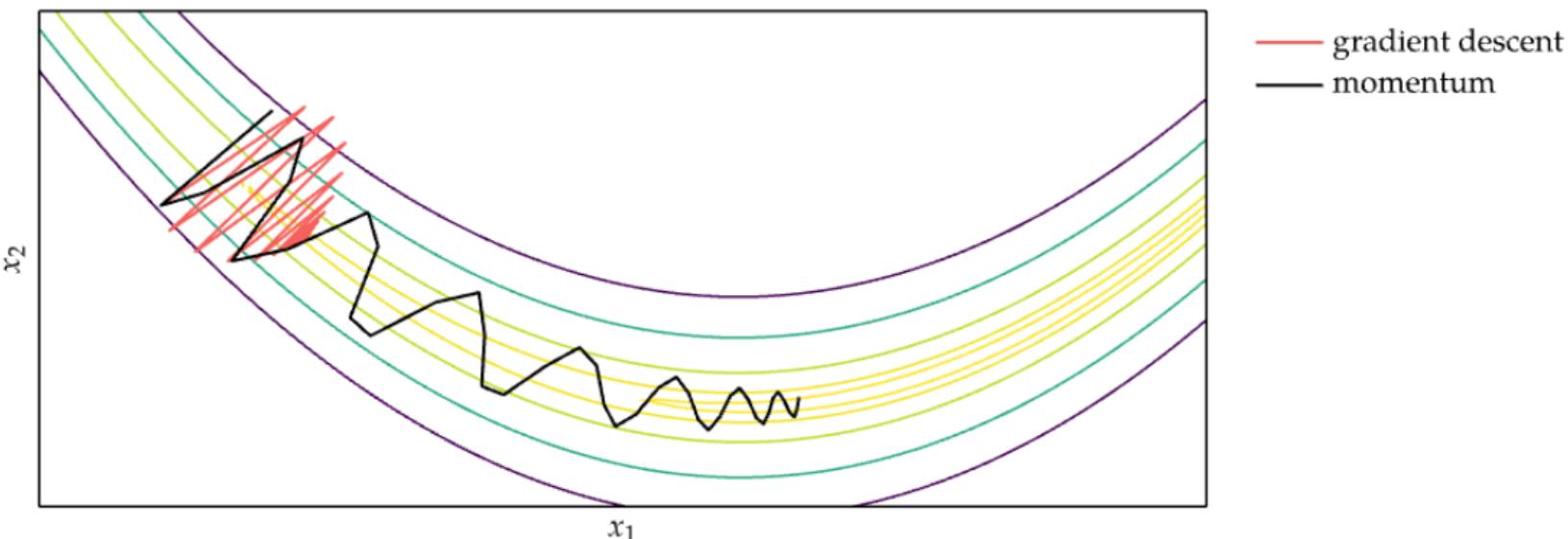
Accelerated Descents

- Addresses common convergence issues
- Some functions have regions with very small gradients (flat surface) where gradient descent gets stuck



Momentum

Rosenbrock function with $b = 100$



Momentum overcomes these issues by replicating the effect of physical momentum

Momentum

Momentum update equations:

$$\begin{aligned}\mathbf{v}_{k+1} &= \beta \mathbf{v}_k - \alpha \nabla f(\mathbf{x}_k) \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{v}_{k+1}\end{aligned}$$

```
import numpy as np

class Momentum(DescentMethod):
    alpha: float # learning rate
    beta: float # momentum decay
    v: np.array # momentum

    def __init__(self, alpha, beta, f, grad, x):
        self.alpha = alpha
        self.beta = beta
        self.v = np.zeros_like(x)

    def step(self, grad, x):
        self.v = self.beta * self.v - self.alpha * ↵
            ↵grad(x)
        return x + self.v
```

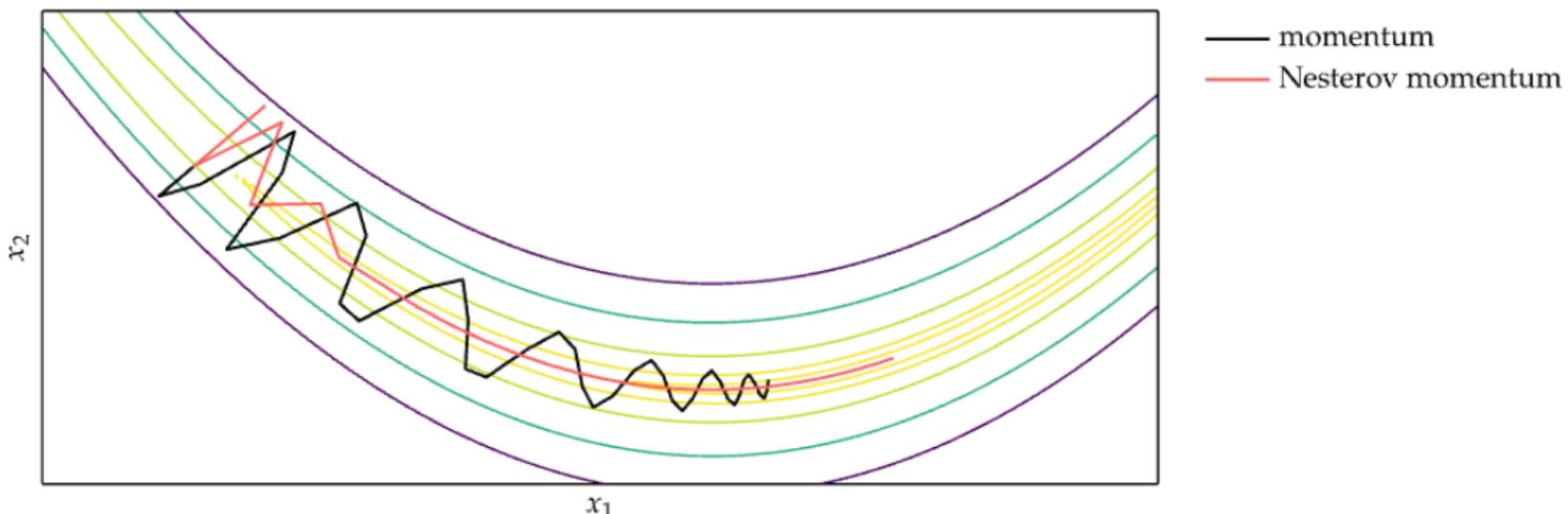
Nesterov Momentum

Issue of momentum: steps do not slow down enough at the bottom of a valley, overshoot.

Nesterov Momentum update equations:

$$\mathbf{v}_{k+1} = \beta \mathbf{v}_k - \alpha \nabla f(\mathbf{x}_k + \beta \mathbf{v}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{v}_{k+1}$$



Adagrad

- Instead of using the same learning rate for all components of \mathbf{x} ,
Adaptive Subgradient method (Adagrad) adapts the learning rate for each component of \mathbf{x} .
For each component of \mathbf{x} , the update equation is

$$x_{i,k+1} = x_{i,k} - \frac{\alpha}{\epsilon + \sqrt{s_{i,k}}} \nabla f_i(\mathbf{x}_k)$$

where

$$s_{i,k} = \sum_{j=1}^k (\nabla f_i(\mathbf{x}_j))^2$$

$$\epsilon \approx 1 \times 10^{-8}, \alpha = 0.01$$

- components of s are strictly nondecreasing, hence learning rate decreases over time

RMSProp

- Extends Adagrad to avoid monotonically decreasing learning rate by maintaining a decaying average of squared gradients

$$\hat{s}_{k+1} = \gamma \hat{s}_k + (1 - \gamma) (\nabla f(\mathbf{x}_k) \odot \nabla f(\mathbf{x}_k)), \quad \gamma \in [0, 1], \quad \odot \text{ element-wise product}$$

Update Equation

$$\begin{aligned}x_{i,k+1} &= x_{i,k} - \frac{\alpha}{\epsilon + \sqrt{\hat{s}_{i,k}}} \nabla f_i(\mathbf{x}_k) \\&= x_{i,k} - \frac{\alpha}{\epsilon + RMS(\nabla f_i(\mathbf{x}_k))} \nabla f_i(\mathbf{x}_k)\end{aligned}$$

root mean square: For n values $\{x_1, x_2, \dots, x_n\}$

$$x_{\text{RMS}} = \sqrt{\frac{1}{n} (x_1^2 + x_2^2 + \dots + x_n^2)}.$$

AdaDelta

Also extends Adagrad to avoid monotonically decreasing learning rate
Modifies RMSProp to eliminate learning rate parameter entirely

$$x_{i,k+1} = x_{i,k} - \frac{RMS(\Delta x_i)}{\epsilon + RMS(\nabla f_i(x))} \nabla f_i(x_k)$$

Adam

- The **adaptive moment estimation method** (Adam), adapts the learning rate to each parameter.
- stores both an exponentially decaying gradient like momentum and an exponentially decaying squared gradient like RMSProp and Adadelta
- At each iteration, a sequence of values are computed

Biased decaying momentum

$$\mathbf{v}_{k+1} = \beta \mathbf{v}_k - \alpha \nabla f(\mathbf{x}_k)$$

Biased decaying squared gradient

$$\mathbf{s}_{k+1} = \gamma \mathbf{s}_k + (1 - \gamma) (\nabla f(\mathbf{x}_k) \odot \nabla f(\mathbf{x}_k))$$

Corrected decaying momentum

$$\hat{\mathbf{v}}_{k+1} = \mathbf{v}_{k+1} / (1 - \gamma_{v,k})$$

Corrected decaying squared gradient

$$\hat{\mathbf{s}}_{k+1} = \mathbf{s}_{k+1} / (1 - \gamma_{s,k})$$

Next iterate

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \hat{\mathbf{v}}_{k+1} / (\epsilon + \sqrt{\hat{\mathbf{s}}_{k+1}})$$

- Defaults: $\alpha = 0.001$, $\gamma_v = 0.9$, $\gamma_s = 0.999$, $\epsilon = 1 \times 10^{-8}$

Adamax

Same as Adam, but based on the max-norm L_∞ .

$$\begin{aligned}\mathbf{s}_{k+1} &= \gamma^\infty \mathbf{s}_k + (1 - \gamma^\infty) (\|\nabla f(\mathbf{x}_k)\|_\infty) \\ &= \max(\gamma \mathbf{s}_k, \|\nabla f(\mathbf{x}_k)\|_\infty)\end{aligned}$$

Nadam

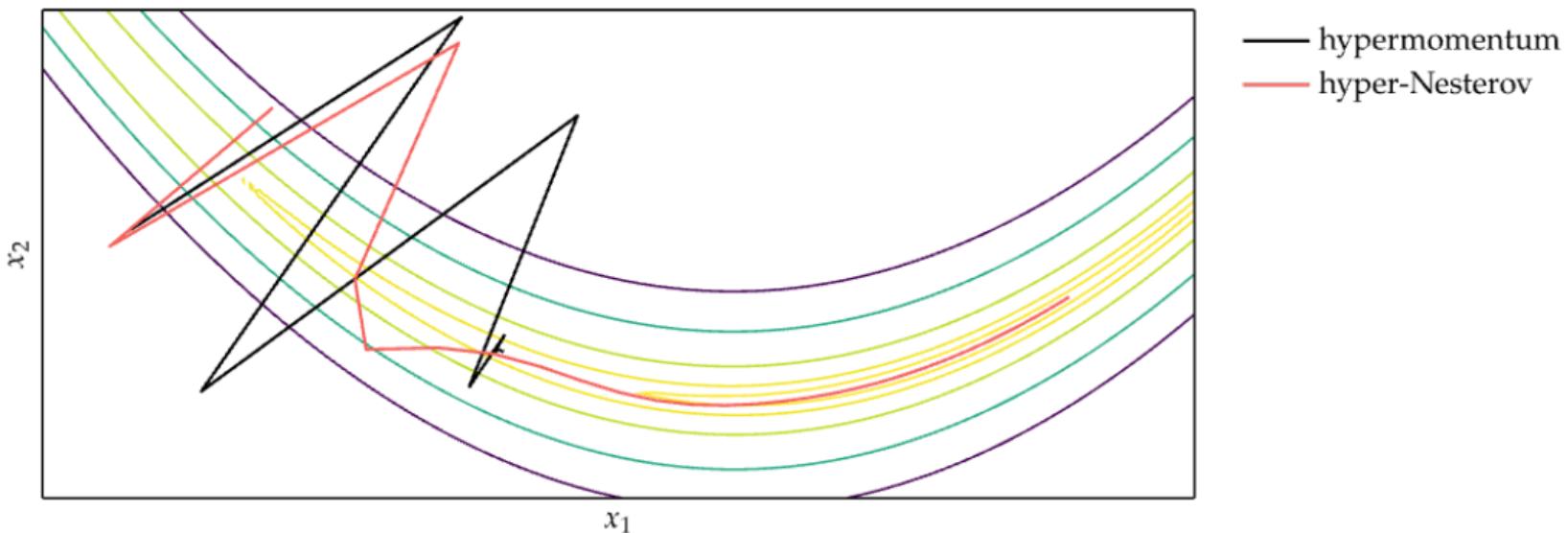
Nadam

- Nesterov-accelerated Adaptive Moment Estimation
- Adam is basically RMSProp with momentum
- We have seen that Nesterov is often more efficient
- Welcome to Nadam: Adam which uses the Nesterov momentum.

Hypergradient Descent

- Learning rate determines how sensitive the method is to the gradient signal.
- Many accelerated descent methods are highly sensitive to hyperparameters such as learning rate.
- Applying gradient descent to a hyperparameter of an underlying descent method is called hypergradient descent
- Requires computing the partial derivative of the objective function with respect to the hyperparameter

Hypergradient Descent



Summary

- Gradient descent follows the direction of steepest descent.
- The conjugate gradient method can automatically adjust to local valleys.
- Descent methods with momentum build up progress in favorable directions.
- A wide variety of accelerated descent methods use special techniques to speed up descent.
- Hypergradient descent applies gradient descent to the learning rate of an underlying descent method.

6. Second-Order Methods

Descent Direction Methods

How to select the descent direction?

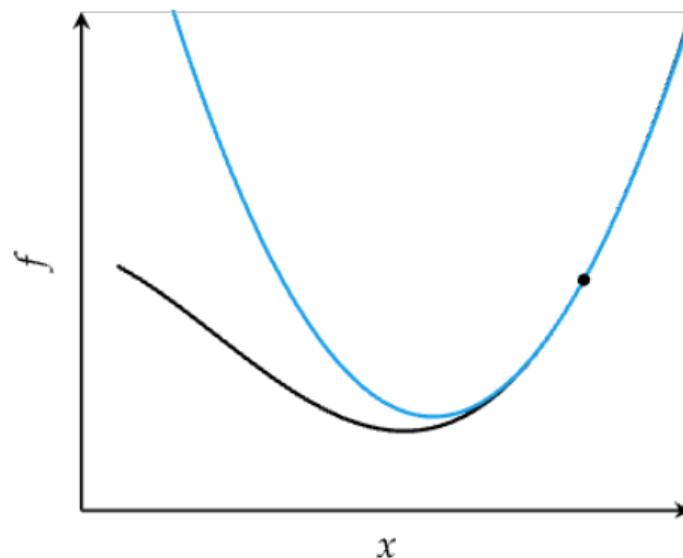
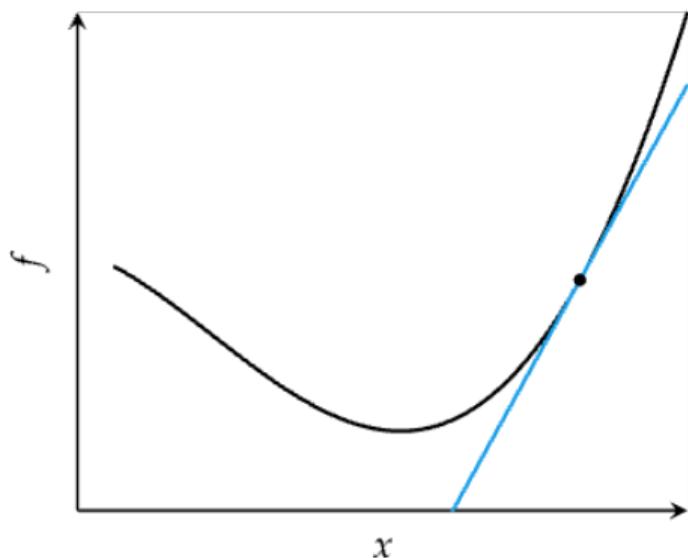
- first-order methods that rely on gradient
- second-order methods that rely on Hessian information

Advantages of second order methods in descent algorithms:

- way of accelerating the iteration [Davidon mid 1950s]
- additional information that can help improve the local model for informing the selection of
 - directions and
 - step lengths

Second-Order Methods

- Locally approximate function as quadratic
- Comparison of first-order and second order approximations



Outline

10. Newton Method

11. Secant Method

12. Quasi-Newton Method

Newton's Method – Univariate

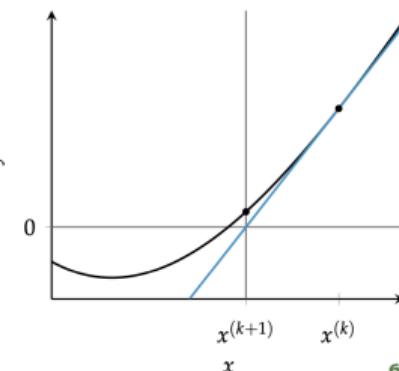
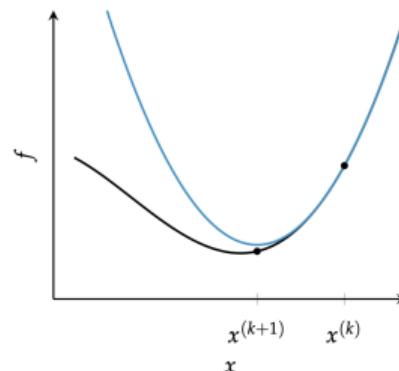
- Approximate a function using second-order Taylor series expansion
- analytically obtain the location where a quadratic approximation has a zero gradient.
- use that location as the next iteration to approach a local minimum.
- Univariate function

$$q(x) = f(x_k) + (x - x_k)f'(x_k) + \frac{(x - x_k)^2}{2}f''(x_k)$$

≡ finding roots of derivative function

$$\frac{dq(x)}{dx} = f'(x_k) + (x - x_k)f''(x_k) = 0$$

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



Newton's Method - Multivariate

- Multivariate function

$$f(\mathbf{x}) \approx q(\mathbf{x}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T H_k (\mathbf{x} - \mathbf{x}_k)$$

- H is the Hessian matrix
- Evaluate the gradient and set it to zero:

$$\nabla q(\mathbf{x}) = \nabla f(\mathbf{x}_k) + H(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) = 0$$

- Multivariate update rule

$$\boxed{\mathbf{x}_{k+1} = \mathbf{x}_k - H_k^{-1} \nabla f(\mathbf{x}_k)}$$

- (If f is quadratic and its Hessian is positive definite, then the update converges to the global minimum in one step.)

Algorithm

Input: $\nabla f, H, \mathbf{x}_0, \epsilon, k_{max}$

Output: \mathbf{x}^*

Set $k = 0, \Delta = 1, \mathbf{x} = \mathbf{x}_0$;

while $\|\Delta\| > \epsilon$ and $k \leq k_{max}$ **do**

$\Delta = H(\mathbf{x})^{-1}\nabla f(\mathbf{x})$;

$\mathbf{x} = \mathbf{x} - \Delta$;

$k = k + 1$;

It can be modified to only give a descent direction $\mathbf{d} = -H(\mathbf{x})^{-1}\nabla f(\mathbf{x})$ and leave the step size to be determined with line search.

Newton's method – Example

Minimize **Booth's function**:

$$f(x) = (x_1 + 2x_2 - 7)^2 + (2x_1 + x_2 - 5)^2$$

- $x_0 = [9, 8]$
- The gradient of Booth's function is:

$$\nabla f(x) = [10x_1 + 8x_2 - 34, 8x_1 + 10x_2 - 38]$$

- The Hessian of Booth's function is:

$$H(x) = \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix}$$

- The first iteration of Newton's method yields:

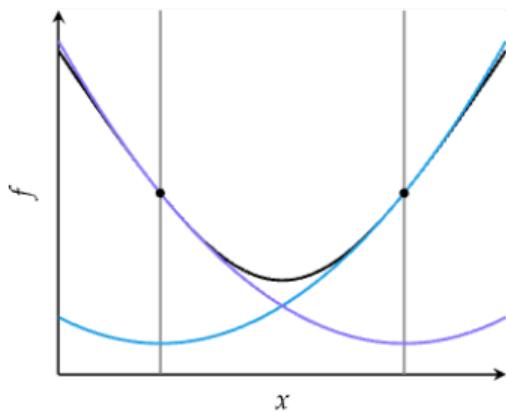
$$\begin{aligned} \mathbf{x}_1 &= \mathbf{x}_0 - H(\mathbf{x}_0)^{-1} \nabla f(\mathbf{x}_0) \\ &= \begin{bmatrix} 9 \\ 8 \end{bmatrix} - \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 10 \cdot 9 + 8 \cdot 8 - 34 \\ 8 \cdot 9 + 10 \cdot 8 - 38 \end{bmatrix} = \begin{bmatrix} 9 \\ 8 \end{bmatrix} - \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 120 \\ 114 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \end{aligned}$$

- Second iteration: The gradient at \mathbf{x}_1 is zero, so we have converged after a single iteration. The Hessian is positive definite everywhere, so \mathbf{x}_1 is the global minimum.

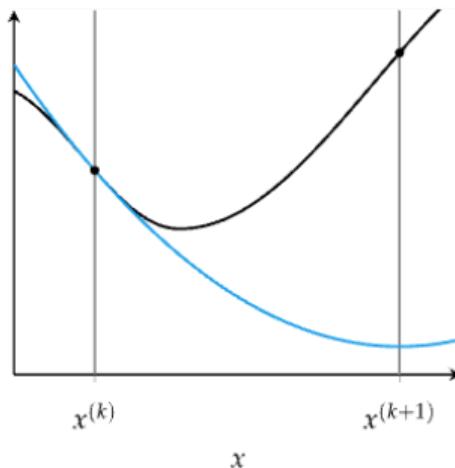
Newton's Method

Common causes of error in Newton's method

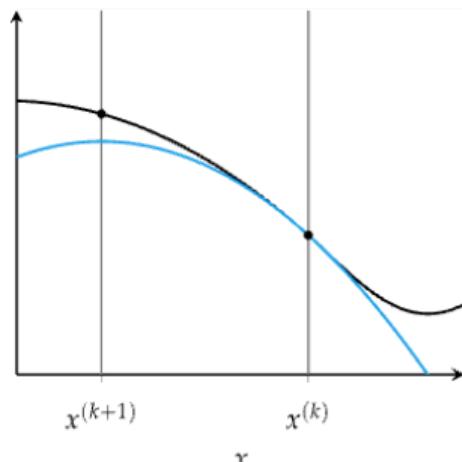
Oscillation



Overshoot



Negative f''



Newton's Method

- has quadratic convergence, meaning the difference between the minimum and the iterate is approximately squared with every iteration.
- This rate of convergence holds for Newton's method starting from x_0 within an interval $I = [x^* - \delta, x^* + \delta]$, for a root x^* , if
 1. $f''(x) \neq 0$ for all points in I ,
 2. $f'''(x)$ is continuous on I , and
 3. $\frac{1}{2} \left| \frac{f'''(x_0)}{f''(x_0)} \right| < c \left| \frac{f'''(x^*)}{f''(x^*)} \right|$ for some $c < \infty$

sufficient closeness condition, ensuring that the function is sufficiently approximated by the Taylor expansion and no overshoot.

Outline

10. Newton Method

11. Secant Method

12. Quasi-Newton Method

Secant Method – Univariate

- For univariate functions, if the second derivative is unknown, it can be approximated using the secant method

$$f''(x_k) = \frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}}$$

- Update equation

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

- It requires an additional initial design point and suffers from the same problems as Newton's method and may take more iterations to converge due to approximating the second derivative.

Outline

10. Newton Method

11. Secant Method

12. Quasi-Newton Method

Quasi-Newton Methods – Multivariate

- Automatic differentiation tools may not be applicable in many situations, and it may be much more costly to work with second derivatives in automatic differentiation software than with the gradient.
- Quasi-Newton methods, like steepest descent, require only the gradient of the objective function to be supplied at each iterate.
- By measuring the changes in gradients, they construct a model of the objective function that is good enough to produce superlinear convergence.
- The improvement over steepest descent is dramatic, especially on difficult problems.

Quasi-Newton Methods – Multivariate

Use an approximation $Q_k \approx H^{-1}(\mathbf{x}_k)$

Input: \mathbf{x}_0 , convergence tolerance $\epsilon > 0$, Q_0 (typically the $n \times n$ identity matrix)

Output: \mathbf{x}^*

Set $k \leftarrow 0$;

while $\|\nabla f(\mathbf{x}_k)\| > \epsilon$ **do**

Compute search direction $\mathbf{d}(\mathbf{x}_k) = -Q_k \nabla f(\mathbf{x}_k)$;

Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}(\mathbf{x}_k)$ where α_k is computed from a line search procedure to satisfy the Wolfe conditions;

Define $\delta_{k+1} \stackrel{\text{def}}{=} \mathbf{x}_{k+1} - \mathbf{x}_k$ and $\gamma_{k+1} \stackrel{\text{def}}{=} \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$;

Compute Q_{k+1} ;

$k \leftarrow k + 1$;

- Davidon-Fletcher-Powell (DFP) method
- Broyden-Fletcher-Goldfarb-Shanno (BFGS) method
- Limited-memory BFGS (L-BFGS) method

Davidon-Fletcher-Powell (DFP) method

$$Q_{k+1} = Q_k - \frac{Q_k \gamma_k \gamma_k^T Q_k}{\gamma_k^T Q_k \gamma_k} + \frac{\delta_k \delta_k^T}{\delta_k^T \gamma_k}$$

where all terms on the right hand side are evaluated at the same iteration k .

The update for Q in the DFP method has three properties:

- Q remains symmetric and positive definite.
- If $f(x) = \frac{1}{2}x^T A x + b^T x + c$, then $Q = A^{-1}$. Thus the DFP has the same convergence properties as the conjugate gradient method.
- For high-dimensional problems, storing and updating Q can be significant compared to other methods like the conjugate gradient method.

Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

$$Q_{k+1} = Q_k - \left(\frac{\delta_k \gamma_k^T Q_k + Q_k \gamma_k \delta_k^T}{\delta_k^T \gamma_k} \right) + \left(1 + \frac{\gamma_k^T Q_k \gamma_k}{\delta_k^T \gamma_k} \right) \frac{\delta_k \delta_k^T}{\delta_k^T \gamma_k}$$

BFGS better than DFP with approximate line search but still uses an $n \times n$ dense matrix.

Theorem: Suppose that f is twice continuously differentiable and that the iterates generated by the BFGS algorithm converge to a minimizer \mathbf{x}^* at which the Hessian matrix G is Lipschitz continuous. Suppose also that the sequence $\|\mathbf{x}_k - \mathbf{x}^*\|$ converges to zero rapidly enough that $\sum_{k=1}^{\infty} \|\mathbf{x}_k - \mathbf{x}^*\| < \infty$. Then \mathbf{x}_k converges to \mathbf{x}^* at a superlinear rate (ie, faster than linear).

Limited-memory BFGS (L-BFGS) method

For **large-scale unconstrained optimization**

It stores the last m values for δ and γ rather than the full inverse Hessian ($i = 1$ oldest, $i = m$ last).

Compute d at x as $d = -z_m$ using:

$$q_m = \nabla f(x_k) \quad q_i = q_{i+1} - \frac{\delta_{i+1}^T q_{i+1}}{\gamma_{i+1}^T \delta_{i+1}} \gamma_{i+1}, \quad i = m-1, \dots, 1$$

$$z_0 = \frac{\delta_m \odot \delta_m \odot q_m}{\gamma_m^T \gamma_m} \quad z_i = z_{i-1} + \delta_{i-1} \left(\frac{\delta_{i-1}^T q_{i-1}}{\gamma_{i-1}^T \delta_{i-1}} - \frac{\gamma_{i-1}^T z_{i-1}}{\gamma_{i-1}^T \gamma_{i-1}} \right), \quad i = 1, \dots, m$$

For minimization, the inverse Hessian Q must remain positive definite.

The initial Hessian is often set to the diagonal of

$$Q_0 = \frac{\gamma_0 \delta_0^T}{\gamma_0^T \gamma_0}$$

Computing the diagonal for the above expression and substituting the result into $z_0 = Q_0 q_0$ results in the equation for z_0 .

