# PV027 Optimization

Tomáš Brázdil

### Resources & Prerequisities

Resources:

- Lectures & tutorials (the main resources)
- Books:

Joaquim R. R. A. Martins and Andrew Ning. Engineering Design Optimization. Cambridge University Press, 2021. ISBN: 9781108833417.

Jorge Nocedal and Stephen J. Wright. Numerical optimization. Springer, 2006. ISBN: 0387303030.

### Resources & Prerequisities

Resources:

- Lectures & tutorials (the main resources)
- Books:

Joaquim R. R. A. Martins and Andrew Ning. Engineering Design Optimization. Cambridge University Press, 2021. ISBN: 9781108833417.

Jorge Nocedal and Stephen J. Wright. Numerical optimization. Springer, 2006. ISBN: 0387303030.

We shall need elementary knowledge and understanding of

- Linear algebra in R<sup>n</sup>
   Operations with vectors and matrices, bases, diagonalization.
- Multi-variable calculus (i.e., in R<sup>n</sup>)
   Partial derivatives, gradients, Hessians, Taylor's theorem.

We will refresh our memories during lectures and tutorials.

### Evaluation

 $\mbox{Oral exam}$  - You will get a manual describing the knowledge necessary for  $\mbox{E}$  and better.

There might be homework assignments that you may discuss at tutorials, but (for this year) there is no mandatory homework.

Please be aware that

This is a **difficult math-based course**.



### What is Optimization

#### Merriam Webster:

An act, process, or methodology of making something (such as a design, system, or decision) as perfect, functional, or effective as possible.

*specifically*: the mathematical procedures (such as finding the maximum of a function) involved in this.

## What is Optimization

#### Merriam Webster:

An act, process, or methodology of making something (such as a design, system, or decision) as perfect, functional, or effective as possible.

*specifically*: the mathematical procedures (such as finding the maximum of a function) involved in this.

#### Britannica

Collection of mathematical principles and methods for solving quantitative problems in many disciplines, including physics, biology, engineering, economics, and business.

Historically, (mathematical/numerical) optimization is called *mathematical programming*.

People optimize in

#### scheduling

- transportation,
- education,
- ▶ ...

People optimize in

- scheduling
  - transportation,
  - education,
  - • •
- investments
  - portfolio management,
  - utility maximization,
  - ▶ ...

People optimize in

- scheduling
  - transportation,
  - education,
  - ...
- investments
  - portfolio management,
  - utility maximization,
  - • •
- industrial design
  - aerodynamics,
  - electrical engineering,
  - • •

People optimize in

- scheduling
  - transportation,
  - education,
  - ...
- investments
  - portfolio management,
  - utility maximization,
  - • •
- industrial design
  - aerodynamics,
  - electrical engineering,
  - **>** ...
- sciences
  - molecular modeling,
  - computational systems biology,

...

People optimize in

- scheduling
  - transportation,
  - education,
  - ...
- investments
  - portfolio management,
  - utility maximization,
  - • •
- industrial design
  - aerodynamics,
  - electrical engineering,
  - **>** ...
- sciences
  - molecular modeling,
  - computational systems biology,
  - **>** ...
- machine learning

## **Optimization Algorithms**

scipy.optimize.minimize

scipy.optimize.minimize(fun, xθ, args=(), method=None, jac=None, hess=None, hessp=None, bounds=None, constraints=(), tol=None, callback=None, options=None)

#### method : str or callable, optional

Type of solver. Should be one of

- 'Nelder-Mead' (see here)
- 'Powell' (see here)
- 'CG' (see here)
- 'BFGS' (see here)
- 'Newton-CG' (see here)
- 'L-BFGS-B' (see here)

## **Optimization Algorithms**

#### sklearn.linear\_model.LogisticRegression

class sklearn.linear\_model.LogisticRegression(penalty='12', \*, dual=False, tol=0.0001, C=1.0, fit\_intercept=True, intercept\_scaling=1, class\_weight=None, random\_state=None, solver='lbfgs', max\_iter=100, multi\_class='auto', verbose=0, warm\_start=False, n\_jobs=None, l1\_ratio=None)

solver : {'lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky', 'sag', 'saga'}, default='lbfgs' Algorithm to use in the optimization problem. Default is 'lbfgs'. To choose a solver,

## Design Optimization Process



## Design Optimization Process



- Consider a company with several plants producing a single product but with different efficiency.
- The goal is to set the production of each plant so that demand for goods is satisfied, but overproduction is minimized.

- Consider a company with several plants producing a single product but with different efficiency.
- The goal is to set the production of each plant so that demand for goods is satisfied, but overproduction is minimized.
- First try: Model each plant's production and maximize the total production efficiency.

This would lead to a solution where only the most efficient plant will produce.

- Consider a company with several plants producing a single product but with different efficiency.
- The goal is to set the production of each plant so that demand for goods is satisfied, but overproduction is minimized.
- First try: Model each plant's production and maximize the total production efficiency.

This would lead to a solution where only the most efficient plant will produce.

► However, after a certain level of demand, no single plant can satisfy the demand ⇒, introducing constraints on the maximum production of the plants.

This would maximize production of the most efficient plant and then the second one, etc.

- Consider a company with several plants producing a single product but with different efficiency.
- The goal is to set the production of each plant so that demand for goods is satisfied, but overproduction is minimized.
- First try: Model each plant's production and maximize the total production efficiency.

This would lead to a solution where only the most efficient plant will produce.

► However, after a certain level of demand, no single plant can satisfy the demand ⇒, introducing constraints on the maximum production of the plants.

This would maximize production of the most efficient plant and then the second one, etc.

- Then you notice that all plant employees must work.
- Then you start solving transportation problems depending on the location of the plants.

- 1. Describe the problem
  - Problem formulation is vital since the optimizer exploits any weaknesses in the model formulation.
  - You might get the "right answer to the wrong question."
  - The problem description is typically informal at the beginning.
- 2. Gather information
  - Identify possible inputs/outputs.
  - Gather data and identify the analysis procedure.



- 3. Define the design variables
  - Identify the quantities that describe the system:

 $x \in \mathbb{R}^n$ 

(i.e., certain characteristics of the system, such as position, investments, etc.)

- The variables are supposed to be independent; the optimizer must be free to choose the components of x independently.
- The choice of variables is typically not unique (e.g., a square can be described by its side or area).
- The variables may affect the functional form of the objective and constraints (e.g., linear vs non-linear).



4. Define the **objective** 

- The function determines if one design is better than another.
- Must be a scalar computable from the variables:

 $f:\mathbb{R}^n\to\mathbb{R}$ 

(e.g., profit, time, potential energy, etc.)
The objective function is either maximized or minimized depending on the application.
The choice is not always obvious: E.g., minimizing just the weight of a vehicle might result in a vehicle being too expensive to be manufactured.



#### 5. Define the constraints

- Prescribe allowed values of the variables.
- May have a general form

$$c(x) \leq 0$$
 or  $c(x) \geq 0$  or  $c(x) = 0$ 

(e.g., time cannot be negative, bounded amount of money to invest) Where  $c : \mathbb{R}^n \to \mathbb{R}$  is a function depending on the variables.



The Optimization Problem consists of

- variables
- ► objective
- constraints

The above components constitute a model.

#### The Optimization Problem consists of

- variables
- objective
- constraints

The above components constitute a model.

**Modelling** is concerned with model building, **optimization** with maximization/minimization of the objective for a given model.

We concentrate on the optimization part but keep in mind that it is intertwined with modeling.

### The Optimization Problem consists of

- variables
- objective
- constraints

The above components constitute a model.

**Modelling** is concerned with model building, **optimization** with maximization/minimization of the objective for a given model.

We concentrate on the optimization part but keep in mind that it is intertwined with modeling.

The **Optimization Problem (OP):** Find settings of variables so that the objective is maximized/minimized while satisfying the constraints.

### The Optimization Problem consists of

- variables
- objective
- constraints

The above components constitute a model.

**Modelling** is concerned with model building, **optimization** with maximization/minimization of the objective for a given model.

We concentrate on the optimization part but keep in mind that it is intertwined with modeling.

The **Optimization Problem (OP):** Find settings of variables so that the objective is maximized/minimized while satisfying the constraints.

An **Optimization Algorithm (OA)** solves the above problem and provides a **solution**, some setting of variables satisfying the constraints and minimizing/maximizing the objective.

# **Optimization Problems**

## **Optimization Problem Formally**

Denote by

- $f: \mathbb{R}^n \to \mathbb{R}$  an objective function,
- x a vector of real variables,

 $g_1, \ldots, g_{n_g}$  inequality constraint functions  $g_i : \mathbb{R}^n \to \mathbb{R}$ .