L-BFGS

(L-BFGS two-loop recursion)

Input:

Output: $Q_k \nabla f_k = z$

Set $q \leftarrow \nabla f_k$;

for $i = k - 1, k - 2, \dots, k - m$ **do**

$$\alpha_i \leftarrow \frac{\delta_i^T q}{\gamma_i \cdot \delta_i^T};$$

$$q \leftarrow q - \alpha_i \cdot \gamma_i;$$

$z \leftarrow Q_0 q$;

for $i = k - m, k - m + 1, \dots, k - 1$
do

$$\beta \leftarrow \frac{\gamma_i^T \cdot r}{\gamma_i \cdot \delta_i^T};$$

$$z \leftarrow z + \delta_i(z_i - \beta);$$

Input:

Output: x^*

Choose starting point x_0 , integer $m > 0$;
 $k \leftarrow 0$;

while not convergence **do**

Set Q_0 ;

Compute $d_k \leftarrow -Q_k \nabla f_k$ from Algorithm
on the left;

Compute $x_{k+1} \leftarrow x_k + \alpha_k d_k$, where α_k is
chosen to satisfy the Wolfe conditions;

if $k > m$ **then**

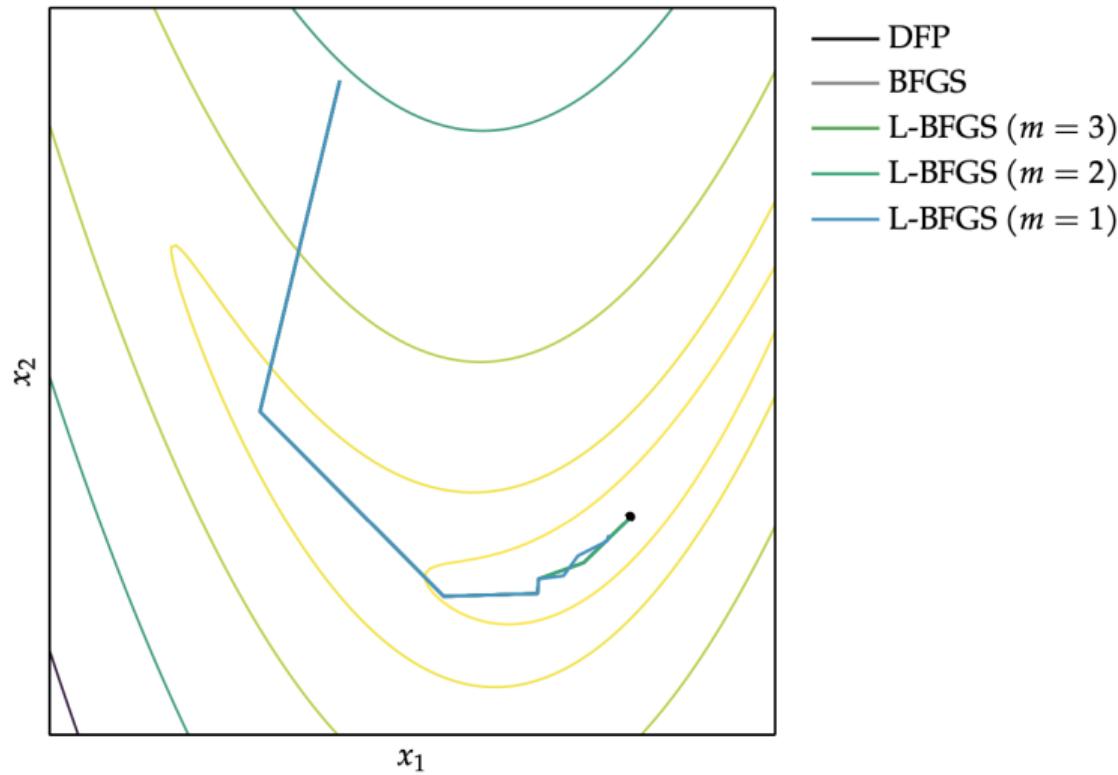
Discard the vector pair $\{\delta_{k-m}, \gamma_{k-m}\}$
from storage;

Compute and save $\delta_k = x_{k+1} - x_k$,

$$\gamma_k = \nabla f_{k+1} - \nabla f_k;$$

$$k \leftarrow k + 1;$$

BFGS Methods - Comparison



Summary

- Incorporating second-order information in descent methods often speeds convergence.
- Newton's method is a root-finding method that leverages second-order information to quickly descend to a local minimum.
- The secant method and quasi-Newton methods approximate Newton's method when the second-order information is not directly available.
- In Python, methods implemented in the module `scipy`
`https://docs.scipy.org/doc/scipy/tutorial/optimize.html`

7. Direct Methods

Derivative-Free Methods

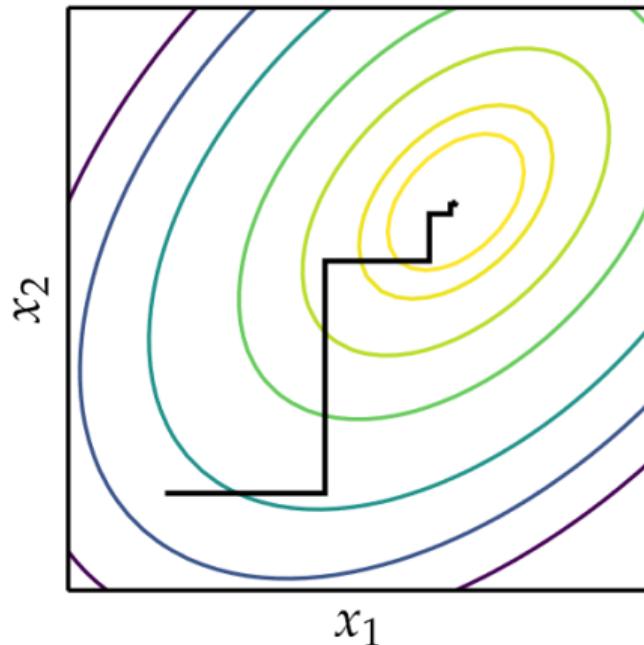
- Also called **direct search methods**, zero-order, black box, pattern search
- Direct method search using function evaluations only

Cyclic Coordinate Search

- Also known as coordinate descent, or taxicab search
- Performs line search in alternating coordinate directions

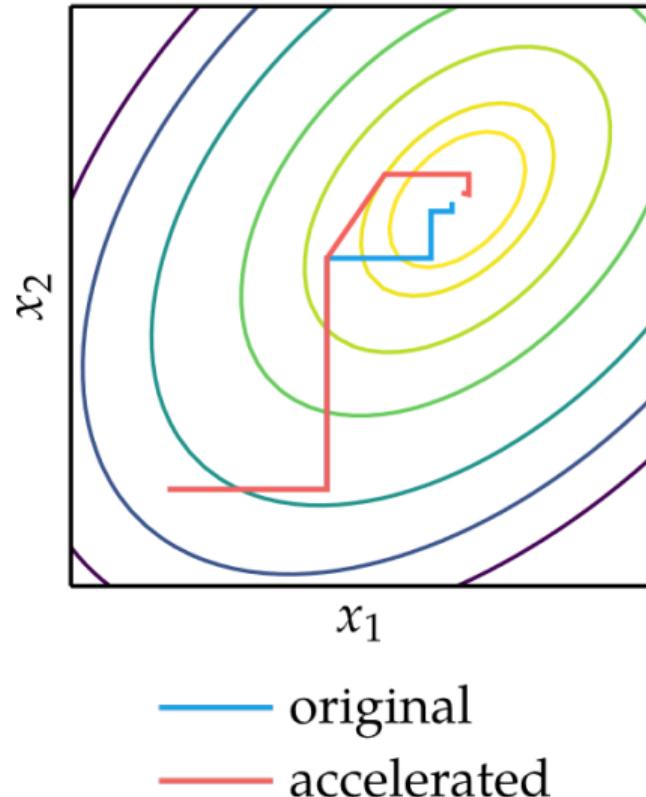
$$x_{1,1} = \operatorname{argmin}_{x_1} f(x_1, x_{2,0}, x_{3,0}, \dots, x_{n,0})$$

$$x_{2,1} = \operatorname{argmin}_{x_2} f(x_{1,1}, x_2, x_{3,1}, \dots, x_{n,1})$$



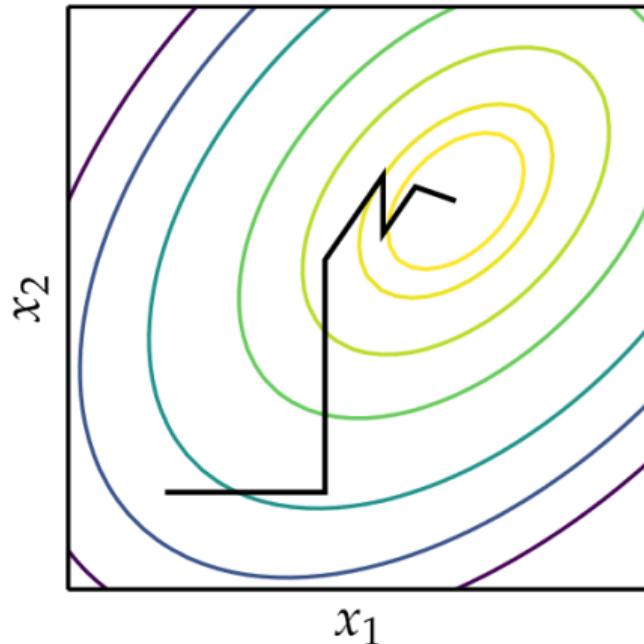
Cyclic Coordinate Search

- Can be augmented to accelerate convergence
- For every full cycle starting with optimizing x_1 along $[1, 0, \dots, 0]$ and ending with x_{n+1} after optimizing along $[0, 0, \dots, 1]$, an additional line search is conducted along the direction $x_{n+1} - x_1$.



Powell's Method

- Similar to Cyclic Coordinate Search, but can search in non-orthogonal directions
- Drops the oldest search direction in favor of the overall direction of progress
- It can lead the search directions to become linearly dependent and the search directions can no longer cover the full design space, and the method may not be able to find the minimum



Powell's Method

Input: f, \mathbf{x}_0

Output: \mathbf{x}^*

search directions $\mathbf{u}_1 = \mathbf{e}_1, \dots, \mathbf{u}_n = \mathbf{e}_i$;

while not converged **do**

for i in $\{1, \dots, n\}$ **do**

$\mathbf{x}_{i+1} \leftarrow \text{line_search}(f, \mathbf{x}_i, \mathbf{u}_i);$

for i in $\{1, \dots, n-1\}$ **do**

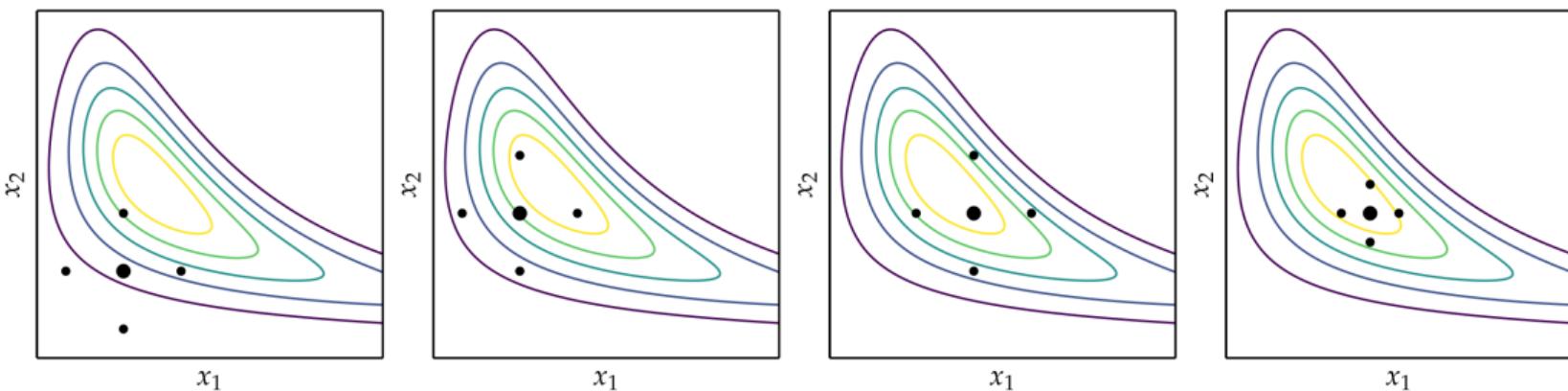
$\mathbf{u}_i \leftarrow \mathbf{u}_{i+1};$

$\mathbf{u}_n \leftarrow \mathbf{x}_{n+1} - \mathbf{x}_1;$

- search directions can become linearly dependent and no longer cover the full design space.
- periodically reset the directions to the canonical basis.

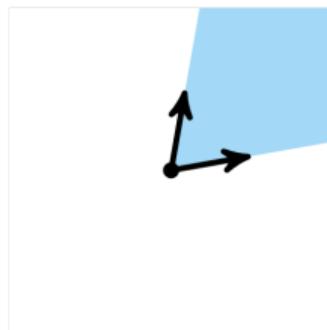
Hooke-Jeeves

- Evaluate $f(\mathbf{x})$ and $f(\mathbf{x} \pm \alpha \mathbf{e}_i)$ for a given step size α in every coordinate direction from an anchoring point \mathbf{x} .
- It accepts any improvement it may find.
- If no improvements are found, it decreases the step size.
- The process repeats until the step size is sufficiently small.
- $2n$ evaluations for an n -dimensional problem

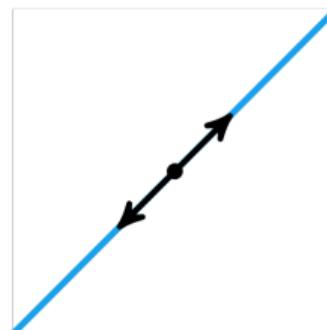


Generalized Pattern Search

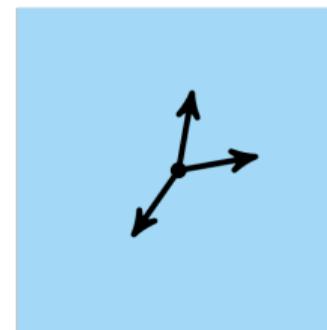
- Generalization of Hooke-Jeeves method
- A pattern P can be constructed from a set of directions D about an anchoring point x with a step size α according to: $P = \{x + \alpha d \text{ for each } d \in D\}$
- Searches in set of directions that **positively spans** (nonnegative linear combination) search space. (if D has full row rank and if $Dx = -D1$ with $x \geq 0$)



only positively spans the cone

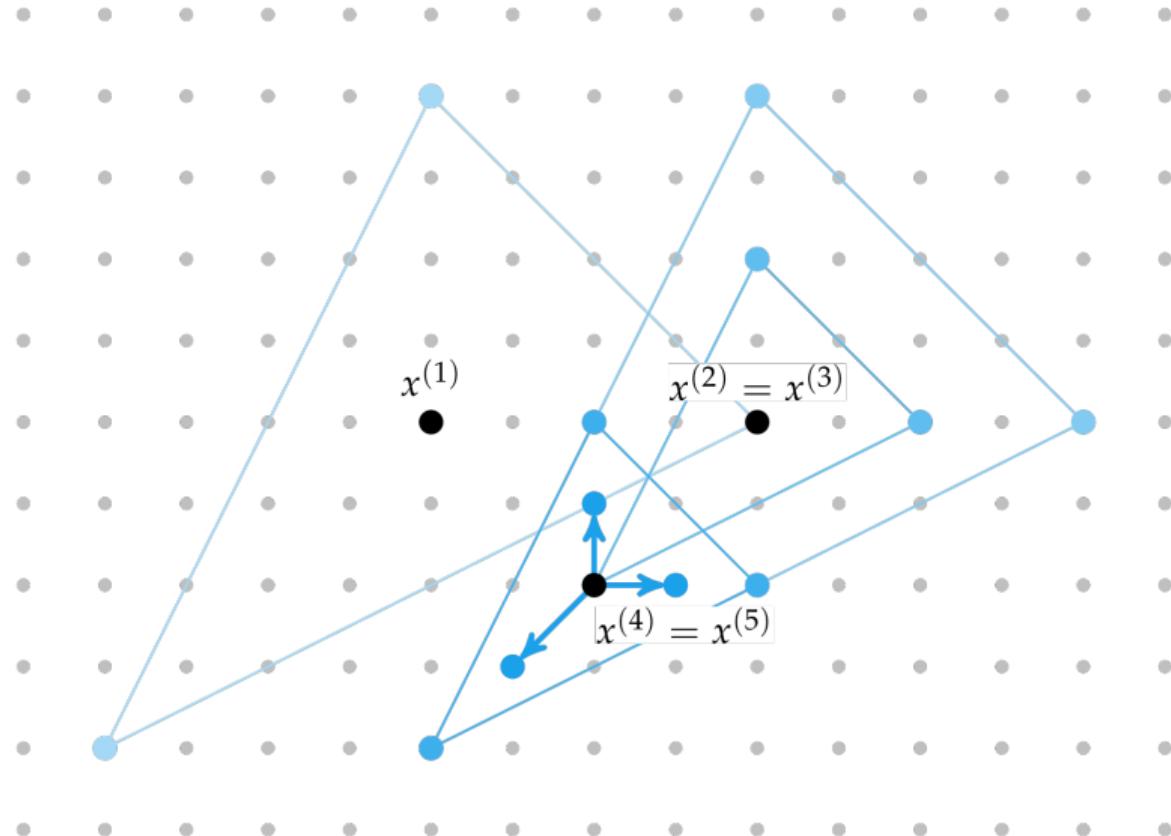


only positively spans 1d space



positively spans \mathbb{R}^2

Generalized Pattern Search



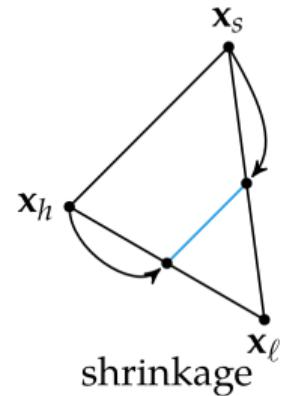
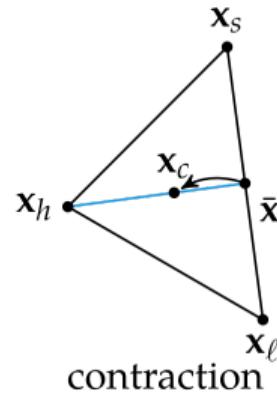
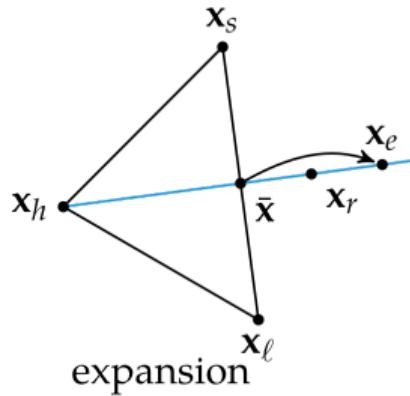
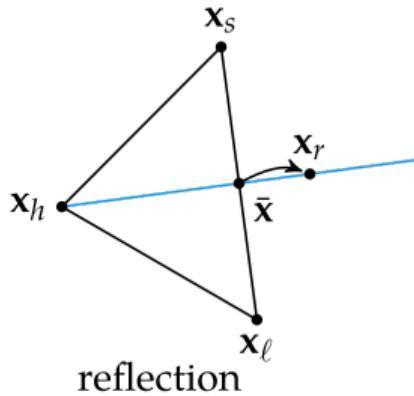
Outline

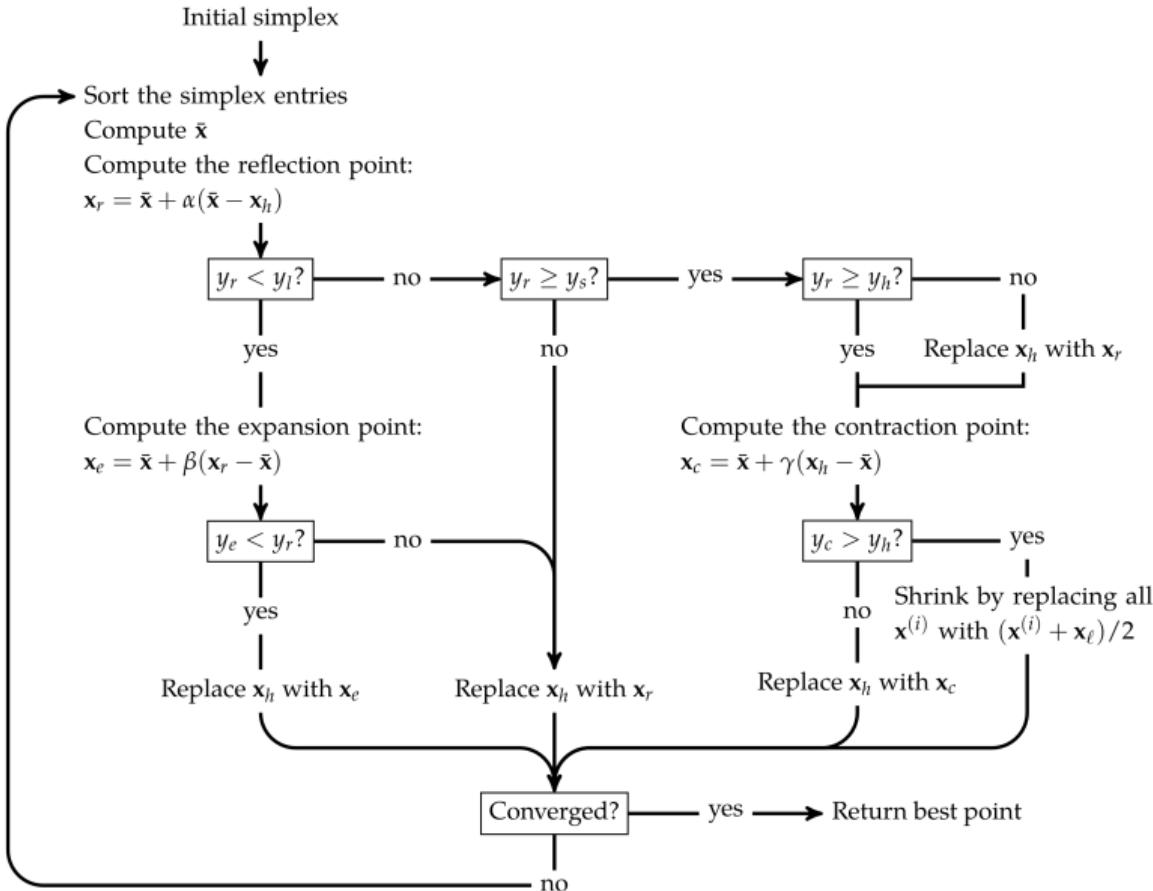
13. Nelder-Mead Simplex Method

14. Divided Rectangles

Nelder-Mead Simplex Method

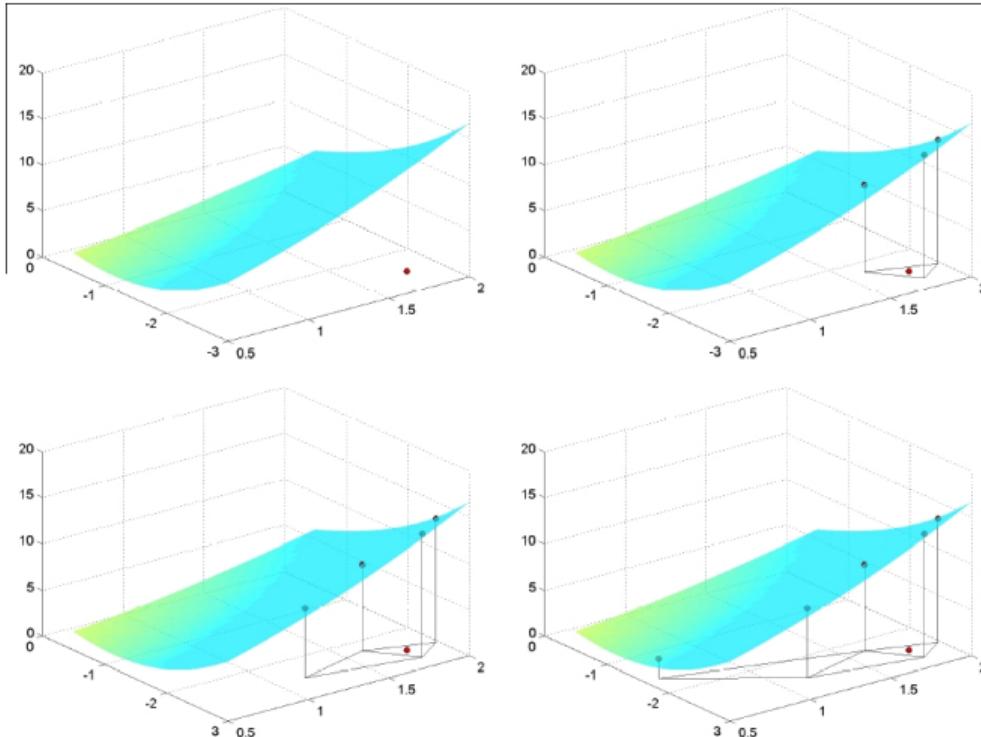
Uses simple algorithm to traverse search space using set of points defining a simplex





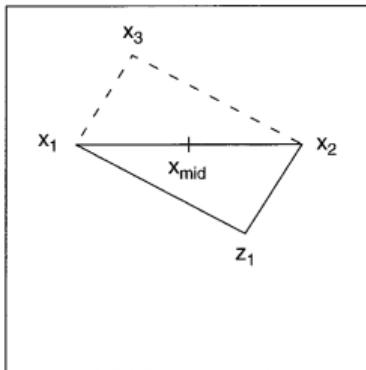
Nelder-Mead

Simplex based method [Spendley et al. (1962)]

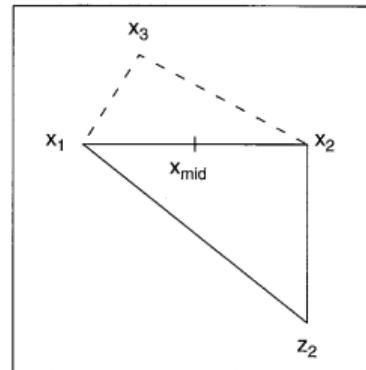


Nelder-Mead (cont.)

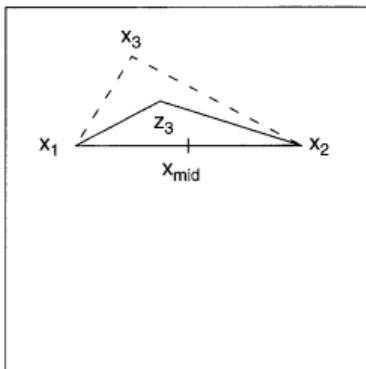
Reflection



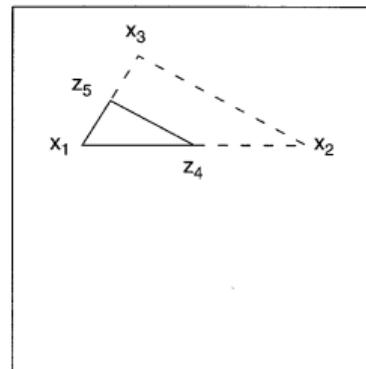
Reflection and expansion



Contraction 1

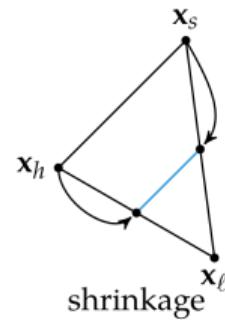
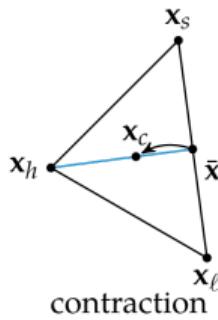
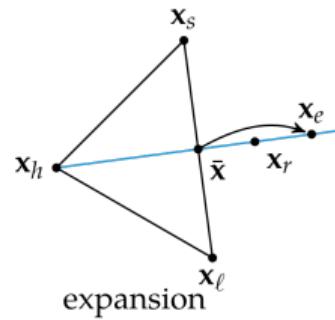
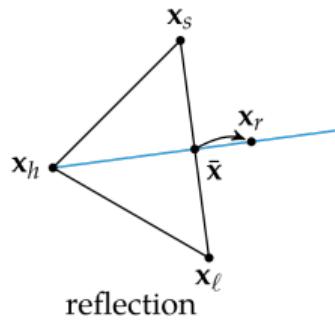


Contraction 2



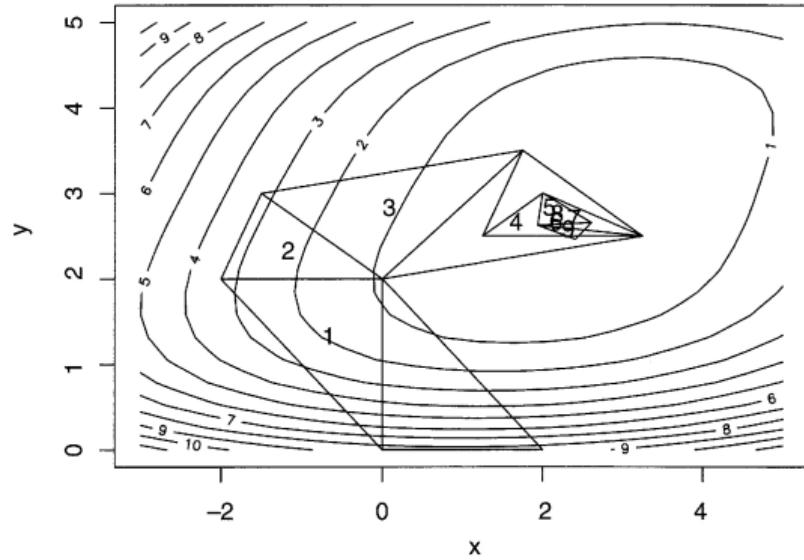
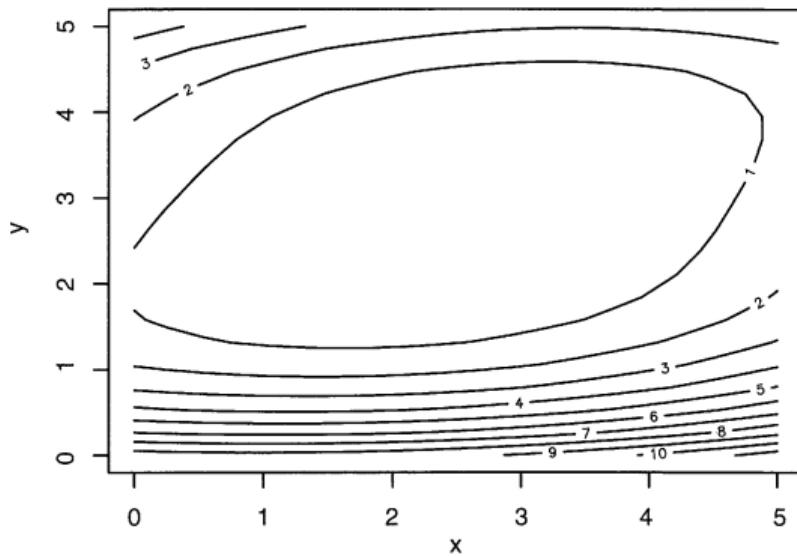
Nelder-Mead (cont.)

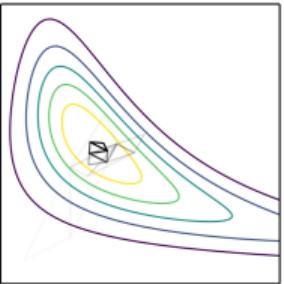
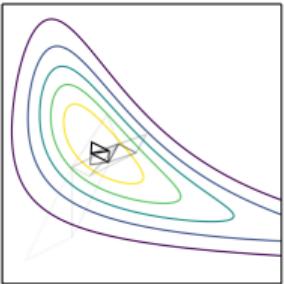
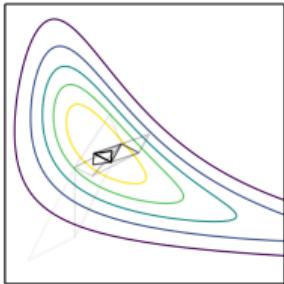
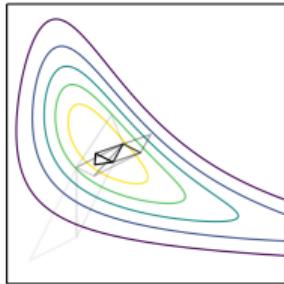
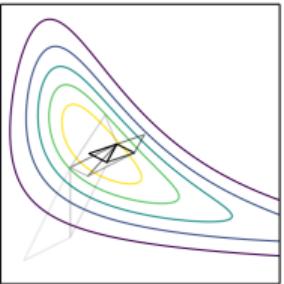
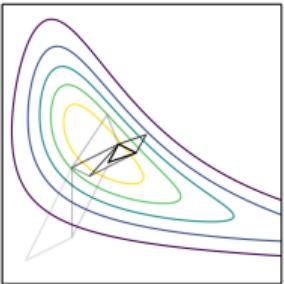
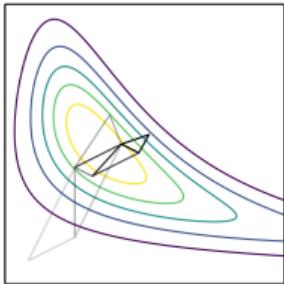
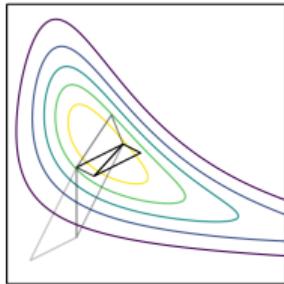
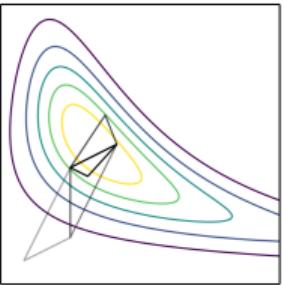
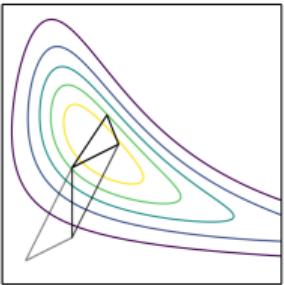
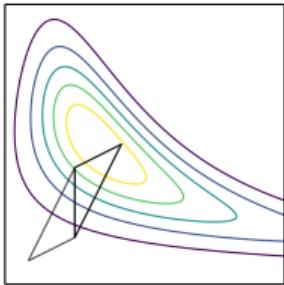
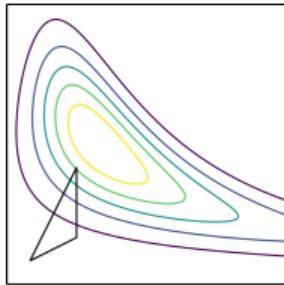
Nelder-Mead simplex method [Nelder and Mead, 1965]:



Nelder-Mead (cont.)

Example:





Nelder-Mead (cont.)

Algorithm: Simplex search

Let x_1, \dots, x_{n+1} be vertices of a simplex

Let h : $y_h = \max_i y_i = \max f(x_i)$ and

l : $y_l = \min_i y_i = \min f(x_i)$

Let \bar{x} be the centroid of points with $i \neq h$ and
 $d(x_i, x_j)$ the distance between two points x_i and x_j
($k \leftarrow 0$)

Reflect iter. k : ($k \leftarrow k + 1$) Generate the reflection x_R of x_h

Case 1 **if** $y_l \leq y_R < y_h$ then $x_h \leftarrow x_R$ and go to **Reflect**

Case 2 **else if** $y_R < y_l$ then

generate the expansion x_E of x_h

if $y_E < y_l$ then $x_h \leftarrow x_E$ and go to **Reflect**

else $x_h \leftarrow x_R$ and go to **Reflect**

Case 3 **else if** $y_R > y_i, \forall i \neq h$ **then** $x_h \leftarrow \min\{y_h, y_R\}$ and generate the contraction 1

if $y_C \leq \min\{y_h, y_R\}$ **then** $x_h \leftarrow x_C$ and go to **Reflect**

else contraction 2 $x_i \leftarrow (x_i + x_l)/2$

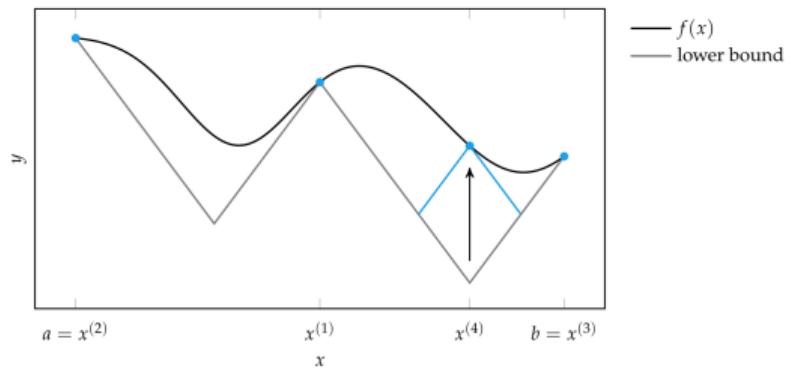
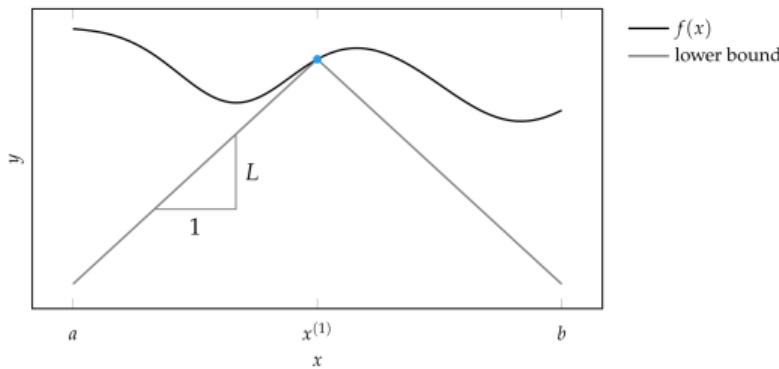
Outline

13. Nelder-Mead Simplex Method

14. Divided Rectangles

DIRECT – Divided Rectangles

- Also called DIRECT for Divided RECTangles
- Recall from Shubert-Piyavskii, a Lipschitz constant is used to provide a lower bound on the function, and a function evaluation is made where this bound is lowest



DIRECT – Divided Rectangles

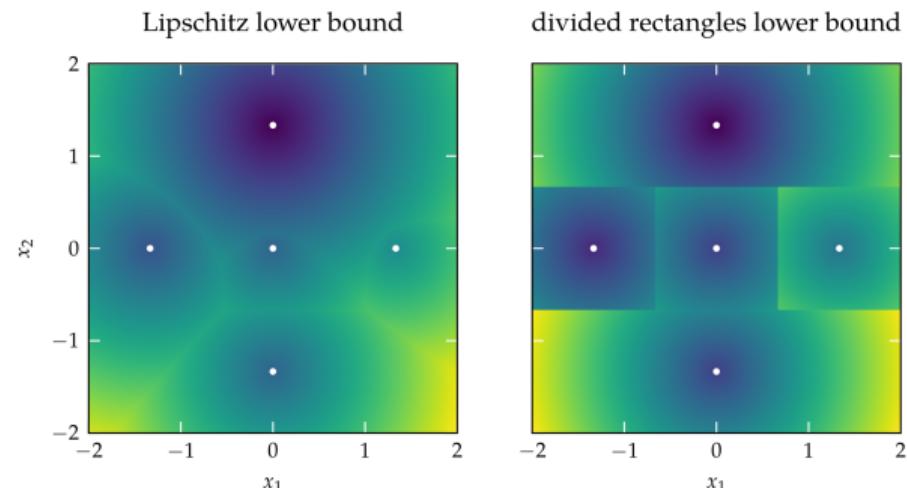
- The notion of Lipschitz continuity can be extended to multiple dimensions.
If f is Lipschitz continuous over a domain X with Lipschitz constant $\ell > 0$, then for a given design x_0 and $y = f(x_0)$, the circular cone

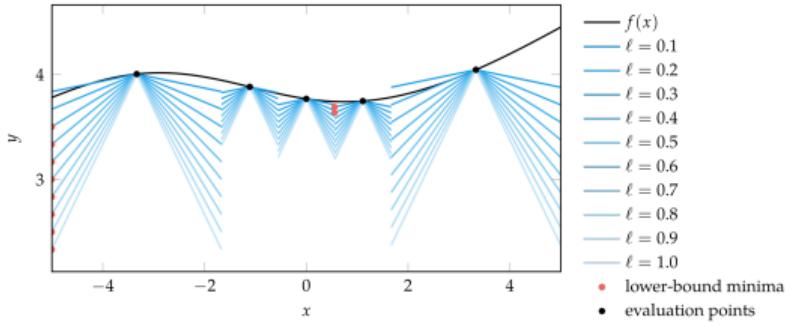
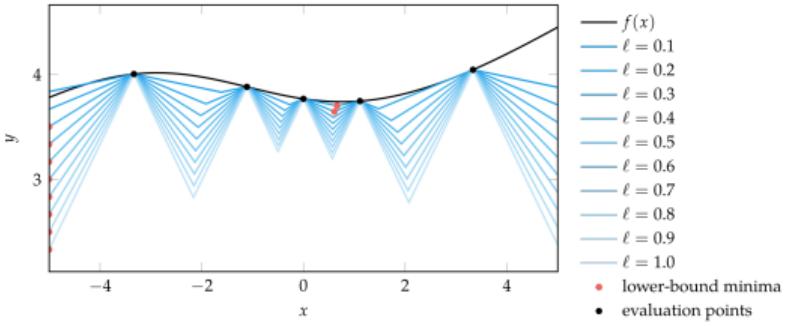
$$f(x_0) - \ell \|x - x_0\|_2$$

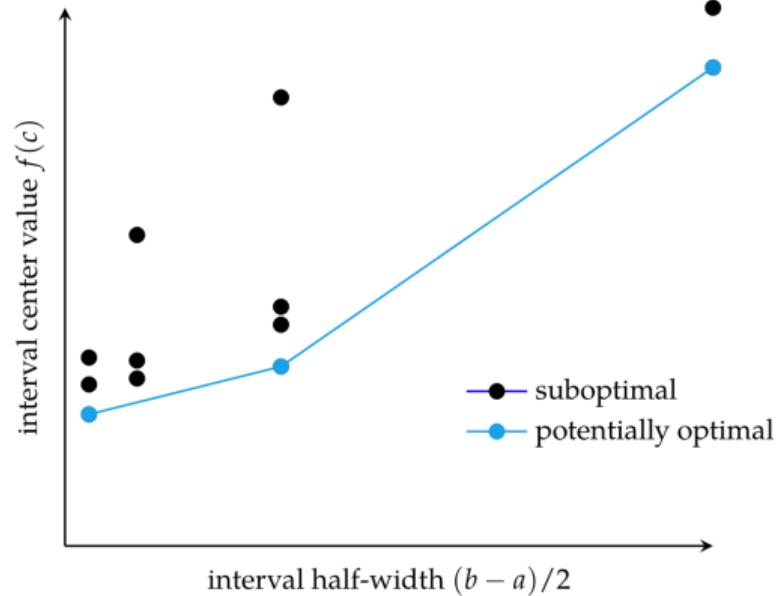
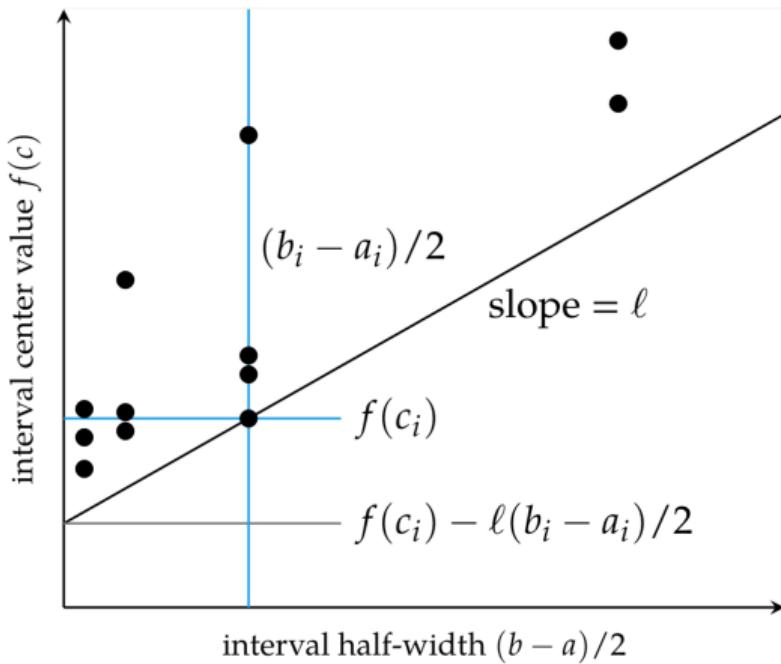
forms a lower bound of f

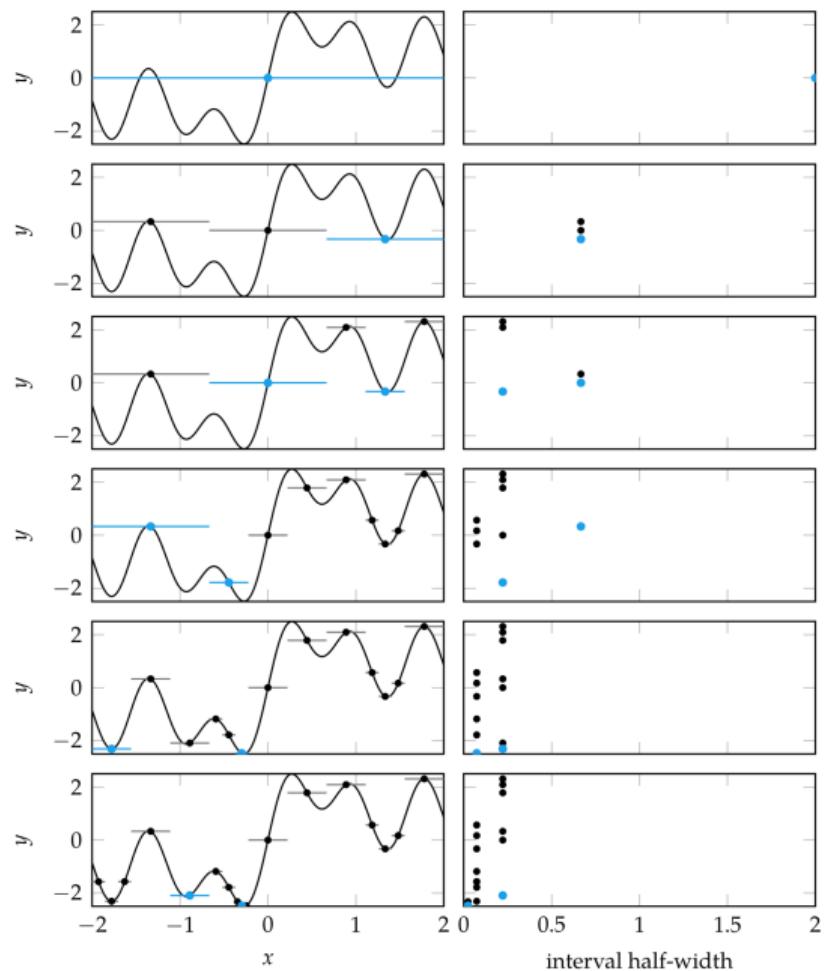
- Given m function evaluations with design points $\{x_1, x_2, \dots, x_m\}$, we can construct a superposition of these lower bounds by taking their maximum:

$$\max_i f(x_i) - \ell \|x - x_i\|_2$$





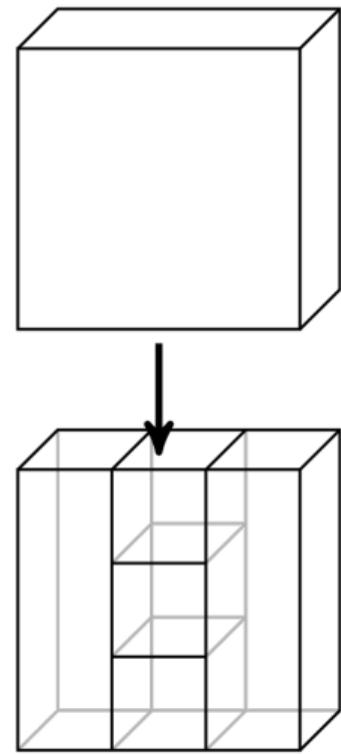
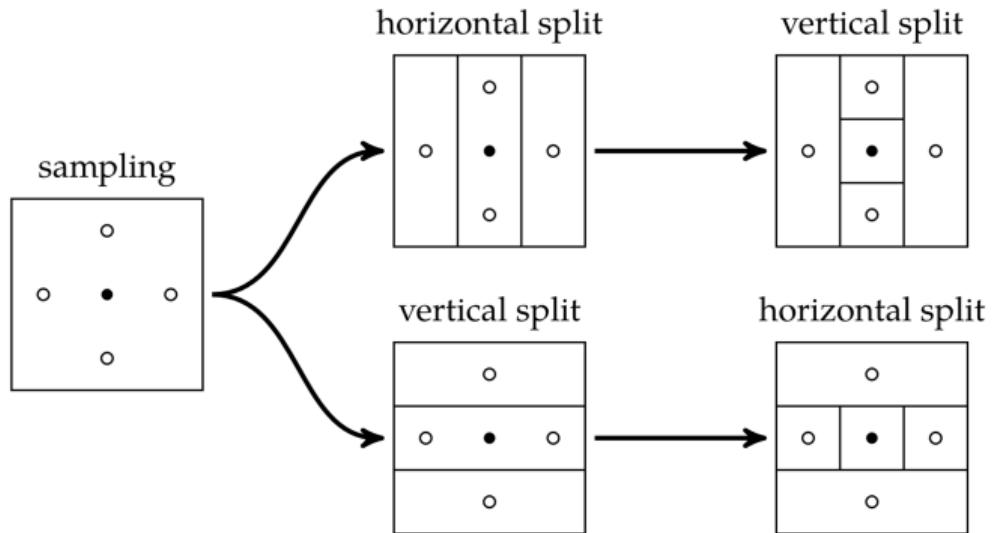




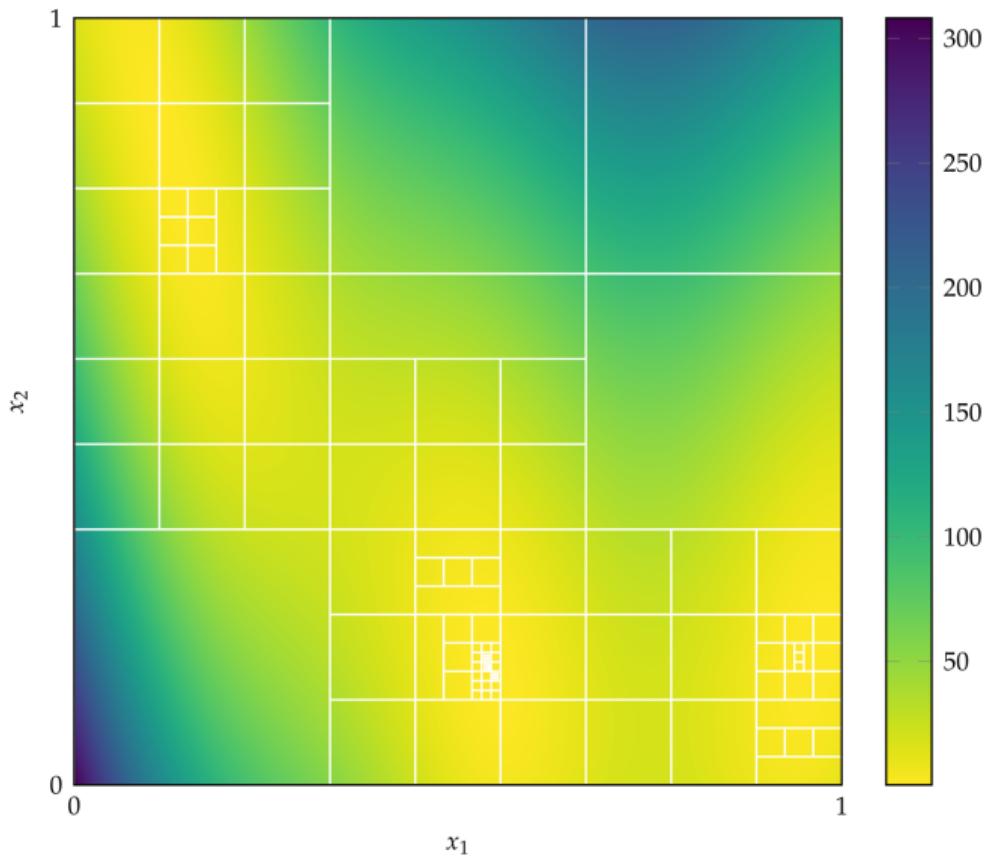
Multivariate DIRECT

- intervals → hyper-rectangles
- normalizes the search space to be the unit hypercube
- divide the rectangles into thirds along the axis directions
- larger rectangles for the points with lower function evaluations
- larger rectangles are prioritized for additional splitting
- when splitting a region without equal side lengths, only the longest dimensions are split

Multivariate DIRECT



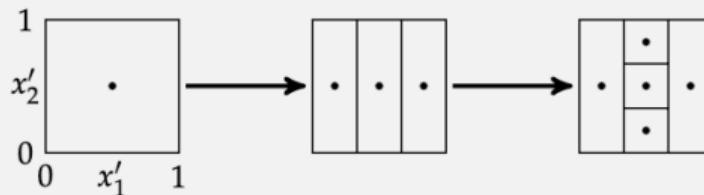
Multivariate DIRECT



Consider using DIRECT to optimize the flower function (appendix B.4) over $x_1 \in [-1, 3]$, $x_2 \in [-2, 1]$. The function is first normalized to the unit hypercube such that we optimize $x'_1, x'_2 \in [0, 1]$:

$$f(x'_1, x'_2) = \text{flower}(4x'_1 - 1, 3x'_2 - 2)$$

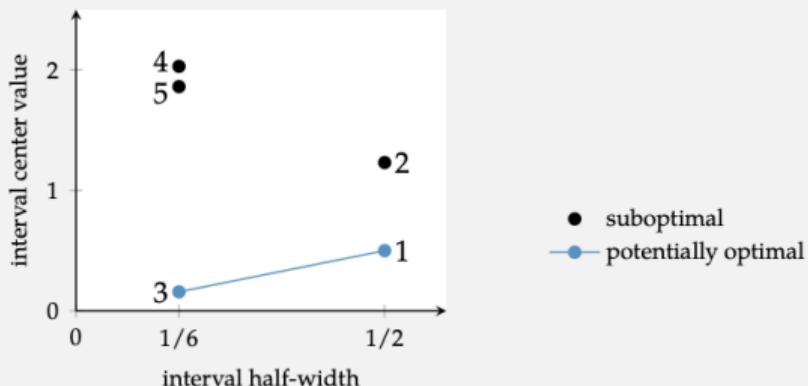
The objective function is sampled at $[0.5, 0.5]$ to obtain 0.158. We have a single interval with center $[0.5, 0.5]$ and side lengths $[1, 1]$. The interval is divided twice, first into thirds in x'_1 and then the center interval is divided into thirds in x'_2 .



the interval width can only take on powers of one-third, hence the interval half-width is
 $\left\| \frac{a-b}{2} \right\|_2 = \left\| \frac{3^{-h}}{2} \right\|_2$ where h is the depth of the rectangle

We now have five intervals:

interval	center	side lengths	vertex distance	center value
1	[1/6, 3/6]	[1/3, 1]	0.527	0.500
2	[5/6, 3/6]	[1/3, 1]	0.527	1.231
3	[3/6, 3/6]	[1/3, 1/3]	0.236	0.158
4	[3/6, 1/6]	[1/3, 1/3]	0.236	2.029
5	[3/6, 5/6]	[1/3, 1/3]	0.236	1.861



We next split on the two intervals centered at [1/6, 3/6] and [3/6, 3/6].

Summary

- Direct methods rely solely on the objective function and do not use derivative information.
- Cyclic coordinate search optimizes one coordinate direction at a time.
- Powell's method adapts the set of search directions based on the direction of progress.
- Hooke-Jeeves searches in each coordinate direction from the current point using a step size that is adapted over time.
- Generalized pattern search is similar to Hooke-Jeeves, but it uses fewer search directions that positively span the design space.
- The Nelder-Mead simplex method uses a simplex to search the design space, adaptively expanding and contracting the size of the simplex in response to evaluations of the objective function.
- The divided rectangles algorithm extends the Shubert-Piyavskii approach to multiple dimensions and does not require specifying a valid Lipschitz constant.

8. Stochastic Methods

Outline

15. Benchmarking

16. Stochastic Methods

Benchmarking in the COCO Platform

- Functions divided in suites.
- Functions, f_i , within suites are distinguished by their identifier $i = 1, 2, \dots$.
- parametrized by the (input) dimension, n , and
- instance number, j . (j as an index to a continuous parameter vector setting, eg, search space translations and rotations).

$$f_i^j \equiv f[n, i, j] : \mathbb{R}^n \rightarrow \mathbb{R} \quad \mathbf{x} \mapsto f_i^j(\mathbf{x}) = f[n, i, j](\mathbf{x}).$$

- Varying n or j leads to a variation of the same function i of a given suite.
- Fixing n and j of function f_i defines an **optimization problem instance** $(n, i, j) \equiv (f_i, n, j)$ that can be presented to the solver.

Why?

Varying the instance parameter j represents a natural randomization for experiments in order to:

- generate repetitions on a single function for deterministic solvers, making deterministic and non-deterministic solvers directly comparable (both are benchmarked with the same experimental setup)
- average away irrelevant aspects of the function definition
- alleviate the problem of overfitting, and
- prevent exploitation of artificial function properties

BBOB Functions

- All benchmark functions are scalable with the dimension.
- Most functions have no specific value of their optimal solution (they are randomly shifted in x -space).
- All functions have an artificially chosen optimal function value (they are randomly shifted in f -space).

Runtime and Target Values

- Runtime of a solver on a problem is the hitting time condition.
- define a non-increasing quality indicator measure and prescribe a set of **target values**, t .
- target values are compared with the best so-far-seen f -value.
- For a single run, the solver run is **successful** on the problem instance (f_i, n, j) when the best-so-far f -value reaches the target value t .
- COCO collects hundreds of different target values from each single run.
- targets $t(i, j)$ depend on the problem instance in a way to make problems comparable
- typically, target values are set to known or estimated optimal solution plus an added **precision**
- **runtime** is the number of f -evaluations needed to **solve** the problem $(f_i, n, j, t(i, j))$.
- only runtimes to comparable target values can be aggregated among problem instances.

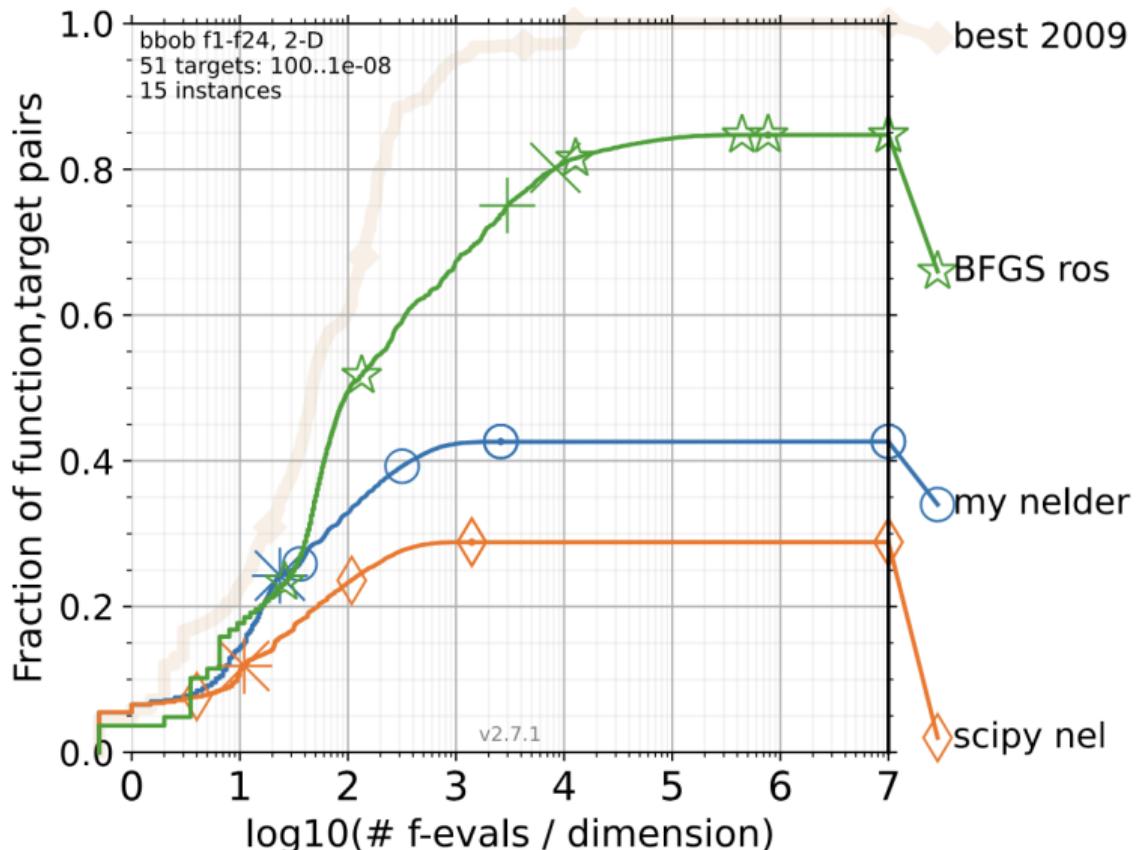
Simulated Restarts

- If a solver does not hit the target t in a given single run, the run is considered to be **unsuccessful**.
- The runtime of this single run remains undefined but is bounded from below by the number of evaluations conducted during the run $\tau \in [T, \infty]$
- T depends on the termination condition encountered. It can be the budget of evaluations.
- For hard problem instances COCO uses **budget-based target values**:
For any given budget, COCO selects from the finite set of recorded target values the easiest (i.e., largest) target for which the expected runtime of all solvers (ERT) exceeds the budget.
- With unsuccessful runs: draw further runs from the set of tried problem instances, uniformly at random with replacement, until find an instance, j , for which $(f_i, n, j, t(i, j))$ is solved.
the runtime is then the sum of the overall time spend and associated to the initially unsolved problem instance.

```
print: '!' if problem.final_target_hit, ':' if restarted and '.' otherwise.
```

Aggregation

- Aggregation is to compute a statistical summary over a set or subset of problem instances over which we assume a uniform distribution
- If we can distinguish between problems easily, for example, according to their input dimension, we can use the information to select the solver, hence not worth aggregating data
- Empirical cumulative distribution functions of runtimes (**runtime ECDFs**)
 - Absolute distributions vs Performance profiles (ECDFs of runtimes relative to the respective best solver)
 - aggregate runtimes from several targets per function (!?)
- arithmetic average, as an estimator of the expected runtime. The estimated expected runtime of the restarted solver, ERT, is often plotted against dimension to indicate scaling with dimension.
alternatives: average of log-runtimes \equiv geometric average or shifted geometric mean



Reference

Nikolaus Hansen, Anne Auger, Raymond Ros, Olaf Mersmann, Tea Tušar & Dimo Brockhoff
(2021) COCO: a platform for comparing continuous optimizers in a black-box setting, Optimization
Methods and Software, 36:1, 114-144, DOI: 10.1080/10556788.2020.1808977

Outline

15. Benchmarking

16. Stochastic Methods

- Employ randomness strategically to help explore design space
- Randomness can help escape local minima
- Increases chance of searching near the global minimum
- Typically rely on pseudo-random number generators to ensure repeatability
- Control over randomness and the exploration vs exploitation trade off.

Noisy Descent

- Saddle points, where the gradient is very close to zero, can cause descent methods to select step sizes that are too small to be useful
- add Gaussian noise at each descent step

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha \nabla f(\mathbf{x}_k) + \epsilon_k$$

$$\epsilon_k \sim \mathcal{N}(0, \sigma_k^2)$$

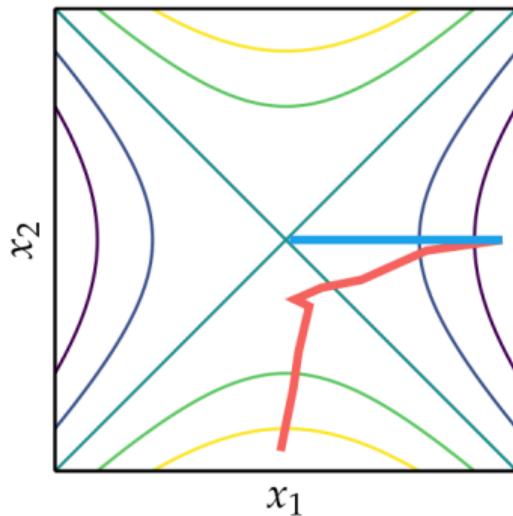
- $\sigma_k = \frac{1}{k}$

Stochastic Gradient Descent

- evaluates gradients using randomly chosen subsets of the training data (**batches**)
- significantly less expensive computationally than calculating the true gradient at every iteration and yields same effect as noisy gradient approximation
- helping traverse past saddle points Convergence guarantees for stochastic gradient descent require that the positive step sizes be chosen such that:

$$\sum_{k=1}^{\infty} \alpha_k = \infty \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty$$

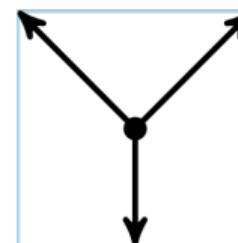
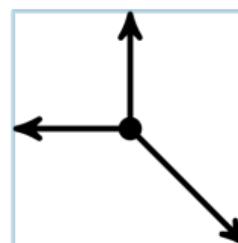
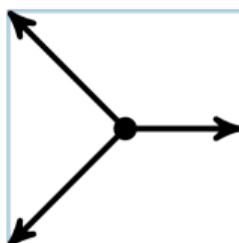
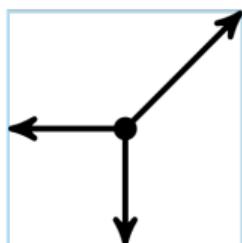
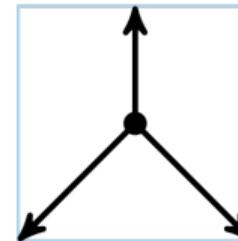
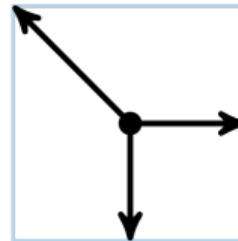
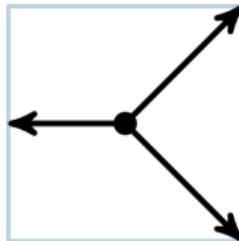
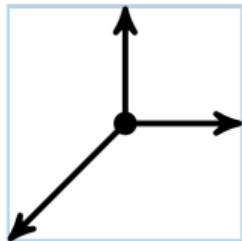
- ensure that the step sizes decrease and allow the method to converge, but not too quickly so as to become stuck away from a local minimum



- stochastic gradient descent
- steepest descent

Mesh Adaptive Direct Search

- Similar to generalized pattern search but uses **random** positive spanning directions
- Example: set of positive spanning sets constructed from nonzero directions $d_1, d_2 \in \{-1, 0, 1\}$.