 $h_1, \ldots, h_{n_h}$  equality constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

## **Optimization Problem Formally**

Denote by

- $f: \mathbb{R}^n \to \mathbb{R}$  an objective function,
- x a vector of real variables,

 $g_1, \ldots, g_{n_g}$  inequality constraint functions  $g_i : \mathbb{R}^n \to \mathbb{R}$ .  $h_1, \ldots, h_{n_h}$  equality constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

The optimization problem is to

 $\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i = 1, \dots, n_g \\ & h_j(x) = 0 \quad j = 1, \dots, n_h \end{array}$ 

## Optimization Problem - Example

$$f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 1)^2$$
  

$$g_1(x_1, x_2) = x_1^2 - x_2$$
  

$$g_2(x_1, x_2) = x_1 + x_2 - 2$$

The optimization problem is

minimize 
$$(x_1-2)^2+(x_2-1)^2$$
 subject to  $\begin{cases} x_1^2-x_2\leq 0,\\ x_1+x_2-2\leq 0. \end{cases}$ 

**Optimization Problem - Example** 

$$f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 1)^2$$
  

$$g_1(x_1, x_2) = x_1^2 - x_2$$
  

$$g_2(x_1, x_2) = x_1 + x_2 - 2$$



A contour of f is defined, for some  $c \in \mathbb{R}$ , by  $\{x \in \mathbb{R}^n \mid f(x) = c\}$  15

Consider the constraints

$$g_i(x) \leq 0$$
  $i = 1, \dots, n_g$   
 $h_j(x) = 0$   $j = 1, \dots, n_h$ 

Consider the constraints

$$g_i(x) \leq 0$$
  $i = 1, \dots, n_g$   
 $h_j(x) = 0$   $j = 1, \dots, n_h$ 

Define the *feasibility region* by

$$\mathcal{F} = \{x \mid g_i(x) \leq 0, h_j(x) = 0, i = 1, \dots, n_g, j = 1, \dots, n_h\}$$

 $x \in \mathcal{F}$  is feasible,  $x \notin \mathcal{F}$  is infeasible.

Consider the constraints

$$g_i(x) \leq 0 \quad i = 1, \dots, n_g$$
  
$$h_j(x) = 0 \quad j = 1, \dots, n_h$$

Define the *feasibility region* by

$$\mathcal{F} = \{x \mid g_i(x) \leq 0, h_j(x) = 0, i = 1, \dots, n_g, j = 1, \dots, n_h\}$$

 $x \in \mathcal{F}$  is *feasible*,  $x \notin \mathcal{F}$  is *infeasible*.

Note that constraints of the form  $g_i(x) \ge 0$  can be easily transformed to the inequality contraints  $-g_i(x) \le 0$ 

Consider the constraints

$$g_i(x) \leq 0 \quad i = 1, \dots, n_g$$
  
$$h_j(x) = 0 \quad j = 1, \dots, n_h$$

Define the *feasibility region* by

$$\mathcal{F} = \{x \mid g_i(x) \leq 0, h_j(x) = 0, i = 1, \dots, n_g, j = 1, \dots, n_h\}$$

#### $x \in \mathcal{F}$ is *feasible*, $x \notin \mathcal{F}$ is *infeasible*.

Note that constraints of the form  $g_i(x) \ge 0$  can be easily transformed to the inequality contraints  $-g_i(x) \le 0$ 

 $x^* \in \mathcal{F}$  is now a *constrained minimizer* if

$$f(x^*) \leq f(x)$$
 for all  $x \in \mathcal{F}$
#### Constraints

Inequality constraints  $g_i(x) \leq 0$  can be active or inactive.



The problem formulation:

- ► A company has two chemical factories F<sub>1</sub> and F<sub>2</sub>, and a dozen retail outlets R<sub>1</sub>,..., R<sub>12</sub>.
- Each F<sub>i</sub> can produce (maximum of) a<sub>i</sub> tons of a chemical each week.
- Each retail outlet  $R_j$  demands at least  $b_j$  tons.
- The cost of shipping one ton from F<sub>i</sub> to R<sub>j</sub> is c<sub>ij</sub>.

The problem formulation:

- ► A company has two chemical factories F<sub>1</sub> and F<sub>2</sub>, and a dozen retail outlets R<sub>1</sub>,..., R<sub>12</sub>.
- Each F<sub>i</sub> can produce (maximum of) a<sub>i</sub> tons of a chemical each week.
- Each retail outlet  $R_j$  demands at least  $b_j$  tons.
- The cost of shipping one ton from  $F_i$  to  $R_j$  is  $c_{ij}$ .

**The problem:** Determine how much each factory should ship to each outlet to satisfy the requirements and minimize cost.

Variables:  $x_{ij}$  for i = 1, 2 and j = 1, ..., 12. Each  $x_{ij}$  (intuitively) corresponds to tons shipped from  $F_i$  to  $R_i$ .

The objective:



Variables:  $x_{ij}$  for i = 1, 2 and j = 1, ..., 12. Each  $x_{ij}$  (intuitively) corresponds to tons shipped from  $F_i$  to  $R_i$ .

The objective:

$$\min\sum_{ij}c_{ij}x_{ij}$$

subject to

$$\sum_{j=1}^{12} x_{ij} \le a_i, \quad i = 1, 2$$
$$\sum_{i=1}^{2} x_{ij} \ge b_j, \quad j = 1, \dots, 12,$$
$$x_{ij} \ge 0, \quad i = 1, 2, \quad j = 1, \dots, 12.$$

Variables:  $x_{ij}$  for i = 1, 2 and j = 1, ..., 12. Each  $x_{ij}$  (intuitively) corresponds to tons shipped from  $F_i$  to  $R_i$ .

The objective:

$$\min\sum_{ij}c_{ij}x_{ij}$$

subject to

$$\sum_{j=1}^{12} x_{ij} \le a_i, \quad i = 1, 2$$
$$\sum_{i=1}^{2} x_{ij} \ge b_j, \quad j = 1, \dots, 12,$$
$$x_{ij} \ge 0, \quad i = 1, 2, \quad j = 1, \dots, 12.$$

The above is *linear programming* problem since both the objective and constraint functions are linear.

## **Discrete Optimization**

In our original optimization problem definition, we consider real (continuous) variables.

Sometimes, we need to assume discrete values. For example, in the previous example, the factories may produce tractors. In such a case, it does not make sense to produce 4.6 tractors.

# **Discrete Optimization**

In our original optimization problem definition, we consider real (continuous) variables.

Sometimes, we need to assume discrete values. For example, in the previous example, the factories may produce tractors. In such a case, it does not make sense to produce 4.6 tractors.

Usually, an integer constraint is added, such as

 $x_i \in \mathbb{Z}$ 

It constrains  $x_i$  only to integer values. This leads to so-called *integer programming*.

Discrete optimization problems have discrete and finite variables.

Our goal is to design the wing shape of an aircraft.

Assume a rectangular wing.

c b

The parameters are called *span b* and *chord c*.

Our goal is to design the wing shape of an aircraft.

Assume a rectangular wing.

b c

The parameters are called *span b* and *chord c*.

However, two other variables are often used in aircraft design: Wing area S and wing aspect ratio AR. It holds that



What exactly are the objectives and constraints?

What exactly are the objectives and constraints?

Our objective function is the power required to keep level flight:

$$f(b,c) = \frac{Dv}{\eta}$$

Here,

D is the drag

That is the aerodynamic force that opposes an aircraft's motion through the air.

•  $\eta$  is the propulsive efficiency

That is the efficiency with which the energy contained in a vehicle's fuel is converted into kinetic energy of the vehicle.

v is the lift velocity

That is the velocity needed to lift the aircraft, which depends on its weight.

For illustration, let us look at the lift velocity v.

For illustration, let us look at the lift velocity v.

In level flight, the aircraft must generate enough lift L to equal its weight W, that is L = W.

For illustration, let us look at the lift velocity v.

In level flight, the aircraft must generate enough lift L to equal its weight W, that is L = W.

The weight partially depends on the wing area:

 $W = W_0 + W_S S$ 

Here S = bc is the wing area, and  $W_0$  is the payload weight.

For illustration, let us look at the lift velocity v.

In level flight, the aircraft must generate enough lift L to equal its weight W, that is L = W.

The weight partially depends on the wing area:

 $W = W_0 + W_S S$ 

Here S = bc is the wing area, and  $W_0$  is the payload weight.

The lift can be approximated using the following formula.

$$L = q \cdot C_L \cdot S$$

Where  $q = \frac{1}{2}\rho v^2$  is the fluid dynamic pressure, here  $\rho$  is the air density,  $C_L$  is a lift coefficient (depending on the wing shape).

For illustration, let us look at the lift velocity v.

In level flight, the aircraft must generate enough lift L to equal its weight W, that is L = W.

The weight partially depends on the wing area:

 $W = W_0 + W_S S$ 

Here S = bc is the wing area, and  $W_0$  is the payload weight.

The lift can be approximated using the following formula.

$$L = q \cdot C_L \cdot S$$

Where  $q = \frac{1}{2}\rho v^2$  is the fluid dynamic pressure, here  $\rho$  is the air density,  $C_L$  is a lift coefficient (depending on the wing shape).

Thus, we may obtain the lift velocity as

$$v = \sqrt{2W/\varrho C_L S} = \sqrt{2(W_0 + W_S bc)/\varrho C_L bc}$$

Similarly, various physics-based arguments provide approximations of the drag D and the propulsion efficiency  $\eta$ .

The drag  $D = D_i + D_f$  is the sum of the induced and viscous drag.

The drag  $D = D_i + D_f$  is the sum of the induced and viscous drag.

The induced drag can be approximated by

$$D_i = W^2/q \pi b^2 \epsilon$$

Here, e is the Oswald efficiency factor, a correction factor that represents the change in drag with the lift of a wing, as compared with an ideal wing having the same aspect ratio.

The drag  $D = D_i + D_f$  is the sum of the induced and viscous drag.

The induced drag can be approximated by

$$D_i = W^2/q \pi b^2 e$$

Here, e is the Oswald efficiency factor, a correction factor that represents the change in drag with the lift of a wing, as compared with an ideal wing having the same aspect ratio.

The viscous drag can be approximated by

 $D_f = k C_f q 2.05 S$ 

Here, k is the form factor (accounts for the pressure drag), and  $C_f$  is the skin friction coefficient that can be approximated by

$$C_f = 0.074/Re^{0.2}$$

Where Re is the Reynolds number that somewhat characterizes air flow patterns around the wing and is defined as follows:

$$\textit{Re} = \rho\textit{vc}/\mu$$

Here  $\mu$  is the air dynamic viscosity.