- Construct lower triangular matrix L sampling from:

$$\{-1/\sqrt{\alpha_k} + 1, -1/\sqrt{\alpha_k} + 2, \dots, 1/\sqrt{\alpha_k} - 1\}$$

- permute rows and columns of L randomly to obtain a matrix D whose columns correspond to n directions that linearly span \mathbb{R}^n . The maximum magnitude among these directions is $1/\sqrt{\alpha_k}$
- add one additional direction $\mathbf{d}_{n+1} = -\sum_{i=1}^n \mathbf{d}_i$ or add n additional directions $\mathbf{d}_{n+j} = -\mathbf{d}_j$
-

$$\alpha_{k+1} \leftarrow \begin{cases} \alpha_k/4 & \text{if no improvement was found in this iteration} \\ \min(1, 4\alpha_k) & \text{otherwise} \end{cases}$$

- If $f(\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha \mathbf{d}) < f(\mathbf{x}_{k-1})$, then the queried point is $\mathbf{x}_{k-1} + 4\alpha \mathbf{d} = \mathbf{x}_k + 3\alpha \mathbf{d}$

Simulated Annealing

- often used on functions with many local minima due to its ability to escape local minima.
- a candidate transition from \mathbf{x} to \mathbf{x}' is sampled from a transition distribution T , eg, multivariate Gaussian

$$\mathbf{x}' = \mathbf{x} + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim T$$

- **Metropolis acceptance criterion:**

$$p(\mathbf{x}, \mathbf{x}') = \begin{cases} 1 & \text{if } \Delta \leq 0 \\ e^{-\frac{\Delta}{t}} & \text{if } \Delta > 0 \end{cases}$$

$$\Delta = f(\mathbf{x}') - f(\mathbf{x})$$

Annealing Plan

- a logarithmic annealing schedule

$$t_k = t_0 \frac{\ln(2)}{\ln(k+1)}$$

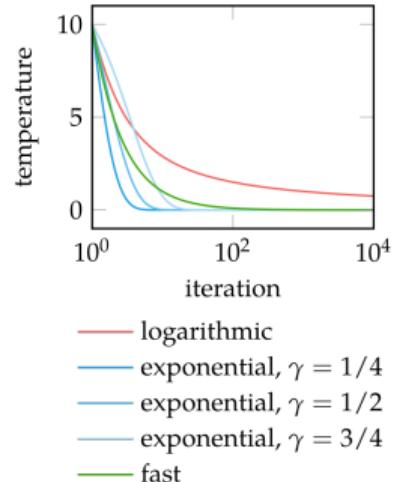
guaranteed to asymptotically reach the global optimum under certain conditions, but it can be slow in practice.

- exponential annealing schedule, more common, uses a simple decay factor:

$$t_{k+1} = \gamma t_k$$

- fast annealing

$$t_k = \frac{t_0}{k}$$



Simulated Annealing

- Corana et al 1987 introduced variable step-size v (separate directional components)
- cycle of random moves, one in each direction

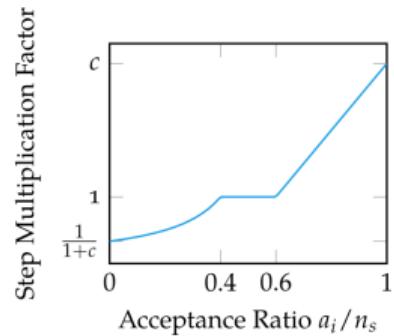
$$\mathbf{x}' = \mathbf{x} + rv_i \mathbf{e}_i$$

where r is randomly sampled from $\{-1, 1\}$

- after n_s cycles, step size is adjusted according to

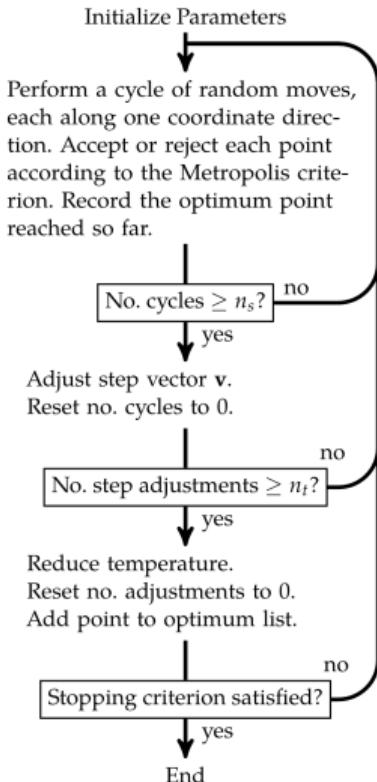
$$v_i = \begin{cases} v_i \left(1 + c_i \frac{a_i/n_s - 0.6}{0.4}\right) & \text{if } a_i > 0.6n_s \\ v_i \left(1 + c_i \frac{0.4 - a_i/n_s}{0.4}\right)^{-1} & \text{if } a_i < 0.4n_s \\ v_i & \text{otherwise} \end{cases}$$

a : accepted steps in each direction; c : typically 2.



regulates the ratio of accepted-to-rejected points to about 50%.

Simulated Annealing



- Temperature reduction occurs every n_t step adjustments, which is every $n_s \cdot n_t$ cycles
- termination when the temperature sank low and no improvement expected or when no movement more than ϵ in last n_ϵ iterations

Cross-Entropy Method

- Maintains explicit probability distribution over design space often called a **proposal distribution**
- Requires choosing a family of parameterized distributions
- At each iteration, a set of design points are **conditionally independently** sampled from the proposal distribution; these are evaluated and ranked
- The best-performing subset of samples, called **elite samples**, are retained
- The proposal distribution parameters are then **updated** based on the elite samples, and the next iteration begins

Cross-Entropy Method

- Cross-entropy is a measure of divergence between two probability distributions p and q (related to Kullback-Leibler divergence)
- Here we measure cross-entropy in a case where one distribution (the one of optimal solutions) is unknown.
- A model is created and then its cross-entropy is measured on the elite set to assess how accurate the model is in predicting this set.
- Let q be the true distribution of the optimal solutions, and p the distribution of solutions as predicted by the model. Since the true distribution is unknown, cross-entropy cannot be directly calculated. Instead, an estimate of cross-entropy is:

$$H(T, p) = - \sum_{i=1}^N \frac{1}{N} \log_2 p(x_i)$$

where N is the size of the elite set, and $p(x)$ is the probability of solution x estimated from the training set T .

Cross-Entropy Method

cross-entropy \equiv Maximum likelihood estimation

- A widely used frequentist estimator is maximum likelihood, in which θ is set to the value that maximizes the **likelihood function** $p(x | \theta)$.
- This corresponds to choosing the value of θ for which the probability of the observed data set is maximized.
- In the machine learning literature, the **negative log of the likelihood function** is called an error function. Because the negative logarithm is a monotonically decreasing function, maximizing the likelihood is equivalent to minimizing the error.
- Suppose our data set consists of N data points $x = \{x_1, \dots, x_N\}$:

$$\mathcal{L}(x | \theta) = p(x | \theta) = \prod_{i=1}^N p(x_i | \theta) \quad \text{likelihood}$$

$$\mathcal{E}(x | \theta) = -\log \mathcal{L}(x) \quad \text{error function}$$

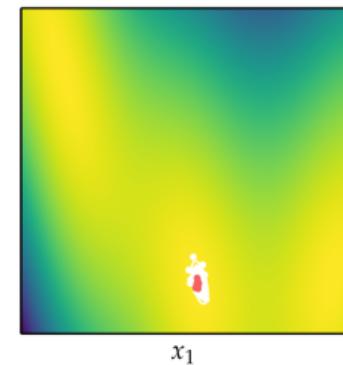
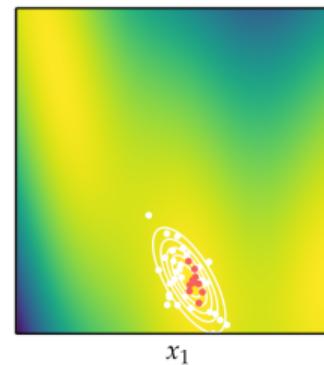
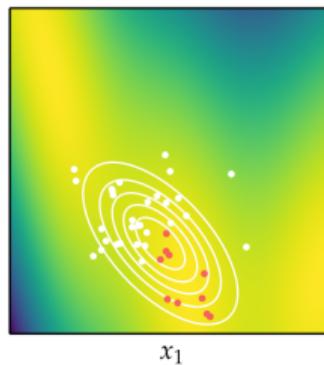
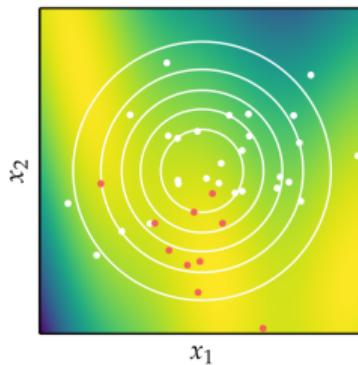
Cross-Entropy Method

-

$$\min_{\theta} \mathcal{E}(x | \theta) = \min_{\theta} (-\log \mathcal{L}(x)) = -\max_{\theta} \log \mathcal{L}(x) \quad \text{maximum log-likelihood}$$

$$\min \left(- \sum_{i=1}^N \log p(x_i | \theta) \right)$$

- hence minimizing the negative of the log-likelihood is equivalent to minimizing the entropy



Multivariate normal distribution

Probability density function

- 1-dimensional:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

- 2-dimensional:

$$f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2[1-\rho^2]}\left[\left(\frac{x - \mu_X}{\sigma_X}\right)^2 - 2\rho\left(\frac{x - \mu_X}{\sigma_X}\right)\left(\frac{y - \mu_Y}{\sigma_Y}\right) + \left(\frac{y - \mu_Y}{\sigma_Y}\right)^2\right]\right)$$

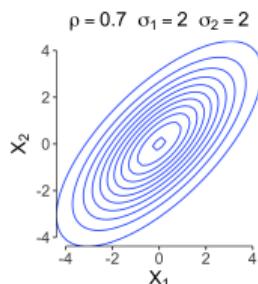
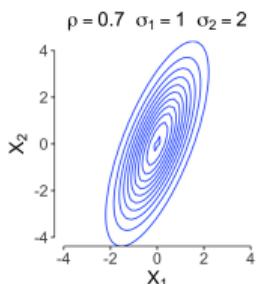
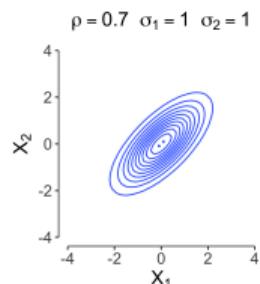
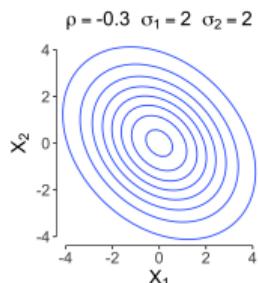
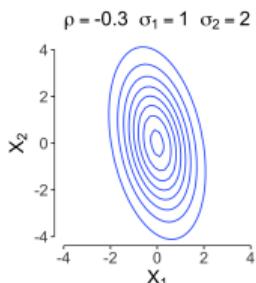
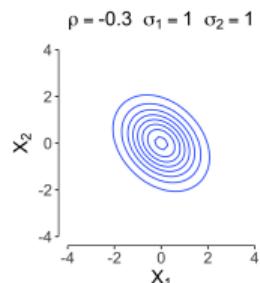
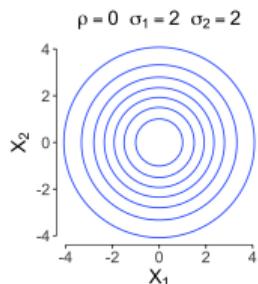
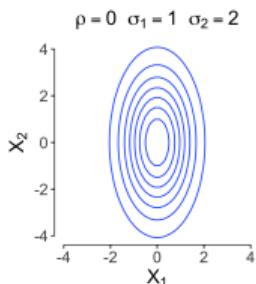
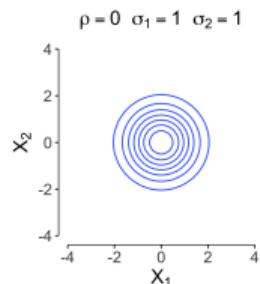
where ρ is the correlation between X and Y and where $\sigma_X > 0$ and $\sigma_Y > 0$. In this case,

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix}.$$

- d dimensional:

$$f_{\mathbf{X}}(x_1, \dots, x_d) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

with symmetric covariance matrix $\boldsymbol{\Sigma}$ positive definite.



(Disadvantage: it is unimodal)

Natural Evolution Strategies

- Similar to cross-entropy method, except instead of parameterizing distribution based on elite samples, it is optimized using **gradient descent**
- The aim is to minimize the expectation

$$E_{\mathbf{x} \sim p(\cdot | \theta)}[f(\mathbf{x})].$$

- The distribution parameter gradient is estimated from the set of function evaluations

Input: $f, \theta, k_{MAX}, N = 100, \alpha = 0.01$

Output: θ

for k in $1, \dots, k_{MAX}$ **do**

 Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ be a conditionally independent sample of size N from $p(\theta)$;

$\theta_{k+1} = \theta_k - \alpha \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i) \nabla_{\theta} \log p(\mathbf{x}_i, \theta_k);$

Natural Evolution Strategies

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{x \sim p(\cdot | \theta)}[f(x)] &= \int \nabla_{\theta} p(x | \theta) f(x) dx \\&= \int \frac{p(x | \theta)}{p(x | \theta)} \nabla_{\theta} p(x | \theta) f(x) dx \\&= \int p(x | \theta) \nabla_{\theta} \log p(x | \theta) f(x) dx \\&= \mathbb{E}_{x \sim p(\cdot | \theta)}[f(x) \nabla_{\theta} \log p(x | \theta)] \\&\approx \frac{1}{m} \sum_{i=1}^m f(x^{(i)}) \nabla_{\theta} \log p(x^{(i)} | \theta)\end{aligned}$$

Natural Evolution Strategies

The multivariate normal distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ is a popular distribution family due to having analytic solutions. The likelihood in d dimensions has the form

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

where $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$. The log likelihood is

$$\log p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

The parameters can be updated using their log likelihood gradients:

$$\nabla_{(\boldsymbol{\mu})} \log p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

$$\nabla_{(\boldsymbol{\Sigma})} \log p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} - \frac{1}{2} \boldsymbol{\Sigma}^{-1}$$

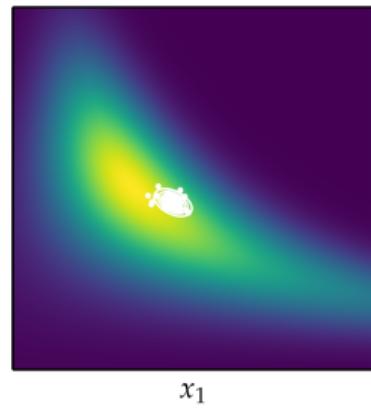
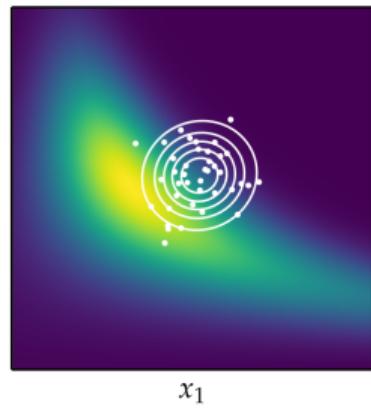
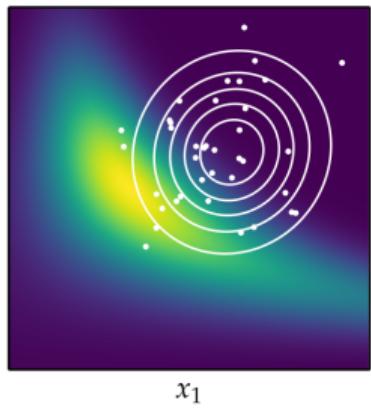
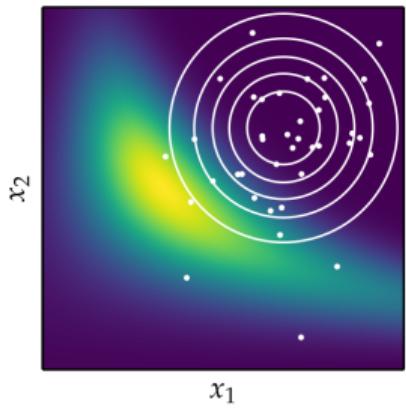
The term $\nabla_{(\boldsymbol{\Sigma})}$ contains the partial derivative of each entry of $\boldsymbol{\Sigma}$ with respect to the log likelihood.

Natural Evolution Strategies

Directly updating Σ may not result in a positive definite matrix, as is required for covariance matrices. One solution is to represent Σ as a product $\mathbf{A}^\top \mathbf{A}$, which guarantees that Σ remains positive semidefinite, and then update \mathbf{A} instead. Replacing Σ by $\mathbf{A}^\top \mathbf{A}$ and taking the gradient with respect to \mathbf{A} yields:

$$\nabla_{(\mathbf{A})} \log p(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{A}) = \mathbf{A} \left[\nabla_{(\Sigma)} \log p(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma) + \nabla_{(\Sigma)} \log p(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma)^\top \right]$$

Natural Evolution Strategies



Covariance Matrix Adaptation Evolutionary Strategy (CMA-ES)

- Same approach as natural evolution strategy and cross entropy method, but the proposal distribution is a multivariate Gaussian parameterized by a covariance matrix.
- At every iteration, m designs are sampled from the multivariate Gaussian:

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma})$$

parameters: mean vector $\boldsymbol{\mu}$, a covariance matrix $\boldsymbol{\Sigma}$, and an additional step-size scalar σ .

- The covariance matrix only increases or decreases in a single direction with every iteration, whereas σ is adapted to control the overall spread of the distribution.
- Design points are sorted $f(\mathbf{x}^{(1)}) \leq f(\mathbf{x}^{(2)}) \leq \dots \leq f(\mathbf{x}^{(m)})$.
- A new mean vector $\boldsymbol{\mu}_{k+1}$ is formed using a weighted average of the first m_e -elite sampled designs:

$$\boldsymbol{\mu}(k+1) \leftarrow \sum_{i=1}^{m_e} w_i \mathbf{x}^{(i)}$$

CMA-ES

- the first m_e elite weights sum to 1, and all the weights approximately sum to 0 and are ordered largest to smallest
- positive and negative weights, more aggressive shift
- The step size σ is updated using a cumulative vector p_1 that tracks steps over time: Comparing the length of p_1 to its expected length under random selection provides the mechanism by which σ is increased or decreased.
- covariance matrix is updated using a cumulative vector p_2 and adjusted weights; the update consists of three components: the previous covariance matrix Σ_k , a rank-one update, and a rank- μ update
Rank-one updates using the cumulation vector allow for correlations between consecutive steps to be exploited
- covariance estimated around original mean μ_k (cross-entropy did it around new mean μ_{k+1})

Summary

- Stochastic methods employ random numbers during the optimization process
- Simulated annealing uses a temperature that controls random exploration and which is reduced over time to converge on a local minimum
- The cross-entropy method and evolution strategies maintain proposal distributions from which they sample in order to inform updates
- Natural evolution strategies uses gradient descent with respect to the log likelihood to update its proposal distribution
- Covariance matrix adaptation is a robust and sample-efficient optimizer that maintains a multivariate Gaussian proposal distribution with a full covariance matrix

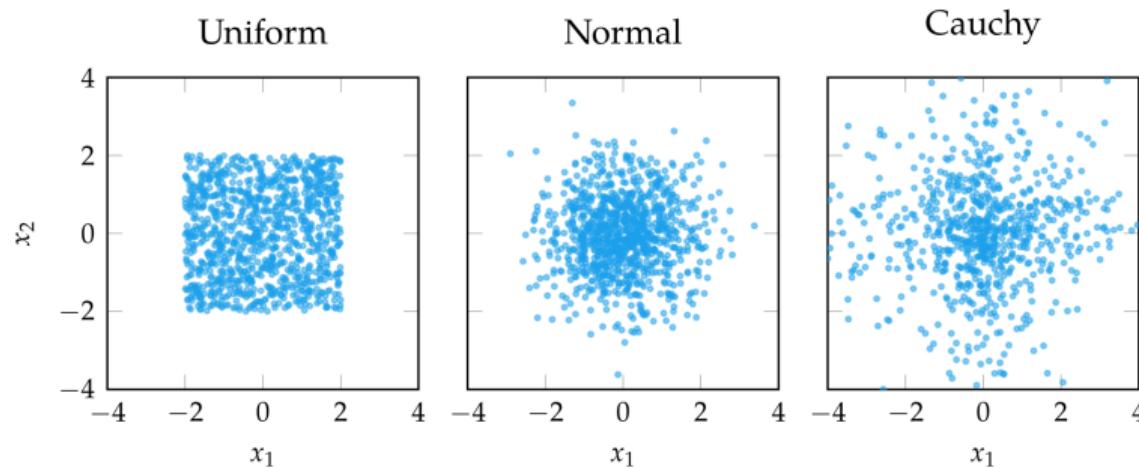
9. Population-based Methods

Population Methods

- Instead of optimizing a single design point, population methods optimize a collection of **individuals**
- A large number of individuals prevents algorithm from being stuck in a local minimum
- Useful information can be shared between individuals
- Stochastic in nature
- Easy to parallelize

Initialization

- Population methods begin with an initial population
- Common initializations are uniform, normal distribution, and Cauchy distribution
- But also space filling designs (later)



Genetic Algorithms

- Inspired by biological evolution where the fittest individuals pass their genetic information to the next generation
- Individuals are interpreted as **chromosomes**
- The fittest individuals are determined by **selection**
- The next generation is formed by selecting the fittest individuals and performing **crossover** and **mutation**

Genetic Algorithms: Chromosomes

- Simplest representation is the binary string chromosome



- Chromosomes are more commonly represented as real-valued chromosomes which are simply real-valued vectors
- Typically initialized randomly

Genetic Algorithms: Selection

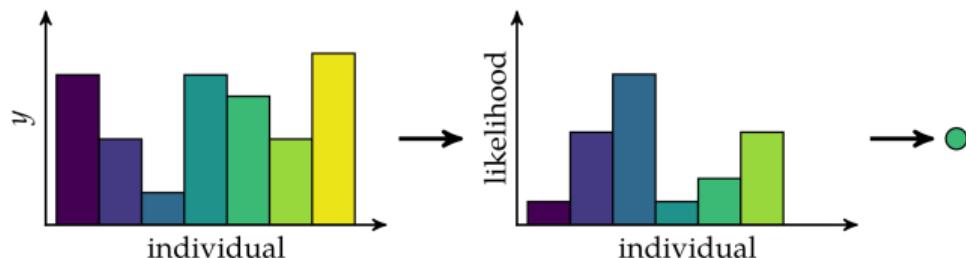
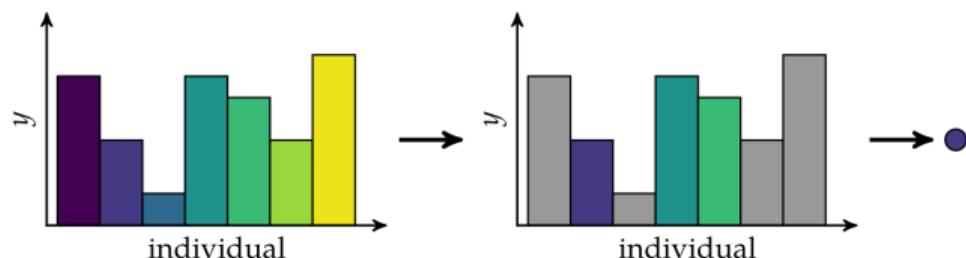
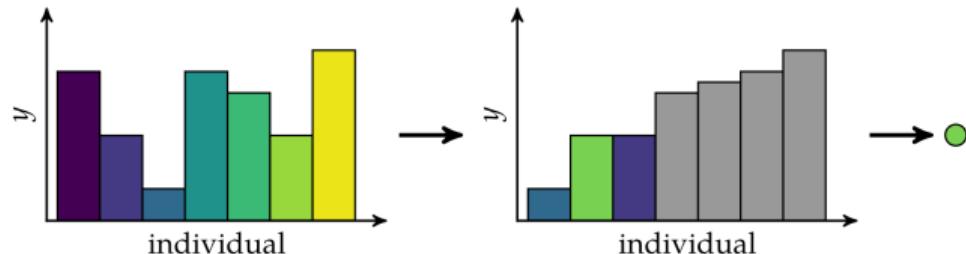
- Determining which individuals pass their genetic information on to the next generation choosing chromosomes to use as parents for the next generation
- Truncation selection: truncate the lowest performers
- Tournament selection: selects fittest out of k randomly chosen individuals
- Roulette Wheel selection: individuals are chosen with probability proportional to their fitness

Genetic Algorithms: Selection

Truncation Selection

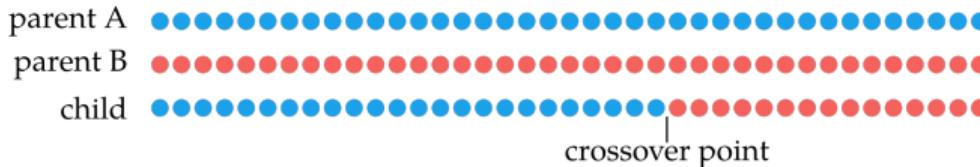
Tournament Selection

Roulette Wheel Selection

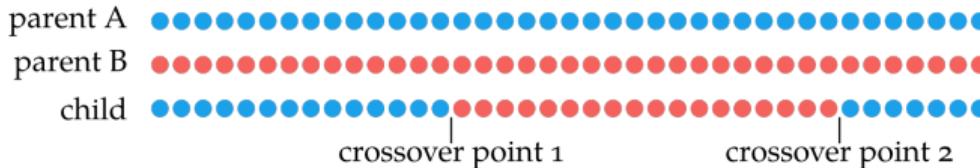


Genetic Algorithms: Crossover

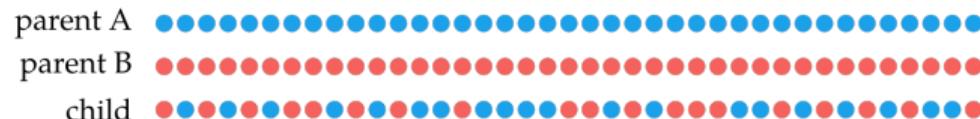
- Combines the chromosomes of the parents to form children
- Single-point crossover: swap occurs after single crossover point



- Two-point crossover: two crossover points



- Uniform crossover: each bit has 50% chance of crossover



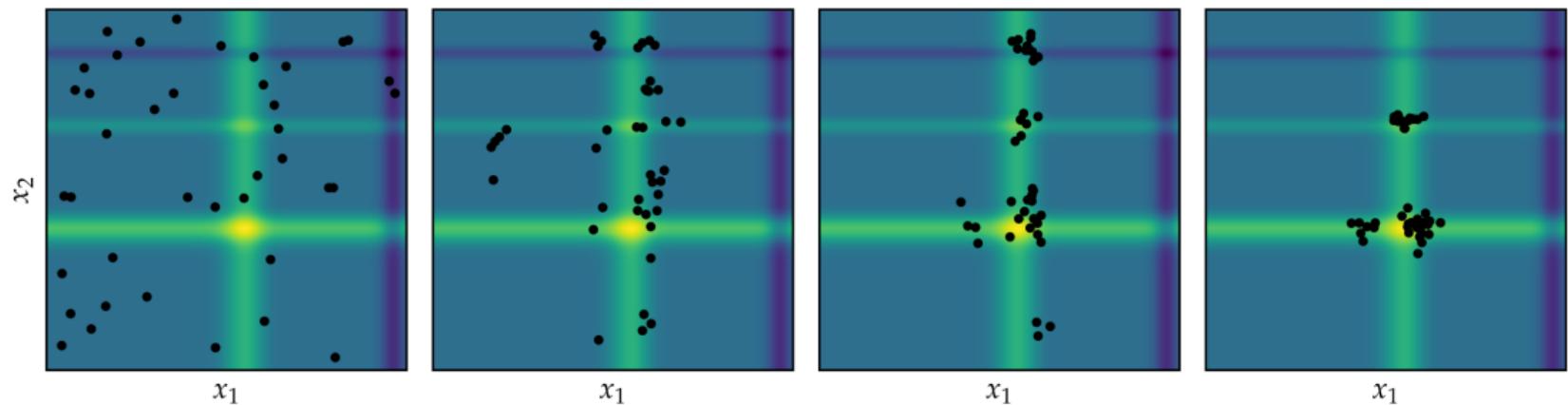
- real values are linearly interpolated between the parents' values x_a and x_b :

Genetic Algorithms: Mutation

- Mutation supports exploration of new areas of design space
- Each bit or real-valued element has a probability of being flipped or modified by noise
- The probability of an element mutating is called **mutation rate**

Genetic Algorithms

Genetic algorithm with truncation selection, single point crossover, and Gaussian mutation applied to Michalewicz function



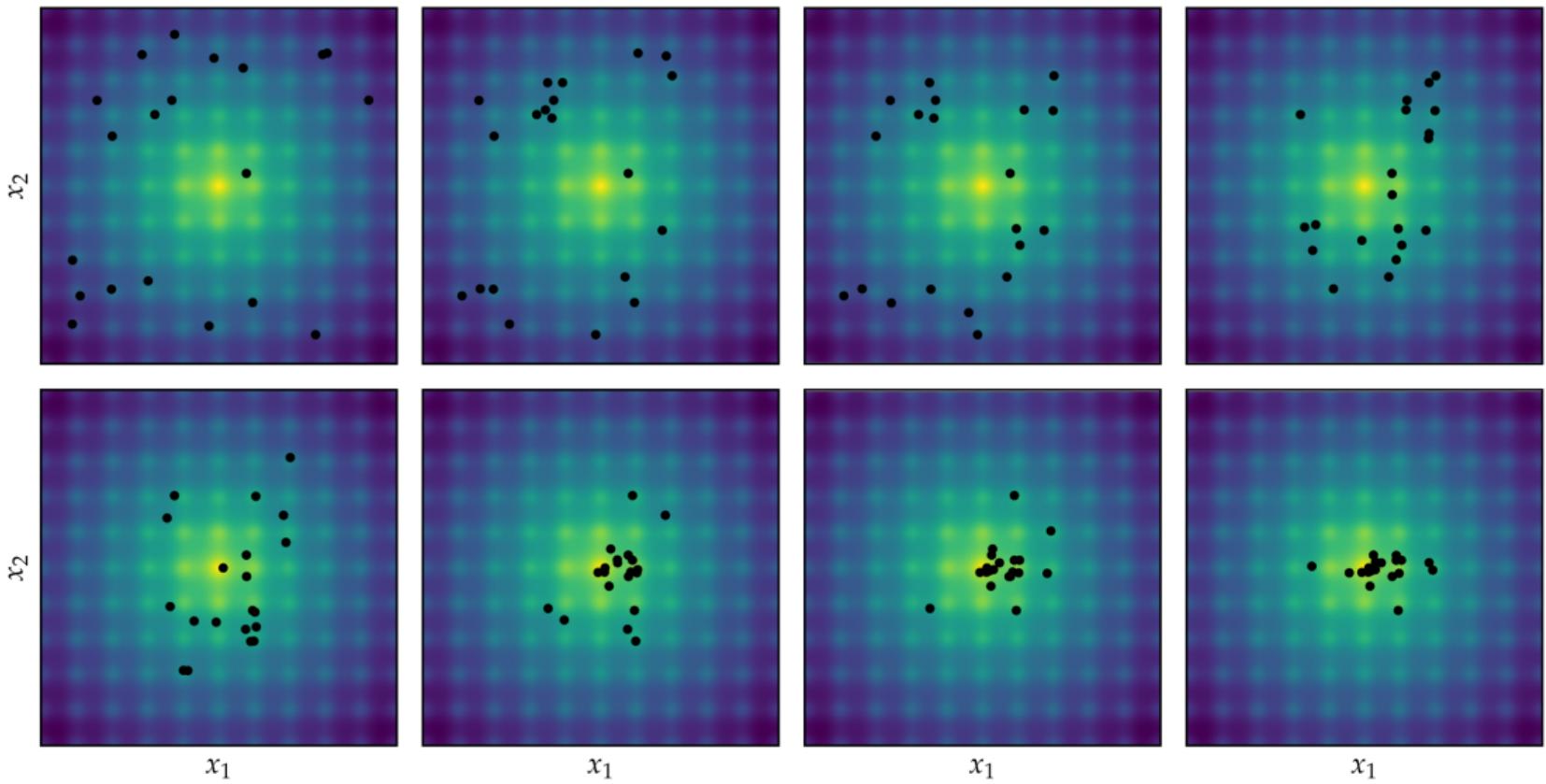
Differential Evolution

Improves each individual \mathbf{x} by recombining other individuals according to a simple formula

1. Choose three random, distinct individuals \mathbf{a} , \mathbf{b} , and \mathbf{c}
2. Construct interim design $\mathbf{z} = \mathbf{a} + w(\mathbf{b} - \mathbf{c})$
3. Choose a random dimension to optimize in
4. Construct candidate \mathbf{x}' via binary crossover of \mathbf{x}' and \mathbf{z}

$$x'_i = \begin{cases} z_i & \text{if } i = j \text{ or with probability } p \\ x_i & \text{otherwise} \end{cases}$$

5. Insert better design between \mathbf{x} and \mathbf{x}' into next generation



Particle Swarm Optimization

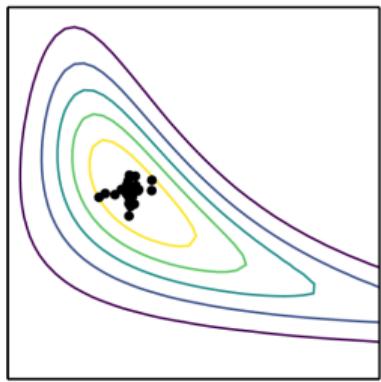
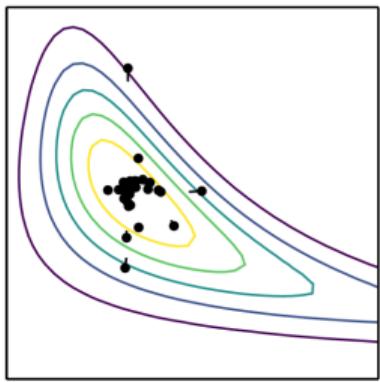
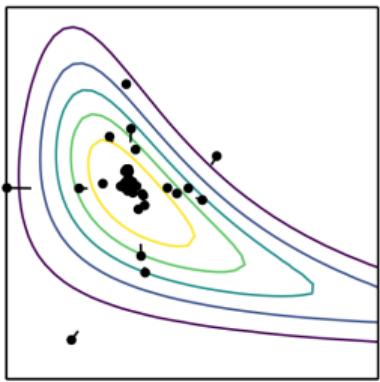
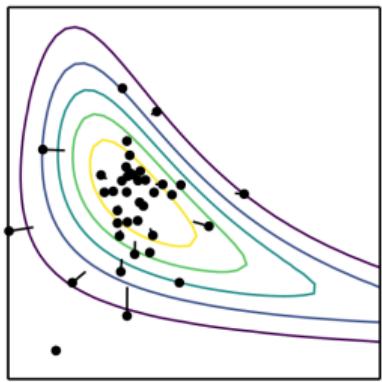
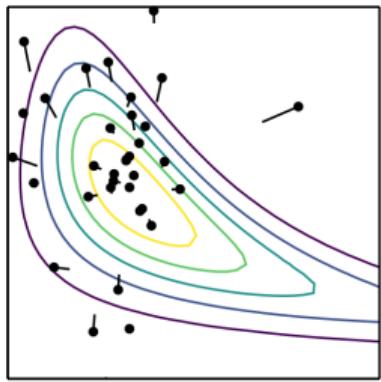
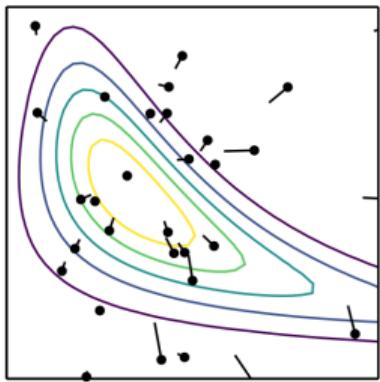
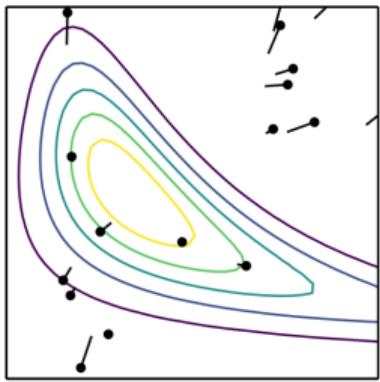
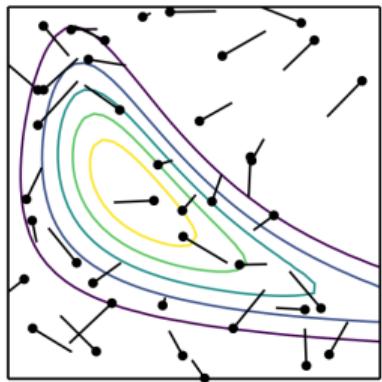
- Each individual, or particle, tracks the following
 - Current position
 - Current velocity
 - Best position seen so far by the particle
 - Best position seen so far by any particle
- At each iteration, these factors produce **force** and **momentum** effects to determine each particle's movement

$$x_i \leftarrow x_i + v_i$$

$$v_i \leftarrow w v_i + c_1 r_1 (x_i^{best} - x_i) + c_2 r_2 (x_i^{best} - x_i)$$

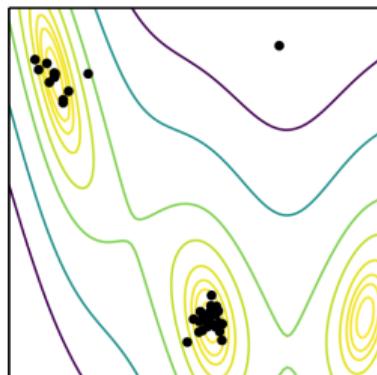
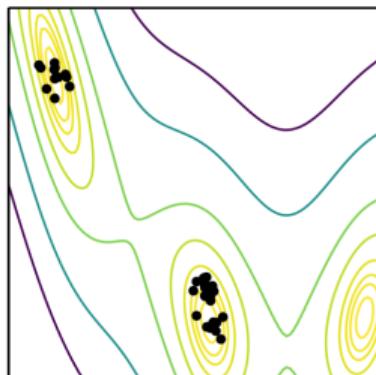
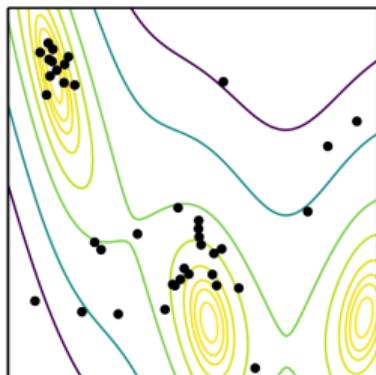
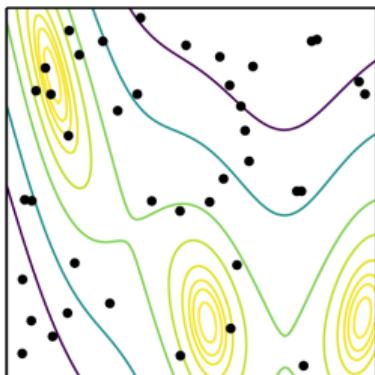
x_{best} is the best location found so far over all particles; w , c_1 , and c_2 are parameters; and r_1 and r_2 are random numbers drawn from $U(0, 1)$

Particle Swarm Optimization



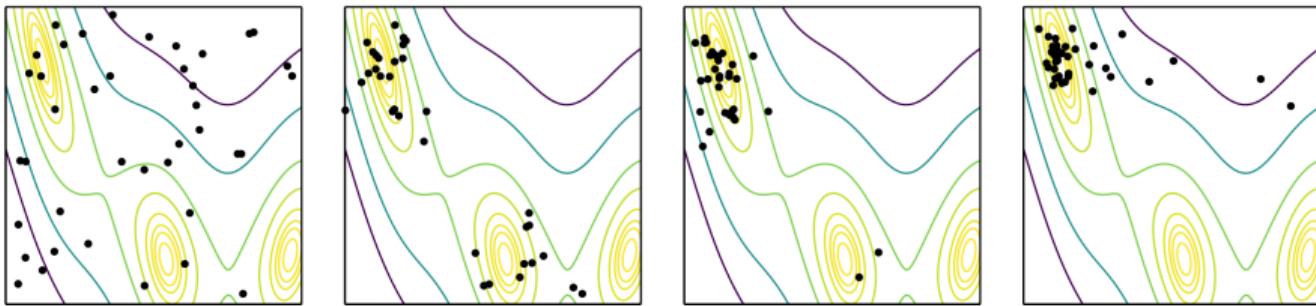
Firefly Algorithm

- Inspired by the way fireflies flash their lights to attract mates
- Attractiveness is determined by low function value
- At each iteration, fireflies move toward the most attractive lights
- Random noise is added to increase exploration



Cuckoo Search

• ...



...

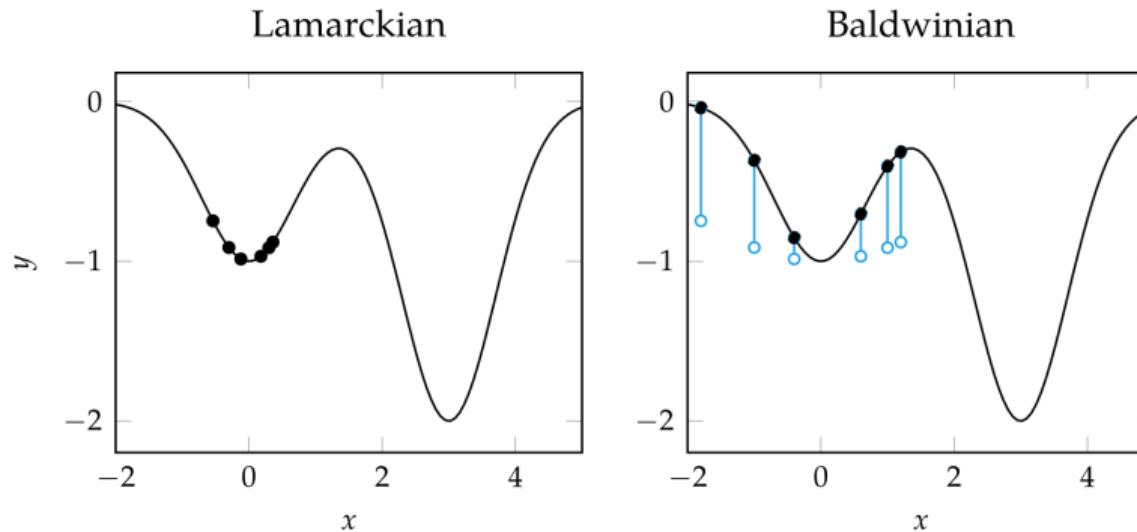
The Evolutionary Computation Bestiary
<http://fcampelo.github.io/EC-Bestiary/>

Hybrid Methods

- Generally, population methods are good at finding the best regions in design space, but do not perform as well as descent methods near the minimizer
- Hybrid methods try to leverage the strength of both methods
- Two hybrid approaches
 - Lamarckian learning
 - Baldwinian learning

Hybrid Methods

- Lamarckian learning
Performs regular descent method update on each individual
- Baldwinian learning
Uses value of descent method update to augment the objective value of each design point



Summary

- Population methods use a collection of individuals in the design space to guide progression toward an optimum
- Genetic algorithms leverage selection, crossover, and mutations to produce better subsequent generations
- Differential evolution, particle swarm optimization, the firefly algorithm, and cuckoo search include rules and mechanisms for attracting design points to the best individuals in the population while maintaining suitable state space exploration
- Population methods can be extended with local search approaches to improve convergence

10. Machine Learning Applications

Introduction

Large-scale machine learning represents a distinctive setting in which traditional nonlinear optimization techniques typically falter

- How do optimization problems arise in machine learning applications and what makes them challenging?
- What have been the most successful optimization methods for large-scale machine learning and why?
- What recent advances have been made in the design of algorithms and what are open questions in this research area?

Two Case Studies

- Logistic regression or support vector machines
convex optimization problems
- deep neural networks
highly nonlinear and nonconvex problems

- batch gradient method
- stochastic gradient method

Enhancements:

- noise reduction methods that attempt to borrow from the strengths of batch methods, such as their fast convergence rates and ability to exploit parallelism;
- methods that incorporate approximate second-order derivative information with the goal of dealing with nonlinearity and ill-conditioning; and
- methods for solving regularized problems designed to avoid overfitting and allow for the use of high-dimensional models.

Text Classification via Convex Optimization

Task: determining whether a text document is one that discusses politics.

- set of examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, where for each $i \in \{1, \dots, n\}$
 \mathbf{x}_i represents the features of a text document (e.g., the words it includes)
 y_i is a label indicating whether the document belongs ($y_i = 1$) or not ($y_i = -1$) to a particular class.
- h prediction function
- measure performance: count how often the program prediction $h(\mathbf{x}_i)$ differs from the correct prediction y_i .
- minimize **empirical risk misclassification**

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h(\mathbf{x}_i) \neq y_i], \quad \text{where } \mathbb{I}[A] = \begin{cases} 1 & \text{if } A \text{ is true,} \\ 0 & \text{otherwise} \end{cases}$$

Text Classification via Convex Optimization

Choosing between prediction functions belonging to a given class by comparing them using cross-validation procedures that involve splitting the examples into three disjoint subsets:

- a training set, optimizing the choice of h by minimizing R_n
- a validation set, generalized performance of each of these remaining candidates is then estimated using the validation set, the best performing of which is chosen as the selected function.
- a testing set, only used to estimate the generalized performance of this selected function

Formalization

- feature vector $\mathbf{x} \in \mathbb{R}^d$ whose components are associated with a prescribed set of vocabulary words; $\|\mathbf{x}\| = 1$
- $h(\mathbf{x}; \mathbf{w}, \tau) = \mathbf{w}^T \mathbf{x} - \tau$, $\mathbf{w} \in \mathbb{R}^d$ and $\tau \in \mathbb{R}^d$
- $\text{sign}(h(\mathbf{x}; \mathbf{w}, \tau))$ discontinuous
- continuous approximation through a loss function that measures a cost for predicting h when the true label is y ;
e.g., one may choose a **log-loss function** of the form

$$L(h, y) = \log(1 + \exp -hy).$$

$$\min_{(\mathbf{w}, \tau) \in \mathbb{R}^d \times \mathbb{R}} L(h(\mathbf{x}_i; \mathbf{w}, \tau), y_i) + \lambda \|\mathbf{w}\|_2^2$$

solve for various λ and choose on the validation set

Deep Neural Networks

Deep Neural Networks: represent hypotheses as computation graphs with tunable weights and compute the gradient of the loss function with respect to those weights in order to fit the training data.

<https://playground.tensorflow.org/>

Perceptual Tasks via Deep Neural Networks

- Prediction function h whose value is computed by applying successive transformations to a given input vector $x_i \in \mathbb{R}^{d_0}$.
- These transformations are made in layers. A canonical fully connected layer performs the computation

$$x_i^{(j)} = s(W_j x_i^{(j-1)} + b_j) \in \mathbb{R}^{d_j}$$

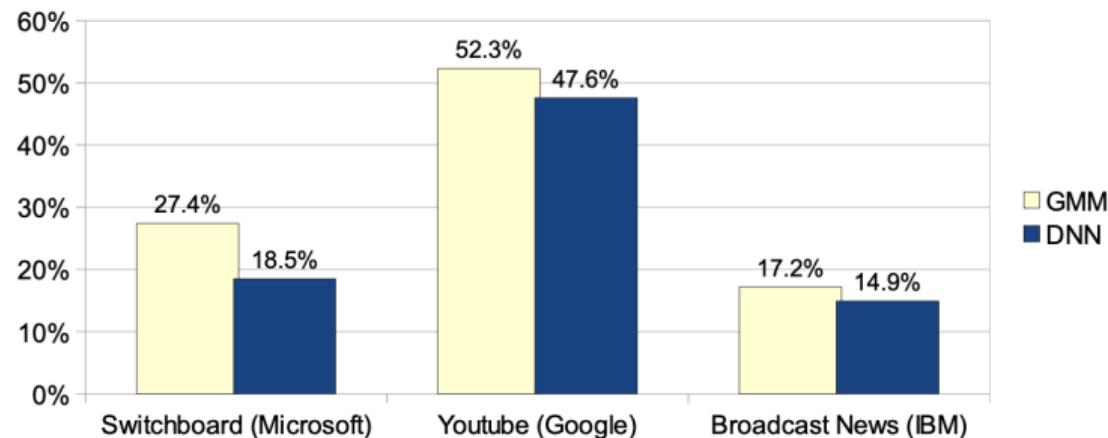
- where $x_i^{(0)} = x_i$, the matrix $W_j \in \mathbb{R}^{d_j \times d_{j-1}}$ and vector $b_j \in \mathbb{R}^{d_j}$ contain the j th layer parameters, and s is a component-wise nonlinear activation function
- $s(x) = 1/(1 + \exp(-x))$ and the hinge function $s(x) = \max(0, x)$ (often called a rectified linear unit (ReLU) in this context)
- $x_i^{(J)}$ leads to the prediction function value $h(x_i; w)$, $w = \{(W_1, b_1), \dots, (W_J, b_J)\}$.
- leads to highly non-linear and non-convex:

$$\min_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^n L(h(x_i; w), y_i)$$

- The gradient with respect to w is made of simple expressions that can be computed by passing information back through the network from the output units.
- the gradient computations for any feedforward computation graph have the same structure as the underlying computation graph.
- gradients can be computed by the chain rule and the algorithmic method of automatic differentiation
- back-propagation in deep learning is simply an application of **reverse mode differentiation**, which applies the chain rule “from the outside in”

Speech recognition

- A contemporary fully connected neural network for speech recognition typically has five to seven layers. This amounts to tens of millions of parameters to be optimized,
- the training may require up to thousands of hours of speech data (representing hundreds of millions of training examples) and weeks of computation on a supercomputer



convolutional neural networks

Convolutional neural networks (CNNs) have proved to be very effective for computer vision and signal processing tasks

ImageNet Large Scale Visual Recognition Competition (ILSVRC) with five convolutional layers and three fully connected layers

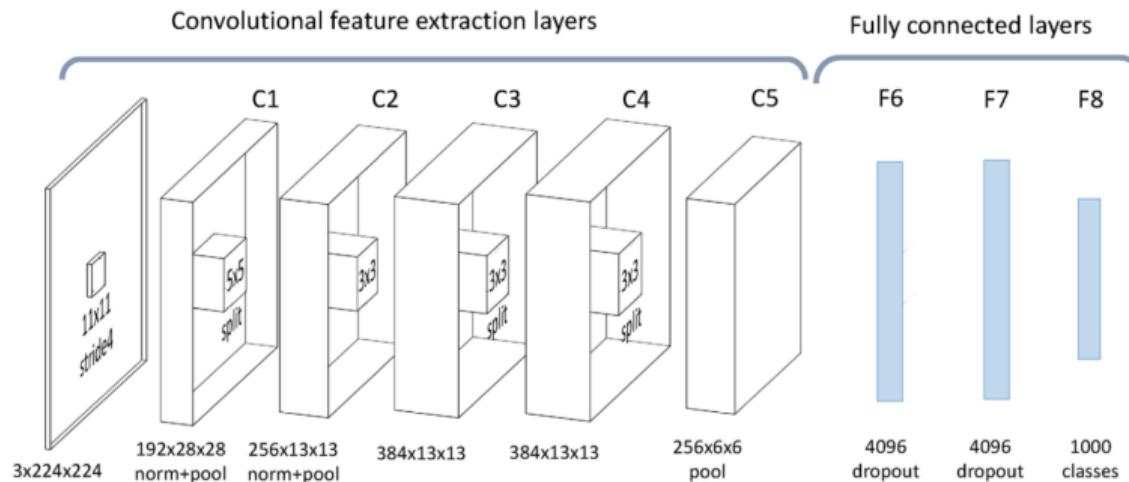


Image Recognition

- input $\mathbf{x}_i^{(j-1)}$ is interpreted as a multichannel image of 224×224 pixels.
- convolutional layers, wherein the parameter matrix \mathbf{W}_j is a circulant matrix
- product $\mathbf{W}_j \mathbf{x}_i^{(j-1)}$ computes the convolution of the image by a trainable filter
- activation functions are piecewise linear functions and can perform more complex operations that may be interpreted as image rectification, contrast normalization, or subsampling.
- output scores represent the odds that the image belongs to each of 1,000 categories.
- 60 million parameters
- training on a few million labeled images takes a few days on a dual GPU workstation.

Fundamentals

- Joint probability distribution function $P(x, y)$ that simultaneously represents the distribution $P(x)$ of inputs as well as the conditional probability $P(y | x)$ of the label y being appropriate for an input x .
- One should seek to find h that yields a small expected risk of misclassification over all possible inputs, i.e., an h that minimizes

$$R(h) = P[h(\mathbf{x}) \neq y] = E[\mathbb{I}[h(\mathbf{x}) \neq y]],$$

which is **variational** since we are optimizing over a set of functions (the h), and is **stochastic** since the objective function involves an expectation.

- without explicit knowledge of P the only tractable option is to construct a surrogate problem that relies solely on the examples $(\mathbf{x}_i, y_i)_{i=1}^n$: minimize the empirical risk

Tasks:

- how to choose the parameterized family of prediction functions \mathcal{H} and
- how to determine (and find) the particular prediction function $h \in \mathcal{H}$ that is optimal.

Choice of Prediction Function

1. \mathcal{H} should contain prediction functions that are able to achieve a low empirical risk over the training set, so as to avoid bias or underfitting the data. (rich family of functions or by using a priori knowledge to select a well-targeted family)
2. the gap between expected risk and empirical risk, namely, $R(h) - R_n(h)$, should be small over all $h \in \mathcal{H}$. (increases with rich family of functions)
3. \mathcal{H} should be efficiently solvable in the corresponding optimization problem (the richer the family of functions and/or the larger training set the more complex the problem becomes)

Choice of Prediction Function

Uniform laws of large numbers and the Hoeffding inequality guarantee that with probability at least $1 - \eta$

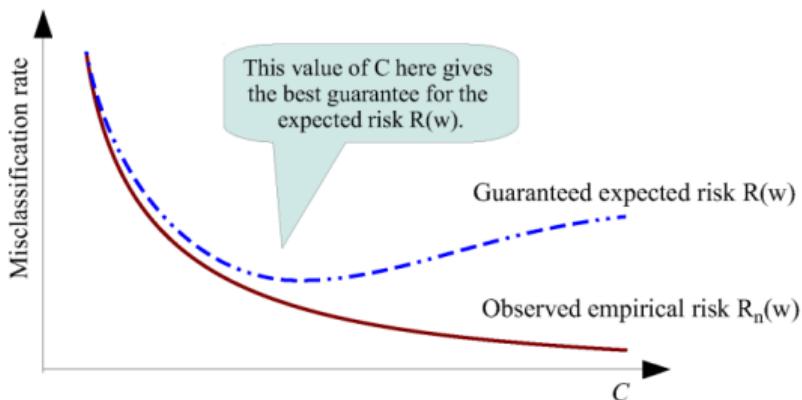
$$\sup_{h \in \mathcal{H}} |R(h) - R_n(h)| \leq \mathcal{O} \left(\sqrt{\frac{1}{2n} \log \left(\frac{2}{\eta} \right)} + \frac{d_{\mathcal{H}}}{n} \log \left(\frac{n}{d_{\mathcal{H}}} \right) \right)$$

- $d_{\mathcal{H}}$ Vapnik-Chervonenkis (VC) dimension (measure of capacity of separating points)
 - not the same as number of parameters
-
- for a fixed $d_{\mathcal{H}}$, uniform convergence is obtained by increasing the number of training points n .
 - for a fixed n , the gap can widen for larger $d_{\mathcal{H}}$.

In practice it is typically easier to estimate with cross-validation experiments.

Structural Risk Minimization

- Rather than choosing a generic family of prediction functions one chooses a structure, i.e., a collection of nested function families.
- structure can be formed as a collection of subsets of a given family \mathcal{H} : given a preference function Ω , choose various values of a hyperparameter C , according to each of which one obtains the subset $\mathcal{H}_C \stackrel{\text{def}}{=} \{h \in \mathcal{H} : \omega(h) \leq C\}$. (C is, eg, degree of a polynomial model function, dimension of an inner layer of a DNN)



$$\begin{aligned} & \min R_n(h) \\ & \text{subject to } \Omega(h) \leq C \end{aligned}$$

Regularized empirical risk

$$R_n(h) + \lambda \Omega(h)$$

validation set is then used to estimate the expected risk corresponding to each C and to choose one.

Structural Risk Minimization

Another approach:

- employ an algorithm for minimizing R_n , but terminate the algorithm early, i.e., before an actual minimizer of R_n is found.
The hyperparameter is played by the training time allowed
- often essential due to computational budget limitations.



Formal Optimization Problem Statements

- We do not consider a variational optimization problem over \mathcal{H} ,
- instead we assume that the prediction function h has a fixed form and is parameterized by a real vector $w \in \mathbb{R}^d$
- for some given $h(\cdot; \cdot) : \mathbb{R}^{d_x} \times \mathbb{R}^d \rightarrow \mathbb{R}^{d_y}$, we consider the family of prediction functions
$$\mathcal{H} \stackrel{\text{def}}{=} \{h(\cdot; w) : w \in \mathbb{R}^d\}$$
- aim to find $h \in \mathcal{H}$ that minimizes a given **loss function** $L : \mathbb{R}^{d_x} \times \mathbb{R}^d \rightarrow \mathbb{R}^{d_y}$, $L(h(x; w), y)$
- Ideally, the expected loss is defined over **any** input-output pair. Assuming probability distribution $P(x, y)$ represents the true input-output relationship:

$$R(w) = \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} L(h(x; w), y) dP(x, y) = E[L(h(x; w), y)] \quad \text{Expected Risk}$$

Formal Optimization Problem Statements

- In practice, one seeks the solution of a problem that involves an estimate of the expected risk R .
- In supervised learning, we have access (either all-at-once or incrementally) to a set of $n \in \mathbb{N}$ independently drawn input-output samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$, with which we define the **empirical risk function** $R_n : \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$R_n(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n L(h(\mathbf{x}_i; \mathbf{w}), y_i) \quad \text{Empirical Risk}$$

(note that before we used misclassification error while now L .)

11. Machine Learning Applications

Simplified Notation

Let ξ be a random seed or the realization of a single (or a set of) sample (x, y) .

For a given (w, ξ) let $f(w; \xi)$ be the composition of the loss function L and the prediction function h

Then:

$$R(w) = \mathbb{E}_\xi[f(w; \xi)] \quad \text{Expected Risk}$$

Let $\{\xi_{[i]}\}_{i=1}^n$ be realizations of ξ corresponding to $\{(x_i, y_i)\}_{i=1}^n$ and $f_i(w) \stackrel{\text{def}}{=} f(w; \xi_{[i]})$

Then:

$$R_n(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w) \quad \text{Empirical Risk}$$

Stochastic vs Batch Optimization Methods

Reduction to minimizing R_n , with $\mathbf{w}_0 \in \mathbb{R}^d$ given (deterministic problem)

Stochastic Approach: Stochastic Gradient (Robbins and Monro, 1951)

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k)$$

i_k is chosen randomly from $\{1, \dots, n\}$, $\alpha_k > 0$.

- very cheap iteration only on one sample.
- $\{\mathbf{w}_k\}$ is a stochastic process determined by the random sequence $\{i_k\}$.
- the direction might not always be a descent but if it is a descent direction in **expectation**, then the sequence $\{\mathbf{w}_k\}$ can be guided toward a minimizer of R_n .

Batch Approach: batch gradient, steepest descent, full gradient method:

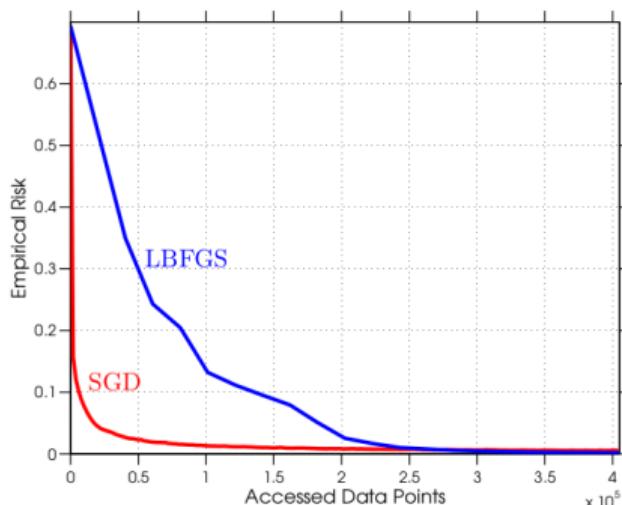
$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \nabla R_n(\mathbf{w}_k) = \mathbf{w}_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_k)$$

- more expensive
- can use all deterministic gradient-based optimization methods
- the sum structure opens up to parallelization

Analogues in simulation: stochastic approximation (SA) and sample average approximation (SAA)

Stochastic Gradient

- In case of redundancy using all the sample data in every iteration is inefficient
- Comparison of the performance of a batch L-BFGS method on number of evaluations of a sample gradient $\nabla f_{i_k}(\mathbf{w}_k)$.
- Each set of n consecutive accesses is called an **epoch**.
- The batch method performs only one step per epoch while SG performs n steps per epoch.



the fast initial improvement achieved by SG, followed by a drastic slowdown after 1 or 2 epochs, is common in practice

SG more sensitive to α_k and starting point

if more epochs, batch may become better

Rate of Convergence

Let $\{\mathbf{x}_k\}$ be a sequence in \mathbb{R}^n that converges to \mathbf{x}^* .

The convergence is said to be **Q-linear** (quotient-linear) if there is a constant $r \in (0, 1)$ such that

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r \quad \text{for all } k \text{ sufficiently large}$$

i.e., the distance to the solution \mathbf{x}^* decreases at each iteration by at least a constant factor bounded away from 1 (i.e., < 1).

Example:

sequence $\{1 + (0.5)^k\}$ converges Q-linearly to 1, with rate $r = 0.5$.

Rate of Convergence

The convergence is said to be **Q-superlinear** if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0$$

Example: the sequence $\{1 + k^{-k}\}$ converges superlinearly to 1.

An even more rapid convergence rate:

The convergence is said to be **Q-quadratic** if

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq M \quad \text{for all } k \text{ sufficiently large}$$

where M is a positive constant, not necessarily less than 1.

Example: the sequence $\{1 + (0.5)^{2^k}\}$.

The values of r and M depend not only on the algorithm but also on the properties of the particular problem.