The propulsion efficiency  $\eta$  can be roughly approximated by the Gaussian efficiency curve.

$$\eta = \eta_{\max} \exp\left(\frac{-(v-\bar{v})^2}{2\sigma^2}\right)$$

Here,  $\bar{\mathbf{v}}$  is the peak propulsive efficiency velocity, and  $\sigma$  is the std of the efficiency function.

The objective function contours:



The objective function contours:



The engineers would refuse the solution: The aspect ratio is much higher than typically seen in airplanes. It adversely affects the structural strength. Add constraints!

Added a constraint on bending stress at the root of the wing:



It looks like a reasonable wing ...







Continuous allows only x<sub>i</sub> ∈ ℝ, discrete allows only x<sub>i</sub> ∈ ℤ, mixed allows variables of both kinds.



- Continuous allows only x<sub>i</sub> ∈ ℝ, discrete allows only x<sub>i</sub> ∈ ℤ, mixed allows variables of both kinds.
- ▶ Single-objective:  $f : \mathbb{R}^n \to \mathbb{R}$ , Multi-objective:  $f : \mathbb{R}^n \to \mathbb{R}^m$



- Continuous allows only x<sub>i</sub> ∈ ℝ, discrete allows only x<sub>i</sub> ∈ ℤ, mixed allows variables of both kinds.
- ▶ Single-objective:  $f : \mathbb{R}^n \to \mathbb{R}$ , Multi-objective:  $f : \mathbb{R}^n \to \mathbb{R}^m$
- Unconstrained: No constraints, just the objective function.



# Smoothness

We consider various classes of problems depending on the smoothness properties of the objective/constraint functions:

 C<sup>0</sup>: Continuous function Continuity allows us to estimate value in small neighborhoods.

Discontinuous functions exist.

C<sup>1</sup>: Continuous first derivatives

The derivatives give information on the slope. If continuous, it changes smoothly, allowing us to estimate the slope locally.

Nondifferentiable continuous functions and differentiable functions with discontinuous derivatives exist.

 C<sup>2</sup>: Continuous second derivatives The second derivatives inform about curvature.

Continuously differentiable functions without second derivatives and twice differentiable functions with discontinuous second derivatives exist.

f(x) = |x| is continuous, f is not differentiable at 0



f(x) = x|x| is differentiable on  $\mathbb{R}$ , f' has no second derivative at 0



$$f(x) = \begin{cases} x^2 \sin(1/x) & \text{if } x \neq 0\\ 0 & \text{if } x = 0 \end{cases}$$



$$f'(x) = \begin{cases} 2x\sin(1/x) - \cos(1/x), & x \neq 0\\ 0, & x = 0 \end{cases}$$



f is differentiable on  $\mathbb{R}$ , f' is not continuous at 0

$$f(x) = \begin{cases} x^4 \sin(1/x) & \text{if } x \neq 0\\ 0 & \text{if } x = 0 \end{cases}$$

f is differentiable on  $\mathbb{R}$ ,

$$f'(x) = \begin{cases} 4x^3 \sin(1/x) - x^2 \cos(1/x), & x \neq 0\\ 0, & x = 0 \end{cases}$$

f' is differentiable on  $\mathbb{R}$ ,

$$f''(x) = \begin{cases} 12x^2 \sin(1/x) - 6x \cos(1/x) - \sin(1/x), & x \neq 0\\ 0, & x = 0 \end{cases}$$

Clearly, f'' does not have a limit at 0 as sin(1/x) oscillates between -1 and 1 and thus is not continuous.

## Linearity

Linear programming: Both the objective and the constraints are linear.



It is possible to solve precisely, efficiently, and in rational numbers (see the linear programming later).

# Multimodality

Denote by  ${\mathcal F}$  the feasibility set.

# $x^*$ is a *(weak) local minimiser* if there is $\varepsilon > 0$ such that $f(x^*) \le f(x)$ for all $x \in \mathcal{F}$ satisfying $||x^* - x|| \le \varepsilon$
# Multimodality

Denote by  ${\mathcal F}$  the feasibility set.

 $x^*$  is a *(weak) local minimiser* if there is  $\varepsilon > 0$  such that

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$  satisfying  $||x^* - x|| \leq \varepsilon$ 

x\* is a (weak) global minimiser if

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$ 

# Multimodality

Denote by  ${\mathcal F}$  the feasibility set.

 $x^*$  is a *(weak) local minimiser* if there is  $\varepsilon > 0$  such that

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$  satisfying  $||x^* - x|| \leq \varepsilon$ 

x\* is a (weak) global minimiser if

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$ 

Global/local minimiser is *strict* if the inequality is strict.

# Multimodality

Denote by  ${\mathcal F}$  the feasibility set.

 $x^*$  is a (weak) local minimiser if there is  $\varepsilon > 0$  such that

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$  satisfying  $||x^* - x|| \leq \varepsilon$ 

x\* is a (weak) global minimiser if

 $f(x^*) \leq f(x)$  for all  $x \in \mathcal{F}$ 

Global/local minimiser is *strict* if the inequality is strict.



Unimodal functions have a single global minimiser in  $\mathcal{F}$ , multimodal have multiple local minimisers in  $\mathcal{F}$ .

Convexity

 $S \subseteq \mathbb{R}^n$  is a *convex set* if the straight line segment connecting any two points in S lies entirely inside S. Formally, for any two points  $x \in S$  and  $y \in S$ , we have  $\alpha x + (1 - \alpha)y \in S$  for all  $\alpha \in [0, 1]$ 

Convexity

 $S \subseteq \mathbb{R}^n$  is a *convex set* if the straight line segment connecting any two points in S lies entirely inside S. Formally, for any two points  $x \in S$  and  $y \in S$ , we have  $\alpha x + (1 - \alpha)y \in S$  for all  $\alpha \in [0, 1]$ 

*f* is a *convex function* if its domain is a convex set and if for any two points *x* and *y* in this domain, the graph of *f* lies below the straight line connecting (x, f(x)) to (y, f(y)) in the space  $\mathbb{R}^{n+1}$ . That is, we have

 $f(\alpha x + (1 - \alpha)y) \le \alpha f(x) + (1 - \alpha)f(y), \quad \text{ for all } \alpha \in (0, 1).$ 

Convexity

 $S \subseteq \mathbb{R}^n$  is a *convex set* if the straight line segment connecting any two points in S lies entirely inside S. Formally, for any two points  $x \in S$  and  $y \in S$ , we have  $\alpha x + (1 - \alpha)y \in S$  for all  $\alpha \in [0, 1]$ 

*f* is a *convex function* if its domain is a convex set and if for any two points *x* and *y* in this domain, the graph of *f* lies below the straight line connecting (x, f(x)) to (y, f(y)) in the space  $\mathbb{R}^{n+1}$ . That is, we have

 $f(\alpha x + (1 - \alpha)y) \le \alpha f(x) + (1 - \alpha)f(y), \quad \text{ for all } \alpha \in (0, 1).$ 

A standard form convex optimization assumes

- convex objective f and convex inequality constraint functions g<sub>i</sub>
- affine equality constraint functions h<sub>j</sub>

#### Implications:

- Every local minimum is a global minimum.
- If the above inequality is strict for all x ≠ y, then there is a unique minimum.

# Stochasticity

Sometimes, the parameters of a model cannot be specified with certainty.

For example, in the transportation model, customer demand cannot be predicted precisely in practice.

However, such parameters may often be statistically estimated and modeled using an appropriate probability distribution.

# Stochasticity

Sometimes, the parameters of a model cannot be specified with certainty.

For example, in the transportation model, customer demand cannot be predicted precisely in practice.

However, such parameters may often be statistically estimated and modeled using an appropriate probability distribution.

*Stochastic optimization* problem is to minimize/maximize the expectation of a statistic parametrized with the variables *x*:

Find x maximizing  $\mathbb{E}f(x; W)$ 

Here, W is a vector of random variables, and the expectation is taken using the probability distribution of these variables.

In this course, we stick with *deterministic optimization*.

# **Optimization Algorithms**

## **Optimization Algorithm**

An optimization algorithm solves the optimization problem, i.e., searches for  $x^*$ , which (in some sense) minimizes the objective f and satisfies the constraints.

Typically, the algorithm computes a set of candidate solutions  $x_0, x_1, \ldots$  and then identifies one resembling a solution.

# **Optimization Algorithm**

An optimization algorithm solves the optimization problem, i.e., searches for  $x^*$ , which (in some sense) minimizes the objective f and satisfies the constraints.

Typically, the algorithm computes a set of candidate solutions  $x_0, x_1, \ldots$  and then identifies one resembling a solution.

The problem is to

compute the candidate solutions,

Complexity of the objective function, difficulties in selection of the candidates, etc.

Select the one closest to a minimum. It is Hard to decide whether a given point is a minimum (even a local one). Example: Neural networks training.

Typically, we are concerned with the following issues:

Typically, we are concerned with the following issues:

 Robustness: OA should perform well on various problems in their class for all reasonable choices of the initial variables.

Typically, we are concerned with the following issues:

- Robustness: OA should perform well on various problems in their class for all reasonable choices of the initial variables.
- Efficiency: OA should not require too much computer time or storage.