Regardless of these values a quadratically convergent sequence will always eventually converge faster than a linearly convergent sequence.

Rate of Convergence

Superlinear convergence (quadratic, cubic, quartic, etc) is regarded as fast and desirable, while sublinear convergence is usually impractical.

- Quasi-Newton methods for unconstrained optimization typically converge Q-superlinearly
- Newton's method converges Q-quadratically under appropriate assumptions.
- Steepest descent algorithms converge only at a Q-linear rate, and when the problem is ill-conditioned the convergence constant r is close to 1.

Rate of Convergence

A slightly weaker form of convergence:

overall rate of decrease in the error, rather than the decrease over each individual step of the algorithm.

We say that convergence is **R-linear** (root-linear) if there is a sequence of nonnegative scalars $\{v_k\}$ such that

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \{v_k\} \text{ for all } k, \text{ and } \{v_k\} \text{ converges Q-linearly to zero.}$$

Theoretical Motivations

- a batch approach can minimize R_n at a fast rate; e.g., if R_n is strongly convex. A batch gradient method, then there exists a constant $\rho \in (0, 1)$ such that, for all $k \in \mathbb{N}$, the training error follows **linear convergence**

$$R_n(\mathbf{w}_k) - R_n^* \leq \mathcal{O}(\rho^k),$$

- rate of convergence of a basic stochastic method is slower than for a batch gradient; e.g., if R_n is strictly convex and each i_k is drawn uniformly from $\{1, \dots, n\}$, then for all $k \in \mathbb{N}$, SG satisfies the **sublinear convergence property**

$$\mathbb{E}[R_n(\mathbf{w}_k) - R_n^*] = \mathcal{O}(1/k).$$

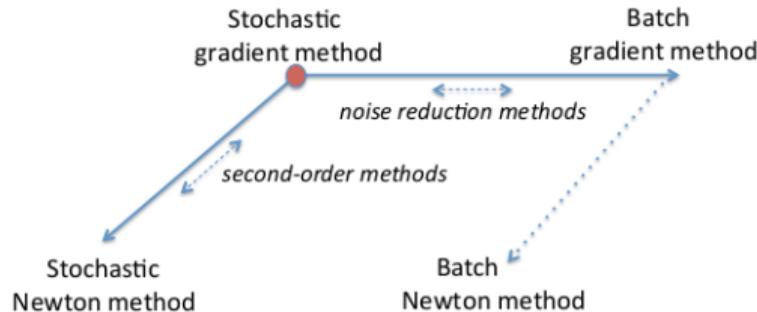
neither the per-iteration cost nor the right-hand side depends on the sample set size n

- in a stochastic optimization setting, SG yields for the expected risk the same convergence rate once substituted $\nabla f_{i_k}(\mathbf{w}_k)$ replaced by $\nabla f(\mathbf{w}_k; \xi_k)$ with each ξ_k drawn independently according to the distribution P

$$\mathbb{E}[R(\mathbf{w}_k) - R^*] = \mathcal{O}(1/k).$$

If $n \gg k$ up to iteration k minimizing R_n same as minimizing R

Beyond SG: Noise Reduction and Second-Order Methods



- on horizontal axis methods that try to improve rate of convergence
- on vertical axis, methods that try to overcome non-linearity and ill-conditioning

Mini-batch Approach small subset of samples, call it $\mathcal{S}_k \subseteq \{1, \dots, n\}$, chosen randomly in each iteration:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \frac{\alpha_k}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{w}_k)$$

due to the reduced variance of the stochastic gradient estimates, the method is easier to tune in terms of choosing the stepsizes $\{\alpha_k\}$.

dynamic sample size and gradient aggregation methods, both of which aim to improve the rate of convergence from sublinear to linear

Outline

17. Analysis of SG

Theoretical Analysis — Preliminaries

convergence properties and worst-case iteration complexity bounds.

$$F(\mathbf{w}) = \begin{cases} R_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) & \text{Empirical Risk} \\ R(\mathbf{w}) = \mathbb{E}_\xi[f(\mathbf{w}; \xi)] & \text{Expected Risk} \end{cases}$$

sampling uniformly with replacement from training set $\rightsquigarrow R_n$
sampling with $P(\xi)$ with replacement from training set $\rightsquigarrow R$.

Procedure SG(...);

Choose an initial iterate \mathbf{w}_0 ;

for $k = 0, 1, \dots$ **do**

Generate a realization of the random variable \mathbf{x}_{ik} ;

Compute a stochastic vector $\mathbf{g}(\mathbf{w}_k, \xi_k)$;

Choose a stepsize $\alpha_k > 0$;

Set the new iterate as $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \mathbf{g}(\mathbf{w}_k, \xi_k)$;

Theoretical Analysis — Preliminaries

ξ_k may represent a single sample or a mini-batch

g may represent a stochastic gradient (biased estimator of $\nabla F(\mathbf{w}_k)$ or a stochastic Netwon or quasi-Newton direction).

$$g(\mathbf{w}_k, \xi_k) = \begin{cases} \nabla f(\mathbf{w}_k; \xi_k) \\ \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(\mathbf{w}_k; \xi_{k,i}) \\ H_k \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(\mathbf{w}_k; \xi_{k,i}) \end{cases}$$

H_k a symmetric positive definite scaling matrix

α_k fixed stepsize or diminishing stepsizes

\mathbf{w}_k can have influence on the sample selection (active learning)

Convergence Analysis – Assumptions

- Assumption 4.1 Lipschitz-continuous objective gradients
- Assumption 4.3 First and second moment limits. The objective function and SG (Algorithm 4.1) satisfy the following:
 - objective function to be bounded below by a scalar F_{\inf} over the region explored by the algorithm.
 - in expectation, the vector $-g(\mathbf{w}_k, \xi_k)$ is a direction of sufficient descent for F from \mathbf{w}_k with a norm comparable to the norm of the gradient
 - the variance of $g(\mathbf{w}_k, \xi_k)$ is restricted, but in a relatively minor manner.

$$\text{Var}_{\xi_k}[g(\mathbf{w}_k, \xi_k)] \leq M + M_V \|\nabla F(\mathbf{w}_k)\|_2^2, \quad M > 0, M_V > 0 \text{ for all } k \in \mathbb{N}$$

- Lemma: Markovian manner in the sense that \mathbf{w}_{k+1} is a random variable that depends only on the iterate \mathbf{w}_k , the seed ξ_k , and the stepsize α_k and not on any past iterates.

Convergence Analysis – Assumptions

- Assumption 4.5 Strong convexity.

The objective function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ is **strongly convex** in that there exists a constant $c > 0$ such that

$$F(\bar{\mathbf{w}}) \geq F(\mathbf{w}) + \nabla F(\mathbf{w})^T (\bar{\mathbf{w}} - \mathbf{w}) + \frac{1}{2}c \|\bar{\mathbf{w}} - \mathbf{w}\|_2^2 \quad \text{for all } (\bar{\mathbf{w}}, \mathbf{w}) \in \mathbb{R}^d \times \mathbb{R}^d$$

or equivalently if there exists $c > 0$:

$$\nabla^2 F(\mathbf{w}) \succeq c$$

(for univariate case: $f''(w) \geq c$), ie, grows at least quadratically.

Hence, F has a unique minimizer, denoted as $\mathbf{w}^* \in \mathbb{R}^d$ with $F^* \stackrel{\text{def}}{=} F(\mathbf{w}^*)$.

Convergence Analysis – Results

- Theorem 4.6 (Strongly Convex Objective, Fixed Stepsize).
- Theorem 4.7 (Strongly Convex Objective, Diminishing Stepsizes)
SG with diminishing step size converges in expectation.
 - role of strong convexity
 - role of initial point
 - trade-offs of mini batches
- Theorem 4.8 (Nonconvex Objective, Fixed Stepsize)
 - While one cannot bound the expected optimality gap as in the convex case, inequality (4.28b) bounds the average norm of the gradient of the objective function observed on $\{\mathbf{w}_k\}$ visited during the first K iterations.
 - classical result for the full gradient method applied to nonconvex functions, namely, that the sum of squared gradients remains finite, implying that

$$\{\|\nabla F(\mathbf{w}_k)\|_2\} \rightarrow 0.$$

- Theorem 4.9 (Nonconvex Objective, Diminishing Stepsizes)
 - for the SG method with diminishing stepsizes, the expected gradient norms cannot stay bounded away from zero
 - the weighted average norm of the squared gradients converges to zero even if the gradients are noisy, (i.e., if $M > 0$ in the Variance upper bounding assumption) one can still conclude that the expected gradient norms cannot asymptotically stay far from zero.

Computational Complexity Analysis

- consider a big data scenario with an infinite supply of training examples, but a **limited computational time budget**. what type of algorithm — e.g., a simple SG or batch gradient method — would provide the best guarantees in terms of achieving a low expected risk?
- $\mathbf{w}^* \in \operatorname{argmin} R(\mathbf{w})$; $\mathbf{w}_n \in \operatorname{argmin} R_n(\mathbf{w})$, $\tilde{\mathbf{w}}_n$ approximate empirical risk minimizer returned by a given optimization algorithm at \mathcal{T}_{\max}
- The tradeoffs associated with this scenario can be formalized as choosing the family of prediction functions \mathcal{H} , the number of examples n , and the optimization accuracy $\epsilon \stackrel{\text{def}}{=} E[R_n(\tilde{\mathbf{w}}_n) - R_n(\mathbf{w}_n)]$ in order to minimize the total error:

$$\underset{\mathcal{H}, n \in \mathbb{N}, \epsilon}{\text{minimize}} E[R(\tilde{\mathbf{w}}_n)] = \overbrace{R(\mathbf{w}^*)}^{\mathcal{E}_{app}(\mathcal{H})} + \overbrace{E[R(\mathbf{w}_n) - R(\mathbf{w}^*)]}^{\mathcal{E}_{est}(\mathcal{H}, n)} + \overbrace{E[R(\tilde{\mathbf{w}}_n) - R(\mathbf{w}_n)]}^{\mathcal{E}_{opt}(\mathcal{H}, n, \epsilon)}$$

subject to $\mathcal{T}(n, \epsilon) \leq \mathcal{T}_{\max}$

Computational Complexity Analysis

- SG, with its sublinear rate of convergence, is more efficient for large-scale learning than (full, batch) gradient-based methods that have a linear rate of convergence.
- reducing the optimization error $\mathcal{E}_{opt}(H, n, \epsilon)$ (evaluated with respect to R rather than R_n) one might need to make up for the additional computing time by: (i) reducing the sample size n , potentially increasing the estimation error $\mathcal{E}_{est}(H, n)$; or (ii) simplifying the function family \mathcal{H} , potentially increasing the approximation error $\mathcal{E}_{app}(\mathcal{H})$.

Computational Complexity Analysis

Keep fixed \mathcal{H} carrying out a worst-case analysis on the influence of the sample size n and optimization tolerance ϵ , which together only influence the estimation and optimization errors.

	Batch	Stochastic
$\mathcal{T}(n, \epsilon)$	$\sim n \log\left(\frac{1}{\epsilon}\right)$	$\frac{1}{\epsilon}$
\mathcal{E}^*	$\sim \frac{\log(\mathcal{T}_{\max})}{\mathcal{T}_{\max}} + \frac{1}{\mathcal{T}_{\max}}$	$\frac{1}{\mathcal{T}_{\max}}$

A stochastic optimization algorithm performs better than batch stochastic in terms of expected error

Large gap between asymptotical behavior and practical realities.

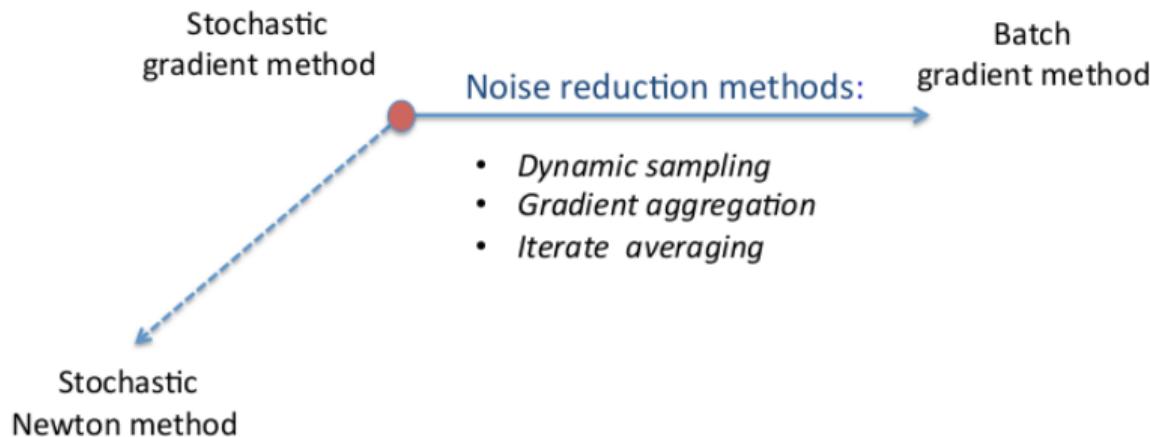
Remarks

- Fragility of the Asymptotic Performance of SG
ok if objective function it includes a squared L_2 -norm regularizer (related to constant c) but regularization parameter should be lowered when the number of samples increases.
- SG good for GPUs but ill-conditioning erodes efficiency of SG
- Distributed computing not working with basic SG because of too frequent updates of w , more promising with mini-batch.
- Alternatives with Faster Convergence: minimizing empirical risk R_n there is information from previous gradients.
 - gradient aggregation methods
 - dynamic sampling approach

12. Machine Learning Applications

Noise Reduction Methods

- SG as the ideal optimization approach for large-scale applications.
- SG suffers from the adverse effect of **noisy gradient estimates**.
 - when fixed stepsizes are used it prevents SG from converging to the solution
 - when a diminishing stepsize sequence $\{\alpha_k\}$ is employed it leads to a slow, sublinear rate of convergence.
- Remedies:



Overview

Achieve linear rate of convergence to the optimal value using a fixed stepsize.

- **Dynamic sampling methods** achieve noise reduction by gradually increasing the mini-batch size used in the gradient computation, thus employing increasingly more accurate gradient estimates as the optimization process proceeds.
- **Gradient aggregation methods** improve the quality of the search directions by storing gradient estimates corresponding to samples employed in previous iterations, updating one (or some) of these estimates in each iteration, and defining the search direction as a weighted average of these estimates.

Rate of convergence remains sublinear but reduces variance of iterates

- **iterate averaging methods** maintain an average of iterates computed during the optimization process and employs a more aggressive stepsize sequence—of order $O(1/\sqrt{k})$ rather than $O(1/k)$.

Reducing Noise at a Geometric Rate

rate of decrease in noise that allows a stochastic-gradient-type method to converge at a linear rate.

Consequence of Lipschitz assumption with ℓ constant:

$$\mathbb{E}_{\xi_k}[F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\alpha_k \nabla F(\mathbf{w}_k)^T \mathbb{E}_{\xi_k}[g(\mathbf{w}_k, \xi_k)] + \frac{1}{2} \alpha_k^2 \ell \mathbb{E}_{\xi_k}[\|g(\mathbf{w}_k, \xi_k)\|_2^2]$$

We want to make the left hand side small (sequence of expected optimality gaps).

Theorem 5.1 (Strongly Convex Objective, Noise Reduction)

The SG method with a fixed stepsize $\bar{\alpha}$ and previous assumptions plus a variance of the stochastic vectors that decreases geometrically

$$\text{Var}_{\xi_k}[g(\mathbf{w}_k, \xi_k)] \leq M\zeta^{k-1}$$

has a sequence of expected optimality gaps that vanishes at a linear rate:

$$\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \omega\rho^{k-1}$$

Dynamic Sample Size Methods

Can we design efficient optimization methods attaining the critical bound on the variance?

Mini-batch stochastic gradient:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \bar{\alpha} g(\mathbf{w}_k, \xi_k)$$

where the stochastic directions are computed for some $\tau > 1$ as

$$g(\mathbf{w}_k, \xi_k) \stackrel{\text{def}}{=} \frac{1}{n_k} \sum_{i \in \mathcal{S}_k} \nabla f(\mathbf{w}_k; \xi_{k,i}) \quad \text{with } n_k \stackrel{\text{def}}{=} |\mathcal{S}_k| = \lceil \tau^{k-1} \rceil.$$

the mini-batch size increases geometrically as a function of the iteration counter k

Corollary 5.2. Let $\{\mathbf{w}_k\}$ be the iterates generated with unbiased gradient estimates, i.e., $\mathbb{E}_{\xi_{k,i}}[\nabla f(\mathbf{w}_k; \xi_{k,i})] = \nabla F(\mathbf{w}_k)$ for all $k \in \mathbb{N}$ and $i \in \mathcal{S}_k$. Then, the variance condition is satisfied, and if all other assumptions of Theorem 5.1 hold, then the expected optimality gap vanishes linearly.

Dynamic Sample Size Methods

Note: we described a method as linearly convergent but the per-iteration cost increases without bound.

Recall that SG method needs $\mathcal{T}(n, \epsilon) \leq 1/\epsilon$ evaluations to guarantee $\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \epsilon$

Theorem 5.3 Suppose that the dynamic sampling SG method is run with a stepsize $\bar{\alpha}$ satisfying “some” bounds and some τ . In addition, suppose that all previous Assumptions hold. Then, the total number of evaluations of a stochastic gradient of the form $\nabla f(\mathbf{w}_k; \xi_{k,i})$ required to obtain $\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \epsilon$ is $O(\epsilon^{-1})$.

Dynamic Sample Size Guidelines

Given the rate of convergence of a batch optimization algorithm on strongly convex functions (i.e., linear, superlinear, etc.), what should be the sampling rate so that the overall algorithm is **efficient** in the sense that it results in the lowest computational complexity?

- if the optimization method has a sublinear rate of convergence, then there is no sampling rate that makes the algorithm “efficient”;
- if the optimization algorithm is linearly convergent, then the sampling rate must be geometric (with restrictions on the constant in the rate) for the algorithm to be “efficient”;
- for superlinearly convergent methods, increasing the sample size at a rate that is slightly faster than geometric will yield an “efficient” method.

Design in Practice

- presetting the sampling rate, ie, $\tau > 1$ before running the optimization algorithm, requires some experimentation. Care must be put in preventing the full sample set from being employed too soon
- adaptive mechanisms to produce descent directions sufficiently often
 - any direction $g(\mathbf{w}_k, \xi_k)$ is a descent direction for F at \mathbf{w}_k if, for some $\chi \in [0, 1]$, one has

$$\delta(\mathbf{w}_k, \xi_k) \stackrel{\text{def}}{=} \|g(\mathbf{w}_k, \xi_k) - \nabla F(\mathbf{w}_k)\|_2 \leq \chi \|g(\mathbf{w}_k, \xi_k)\|_2$$

verifying the inequality may be costly because involves the evaluation of $\nabla F(\mathbf{w}_k)$, one can estimate the left-hand side $\delta(\mathbf{w}_k, \xi_k)$, and then choose n_k so it holds sufficiently often.

- The sample variance obtained by sampling without replacements is bounded above by $\chi^2 \|g(\mathbf{w}_k, \xi_k)\|_2^2$
- If this condition is not satisfied, then increase the sample size to a size that one might predict would satisfy such a condition.
- no guarantee that the size n_k increases at a geometric rate. Remedy: if the adaptive increases the sampling rate more slowly than a preset geometric sequence, then a growth in the sample size is imposed.

Gradient Aggregation

- Rather than compute increasingly more **new** stochastic gradient information in each iteration, achieve a lower variance by **reusing and/or revising** previously computed information
- achieve a linear rate of convergence on strongly convex problems.
- improved rate is achieved primarily by either an increase in computation or an increase in storage.
- works on finite sums like R_n

Procedure SVRG ; # Methods for Minimizing an Empirical Risk R_n

Choose an initial iterate $\mathbf{w}_1 \in \mathbb{R}^d$, stepsize $\alpha > 0$ and a positive integer m ;

for $k = 1, 2, \dots$ **do**

 Compute the batch gradient $\nabla R_n(\mathbf{w}_k)$;

 Initialize $\tilde{\mathbf{w}}_1 \leftarrow \mathbf{w}_k$;

for $j = 1, \dots, m$ **do**

$\tilde{\mathbf{g}}_j \leftarrow \nabla f_{i_j}(\tilde{\mathbf{w}}_j) - (\nabla f_{i_j}(\mathbf{w}_k) - \nabla R_n(\mathbf{w}_k))$; # $\nabla R_n(\mathbf{w}_k)$ from batch gradient

$\tilde{\mathbf{w}}_{j+1} \leftarrow \tilde{\mathbf{w}}_j - \alpha \tilde{\mathbf{g}}_j$;

 Option (a): Set $\mathbf{w}_{k+1} = \tilde{\mathbf{w}}_{m+1}$;

 Option (b): Set $\mathbf{w}_{k+1} = \frac{1}{m} \sum_{j=1}^m \tilde{\mathbf{w}}_{j+1}$;

 Option (c): Choose j uniformly from $\{1, \dots, m\}$ and set $\mathbf{w}_{k+1} = \tilde{\mathbf{w}}_{j+1}$;

- since $E_{i_j}[\nabla f_{i_j}(\mathbf{w}_k)] = \nabla R_n(\mathbf{w}_k)$, one can view $\nabla f_{i_j}(\mathbf{w}_k) - \nabla R_n(\mathbf{w}_k)$ as the bias in the gradient estimate $\nabla f_{i_j}(\mathbf{w}_k)$.
- sampled gradient $\nabla f_{i_j}(\tilde{\mathbf{w}}_j)$ is corrected based on a perceived bias. Overall, $\tilde{\mathbf{g}}_j$ represents an unbiased estimator of $\nabla R_n(\tilde{\mathbf{w}}_j)$, but with a variance that one can expect to be smaller than as in simple SG

SAGA

in each iteration, it computes a stochastic vector \mathbf{g}_k as the average of stochastic gradients evaluated at previous iterates.

Procedure SAGA ; # Method for Minimizing an Empirical Risk R_n
Choose an initial iterate $\mathbf{w}_1 \in \mathbb{R}^d$ and stepsize $\alpha > 0$;
for $i = 1, \dots, n$ **do**
 Compute $\nabla f_i(\mathbf{w}_1)$;
 Store $\nabla f_i(\mathbf{w}_{[i]}) \leftarrow \nabla f_i(\mathbf{w}_1)$; # $\mathbf{w}_{[i]}$ represents the latest iterate at which ∇f_i
for $k = 1, 2, \dots$ **do**
 Choose j uniformly in $\{1, \dots, n\}$;
 Compute $\nabla f_j(\mathbf{w}_k)$;
 Set $\mathbf{g}_k \leftarrow \nabla f_j(\mathbf{w}_k) - \nabla f_j(\mathbf{w}_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_{[i]})$;
 Store $\nabla f_j(\mathbf{w}_{[j]}) \leftarrow \nabla f_j(\mathbf{w}_k)$;
 Set $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha \mathbf{g}_k$;

As in SVRG, the method employs unbiased gradient estimates, but with variances that are expected to be less than the stochastic gradients that would be employed in a basic SG routine

SAGA

- Same per-iteration costs as basic SG
- on strongly convex R_n can achieve a linear rate of convergence but needs knowledge of at least ℓ .
- More effective initialization instead of evaluating all the gradients $\{\nabla f_i\}_{i=1}^n$ at the initial point. For example, one could perform one epoch of simple SG, or one can assimilate iterates one-by-one and compute \mathbf{g}_k only using the gradients available up to that point.
- SAGA needs to store n stochastic gradient vectors
- for very large n , gradient aggregation methods are comparable to batch algorithms and therefore cannot beat SG in this regime

Iterated Averaging Methods

- for minimizing a continuously differentiable F with unbiased gradient estimates, the idea is to employ the iteration:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha g(\mathbf{w}_k, \xi_k)$$

$$\tilde{\mathbf{w}}_{k+1} \leftarrow \frac{1}{k+1} \sum_{j=1}^{k+1} \mathbf{w}_j$$

where $\{\tilde{\mathbf{w}}_k\}$ has no effect on the computation of the SG iterate sequence $\{\mathbf{w}_k\}$

- with stepsizes diminishing at a slow rate of $\mathcal{O}(1/(k^a))$ for some $a \in (\frac{1}{2}, 1)$ on strongly convex objectives, yields that $\mathbb{E}[\|\mathbf{w}_k - \mathbf{w}^*\|_2^2] = \mathcal{O}(1/(k^a))$ while $\mathbb{E}[\|\tilde{\mathbf{w}}_k - \mathbf{w}^*\|_2^2] = \mathcal{O}(1/k)$.
- in certain cases this combination of long steps and averaging yields an optimal constant in $\mathbb{E}[\|\tilde{\mathbf{w}}_k - \mathbf{w}^*\|_2^2]$ in the sense that no rescaling of the steps—through multiplication with a positive definite matrix (second order methods) can improve the asymptotic rate or constant.

Second Order Methods

- Address the adverse effects of high nonlinearity and ill-conditioning of the objective function through the use of second-order information.
- improve convergence rates of batch methods or the constants involved in the sublinear convergence rate of stochastic methods
- First-order methods are **not scale invariant**. Consider:
 F continuously differentiable function $F : \mathbb{R}^d \rightarrow \mathbb{R}$

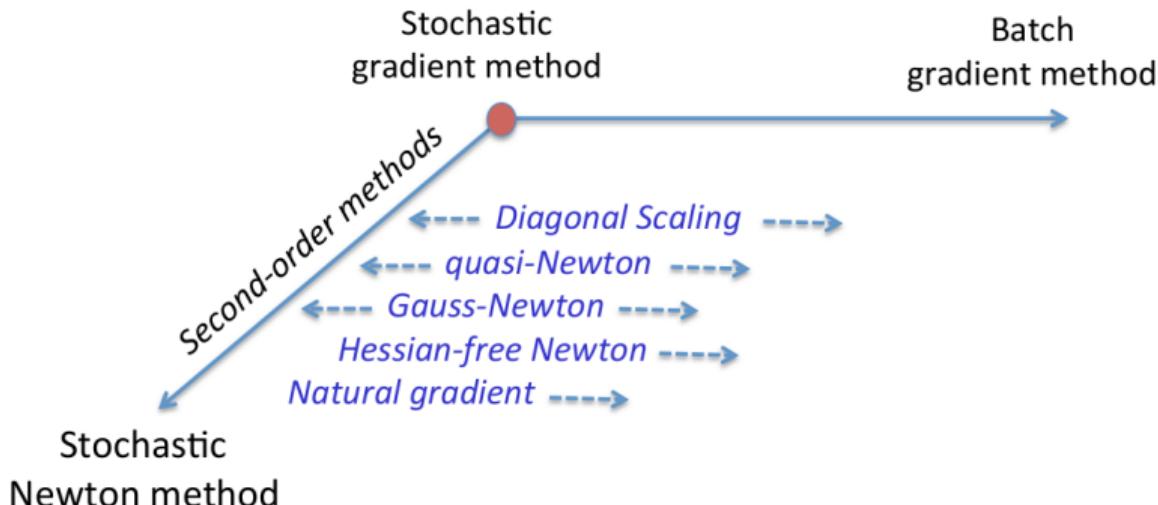
$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \nabla F(\mathbf{w}_k)$$

linear transformation of the variables $\{\mathbf{w}_k\} = \{B\tilde{\mathbf{w}}_k\}$. $\min_{\tilde{\mathbf{w}}} F(B\tilde{\mathbf{w}}_k)$

$$\tilde{\mathbf{w}}_{k+1} \leftarrow \tilde{\mathbf{w}}_k - \alpha_k B \nabla F(B\tilde{\mathbf{w}}_k) \quad \Rightarrow \quad \mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k B^2 \nabla F(\mathbf{w}_k)$$

They will perform differently. With $\alpha = 1$ and $B = (\nabla^2 F(\mathbf{w}_1))^{-1/2}$ we get Newton's method

- Newton's method achieves a quadratic rate of convergence if w_1 is sufficiently close to a strong minimizer. On the other hand, stochastic methods like the SG method cannot achieve a convergence rate that is faster than sublinear, regardless of the choice of B .
- careful use of successive re-scalings based on (approximate) second-order derivatives can be beneficial between the stochastic and batch regimes.



Hessian-Free Inexact Newton Methods

- Newton's method:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + \alpha_k \mathbf{s}_k$$

where \mathbf{s}_k satisfies $\nabla^2 F(\mathbf{w}_k) \mathbf{s}_k = -\nabla F(\mathbf{w}_k)$.

- one can solve the linear system inexactly through an iterative approach such as the conjugate gradient (CG) method.
- By ensuring that the linear solves are accurate enough, such an inexact Newton-CG method can enjoy a superlinear rate of convergence
- For a smooth objective function F , one can compute $\nabla^2 F(\mathbf{w}) \mathbf{d}$ at a cost that is a small multiple of the cost of evaluating ∇F , and without forming the Hessian, which would require $\mathcal{O}(d^2)$ storage
- exploit structure of risk measures

- iterations are more tolerant to noise in the hessian estimate than it is to noise in the gradient estimate
- employs a smaller, conditionally (given w_k) uncorrelated, sample for defining the Hessian than for the stochastic gradient estimate
- can be combined with a backtracking (Armijo) line search or trust region
- (subsampled) Hessian-vector products can be computed efficiently in ML tasks

Example 6.2. Consider a binary classification problem where the training function is given by the logistic loss with an ℓ_2 -norm regularization parameterized by $\lambda > 0$:

$$R_n(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \frac{\lambda}{2} \|w\|^2. \quad (6.8)$$

A (subsampled) Hessian-vector product can be computed efficiently by observing that

$$\nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) d = \frac{1}{|\mathcal{S}_k^H|} \sum_{i \in \mathcal{S}_k^H} \frac{\exp(-y_i w_k^T x_i)}{(1 + \exp(-y_i w_k^T x_i))^2} (x_i^T d) x_i + \lambda d.$$

- Stochastic Quasi-Newton Methods:
Like BFGS
- Gauss-Newton Methods
constructs an approximation to the Hessian using only first-order information, and this approximation is guaranteed to be positive semidefinite, even when the full Hessian itself may be indefinite.
The price to pay for this convenient representation is that it ignores second-order interactions between elements of the parameter vector w , which might mean a loss of curvature information that could be useful for the optimization process.

Summary

- Ways to cope with the problems in machine learning
- SG might not be the best choice for parallelization
- How about other methods like CMA-ES?