Typically, we are concerned with the following issues:

- Robustness: OA should perform well on various problems in their class for all reasonable choices of the initial variables.
- *Efficiency*: OA should not require too much computer time or storage.
- Accuracy: OA should be able to identify a solution with precision without being overly sensitive to
  - errors in the data/model
  - the arithmetic rounding errors



# Order and Search

#### Order

- Zeroth = gradient-free: no info about derivatives is used
- First = gradient-based: use info about first derivatives (e.g., gradient descent)
- Second = use info about first and second derivatives (e.g., Newton's method)

# Order and Search

#### Order

- Zeroth = gradient-free: no info about derivatives is used
- First = gradient-based: use info about first derivatives (e.g., gradient descent)
- Second = use info about first and second derivatives (e.g., Newton's method)

Search

- Local search = start at a point and search for a solution by successively updating the current solution (e.g., gradient descent)
- Global search tries to span the whole space (e.g., grid search)

For some algorithms and under specific assumptions imposed on the optimization problem, we can do the following:

Prove that the algorithm converges to an optimum/minimum.

For some algorithms and under specific assumptions imposed on the optimization problem, we can do the following:

- Prove that the algorithm converges to an optimum/minimum.
- Determine the rate of convergence.

For some algorithms and under specific assumptions imposed on the optimization problem, we can do the following:

- Prove that the algorithm converges to an optimum/minimum.
- Determine the rate of convergence.
- Decide whether we are at (or close to) an optimum/minimum.

For some algorithms and under specific assumptions imposed on the optimization problem, we can do the following:

- Prove that the algorithm converges to an optimum/minimum.
- Determine the rate of convergence.
- Decide whether we are at (or close to) an optimum/minimum.

For example, for linear optimization problems, the simplex algorithm converges to a minimum (or says that there is no minimum) in, at most, exponentially many steps, and we may efficiently decide whether we have reached a minimum.

For some algorithms and under specific assumptions imposed on the optimization problem, we can do the following:

- Prove that the algorithm converges to an optimum/minimum.
- Determine the rate of convergence.
- Decide whether we are at (or close to) an optimum/minimum.

For example, for linear optimization problems, the simplex algorithm converges to a minimum (or says that there is no minimum) in, at most, exponentially many steps, and we may efficiently decide whether we have reached a minimum.

We may prove only some or none of the properties for some algorithms.

There are (almost) infinitely many heuristic algorithms without provable convergence, often motivated by the behaviors of various animals.

Deterministic vs Stochastic and Static vs Dynamic

*Stochastic optimization* is based on a random selection of candidate solutions.

Evolutionary algorithms contain some randomness (e.g., in the form of random mutations).

Also, various variants of the gradient-based methods are often randomized (e.g., variants of the stochastic gradient descent).

Deterministic vs Stochastic and Static vs Dynamic

*Stochastic optimization* is based on a random selection of candidate solutions.

Evolutionary algorithms contain some randomness (e.g., in the form of random mutations).

Also, various variants of the gradient-based methods are often randomized (e.g., variants of the stochastic gradient descent).

In this course, we stick to *static* optimization problems where we solve the optimization problem only once.

In contrast, the *dynamic* optimization, a sequence of (usually) dependent optimization problems are solved sequentially.

For example, consider driving a car where the driver must react optimally to changing situations several times per second.

Dynamic optimization problems are usually defined using a kind of (Markov) decision process.

# Summary

The course consists of the following main parts:

- Unconstrained optimization
  - Non-linear objectives, (twice) differentiable
  - Second-order methods (quasi-Newton)
- Constrained optimization
  - Non-linear objectives and constraints, (twice) differentiable
  - Lagrange multipliers, Newton-Lagrange method
  - Quadratic programming (a little bit)
- Linear programming
  - Linear objectives and constraints
  - Simplex algorithm deep dive (including the degenerate case)
- Integer linear programming
  - Linear objectives and mixed integer linear constraints
  - Branch-and-bound, Gomory cuts algorithms
- A little bit on non-differentiable algorithms.

You will need to understand: Calculus in  $\mathbb{R}^n$  (gradient, Hessian) and linear algebra in  $\mathbb{R}^n$  (vectors, matrices, geometry)

# Single-variable Objectives

## Unconstrained Single Variable Optimization Problem

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable x

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

# Unconstrained Single Variable Optimization Problem

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable x

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

We consider

- f continuously differentiable
- f twice continuously differentiable

Present the following methods:

- Gradient descent
- Newton's method
- Secant method

# Gradient Based Methods

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable  $x \in \mathbb{R}$ 

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

### Gradient Based Methods

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable  $x \in \mathbb{R}$ 

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

Assume that

$$f'(x) = \lim_{h \to 0} rac{f(x+h) - f(x)}{h}$$
 for  $x \in \mathbb{R}$ 

is continuous on  $\mathbb{R}$ .

Denote by  $C^1$  the set of all continuously differentiable functions.

### Gradient Descent in Single Variable

Gradient descent algorithm for finding a local minimum of a function f, using a variable step length.

**Input:** Function f with first derivative f', initial point  $x_0$ , initial step length  $\alpha_0 > 0$ , tolerance  $\epsilon > 0$ 

**Output:** A point x that approximately minimizes f(x)

- 1: Set  $k \leftarrow 0$
- 2: while  $|f'(x_k)| > \epsilon$  do
- 3: Calculate the derivative:  $y' \leftarrow f'(x_k)$
- 4: Update  $x_{k+1} \leftarrow x_k \alpha_k \cdot y'$
- 5: Update step length  $\alpha_k$  to  $\alpha_{k+1}$  based on a certain strategy
- 6: Increment k
- 7: end while
- 8: **return** *x*<sub>*k*</sub>

Convergence of Single Variable Gradient Descent

#### Theorem 1

Assume that f is

- differentiable, i.e., that f' exists,
- bounded below, i.e., there is B ∈ ℝ such that f(x) ≥ B for all x ∈ ℝ,

► L-smooth, i.e., there is 
$$L > 0$$
 such that  $|f'(x) - f'(x')| \le L|x - x'|$  for all  $x, x' \in \mathbb{R}$ 

Consider a sequence  $x_0, x_1, \ldots$  computed by the gradient descent algorithm for f. Assume a constant step length  $\alpha \leq \frac{1}{L}$ . Then  $\lim_{k\to\infty} |f'(x_k)| = 0$  and, moreover,

$$\min_{0 \le t < T} |f'(x_t)| \le \sqrt{\frac{2L(f(x_0) - B)}{T}}$$

#### Example

Consider the following objective function f

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



#### Example

Consider the objective function f

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

Assume  $x_0 = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-4}$ , i.e., we stop when  $|x_{k+1} - x_k| < \epsilon$ .

Consider the step length  $\alpha = 1$ .

#### Example

Consider the objective function f

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

Assume  $x_0 = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-4}$ , i.e., we stop when  $|x_{k+1} - x_k| < \epsilon$ .

Consider the step length  $\alpha = 1$ .

We compute

$$f'(x) = x - \cos x.$$

Then,

$$\begin{aligned} x_1 &= 0.5 - (0.5 - \cos 0.5) \\ &= 0.5 - (-0.37758) \\ &= 0.87758 \end{aligned}$$
Continuing in the same way:

$x_1 = 0.87758$	$x_{12} = 0.73724$
$x_2 = 0.63901$	$x_{13} = 0.74033$
$x_3 = 0.80269$	$x_{14} = 0.73825$
$x_4 = 0.69478$	$x_{15} = 0.73965$
$x_5 = 0.76820$	$x_{16} = 0.73870$
$x_6 = 0.71917$	$x_{17} = 0.73934$
$x_7 = 0.75236$	$x_{18} = 0.73891$
$x_8 = 0.73008$	$x_{19} = 0.73920$
$x_9 = 0.74512$	$x_{20} = 0.73901$
$x_{10} = 0.73501$	$x_{21} = 0.73914$
$x_{11} = 0.74183$	$x_{22} = 0.73905$

Note that  $|x_{22} - x_{21}| < 10^{-4}$ .

What if we consider the step length 1/k? Then

- $x_1 = 0.50000$
- $x_2 = 0.87758$
- $x_3 = 0.75830$
- $x_4 = 0.74753$
- $x_5 = 0.74399$
- $x_6 = 0.74235$
- $x_7 = 0.74144$
- $x_8 = 0.74087$
- $x_9 = 0.74050$
- $x_{10} = 0.74024$
- $x_{11} = 0.74004$
- $x_{12} = 0.73990$
- $x_{13} = 0.73978$
- $x_{14} = 0.73969$

Note that  $|x_{14} - x_{13}| < 10^{-4}$  but  $x_{14}$  is far from the solution which is 0.7390....

What if we consider the step length 1/k? Then

$x_1 = 0.50000$	$x_{115} = 0.739100605$
$x_2 = 0.87758$	$x_{116} = 0.739100379$
$x_3 = 0.75830$	$x_{117} = 0.739100159$
<i>x</i> <sub>4</sub> = 0.74753	$x_{118} = 0.739099944$
$x_5 = 0.74399$	$x_{119} = 0.739099734$
$x_6 = 0.74235$	$x_{120} = 0.739099529$
$x_7 = 0.74144$	$x_{121} = 0.739099328$
<i>x</i> <sub>8</sub> = 0.74087	$x_{122} = 0.739099132$
$x_9 = 0.74050$	$x_{123} = 0.739098940$
$x_{10} = 0.74024$	$x_{124} = 0.739098752$
$x_{11} = 0.74004$	$x_{125} = 0.739098568$
$x_{12} = 0.73990$	$x_{126} = 0.739098388$
<i>x</i> <sub>13</sub> = 0.73978	$x_{127} = 0.739098212$
$x_{14} = 0.73969$	$x_{128} = 0.739098040$

#### Gradient descent with the step length = 1.0:



#### Gradient descent with the step length = 1/k:



Gradient descent with the step length  $= 1/k^2$ :



It does not seem to converge to the same number as the previous step lengths.

#### Gradient descent with the step length = 1.0:



#### Gradient descent with the step length = 1/k:



- ► The objective must be differentiable, however:
  - Can be extended to functions with few non-linearities by considering differentiable parts or sub-gradients.
  - There are methods for differentiable approximation of non-differentiable functions.

- ► The objective must be differentiable, however:
  - Can be extended to functions with few non-linearities by considering differentiable parts or sub-gradients.
  - There are methods for differentiable approximation of non-differentiable functions.
- GD is sensitive to the initial point: Converges to a local minimum for a small step length (typically) to the closest one.

- ► The objective must be differentiable, however:
  - Can be extended to functions with few non-linearities by considering differentiable parts or sub-gradients.
  - There are methods for differentiable approximation of non-differentiable functions.
- GD is sensitive to the initial point: Converges to a local minimum for a small step length (typically) to the closest one.
- GD is quite sensitive to the step length.

Might be very slow or too fast (even overshoot and diverge).

- ► The objective must be differentiable, however:
  - Can be extended to functions with few non-linearities by considering differentiable parts or sub-gradients.
  - There are methods for differentiable approximation of non-differentiable functions.
- GD is sensitive to the initial point: Converges to a local minimum for a small step length (typically) to the closest one.
- GD is quite sensitive to the step length.
   Might be very slow or too fast (even overshoot and diverge).
- For convex functions, the algorithm converges to a minimum (if it converges).

- ► The objective must be differentiable, however:
  - Can be extended to functions with few non-linearities by considering differentiable parts or sub-gradients.
  - There are methods for differentiable approximation of non-differentiable functions.
- GD is sensitive to the initial point: Converges to a local minimum for a small step length (typically) to the closest one.
- GD is quite sensitive to the step length.
   Might be very slow or too fast (even overshoot and diverge).
- For convex functions, the algorithm converges to a minimum (if it converges).

Straightforward to implement if the derivatives are available.

GD is much more interesting in multiple variables, forming the basis for neural network learning (see later).

Better algorithm for unimodal functions using just derivatives?

# Newton's Method

An objective function  $f : \mathbb{R} \to \mathbb{R}$ A variable  $x \in \mathbb{R}$ 

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

# Newton's Method

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable  $x \in \mathbb{R}$ 

Find  $x^*$  such that

$$f(x^*) \leq \min_{x \in \mathbb{R}} f(x)$$

Assume that

$$f''(x) = \lim_{h o 0} rac{f'(x+h) - f'(x)}{h} \quad ext{for } x \in \mathbb{R}$$

is continuous on  $\mathbb{R}$ .

Denote by  $\mathcal{C}^2$  the set of all twice continuously differentiable functions.

## Taylor Series Approximation

We would need the o-notation: Given functions  $f,g:\mathbb{R}\to\mathbb{R}$  we write f=o(g) if

$$\lim_{x\to 0}\frac{f(x)}{g(x)}=0$$

#### Taylor Series Approximation

We would need the o-notation: Given functions  $f,g:\mathbb{R}\to\mathbb{R}$  we write f=o(g) if

$$\lim_{x\to 0}\frac{f(x)}{g(x)}=0$$

Consider a function  $f : \mathbb{R} \to \mathbb{R}$  and  $x_0 \in \mathbb{R}$ . Assume that f is twice differentiable at  $x_0$ . Then for all  $x \in \mathbb{R}$  we have that

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + o(|x - x_0|^2)$$

#### Taylor Series Approximation

We would need the o-notation: Given functions  $f,g:\mathbb{R}\to\mathbb{R}$  we write f=o(g) if

$$\lim_{x\to 0}\frac{f(x)}{g(x)}=0$$

Consider a function  $f : \mathbb{R} \to \mathbb{R}$  and  $x_0 \in \mathbb{R}$ . Assume that f is twice differentiable at  $x_0$ . Then for all  $x \in \mathbb{R}$  we have that

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + o(|x - x_0|^2)$$

Thus, such f can be reasonably approximated around  $x_0$  with a quadratic function

$$f(x) \approx q(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

## Newton's Method Idea

The method computes successive approximations  $x_0, x_1, \ldots, x_k, \ldots$  as the GD.

## Newton's Method Idea

The method computes successive approximations  $x_0, x_1, \ldots, x_k, \ldots$  as the GD.

To compute  $x_{k+1}$ , a quadratic approximation

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

is considered around  $x_k$ .



#### Newton's Method Idea

The method computes successive approximations  $x_0, x_1, \ldots, x_k, \ldots$  as the GD.

To compute  $x_{k+1}$ , a quadratic approximation

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

is considered around  $x_k$ .



Then  $x_{k+1}$  is set to the extreme point of q(x) (i.e.,  $q'(x_{k+1}) = 0$ ).

Now note that for

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

Now note that for

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

we have

$$q'(x) = f'(x_k) + f''(x_k)(x - x_k)$$

Now note that for

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

we have

$$q'(x) = f'(x_k) + f''(x_k)(x - x_k)$$

and thus

$$q'(x) = 0$$
 iff  $x = x_k - \frac{f'(x_k)}{f''(x_k)}$ 

Now note that for

$$q(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(x_k)(x - x_k)^2$$

we have

$$q'(x) = f'(x_k) + f''(x_k)(x - x_k)$$

and thus

$$q'(x) = 0$$
 iff  $x = x_k - \frac{f'(x_k)}{f''(x_k)}$ 

Newton's method then sets

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

**Input:** A function f with derivative f' and second derivative f'', initial point  $x_0$ , tolerance  $\epsilon > 0$ 

**Output:** A point x that approximately minimizes f(x)

1: Set  $k \leftarrow 0$ 

2: while 
$$|x_{k+1} - x_k| > \epsilon$$
 do

- 3: Calculate the derivative:  $y' \leftarrow f'(x_k)$
- 4: Calculate the second derivative :  $y'' \leftarrow f''(x_k)$
- 5: Update the estimate:  $x_{k+1} \leftarrow x_k \frac{y'}{y''}$
- 6: Increment k
- 7: end while
- 8: **return** *x*<sub>*k*</sub>

Note that the method implicitly assumes that  $f''(x_k) \neq 0$  in every iteration.

Consider the following objective function f

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

Assume  $x_0 = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-5}$ , i.e., we stop when  $|x_{k+1} - x_k| \le \epsilon$ .

Consider the following objective function f

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

Assume  $x_0 = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-5}$ , i.e., we stop when  $|x_{k+1} - x_k| \le \epsilon$ .

We compute

$$f'(x) = x - \cos x$$
,  $f''(x) = 1 + \sin x$ .

Consider the following objective function f

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

Assume  $x_0 = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-5}$ , i.e., we stop when  $|x_{k+1} - x_k| \le \epsilon$ .

We compute

$$f'(x) = x - \cos x$$
,  $f''(x) = 1 + \sin x$ .

Hence,

$$x_1 = 0.5 - \frac{0.5 - \cos 0.5}{1 + \sin 0.5}$$
$$= 0.5 - \frac{-0.3775}{1.479}$$
$$= 0.7552$$

. . .

Proceeding similarly, we obtain

$$x_{2} = x_{1} - \frac{f'(x_{1})}{f''(x_{1})} = x_{1} - \frac{0.02710}{1.685} = 0.7391$$
  

$$x_{3} = x_{2} - \frac{f'(x_{2})}{f''(x_{2})} = x_{2} - \frac{9.461 \times 10^{-5}}{1.673} = 0.7390851339$$
  

$$x_{4} = x_{3} - \frac{f'(x_{3})}{f''(x_{3})} = x_{3} - \frac{1.17 \times 10^{-9}}{1.673} = 0.7390851332$$

Proceeding similarly, we obtain

$$x_{2} = x_{1} - \frac{f'(x_{1})}{f''(x_{1})} = x_{1} - \frac{0.02710}{1.685} = 0.7391$$
  

$$x_{3} = x_{2} - \frac{f'(x_{2})}{f''(x_{2})} = x_{2} - \frac{9.461 \times 10^{-5}}{1.673} = 0.7390851339$$
  

$$x_{4} = x_{3} - \frac{f'(x_{3})}{f''(x_{3})} = x_{3} - \frac{1.17 \times 10^{-9}}{1.673} = 0.7390851332$$

Note that

. . .

$$|x_4 - x_3| < \epsilon = 10^{-5}$$
  
 $f'(x_4) = -8.6 \times 10^{-6} \approx 0$   
 $f''(x_4) = 1.673 > 0$ 

So, we conclude that  $x^* \approx x_4$  is a strict minimizer. However, remember that the above does not have to be true!

#### Convergence

Newton's method works well if f''(x) > 0 everywhere.

However, if f''(x) < 0 for some x, Newton's method may fail to converge to a minimizer (converges to a point x where f'(x) = 0):



If the method converges to a minimizer, it does so *quadratically*. What does this mean?

# Types of Convergence Rates

#### Linear Convergence

An algorithm is said to have linear convergence if the error at each step is proportionally reduced by a constant factor:

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = r, \quad 0 < r < 1$$

# Types of Convergence Rates

#### Linear Convergence

An algorithm is said to have linear convergence if the error at each step is proportionally reduced by a constant factor:

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = r, \quad 0 < r < 1$$

#### Superlinear Convergence

Convergence is superlinear if:

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = 0$$

This often requires an algorithm to utilize second-order information.

## Quadratic Convergence of Newton's Method

#### Quadratic Convergence

Quadratic convergence is achieved when the number of accurate digits roughly doubles with each iteration:

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^2} = C, \quad C > 0$$

## Quadratic Convergence of Newton's Method

#### Quadratic Convergence

Quadratic convergence is achieved when the number of accurate digits roughly doubles with each iteration:

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^2} = C, \quad C > 0$$

Newton's method is a classic example of an algorithm with quadratic convergence.