13. Constrained Optimization

Constrained Optimization

- Minimizing an objective subject to design point restrictions called **constraints**
- A variety of techniques transform constrained optimization problems into unconstrained problems
- New optimization problem statement

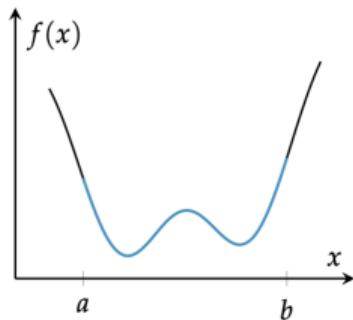
$$\underset{x}{\text{minimize}} \ f(x)$$

subject to $x \in \mathcal{X}$

- The set $\mathcal{X} \subset \mathbb{R}$ is called the **feasible set**

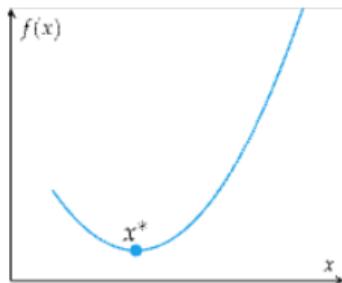
Constrained Optimization

Constraints that bound feasible set can change the optimizer

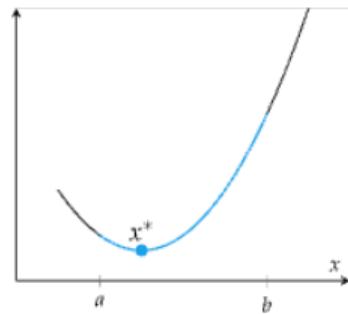


$$\begin{aligned} & \underset{x}{\text{minimize}} \quad f(x) \\ & \text{subject to } x \in [a, b] \end{aligned}$$

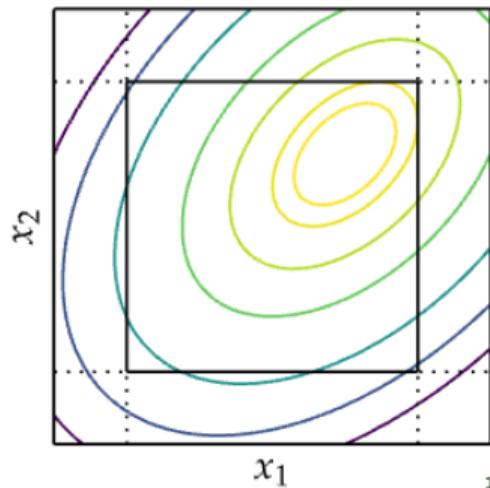
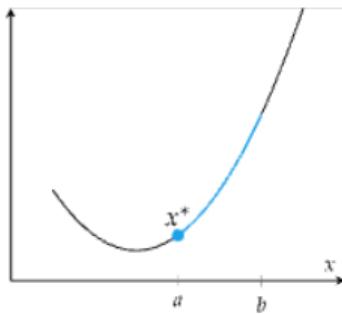
Unconstrained



Constrained, Same Solution



Constrained, New Solution



Constraint Types

- Generally, constraints are formulated using two types:

1. **Equality constraints:** $h(x) = 0$

2. **Inequality constraints:** $g(x) \leq 0$

- Any optimization problem can be written as

$$\underset{x}{\text{minimize}} \ f(x)$$

subject to $g_i(x) \leq 0$ for all i in $\{1, \dots, m\}$

$h_j(x) = 0$ for all j in $\{1, \dots, \ell\}$

$$\underset{x}{\text{minimize}} \ f(x)$$

subject to $g(x) \leq 0$

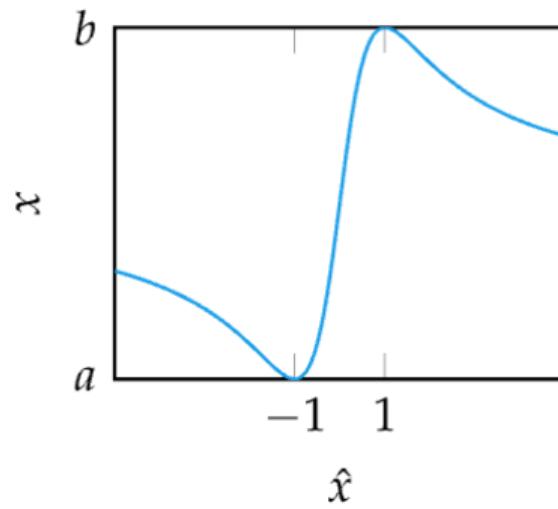
$h(x) = 0$

f and the functions h and g are all smooth, real-valued functions on a subset of \mathbb{R}^n

Transformations to Remove Constraints

- If necessary, some problems can be reformulated to incorporate constraints into the objective function
- If x is constrained between a and b

$$x = t_{a,b}(\hat{x}) = \frac{b+a}{2} + \frac{b-a}{2} \left(\frac{2\hat{x}}{1+\hat{x}^2} \right)$$



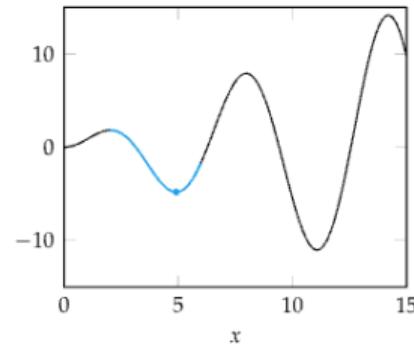
Transformations to Remove Constraints

Example

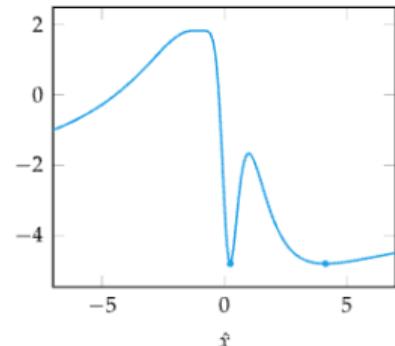
$$\underset{x}{\text{minimize}} \quad x \sin(x)$$

subject to $2 \leq x \leq 6$

$$f(x) = x \sin(x)$$



$$(f \circ T_{2,6})(\hat{x})$$



$$\underset{\hat{x}}{\text{minimize}} \quad t_{2,6}(\hat{x}) \sin(t_{2,6}(\hat{x}))$$

$$\underset{\hat{x}}{\text{minimize}} \quad \left(4 + 2 \left(\frac{2\hat{x}}{1 + \hat{x}^2}\right)\right) \sin\left(4 + 2 \left(\frac{2\hat{x}}{1 + \hat{x}^2}\right)\right)$$

Transformations to Remove Constraints

Example

$$\underset{x}{\text{minimize}} \ f(\mathbf{x})$$

$$\text{subject to } h(\mathbf{x}) = x_1^2 + x_2^2 + \dots + x_n^2 - 1 = 0$$

- Solve for one of the variables to eliminate constraint:

$$x_n = \pm \sqrt{1 - x_1^2 - x_2^2 - \dots - x_{n-1}^2}$$

- Transformed, unconstrained optimization problem:

$$\underset{x}{\text{minimize}} \ \left(\left[x_1, x_2, \dots, x_{n-1}, \pm \sqrt{1 - x_1^2 - x_2^2 - \dots - x_{n-1}^2} \right] \right)$$

Lagrangian Relaxation

- With only equality constraints, critical points (local minima, global minima, or saddle points optimal) where gradient of f and the gradient of h are aligned
- The method of Lagrangian relaxation is used to optimize a function subject to (equality) constraints
- Lagrangian multipliers refer to the variables introduced by the method denoted by λ

1. Form **Lagrangian relaxation**

$$\underset{x}{\text{minimize}} \ f(x)$$

subject to $h(x) = 0$

$$\mathcal{L}(x, \lambda) = f(x) - \lambda h(x)$$

2. Set $\nabla_x \mathcal{L}(x, \lambda) = 0$ and $\nabla_\lambda \mathcal{L}(x, \lambda) = 0$ to get

$$\nabla f(x) = \lambda \nabla h(x) \quad h(x) = 0$$

3. solve for x and λ

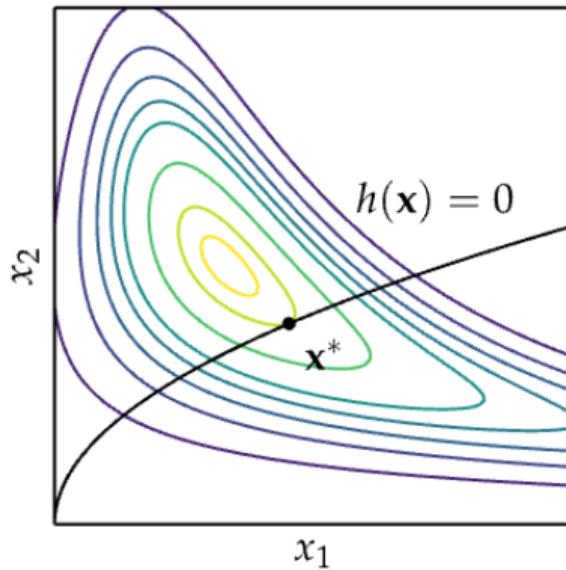
Example

$$\text{minimize } -\exp \left(-\left(x_1 x_2 - \frac{3}{2} \right)^2 - \left(x_2 - \frac{3}{2} \right)^2 \right)$$

$$\text{subject to } x_1 - x_2^2 = 0$$

Lagrangian Relaxation

Intuitively, the method of Lagrange multipliers finds the point \mathbf{x}^* where the constraint function is orthogonal to the gradient



Lagrangian Relaxation with Inequality Constraints

$$\underset{x}{\text{minimize}} \ f(x)$$

subject to $g(x) \leq 0$

- If solution lies at the constraint boundary, the constraint is called **active**, and the Lagrangian condition holds for a non-negative constant μ :

$$\nabla f(x) + \mu \nabla g(x) = 0$$

- If the solution lies within the boundary, the constraint is called **inactive**, and the optimal solution simply lies where

$$\nabla f(x) = 0$$

that is, the Lagrangian condition holds with $\mu = 0$

Lagrangian Relaxation with Inequality Constraints

$$\underset{x}{\text{minimize}} \ f(\mathbf{x})$$

$$\text{subject to } g(\mathbf{x}) \leq 0$$

- We create the Lagrangian relaxation such that it goes to ∞ outside the feasibility set ($g(\mathbf{x}) \not\leq 0$):

$$\mathcal{L}_{\infty}(\mathbf{x}) = f(\mathbf{x}) + \infty(g(\mathbf{x}) > 0)$$

impractical: discontinuous and nondifferentiable.

- Instead, for $\mu > 0$:

$$\mathcal{L}(\mathbf{x}, \mu \geq 0) = f(\mathbf{x}) + \mu g(\mathbf{x})$$

$$\mathcal{L}_{\infty}(\mathbf{x}) = \underset{\mu \geq 0}{\text{maximize}} \ \mathcal{L}(\mathbf{x}, \mu)$$

for \mathbf{x} infeasible, $\mathcal{L}_{\infty}(\mathbf{x}) = \infty$; for \mathbf{x} feasible, $\mathcal{L}_{\infty}(\mathbf{x}) = f(\mathbf{x})$

- The new optimization problem becomes

$$\underset{x}{\text{minimize}} \underset{\mu \geq 0}{\text{maximize}} \ \mathcal{L}(\mathbf{x}, \mu)$$

This is called the **primal problem**

Necessary Conditions – KKT Conditions

minimize_x $f(x)$

subject to $\mathbf{g}(x) \leq 0$
 $\mathbf{h}(x) = 0$

Any critical point \mathbf{x}^* must satisfy the **Karush-Kuhn-Tucker conditions**

1. primal feasibility: $\mathbf{g}(x^*) \leq 0$ and $\mathbf{h}(x^*) = 0$
2. dual feasibility: penalization is towards feasibility $\boldsymbol{\mu} \geq 0$
3. complementary slackness: either μ_i or $g_i(x^*)$ is zero.
$$\mu_i g_i(x^*) = 0, \text{ for } i = 1, \dots, m.$$
4. stationarity: objective function tangent to each active constraint

$$\nabla f(x^*) + \sum_i \mu_i \nabla g_i(x^*) + \sum_j \lambda_j \nabla h_j(x^*) = 0$$

Necessary Conditions – KKT Conditions

Particular cases

- f concave, \mathbf{g} convex: then KKT are also sufficient
- Pathological cases

In vector form:

$$\begin{cases} \nabla f(\mathbf{x}^*) + \mu \cdot \nabla \mathbf{g}(\mathbf{x}^*) + \lambda \cdot \nabla \mathbf{h}(\mathbf{x}^*) = 0 \\ \mu \cdot \mathbf{g}(\mathbf{x}^*) = 0 \\ \mathbf{g}(\mathbf{x}^*) \leq 0, \quad \mathbf{h}(\mathbf{x}^*) = 0 \\ \mu \geq 0 \end{cases}$$

Duality

- Generalized Lagrangian Relaxation:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_i \mu_i g_i(\mathbf{x}) + \sum_j \lambda_j h_j(\mathbf{x})$$

- the primal form is

$$\underset{\mathbf{x}}{\text{minimize}} \underset{\boldsymbol{\mu} \geq 0, \boldsymbol{\lambda}}{\text{maximize}} \mathcal{L}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

- Reversing the order of operations leads to the **dual form**

$$\underset{\boldsymbol{\mu} \geq 0, \boldsymbol{\lambda}}{\text{maximize}} \underset{\mathbf{x}}{\text{minimize}} \mathcal{L}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\lambda})$$

- In some cases, the dual problem is easier to solve computationally than the original problem.
In other cases, the dual can be used to obtain easily a lower bound on the optimal value of the objective for the primal problem. The dual has also been used to design algorithms for solving the primal problem.

Duality

Theorem (Max-min inequality)

For any function $f : Z \times W \rightarrow \mathbb{R}$,

$$\sup_{z \in Z} \inf_{w \in W} f(z, w) \leq \inf_{w \in W} \sup_{z \in Z} f(z, w).$$

Proof: see wikipedia



- When f , W , and Z are convex the inequality becomes equality and we have a strong max–min property (or a saddle-point property).
- For us:

$$\underset{\mu \geq 0, \lambda}{\text{maximize}} \underset{x}{\text{minimize}} \mathcal{L}(x, \mu, \lambda) \leq \underset{x}{\text{minimize}} \underset{\mu \geq 0, \lambda}{\text{maximize}} \mathcal{L}(x, \mu, \lambda)$$

- Therefore, the solution to the dual problem d^* is a lower bound to the primal solution p^*
- The inner part of the dual problem can be used to define the **dual function** or **dual objective**

$$\mathcal{D}(\mu \geq 0, \lambda) = \underset{x}{\text{minimize}} \mathcal{L}(x, \mu, \lambda)$$

Duality

- The dual function is concave. Gradient ascent on a concave function always converges to the global maximum.
- **Dual Problem:** $\max \mathcal{D}(\lambda)$ subject to $\lambda \geq 0$
- Optimizing the dual problem is easy whenever minimizing the Lagrangian with respect to x is easy.
- For any $\mu \geq 0$ and any λ , we have

$$\mathcal{D}(\mu \geq 0, \lambda) \leq p^*$$

- The difference between dual and primal solutions d^* and p^* is called the **duality gap**
- Showing zero-duality gap is a “certificate” of optimality

Penalty methods

- Penalty methods are a way of reformulating a constrained optimization problem as an unconstrained problem by penalizing the objective function value when constraints are violated

Example

$$\underset{x}{\text{minimize}} \ f(x)$$

$$\text{subject to } g(x) \leq 0$$

$$h(x) = 0$$

$$\min_x f(x) + \rho \cdot p_{count}(x)$$

$$\text{s.t. } p_{count}(x) = \sum_i (g_i(x) > 0) + \sum_j (h_j(x) \neq 0)$$

Penalty Methods

Procedure penalty_method;

Input: $f, p, \mathbf{x}, k_{max}; \rho = 1, \gamma = 2$

Output: \mathbf{x} solution

for k in $1, \dots, k_{max}$ **do**

$\mathbf{x} \leftarrow \text{minimize}_{\mathbf{x}} \{f(\mathbf{x}) + \rho \cdot p(\mathbf{x})\};$

$\rho \leftarrow \rho \cdot \gamma;$

if $p(\mathbf{x}) = 0$ **then**

return $\mathbf{x};$

return $\mathbf{x};$

Penalty methods

- Count penalty:

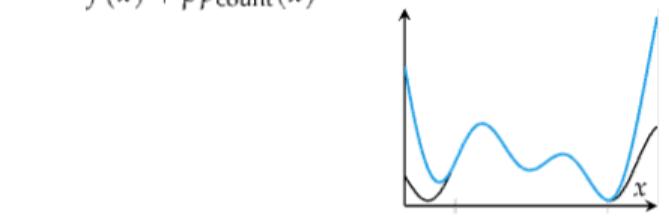
$$p_{\text{count}}(\mathbf{x}) = \sum_i (g_i(\mathbf{x}) > 0) + \sum_j (h_j(\mathbf{x}) \neq 0)$$



— $f(x)$
— $f(x) + \rho p_{\text{count}}(\mathbf{x})$

- Quadratic penalty:

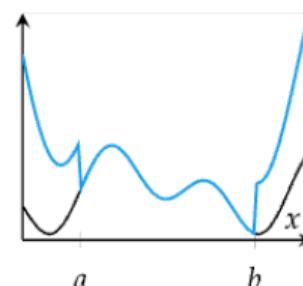
$$p_{\text{quadratic}}(\mathbf{x}) = \sum_i \max(g_i(\mathbf{x}), 0)^2 + \sum_j h_j(\mathbf{x})^2$$



— $f(x)$
— $f(x) + \rho p_{\text{quadratic}}(\mathbf{x})$

- Mixed Penalty:

$$p_{\text{mixed}}(\mathbf{x}) = \rho_1 p_{\text{count}}(\mathbf{x}) + \rho_2 p_{\text{quadratic}}(\mathbf{x})$$



— $f(x)$
— $f(x) + \rho_{\text{mixed}}(\mathbf{x})$

Augmented Lagrange Method

- Adaptation of penalty method for equality constraints

$$p_{\text{Lagrange}}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2}\rho \sum_i (h_i(\mathbf{x}))^2 - \sum_i \lambda_i h_i(\mathbf{x})$$

Procedure augmented_lagrange_method;

Input: $f, h, \mathbf{x}, k_{\max}; \rho = 1, \gamma = 2$)

$\lambda \leftarrow 0;$

for k in $1, \dots, k_{\max}$ **do**

$p \leftarrow (\mathbf{x} \mapsto \rho/2 \cdot \sum_i (h_i(\mathbf{x}))^2 - \lambda \cdot h(\mathbf{x}))$;

$\mathbf{x} \leftarrow \text{minimize}_{\mathbf{x}} \{f(\mathbf{x}) + p(\mathbf{x})\}$;

$\lambda \leftarrow \lambda - \rho \cdot h(\mathbf{x})$;

$\rho \leftarrow \rho \cdot \gamma$;

return \mathbf{x} ;

- λ converges towards the Lagrangian multiplier

Interior Point Methods

- Also called **barrier methods**, interior point methods ensure that each step is feasible
- This allows premature termination to return a nearly optimal, feasible point
- Barrier functions are implemented similar to penalties but must meet the following conditions:
 1. Continuous
 2. Non-negative
 3. Approach infinity as x approaches boundary

Interior Point Methods

- Inverse Barrier:

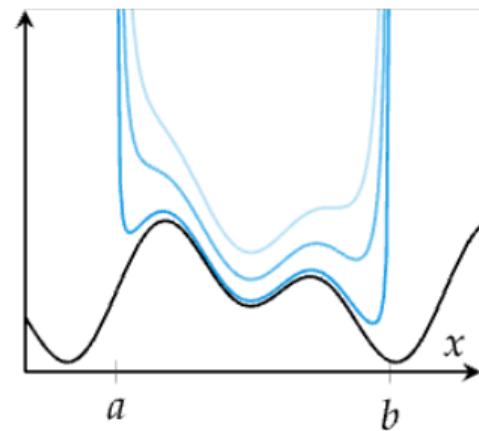
$$p_{\text{barrier}}(\mathbf{x}) = - \sum_i \frac{1}{g_i(\mathbf{x})}$$

- Log Barrier:

$$p_{\text{barrier}}(\mathbf{x}) = - \sum_i \begin{cases} \log(-g_i(\mathbf{x})) & \text{if } g_i(\mathbf{x}) \geq -1 \\ 0 & \text{otherwise} \end{cases}$$

New optimization problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + \frac{1}{\rho} p_{\text{barrier}}(\mathbf{x})$$



- $f(x)$
- $f(x) + p_{\text{barrier}}(x)$
- $f(x) + \frac{1}{2} p_{\text{barrier}}(x)$
- $f(x) + \frac{1}{10} p_{\text{barrier}}(x)$

Interior Point Methods

```
Procedure interior_point_method;  
Input:  $f, p, \mathbf{x}; \rho = 1, \gamma = 2, \epsilon = 0.001$   
 $\Delta \leftarrow \infty;$   
while  $\Delta > \epsilon$  do  
     $\mathbf{x}' \leftarrow \text{minimize}_{\mathbf{x}} \{f(\mathbf{x}) + p(\mathbf{x})/\rho\};$   
     $\Delta \leftarrow \|\mathbf{x}' - \mathbf{x}\|;$   
     $\mathbf{x} \leftarrow \mathbf{x}';$   
     $\rho \leftarrow \rho \cdot \gamma;$   
return  $\mathbf{x};$ 
```

- Line searches $f(\mathbf{x} + \alpha \mathbf{d})$ are constrained to the interval $\alpha = [0, \alpha_u]$, where α_u is the step to the nearest boundary.
In practice, α_u is chosen such that $\mathbf{x} + \alpha \mathbf{d}$ is just inside the boundary to avoid the boundary singularity.
- Needs an initial **feasible** solutions. Typically, found by solving:

$$\underset{\mathbf{x}}{\text{minimize}} \quad p_{\text{quadratic}}(\mathbf{x})$$

Summary

- Constraints are requirements on the design points that a solution must satisfy
- Some constraints can be transformed or substituted into the problem to result in an unconstrained optimization problem
- Analytical methods using Lagrange multipliers yield the generalized Lagrangian and the necessary conditions for optimality under constraints
- A constrained optimization problem has a dual problem formulation that is easier to solve and whose solution is a lower bound of the solution to the original problem
- Penalty methods penalize infeasible solutions and often provide gradient information to the optimizer to guide infeasible points toward feasibility
- Interior point methods maintain feasibility but use barrier functions to avoid leaving the feasible set

14. Linear Constrained Optimization

Problem Formulation

- If an optimization problem has a linear objective and constraints, it is called a **linear programming problem (linear program, LP)**
- The general form of a linear program is:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}$$

$$\text{subject to } A\mathbf{x} \leq \mathbf{b}$$

$$D\mathbf{x} \geq \mathbf{e}$$

$$F\mathbf{x} = \mathbf{g}$$

$$\mathbf{x}, \mathbf{c} \in \mathbb{R}^n,$$

$$A \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m$$

$$D \in \mathbb{R}^{p \times n}, \mathbf{e} \in \mathbb{R}^p$$

$$F \in \mathbb{R}^{q \times n}, \mathbf{g} \in \mathbb{R}^q$$

Numerical Example

$$\underset{x_1, x_2, x_3}{\text{minimize}} \quad 2x_1 - 3x_2 + 7x_3$$

$$\text{subject to} \quad 2x_1 + 3x_2 - 8x_3 \leq 5$$

$$4x_1 + x_2 + 3x_3 \leq 9$$

$$x_1 - 5x_2 - 3x_3 \geq -4$$

$$x_1 + x_2 + 2x_3 = 1$$

Modelling in Linear Programming

Example

Given a set of items I , each item with a price p_i and a value v_i , i in I , select the subset of items that maximizes the total value collected subject to a total expense that does not exceed a given budget B .

$$\max \sum_{i \in I} p_i x_i$$

$$\text{s.t. } \sum_{i \in I} v_i x_i \leq B$$

$$x_i \in \{0, 1\}, \quad \text{for all } i \text{ in } I$$

Modelling in Linear Programming

Many problems can be converted into linear programs that have the same solution.

Example

$$\text{minimize } L_1 = \|Ax - b\|_1$$

$$\min 1^T s$$

$$\text{s.t. } Ax - b \leq s$$

$$-(Ax - b) \leq s$$

Example

$$\text{minimize } L_\infty = \|Ax - b\|_\infty$$

$$\min t$$

$$\text{s.t. } Ax - b \leq t$$

$$-(Ax - b) \leq t$$

Problem Formulation

Every general form linear program can be rewritten more compactly in **standard form**

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}$$

$$\text{subject to } A\mathbf{x} \leq \mathbf{b}$$

$$\mathbf{x} \geq 0$$

$$\mathbf{x}, \mathbf{c} \in \mathbb{R}^n,$$

$$A \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m$$

Example

$$\text{minimize } 5x_1 + 4x_2$$

$$\text{s.t. } 2x_1 + 3x_2 \leq 5$$

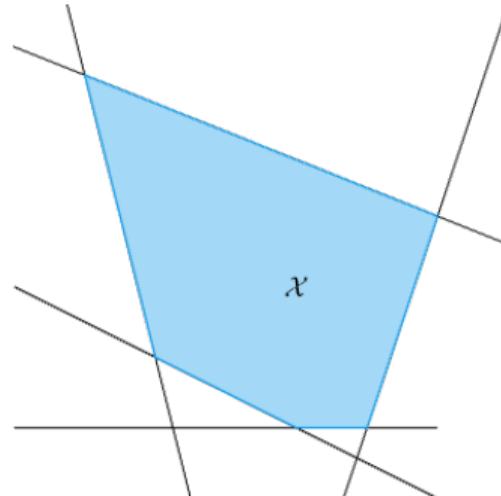
$$4x_1 + x_2 \leq 11$$

Problem Formulation

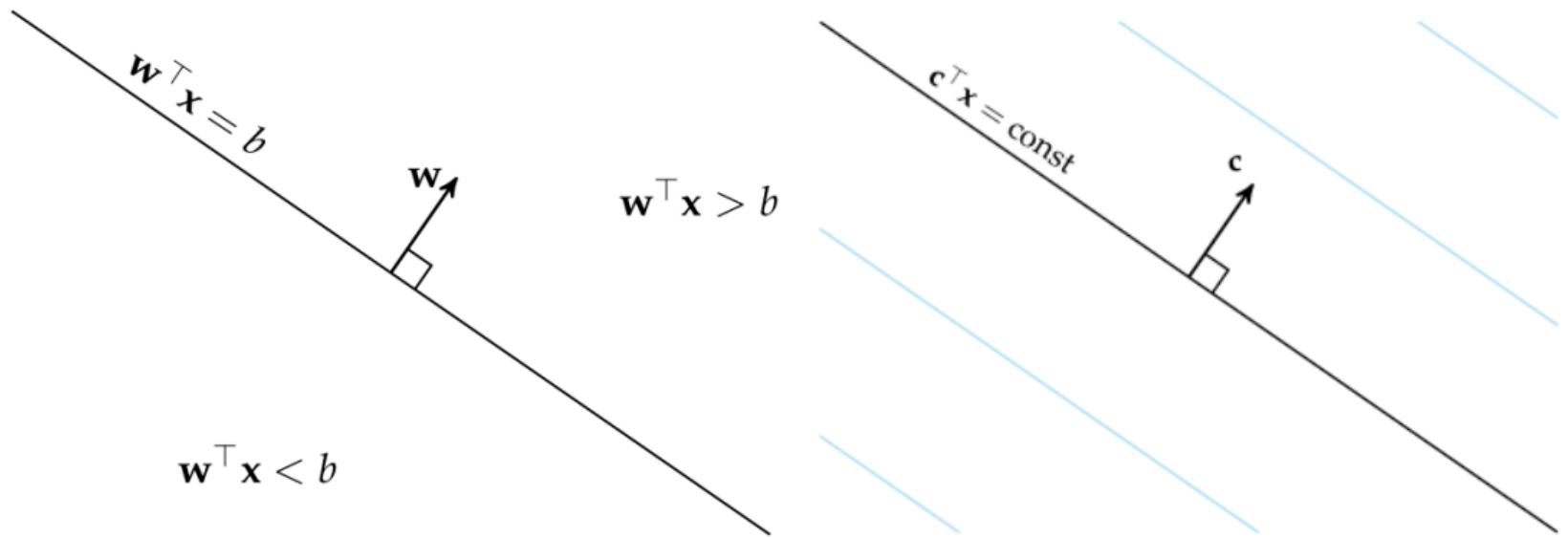
- Each inequality constraint defines a planar boundary of the feasible set called a **half-space**
- The set of inequality constraints define the intersection of multiple half-spaces forming a **convex set**
- Convexity of the feasible set, along with convexity of the objective function, implies that if we find a local feasible minimum, it is also a global feasible minimum.

$$\underset{x}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}$$

$$\begin{aligned} & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \quad \mathbf{x} \geq 0 \end{aligned}$$



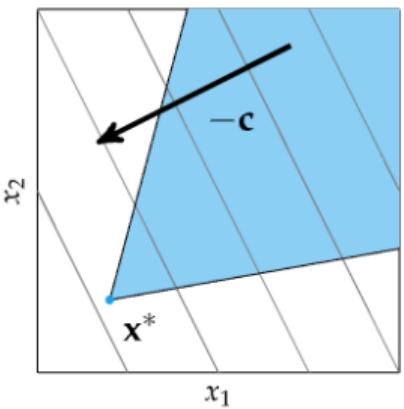
Half-Spaces and Supporting Hyperplanes



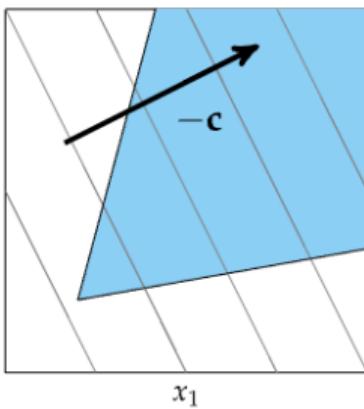
Problem Formulation

- How many solutions are there?

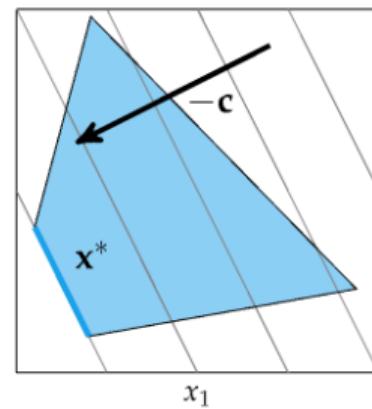
One Solution



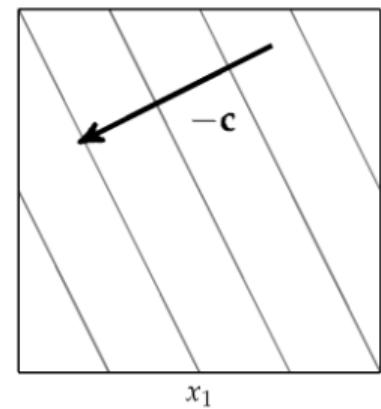
Unbounded Solution



Infinite Solutions



No Solution



Problem Formulation

Linear programs are often solved in **equality form**

$$\underset{x}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}$$

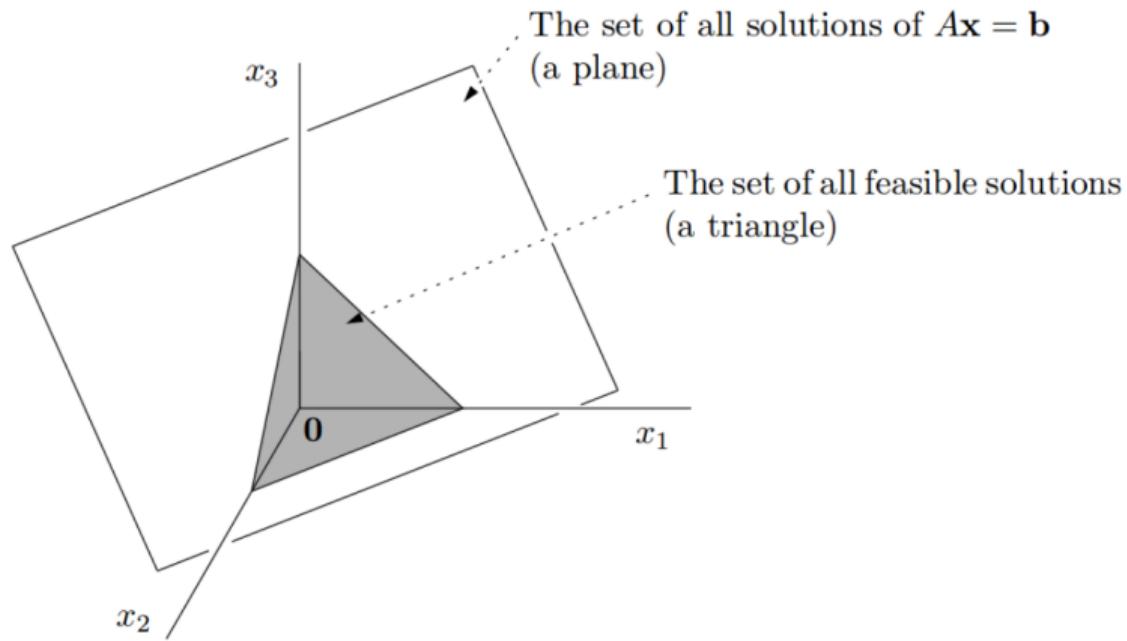
$$\text{subject to } A\mathbf{x} = \mathbf{b}$$

$$\mathbf{x} \geq 0$$

$$\mathbf{x}, \mathbf{c} \in \mathbb{R}^{2n+m},$$

$$A \in \mathbb{R}^{m \times 2n+m},$$

$$\mathbf{b} \in \mathbb{R}^m$$



Simplex Algorithm

- Guaranteed to solve any feasible and bounded linear program
- Works on the equality form
- Assumes that rows of A are linearly independent and $m \leq n'$ ($n' \leq 2n + m$)
- The feasible set of a linear program forms a **polytope** (polyhedra bounded by faces of $n - 1$ dimension)
- The simplex algorithm moves between **vertices** of the polytope until it finds an optimal **vertex**
- Points on faces not perpendicular to c can be improved by sliding along the face in the direction of the projection of $-c$ onto the face.