Theorem 2 (Quadratic Convergence of Newton's Method) Let  $f : \mathbb{R} \to \mathbb{R}$  satisfy  $f \in C^2$  and suppose  $x^*$  is a minimizer of fsuch that  $f''(x^*) > 0$ . Assume Lipschitz continuity of f''. If the initial guess  $x_0$  is sufficiently close to  $x^*$ , then the sequence  $\{x_k\}$ computed by the Newton's method converges quadratically to  $x^*$ .
## Newton's Method of Tangents

Newton's method is also a technique for finding roots of functions. In our case, this means finding a root of f'.

### Newton's Method of Tangents

Newton's method is also a technique for finding roots of functions. In our case, this means finding a root of f'.

Denote g = f'. Then Newton's approximation goes like this:



## Secant Method

What if f'' is unavailable, but we want to use something like Newton's method (with its superlinear convergence)?

### Secant Method

What if f'' is unavailable, but we want to use something like Newton's method (with its superlinear convergence)?

Assume  $f \in C^1$  and try to approximate f'' around  $x_{k-1}$  with

$$f''(x) \approx rac{f'(x) - f'(x_{k-1})}{x - x_{k-1}}$$

Substituting x with  $x_k$ , we obtain

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

### Secant Method

What if f'' is unavailable, but we want to use something like Newton's method (with its superlinear convergence)?

Assume  $f \in C^1$  and try to approximate f'' around  $x_{k-1}$  with

$$f''(x) \approx rac{f'(x) - f'(x_{k-1})}{x - x_{k-1}}$$

Substituting x with  $x_k$ , we obtain

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

Then, we may try to use Newton's step with this approximation:

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

Is the rate of convergence superlinear?

Consider the following objective function f

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

Assume  $x_0 = 0.5$  and  $x_1 = 1.0$ .

Now, we need to initialize the first two values.

Consider the following objective function f

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

Assume  $x_0 = 0.5$  and  $x_1 = 1.0$ .

Now, we need to initialize the first two values.

We have  $f'(x) = x - \cos x$ Hence,

$$egin{aligned} x_2 &= 1.0 - rac{1.0 - 0.5}{(1.0 - \cos 1.0) - (0.5 - \cos 0.5)} (0.5 - \cos 0.5) \ &= 0.7254 \end{aligned}$$

Continuing, we obtain:

 $\begin{aligned} x_0 &= 0.5 \\ x_1 &= 1.0 \\ x_2 &= 0.72548 \\ x_3 &= 0.73839 \\ x_4 &= 0.739087 \\ x_5 &= 0.739085132 \\ x_6 &= 0.739085133 \end{aligned}$ 

Start the secant method with the approximation given by Newton's method:

. . .

$$\begin{array}{l} x_0 = 0.5 \\ x_1 = 0.7552 \\ x_2 = 0.7381 \\ x_3 = 0.739081 \\ x_5 = 0.7390851339 \\ x_6 = 0.7390851332 \end{array}$$

Compare with Newton's method:

. . .

$$\begin{array}{l} x_0 = 0.5 \\ x_1 = 0.7552 \\ x_2 = 0.7391 \\ x_3 = 0.7390851339 \\ x_4 = 0.73908513321516067229 \\ x_5 = 0.73908513321516067229 \end{array}$$

### Superlinear Convergence of Secant Method

Theorem 3 (Superlinear Convergence of Secant Method) Assume  $f : \mathbb{R} \to \mathbb{R}$  twice continuously differentiable and  $x^*$ a minimizer of f. Assume f'' Lipschitz continuous and  $f''(x^*) > 0$ . The sequence  $\{x_k\}$  generated by the Secant method converges to  $x^*$  superlinearly if  $x_0$  and  $x_1$  are sufficiently close to  $x^*$ .

The rate of convergence p of the Secant method is given by the positive root of the equation  $p^2 - p - 1 = 0$ , which is  $p = \frac{1+\sqrt{5}}{2} \approx 1.618$  (the golden ratio). Formally,

$$\lim_{k \to \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^{\frac{1+\sqrt{5}}{2}}} = C, \quad C > 0$$

## Secant Method for Root Finding

As for Newton's method of tangents, the secant method can be seen as a method for finding a root of f'.

### Secant Method for Root Finding

As for Newton's method of tangents, the secant method can be seen as a method for finding a root of f'.

Denote g = f'. Then the secant method approximation is



### General Form

Note that all methods have similar update formula:

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

Different choice of  $a_k$  produce different algorithm:

# Summary

### Newton's method

- Converges quickly to an extremum under rather strict conditions (see Theorem 2)
- The choice of the initial point is critical; the method may diverge to a stationary point, which is not a minimizer. The method may also cycle.
- If the second derivative is very small, close to the minimizer, the method can be very slow (the quadratic convergence is guaranteed only if the second derivative is non-zero at the minimizer and the constants depend on the second derivative).

# Summary

### Newton's method

- Converges quickly to an extremum under rather strict conditions (see Theorem 2)
- The choice of the initial point is critical; the method may diverge to a stationary point, which is not a minimizer. The method may also cycle.
- If the second derivative is very small, close to the minimizer, the method can be very slow (the quadratic convergence is guaranteed only if the second derivative is non-zero at the minimizer and the constants depend on the second derivative).

#### Secant method

- The second derivative is not needed.
- Superlinear (but not quadratic) convergence for an initial point close to a minimum (under rather strict conditions Theorem 3)

## Constrained Single Variable Optimization Problem

An objective function  $f : \mathbb{R} \to \mathbb{R}$ 

A variable x

A constraint

 $a_0 \leq x \leq b_0$ 

Consider the following cases:

- ▶ *f* continuously differentiable on [*a*<sub>0</sub>, *b*<sub>0</sub>]
- f twice continuously differentiable on  $[a_0, b_0]$

**Homework:** Modify the gradient descent and Newton's method to work on the bounded interval (the above definitions guarantee continuous differentiability at  $a_0$  and  $b_0$ ).

# **Unconstrained Optimization Overview**

### Notation

In what follows, we will work with vectors in  $\mathbb{R}^n$ .

The vectors will be (usually) denoted by  $x \in \mathbb{R}^n$ .

We often consider sequences of vectors,  $x_0, x_1, \ldots, x_k, \ldots$ 

The index k will usually indicate that  $x_k$  is the k-the vector in a sequence.

When we talk (relatively rarely) about components of vectors, we use *i* as an index, i.e.,  $x_i$  will be the *i*-th component of  $x \in \mathbb{R}^n$ .

We denote by ||x|| the Euclidean norm of x.

We denote by  $||x||_{\infty}$  the  $\mathcal{L}^{\infty}$  norm giving the maximum of absolute values of components of x.

We ocasionally use the matrix norm ||A||, consistent with the Euclidean norm, defined by

$$||A|| = \sup_{||x||=1} ||Ax|| = \sqrt{\lambda_1}$$

Here  $\lambda_1$  is the largest eigenvalue of  $A^{\top}A$ .

# How to Recognize (Local) Minimum

How do we verify that  $x^* \in \mathbb{R}^n$  is a minimizer of f?



# How to Recognize (Local) Minimum

How do we verify that  $x^* \in \mathbb{R}^n$  is a minimizer of f?



Technically, we should examine *all* points in the immediate vicinity if one has a smaller value (impractical).

Assuming the smoothness of f, we may benefit from the "stable" behavior of f around  $x^*$ .

### Derivatives and Gradients

The gradient of  $f : \mathbb{R}^n \to \mathbb{R}$ , denoted by  $\nabla f(x)$ , is a column vector of first-order partial derivatives of the function concerning each variable:

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

Where each partial derivative is defined as the following limit:

$$\frac{\partial f}{\partial \mathbf{x}_{i}} = \lim_{\varepsilon \to 0} \frac{f(x_{1}, \dots, x_{i} + \varepsilon, \dots, x_{n}) - f(x_{1}, \dots, x_{i}, \dots, x_{n})}{\varepsilon}$$

## Gradient



The gradient is a vector pointing in the direction of the most significant function increase from the current point.

### Gradient

Consider the following function of two variables:

$$f(x_1, x_2) = x_1^3 + 2x_1x_2^2 - x_2^3 - 20x_1.$$

$$\nabla f(x_1, x_2) = \begin{bmatrix} 3x_1^2 + 2x_2^2 - 20 \\ 4x_1x_2 - 3x_2^2 \end{bmatrix}$$



### Directional Derivatives vs Gradient

The rate of change in a direction p is quantified by a directional derivative, defined as

$$abla_{p}f(x) = \lim_{\varepsilon \to 0} \frac{f(x + \varepsilon p) - f(x)}{\varepsilon}$$

We can find this derivative by projecting the gradient onto the desired direction p using the dot product  $\nabla_p f(x) = (\nabla f(x))^\top p$ 



(Here, we assume continuous partial derivatives.)

### Geometry of Gradient

Consider the geometric interpretation of the dot product:

 $abla_p f(x) = (
abla f(x))^\top p = ||\nabla f|| \, ||p|| \cos \theta$ 

Here  $\theta$  is the angle between  $\nabla f$  and p.

### Geometry of Gradient

Consider the geometric interpretation of the dot product:

$$\nabla_p f(x) = (\nabla f(x))^\top p = ||\nabla f|| ||p|| \cos \theta$$

Here  $\theta$  is the angle between  $\nabla f$  and p.

The directional derivative is maximized by  $\theta = 0$ , i.e., when  $\nabla f$  and p point in the same direction.