Fundamental Theorem of LP

Theorem (Fundamental Theorem of Linear Programming)

Given:

$$\min\{\mathbf{c}^T \mathbf{x} \mid \mathbf{x} \in P\} \text{ where } P = \{\mathbf{x} \in \mathbb{R}^n \mid A\mathbf{x} \leq \mathbf{b}\}$$

If P is a bounded polyhedron and not empty and \mathbf{x}^* is an optimal solution to the problem, then:

- \mathbf{x}^* is an extreme point (vertex) of P , or
- \mathbf{x}^* lies on a face $F \subset P$ of optimal solution



Proof:

- assume \mathbf{x}^* not a vertex of P then \exists a ball around it still in P . Show that a point in the ball has better cost
- if \mathbf{x}^* is not a vertex then it is a convex combination of vertices. Show that all points are also optimal.

Simplex Algorithm

- Every vertex for a linear program in equality form can be uniquely defined by $n - m$ components of \mathbf{x} that equal zero.
- choosing m design variables and setting the remaining variables to zero effectively removes $n - m$ columns of A , yielding an $m \times m$ constraint matrix
- the m selected columns of the matrix A are called **basis** and denoted by B : $x_i \geq 0$ for $i \in B$
- the $n - m$ columns not in B are called **not in basis** and are denoted by V : $x_i = 0$ for $i \in V$.

$$A\mathbf{x} = A_B \mathbf{x}_B = \mathbf{b} \implies \mathbf{x}_B = A_B^{-1} \mathbf{b}$$

Simplex Algorithm

- every vertex has an associated partition (B, V) ,
- not every partition corresponds to a vertex.
 A_B might be not invertible or the point x_B might not be ≥ 0 .
- identifying partitions that correspond to vertices corresponds to solving an LP problem as well!

Two phases of the algorithm

1. **Initialization Phase:** finding a feasible starting vertex
2. **Optimization Phase:** finding the optimal vertex

Simplex Algorithm: FONCs

Lagrangian function:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mathbf{c}^T \mathbf{x} - \boldsymbol{\mu}^T \mathbf{x} - \boldsymbol{\lambda}^T (\mathbf{A}\mathbf{x} - \mathbf{b})$$

Conditions for Optimality for linear programs: KKT are also sufficient:

- feasibility: $\mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0$
- dual feasibility: $\boldsymbol{\mu} \geq 0$
- complementary slackness: $\boldsymbol{\mu} \cdot \mathbf{x} = 0$
- stationarity: $\mathbf{A}^T \boldsymbol{\lambda} + \boldsymbol{\mu} = \mathbf{c}$

$$A^T \lambda + \mu = c \implies \begin{cases} A_B^T \lambda + \mu_B = c_B \\ A_V^T \lambda + \mu_V = c_V \end{cases}$$

- We can choose $\mu_B = 0$ to satisfy complementarity slackness (because $x_B \geq 0$)

$$\mu_V = c_V - (A_B^{-1} A_V)^T c_B$$

- Knowing μ_V allows us to assess the optimality of the vertices. If μ_B contains negative components, then dual feasibility is not satisfied and the vertex is sub-optimal.
- maintain a partition (B, V) , which corresponds to a vertex of the feasible set polytope.
- The partition can be updated by swapping indices between B and V . Such a swap equates to moving from one vertex along an edge of the feasible set polytope to another vertex.

Simplex Algorithm: Optimization Phase

Pivoting

- $q \in V$ to enter in B

$$A\mathbf{x}' = A_B\mathbf{x}'_B + A_{\{q\}}\mathbf{x}'_q = A_B\mathbf{x}_B = A\mathbf{x} = \mathbf{b}$$

- $p \in B$ to leave B becomes zero during the transition.

$$\mathbf{x}'_B = \mathbf{x}_B - A_B^{-1}A_{\{q\}}\mathbf{x}'_q \implies (\mathbf{x}'_B)_p = 0 = (\mathbf{x}_B)_p - (A_B^{-1}A_{\{q\}})_p \mathbf{x}'_q$$

- leaving index is obtained using the **minimum ratio test**: compute \mathbf{x}'_q for each potential leaving index p and select the leaving index p that yields the smallest \mathbf{x}'_q .
- Choosing an entering index q decreases the objective function value by

$$\mathbf{c}^T \mathbf{x}' = \mathbf{c}_B^T \mathbf{x}'_B + c_q \mathbf{x}'_q = \mathbf{c}^T \mathbf{x} + \mu_q \mathbf{x}'_q$$

- The objective function decreases only if μ_q is negative.

Simplex Algorithm: Optimization Phase

- In order to progress toward optimality, we must choose an index q in V such that μ_q is negative. If all components of μ_V are non-negative, we have found a global optimum.
- Since there can be multiple negative entries in μ_V , Several possible heuristics to search for optimal vertex (choose next q)
 - **Dantzig's rule:** choose most negative entry in μ_V ; easy to calculate
 - **Greedy heuristic (largest decrease):** maximally reduces objective at each step
 - **Bland's rule:** chooses first vertex found with negative μ_V ; useful for preventing or breaking out of cycles

Simplex Algorithm: Initialization Phase

- The starting vertex of the optimization phase is found by solving an additional **auxiliary linear program** that has a known feasible starting vertex

$$\underset{x,z}{\text{minimize}} \quad \begin{bmatrix} 0^T & 1^T \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}$$

$$\begin{bmatrix} A & Z \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix} = b$$

$$\begin{bmatrix} x \\ z \end{bmatrix} \geq 0$$

- The solution is a feasible vertex in the original linear program

Dual Certificates

- Verification that the solution returned by the algorithm is actually the correct solution
- Recall that the solution to the dual problem, d^* provides a lower bound to the solution of the primal problem, p^*
- If $d^* = p^*$ then p^* is guaranteed to be the unique optimal value because the duality gap is zero
- What happens if one of the two is unbounded or infeasible?

Dual Certificates

Linear programs have a simple dual form:

Primal form (equality)

$$\underset{x}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}$$

$$\text{subject to } A\mathbf{x} = \mathbf{b}$$

$$\mathbf{x} \geq 0$$

Dual form

$$\underset{\lambda}{\text{maximize}} \quad \mathbf{b}^T \lambda$$

$$\text{subject to } A^T \lambda \leq \mathbf{c}$$

Strong Duality Theorem

Due to Von Neumann and Dantzig 1947 and Gale, Kuhn and Tucker 1951.

Theorem (Strong Duality Theorem)

Given:

$$(P) \min\{\mathbf{c}^T \mathbf{x} \mid A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0\}$$

$$(D) \max\{\mathbf{b}^T \boldsymbol{\lambda} \mid A^T \boldsymbol{\lambda} \geq \mathbf{c}\}$$

exactly one of the following occurs:

1. (P) and (D) are both infeasible
2. (P) is unbounded and (D) is infeasible
3. (P) is infeasible and (D) is unbounded
4. (P) has feasible solution, then let an optimal be: $\mathbf{x}^* = [x_1^*, \dots, x_n^*]$
(D) has feasible solution, then let an optimal be: $\boldsymbol{\lambda}^* = [\lambda_1^*, \dots, \lambda_m^*]$, then:

$$p^* = \mathbf{c}^T \mathbf{x}^* = \mathbf{b}^T \boldsymbol{\lambda}^* = d^*$$

Summary

- Linear programs are problems consisting of a linear objective function and linear constraints
- The simplex algorithm can optimize linear programs globally in an efficient manner
- Dual certificates allow us to verify that a candidate primal-dual solution pair is optimal
- Linear programs can be solved to optimality for problems with millions of variables.

15. Sampling Plans

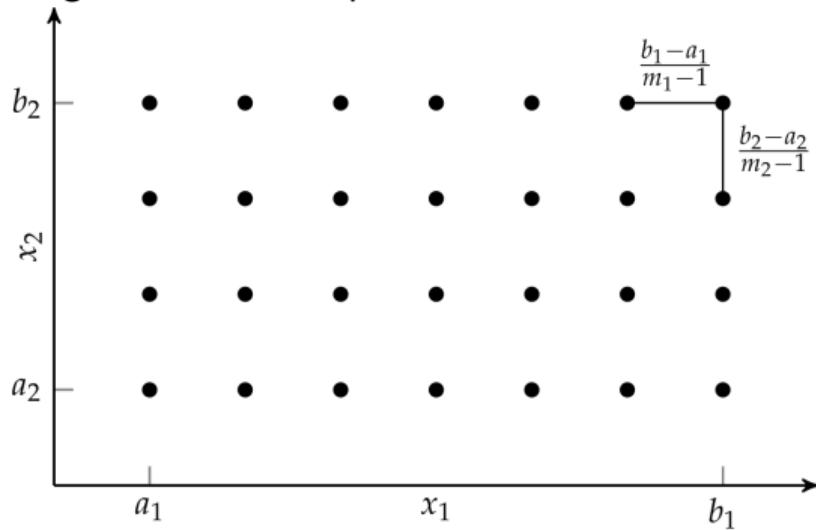
Sampling Plans

- In all nonlinear non convex optimization, to generate good initial design points
- With computationally costly functions, to create an initial set of design points from where to build a **surrogate models** to optimize in place of the original function
- In hyperparameter tuning

Full Factorial Design

- Factors and levels, terms from the field of Experimental Design in Statistics
- Uniform and evenly spaced samples across domain
- Simple, easy to implement, and covers domain
- Optimization over the points known as **grid search**
- Sample count grows exponentially with dimension: n^m
- Can be coarse and miss local features

$a_i \leq x_i \leq b_i$ for each component i .
grid with m_i samples in the i th dimension



Random Sampling

- Uses pseudorandom number generator to define samples according to our desired distribution
- If variable bounds are known, a common choice is independent **uniform distributions** across domains of possible variable values
 $[a_1, b_1] \times \dots \times [a_n, b_n]$
- Ideally, if enough points are sampled and the right distribution is chosen, the design space will be covered

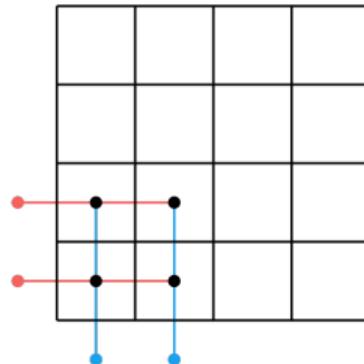
Uniform Projection Plans

- A **uniform projection plan** is a sampling plan over a discrete grid where the distribution over each dimension is uniform.

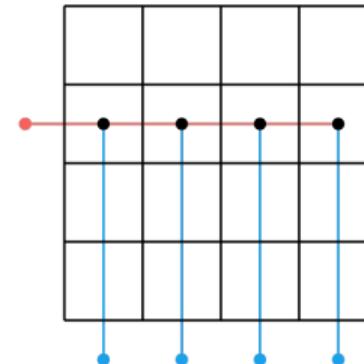
Example

In 2D, $m \times m$ sampling grid (as in full factorial), but, instead of taking all m^2 samples, we want to sample only m positions.

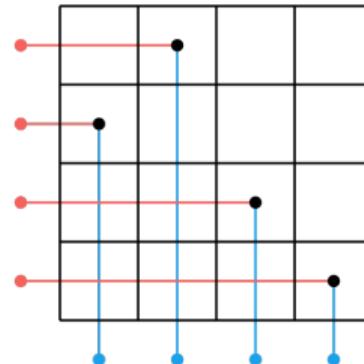
too clustered



no variation in one component

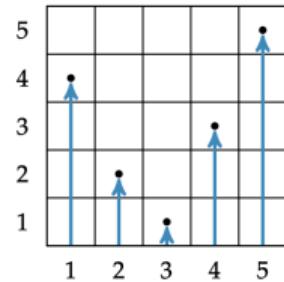


uniform projection



Uniform Projection Plans

Example (Random m -permutations)



$$p = 4 \ 2 \ 1 \ 3 \ 5$$

Example (Latin square)

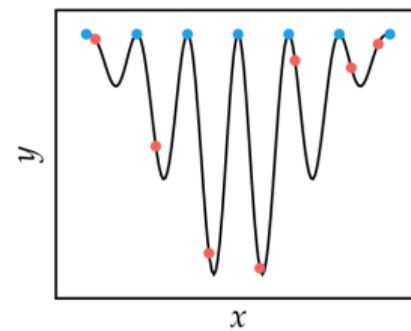
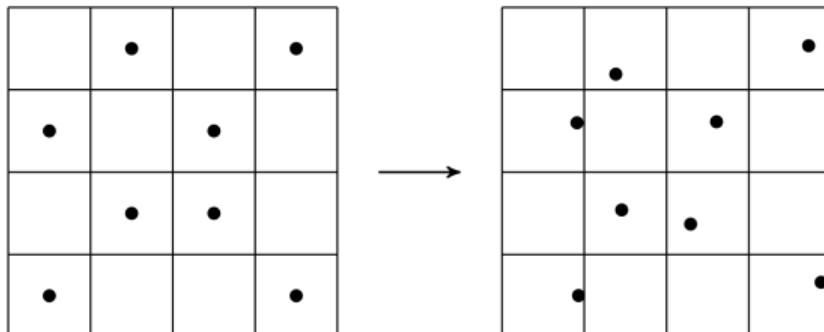
Latin squares are $m \times m$ grids where each row contains each integer 1 through m and each column contains each integer 1 through m .

Latin-hypercubes are a generalization to any number of dimensions (note that the points remain m) N rooks on a chess board without threatening each other

4	1	3	2
1	4	2	3
3	2	1	4
2	3	4	1

Stratified Sampling

- Each point is sampled uniformly at random within each grid cell instead of the center
- Cells decided by Full Factorial or Uniform Projection Plans
- Can capture details that regularly-spaced samples might miss

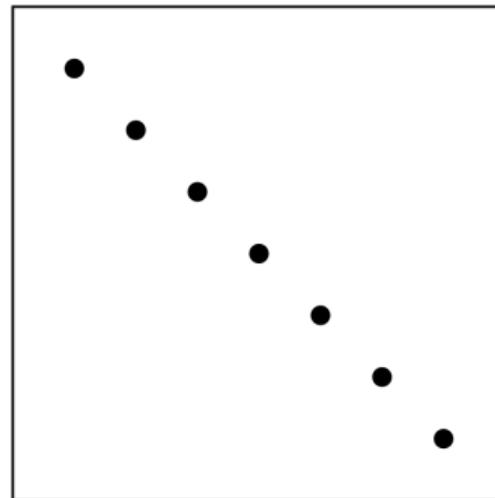


— $f(x)$
● sampling on grid
● stratified sampling

Space Filling Metrics

- A sampling plan may cover a search space fully, but still leave large areas unexplored

Example (Uniform Projection Plan)



- **space-filling metrics** quantify this aspect measuring the degree to which a sampling plan $X \subseteq \mathcal{X}$ fills the design space

Space-Filling Metrics: Discrepancy

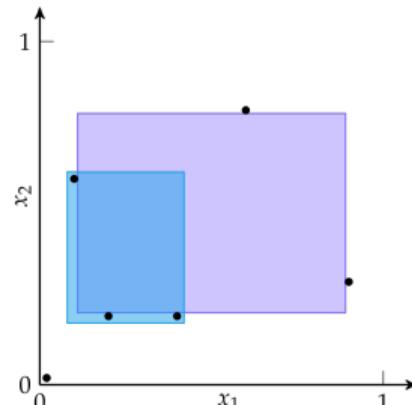
- **Discrepancy**: measure of ability of the sampling plan X to fill a hyper-rectangular design space
- It is given by hyper-rectangular subset \mathcal{H} with the maximum difference between the fraction of samples in \mathcal{H} and the volume of \mathcal{H} 's.

$$d(X) = \sup_{\mathcal{H}} \left| \frac{\#(X \cap \mathcal{H})}{\#X} - \lambda(\mathcal{H}) \right|$$

$\lambda(\mathcal{H})$ is the n -dimensional volume of \mathcal{H} , ie, the product of the side lengths of \mathcal{H}

We wish to have a plan X with low discrepancy

Often very difficult to compute directly



d for the purple rectangle is $>$ than d for the blue rectangle

Space-Filling Metrics: Pairwise Distances

- Method of measuring relative space-filling performance of two m -point sampling plans
- Better spread-out plans will have larger pairwise distances:
 1. compute all pairwise distances between all points within each sampling plan
 2. sort the pairwise distances of each set in ascending order
 3. the plan with the first pairwise distance exceeding the other is considered more space-filling
- Suggests simple algorithm:
 1. produce a set of randomly distributed sampling plans,
 2. pick the one with greatest pairwise distances
- Possible also for uniform projection plans, by mutating them with swaps and simulated annealing.

Space-Filling Metrics: Morris-Mitchell Criterion

- Alternative to previously suggested algorithm that simplifies the optimization problem

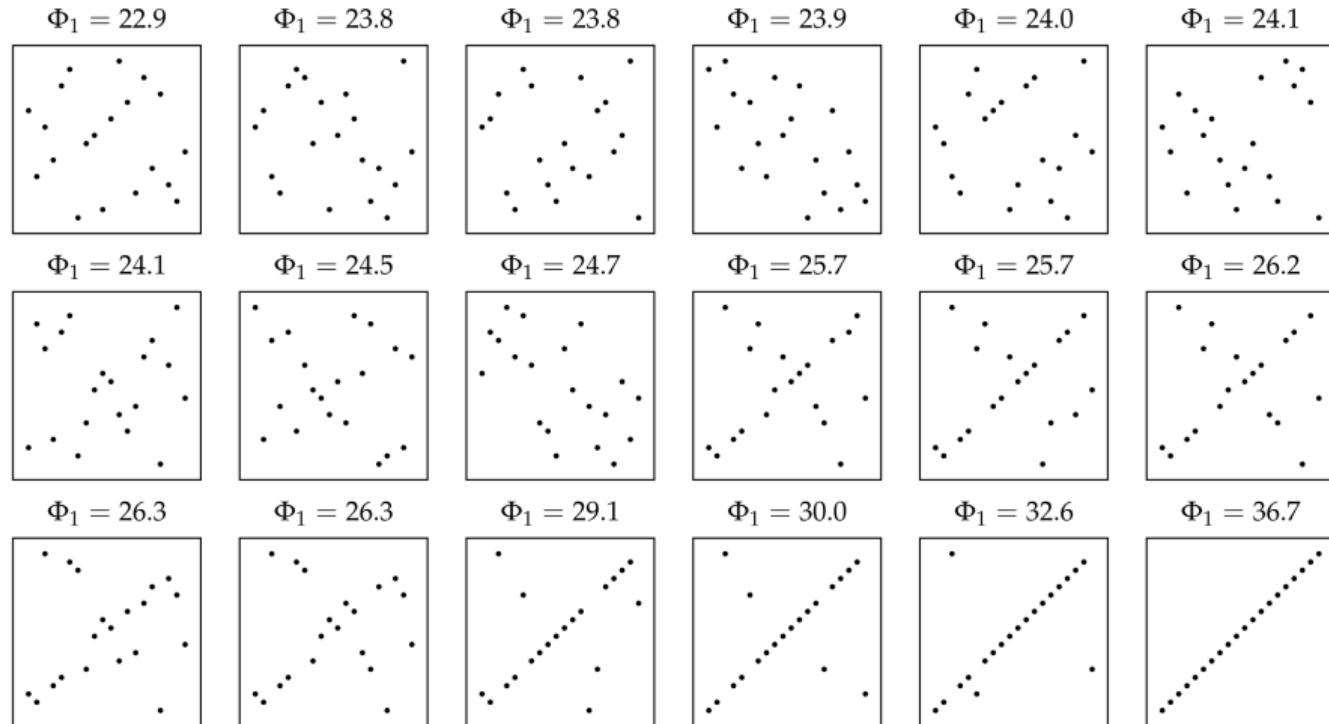
$$\underset{X}{\text{minimize}} \quad \underset{q \in \{1, 2, 3, 10, 20, 50, 100\}}{\text{maximize}} \quad \Phi_q(X)$$

$$\Phi_q(X) = \left(\sum_i d_i^{-q} \right)^{\frac{1}{q}}$$

where d_i is the i th pairwise distance between points in X and $q > 0$ is a tunable parameter. Larger values of q give higher penalties to large distances.

Space-Filling Metrics: Morris-Mitchell Criterion

Uniform projection plans sorted from best to worst according to Φ_1



Space-Filling Subsets

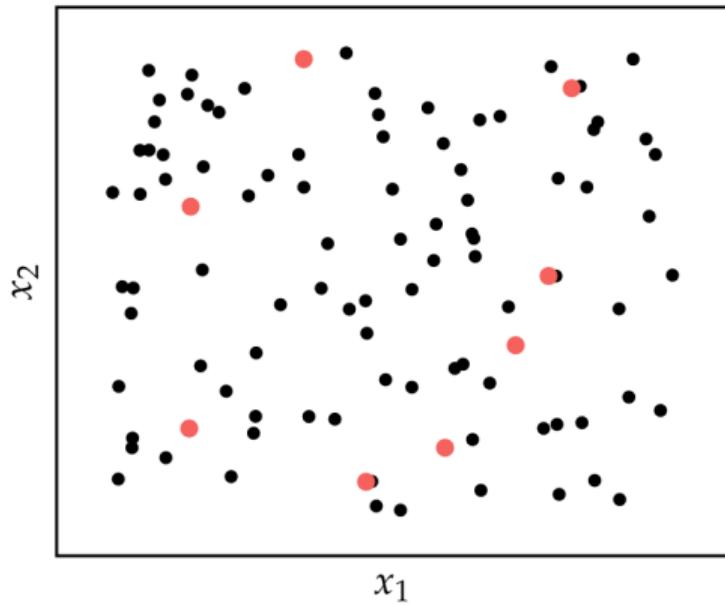
- Often, the set of possible sample points is constrained to be a subset of available choices
- A space-filling metric for a subset S within a finite set X is the maximum distance between a point in X and the closest point in S , using a norm to measure distance

$$d_{\max}(X, S) = \underset{x \in X}{\text{maximize}} \underset{s \in S}{\text{minimize}} \|s - x\|_q$$

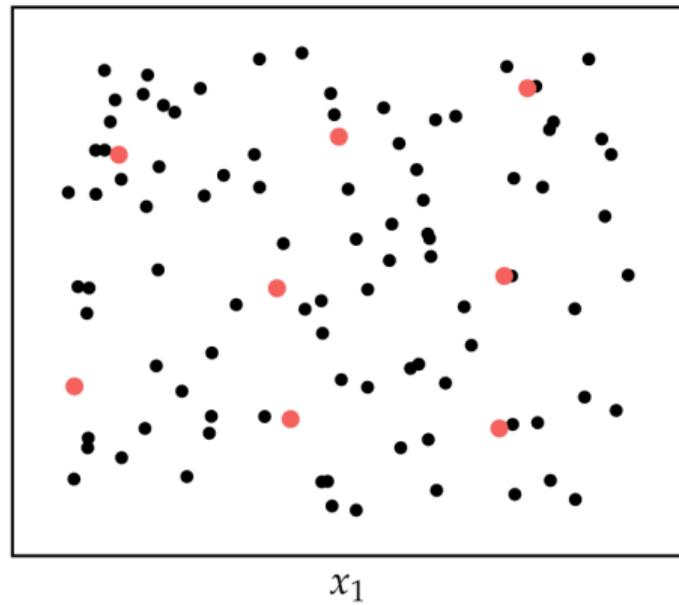
- A space-filling subset minimizes this metric
- Often computationally intractable, but heuristics like (repeated) greedy construction and exchange-search often produce acceptable results

Space-Filling Subsets

greedy local search

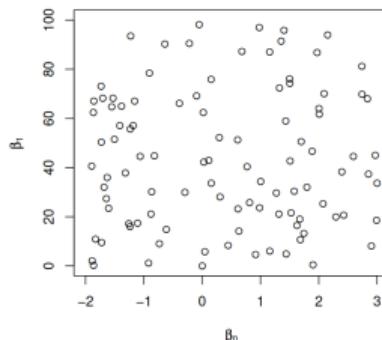
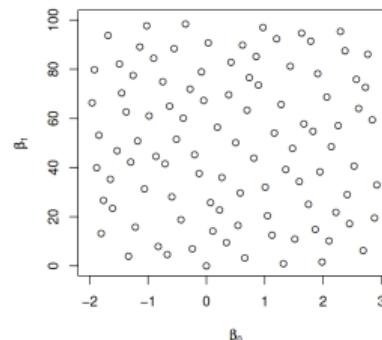


exchange algorithm



Quasi-Random Sequences

- Also called **low-discrepancy sequences**, **quasi-random sequences** are deterministic sequences that systematically fill a space such that their integral over the space converges as fast as possible
- Used for fast convergence in Monte Carlo integration, which approximates an integral by sampling points in a domain
- Quasi-random sequences are typically constructed for the unit n -dimensional hypercube, $[0, 1]^n$. Any multidimensional function with bounds on each variable can be transformed into such a hypercube.



Quasi-Random Sequences

- Additive Recurrence: Recursively adds irrational numbers
- Halton Sequence: sequence of fractions generated with coprime numbers
- Sobol Sequence: recursive XOR operation with carefully chosen numbers

Quasi-Random Sequences: Additive Recurrence

- Recursively adds irrational numbers

$$x_{k+1} = x_k + c \pmod{1}$$

c irrational

$$c = 1 - \varphi = \frac{\sqrt{5} - 1}{2} \approx 0.618034$$

φ is golden ratio

- We can construct a space-filling set over n dimensions using an additive recurrence sequence for each coordinate, each with its own value of c .
- square roots of the primes are known to be irrational, and can thus be used to obtain different sequences for each coordinate:

$$c_1 = \sqrt{2}, c_2 = \sqrt{3}, c_3 = \sqrt{5}, c_4 = \sqrt{7}, c_5 = \sqrt{11}, \dots$$

Quasi-Random Sequences: Halton Sequence

- single-dimensional version, called van der Corput sequences, generates sequences where the unit interval is divided into powers of base b . For example, $b = 2$

$$X = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{3}{4}, \frac{1}{8}, \frac{5}{8}, \frac{3}{8}, \frac{7}{8}, \frac{1}{16}, \dots \right\}$$

whereas $b = 5$

$$X = \left\{ \frac{1}{5}, \frac{2}{5}, \frac{3}{5}, \frac{4}{5}, \frac{1}{25}, \frac{6}{25}, \frac{11}{25}, \dots \right\}$$

- Multi-dimensional space-filling sequences use one van der Corput sequence for each dimension, each with its own base b . The bases, however, must be **coprime** in order to be uncorrelated.
- Two integers are **coprime** if the only positive integer that divides them both is 1, eg, 8 and 9.
- Correlation can be avoided by the leaped Halton method, which takes every p th point, where p is a prime different from all coordinate bases.

Quasi-Random Sequences: Sobol Sequence

- Recursive XOR operation with carefully chosen numbers.
- XOR (\oplus) returns true if and only if both inputs are different
- For n -dimensional hypercube $I^n = [0, 1]^n$, the i th point of the sequence x_i for dimension j is calculated as:

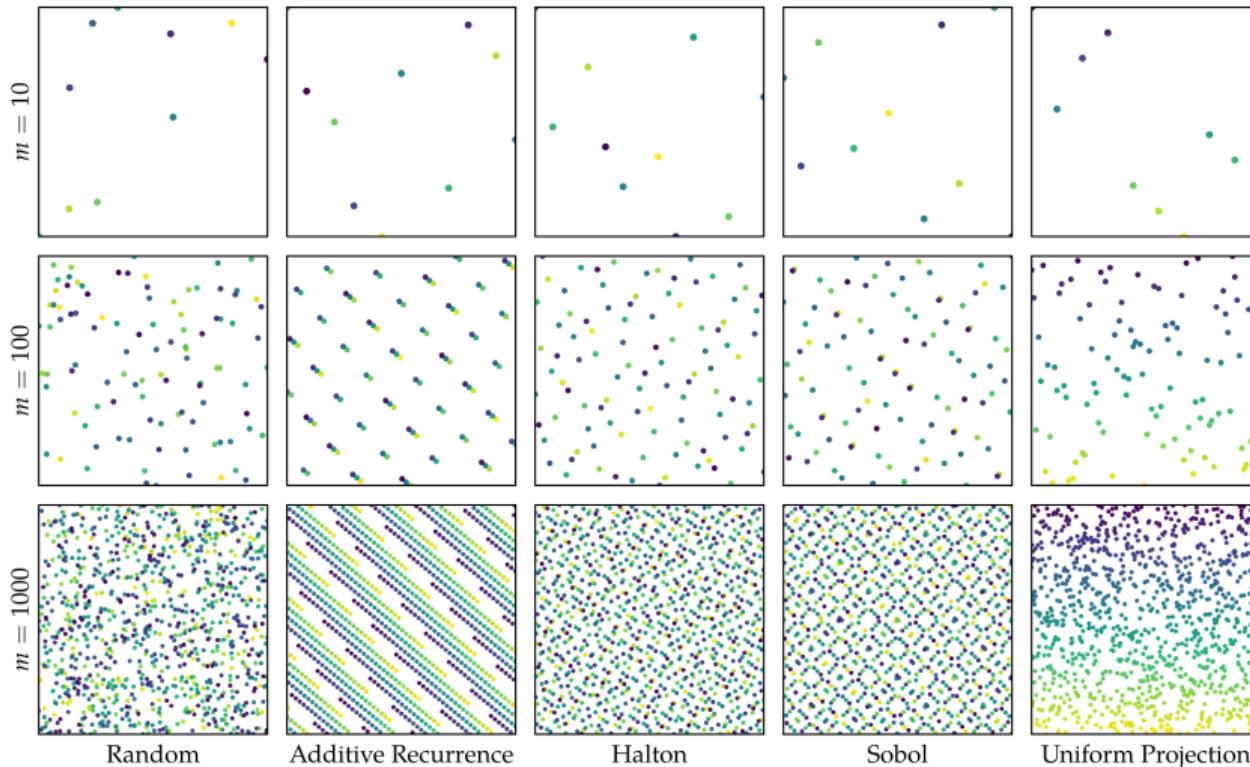
$$x_{i,j} = x_{i-1,j} \oplus v_{k,j}$$

$v_{k,j}$ is the j th dimension of the k th direction number.

- **direction numbers** $v_{k,j} = (0.v_{k,j,1}v_{k,j,2}\dots)_2$ where $v_{k,j,m}$ denotes the m th digit after the binary point.
- Tables of direction numbers with different properties have been proposed.
- Initialization: unit initialisation: ℓ th left most bit set to one $v_{k,j,\ell} = 1$ for all k and j and all others to be zero

Quasi-Random Sequences

space-filling sampling plans in two dimensions. Samples are colored according to the order in which they are sampled. The uniform projection plan was generated randomly and is not optimized.



Summary

- Sampling plans are used to cover search spaces with a limited number of points
- Full factorial sampling, which involves sampling at the vertices of a uniformly discretized grid, requires a number of points exponential in the number of dimensions
- Uniform projection plans, which project uniformly over each dimension, can be efficiently generated and can be optimized to be space-filling
- Greedy construction and the exchange local search algorithm can be used to find a subset of points that maximally fill a space
- Quasi-random sequences are deterministic procedures by which space-filling sampling plans can be generated