Hessian

Taking derivative twice, possibly w.r.t. different variables, gives the Hessian of  $\boldsymbol{f}$ 

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

Note that the Hessian is a function which takes  $x \in \mathbb{R}^n$  and gives a  $n \times n$ -matrix of second derivatives of f.

٠

Hessian

Taking derivative twice, possibly w.r.t. different variables, gives the Hessian of  $\boldsymbol{f}$ 

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

Note that the Hessian is a function which takes  $x \in \mathbb{R}^n$  and gives a  $n \times n$ -matrix of second derivatives of f.

We have

$$H_{ij}=\frac{\partial^2 f}{\partial x_i\partial x_j}.$$

If f has continuous second partial derivatives, then H is symmetric, i.e.,  $H_{ij} = H_{ji}$ .

Let x be fixed and let g(t) = f(x + tp).

What exactly are g'(0) and g''(0)?

Let x be fixed and let g(t) = f(x + tp).

What exactly are g'(0) and g''(0)?

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top} p = \nabla_p f(x+tp)$$

and thus  $g'(0) = \nabla_p f(x)$ .

Let x be fixed and let g(t) = f(x + tp).

What exactly are g'(0) and g''(0)?

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top} p = \nabla_p f(x+tp)$$

and thus  $g'(0) = \nabla_p f(x)$ .

Now note that for all z we have

$$\nabla(\nabla_p f(z)) = \nabla([\nabla f(z)]^\top p) = (\nabla^2 f(z))p = H(z)p$$

Let x be fixed and let g(t) = f(x + tp).

What exactly are g'(0) and g''(0)?

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top} p = \nabla_p f(x+tp)$$

and thus  $g'(0) = \nabla_p f(x)$ .

Now note that for all z we have

$$\nabla(\nabla_p f(z)) = \nabla([\nabla f(z)]^\top p) = (\nabla^2 f(z))p = H(z)p$$

Thus, for our fixed x we have

$$g''(t) = (g'(t))'$$
  
=  $(\nabla_p f(x+tp))' = [\nabla(\nabla_p f(x+tp))]^\top p$   
=  $(H(x+tp)p)^\top p = p^\top H(x+tp)p$ 

Let x be fixed and let g(t) = f(x + tp).

What exactly are g'(0) and g''(0)?

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top} p = \nabla_p f(x+tp)$$

and thus  $g'(0) = \nabla_p f(x)$ .

Now note that for all z we have

$$\nabla(\nabla_p f(z)) = \nabla([\nabla f(z)]^\top p) = (\nabla^2 f(z))p = H(z)p$$

Thus, for our fixed x we have

$$g''(t) = (g'(t))'$$
  
=  $(\nabla_p f(x+tp))' = [\nabla(\nabla_p f(x+tp))]^\top p$   
=  $(H(x+tp)p)^\top p = p^\top H(x+tp)p$ 

Thus  $g''(0) = p^{\top} H(x)p$ .

## Principal Curvature Directions

Fix x and consider H = H(x). Consider unit eigenvectors  $\hat{v}_k$  of H:

 $H\hat{v}_k = \kappa_k\hat{v}_k$ 

For symmetric H, the unit eigenvectors form an orthonormal basis,

### Principal Curvature Directions

Fix x and consider H = H(x). Consider unit eigenvectors  $\hat{v}_k$  of H:

$$H\hat{v}_k = \kappa_k \hat{v}_k$$

For symmetric H, the unit eigenvectors form an orthonormal basis, and there is a rotation matrix R such that

$$H = RDR^{-1} = RDR^{\top}$$

Here *D* is diagonal with  $\kappa_1, \ldots, \kappa_n$  on the diagonal.

If  $\kappa_1 \geq \cdots \geq \kappa_n$ , the direction of  $\hat{v}_1$  is the maximum curvature direction of f at x.



Consider  $f(x) = x^{\top} H x$  where

$$H = \begin{pmatrix} 4/3 & 0 \\ 0 & 1 \end{pmatrix}$$

The eigenvalues are

 $\kappa_1 = 4/3$   $\kappa_2 = 1$ 

Their corresponding eigenvectors are  $(1,0)^{\top}$  and  $(0,1)^{\top}$ .


Consider  $f(x) = x^{\top} H x$  where

$$H = \begin{pmatrix} 4/3 & 0 \\ 0 & 1 \end{pmatrix}$$

The eigenvalues are

 $\kappa_1 = 4/3$   $\kappa_2 = 1$ 

Their corresponding eigenvectors are  $(1,0)^{\top}$  and  $(0,1)^{\top}$ .

Note that

$$f(x) = \kappa_1 x_1^2 + \kappa_2 x_2^2$$

Considering a direction vector p and  $x = (0,0)^{\top}$  we get

$$g(t) = f(x + tp) = f(tp) = t^2 (\kappa_1 p_1^2 + \kappa_2 p_2^2)$$

which is a parabola with  $g'' = 2 \left( \kappa_1 p_1^2 + \kappa_2 p_2^2 \right)$ .



Consider  $f(x) = x^{\top} H x$  where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

Consider  $f(x) = x^{\top} H x$  where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

The eigenvalues are

$$\kappa_1 = \frac{1}{6}(7 + \sqrt{5})$$
  $\kappa_2 = \frac{1}{6}(7 - \sqrt{5})$ 

Their corresponding eigenvectors are

$$\hat{\mathbf{v}}_{1} = \left(rac{1}{2}(1+\sqrt{5}),1
ight) \quad \hat{\mathbf{v}}_{2} = \left(rac{1}{2}(1-\sqrt{5}),1
ight)$$



Consider  $f(x) = x^{\top} Hx$  where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

The eigenvalues are

$$\kappa_1 = \frac{1}{6}(7 + \sqrt{5})$$
  $\kappa_2 = \frac{1}{6}(7 - \sqrt{5})$ 

Their corresponding eigenvectors are

$$\hat{\mathbf{v}_1} = \left(rac{1}{2}(1+\sqrt{5}),1
ight) \quad \hat{\mathbf{v}_2} = \left(rac{1}{2}(1-\sqrt{5}),1
ight)$$

Note that

$$egin{aligned} \mathcal{H} = egin{pmatrix} \hat{v}_1 & \hat{v}_2 \end{pmatrix} egin{pmatrix} \kappa_1 & 0 \ 0 & \kappa_2 \end{pmatrix} egin{pmatrix} \hat{v}_1 & \hat{v}_2 \end{pmatrix}^ op \end{aligned}$$

Here  $(\hat{v}_1 \ \hat{v}_2)$  is a 2 × 2 matrix whose columns are  $\hat{v}_1, \hat{v}_2$ .



## Hessian Visualization Example

Consider

$$f(x_1, x_2) = x_1^3 + 2x_1x_2^2 - x_2^3 - 20x_1.$$

And it's Hessian.

$$H(x_1, x_2) = \begin{bmatrix} 6x_1 & 4x_2 \\ 4x_2 & 4x_1 - 6x_2 \end{bmatrix}.$$



#### Theorem 4 (Taylor)

Suppose that  $f : \mathbb{R}^n \to \mathbb{R}$  is twice continuously differentiable and that  $p \in \mathbb{R}^n$ . Then, we have

$$f(x+p) = f(x) + \nabla f(x)^{T} p + \frac{1}{2} p^{T} H(x) p + o(||p||^{2}).$$

Here  $H = \nabla^2 f$  is the Hessian of f.

## First-Order Necessary Conditions

#### Theorem 5

If  $x^*$  is a local minimizer and f is continuously differentiable in an open neighborhood of  $x^*$ , then  $\nabla f(x^*) = 0$ .





Note that  $\nabla f(x^*) = 0$  does not tell us whether  $x^*$  is a minimizer, maximizer, or a saddle point.

Note that  $\nabla f(x^*) = 0$  does not tell us whether  $x^*$  is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from  $x^*$  might tell us what  $x^*$  is, right?

Note that  $\nabla f(x^*) = 0$  does not tell us whether  $x^*$  is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from  $x^*$  might tell us what  $x^*$  is, right?

Note that  $\nabla f(x^*) = 0$  does not tell us whether  $x^*$  is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from  $x^*$  might tell us what  $x^*$  is, right?

All comes down to the *definiteness* of  $H := H(x^*)$ .

- H is positive definite if p<sup>T</sup>Hp > 0 for all p iff all eigenvalues of H are positive
- H is positive semi-definite if p<sup>⊤</sup>Hp ≥ 0 for all p iff all eigenvalues of H are nonnegative
- H is negative semi-definite if p<sup>T</sup>Hp ≤ 0 for all p iff all eigenvalues of H are nonpositive
- ► H is negative definite if p<sup>T</sup> Hp < 0 for all p iff all eigenvalues of H are negative
- H is indefinite if it is not definite in the above sense iff H has at least one positive and one negative eigenvalue.

#### Definiteness



## Second-Order Necessary Condition

Theorem 6 (Second-Order Necessary Conditions)

If  $x^*$  is a local minimizer of f and  $\nabla^2 f$  is continuous in a neighborhood of  $x^*$ , then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive semidefinite.

Theorem 7 (Second-Order Sufficient Conditions)

Suppose that  $\nabla^2 f$  is continuous in a neighborhood of  $x^*$  and that  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite. Then  $x^*$  is a strict local minimizer of f.



Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

From the second equation, we have that  $x_2 = x_1$ . Substituting this into the first equation yields

$$x_1\left(2x_1^2+6x_1+1\right)=0.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

From the second equation, we have that  $x_2 = x_1$ . Substituting this into the first equation yields

$$x_1\left(2x_1^2+6x_1+1\right)=0.$$

The solution of this equation yields three points:

$$x_{\mathcal{A}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad x_{\mathcal{B}} = \begin{bmatrix} -\frac{3}{2} - \frac{\sqrt{7}}{2} \\ -\frac{3}{2} - \frac{\sqrt{7}}{2} \end{bmatrix}, \quad x_{\mathcal{C}} = \begin{bmatrix} \frac{\sqrt{7}}{2} - \frac{3}{2} \\ \frac{\sqrt{7}}{2} - \frac{3}{2} \end{bmatrix}.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify  $x_A, x_B, x_C$ , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify  $x_A, x_B, x_C$ , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

The Hessian, at the first point, is

$$H(x_{\mathcal{A}}) = \left[\begin{array}{cc} 3 & -2 \\ -2 & 2 \end{array}\right],$$

whose eigenvalues are  $\kappa_1 \approx 0.438$  and  $\kappa_2 \approx 4.561$ . Because both eigenvalues are positive, this point is a local minimum.

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify  $x_A, x_B, x_C$ , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

For the second point,

$$H(x_B) = \left[ egin{array}{cc} 3(3+\sqrt{7}) & -2 \ -2 & 2 \end{array} 
ight].$$

The eigenvalues are  $\kappa_1 \approx 1.737$  and  $\kappa_2 \approx 17.200$ , so this point is another local minimum.

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify  $x_A, x_B, x_C$ , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

For the third point,

$$H(x_C) = \begin{bmatrix} 9 - 3\sqrt{7} & -2\\ -2 & 2 \end{bmatrix}$$

The eigenvalues for this Hessian are  $\kappa_1 \approx -0.523$  and  $\kappa_2 \approx 3.586$ , so this point is a saddle point.



#### Proofs of Some Theorems Optional

## Taylor's Theorem

To prove the theorems characterizing minima/maxima, we need the following form of Taylor's theorem:

#### Theorem 8 (Taylor)

Suppose that  $f : \mathbb{R}^n \to \mathbb{R}$  is continuously differentiable and that  $p \in \mathbb{R}^n$ . Then we have that.

$$f(x+p) = f(x) + \nabla f(x+tp)^T p,$$

for some  $t \in (0, 1)$ . Moreover, if f is twice continuously differentiable, we have that

$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x+tp)p,$$

for some  $t \in (0, 1)$ .

## Proof of Theorem 5 (Optional)

We prove that if  $x^*$  is a local minimizer and f is continuously differentiable in an open neighborhood of  $x^*$ , then  $\nabla f(x^*) = 0$ .

Suppose for contradiction that  $\nabla f(x^*) \neq 0$ . Define the vector  $p = -\nabla f(x^*)$  and note that  $p^T \nabla f(x^*) = - \|\nabla f(x^*)\|^2 < 0$ . Because  $\nabla f$  is continuous near  $x^*$ , there is a scalar T > 0 such that

$$p^T \nabla f(x^* + tp) < 0, \quad \text{ for all } t \in [0, T]$$

For any  $ar{t} \in (0, T]$ , we have by Taylor's theorem that

$$f\left(x^{*}+ar{t}p
ight)=f\left(x^{*}
ight)+ar{t}p^{T}
abla f\left(x^{*}+tp
ight), \hspace{0.5cm} ext{ for some }t\in(0,ar{t}).$$

Therefore,  $f(x^* + \bar{t}p) < f(x^*)$  for all  $\bar{t} \in (0, T]$ . We have found a direction leading away from  $x^*$  along which f decreases, so  $x^*$  is not a local minimizer, and we have a contradiction.

## Proof of Theorem 6 (Optional)

We prove that if  $x^*$  is a local minimizer of f and  $\nabla^2 f$  is continuous in an open neighborhood of  $x^*$ , then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive semidefinite.

We know that  $\nabla f(x^*) = 0$ . For contradiction, assume that  $\nabla^2 f(x^*)$  is not positive semidefinite.

Then we can choose a vector p such that  $p^T \nabla^2 f(x^*) p < 0$ .

As  $\nabla^2 f$  is continuous near  $x^*$ ,  $p^T \nabla^2 f(x^* + tp) p < 0$  for all  $t \in [0, T]$  where T > 0.

By Taylor we have for all  $\overline{t} \in (0, T]$  and some  $t \in (0, \overline{t})$ 

$$f(x^* + \bar{t}p) = f(x^*) + \bar{t}p^T \nabla f(x^*) + \frac{1}{2} \bar{t}^2 p^T \nabla^2 f(x^* + tp) p < f(x^*).$$

Thus,  $x^*$  is not a local minimizer.

## Proof of Theorem 7 (Optional)

We prove the following: Suppose that  $\nabla^2 f$  is continuous in an open neighborhood of  $x^*$  and that  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite. Then  $x^*$  is a strict local minimizer of f.

Because the Hessian is continuous and positive definite at  $x^*$ , we can choose a radius r > 0 so that  $\nabla^2 f(x)$  remains positive definite for all x in the open ball  $\mathcal{D} = \{z \mid ||z - x^*|| < r\}$ . Taking any nonzero vector p with ||p|| < r, we have  $x^* + p \in \mathcal{D}$  and so

$$f(x^* + p) = f(x^*) + p^T \nabla f(x^*) + \frac{1}{2} p^T \nabla^2 f(z) p$$
  
=  $f(x^*) + \frac{1}{2} p^T \nabla^2 f(z) p$ ,

where  $z = x^* + tp$  for some  $t \in (0, 1)$ . Since  $z \in D$ , we have  $p^T \nabla^2 f(z) p > 0$ , and therefore  $f(x^* + p) > f(x^*)$ , giving the result.

# Unconstrained Optimization Algorithms

## Search Algorithms

We consider algorithms that

- Start with an initial guess x<sub>0</sub>
- ▶ Generate a sequence of points *x*<sub>0</sub>, *x*<sub>1</sub>,...
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute  $x_{k+1}$  the algorithms use the information about f at the previous iterates  $x_0, x_1, \ldots, x_k$ .

## Search Algorithms

We consider algorithms that

- Start with an initial guess x<sub>0</sub>
- Generate a sequence of points  $x_0, x_1, \ldots$
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute  $x_{k+1}$  the algorithms use the information about f at the previous iterates  $x_0, x_1, \ldots, x_k$ .

The *monotone* algorithms satisfy  $f(x_{k+1}) < f(x_k)$ .

## Search Algorithms

We consider algorithms that

- Start with an initial guess x<sub>0</sub>
- Generate a sequence of points  $x_0, x_1, \ldots$
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute  $x_{k+1}$  the algorithms use the information about f at the previous iterates  $x_0, x_1, \ldots, x_k$ .

The *monotone* algorithms satisfy  $f(x_{k+1}) < f(x_k)$ .

There are two overall strategies:

- Line search
- Trust region

## Line Search Overview

To compute  $x_{k+1}$ , a line search algorithm chooses

- $\blacktriangleright$  direction  $p_k$
- step size  $\alpha_k$

and computes

$$x_{k+1} = x_k + \alpha_k p_k$$

## Line Search Overview

To compute  $x_{k+1}$ , a line search algorithm chooses

- direction p<sub>k</sub>
- step size  $\alpha_k$

and computes

 $x_{k+1} = x_k + \alpha_k p_k$ 

The vector  $p_k$  should be a *descent* direction, i.e., a direction in which f decreases locally.

## Line Search Overview

To compute  $x_{k+1}$ , a line search algorithm chooses

direction p<sub>k</sub>

• step size  $\alpha_k$ 

and computes

 $x_{k+1} = x_k + \alpha_k p_k$ 

The vector  $p_k$  should be a *descent* direction, i.e., a direction in which f decreases locally.

 $\alpha_k$  is selected to approximately solve

$$\min_{\alpha>0}f(x_k+\alpha p_k)$$

However, typically, an exact solution is expensive and unnecessary. Instead, line search algorithms inspect a limited number of trial step lengths and find one that decreases f appropriately (see later).



A descent direction does not have to be followed to the minimum.


# Trust Region

To compute  $x_{k+1}$ , a trust region algorithm chooses

• model function  $m_k$  whose behavior near  $x_k$  is similar to f

▶ a trust region  $R \subseteq \mathbb{R}^n$  around  $x_k$ . Usually R is the ball defined by  $||x - x_k|| \leq \Delta$  where  $\Delta > 0$  is trust region radius.

and finds  $x_{k+1}$  solving

 $\min_{x\in R}m_k(x)$ 

# Trust Region

To compute  $x_{k+1}$ , a trust region algorithm chooses

• model function  $m_k$  whose behavior near  $x_k$  is similar to f

▶ a trust region  $R \subseteq \mathbb{R}^n$  around  $x_k$ . Usually R is the ball defined by  $||x - x_k|| \leq \Delta$  where  $\Delta > 0$  is trust region radius.

and finds  $x_{k+1}$  solving

 $\min_{x\in R}m_k(x)$ 

If the solution does not sufficiently decrease f, we shrink the trust region and re-solve.

### **Trust Region**

To compute  $x_{k+1}$ , a trust region algorithm chooses

• model function  $m_k$  whose behavior near  $x_k$  is similar to f

▶ a *trust region*  $R \subseteq \mathbb{R}^n$  around  $x_k$ . Usually R is the ball defined by  $||x - x_k|| \leq \Delta$  where  $\Delta > 0$  is *trust region radius*.

and finds  $x_{k+1}$  solving

 $\min_{x\in R}m_k(x)$ 

If the solution does not sufficiently decrease f, we shrink the trust region and re-solve.

The model  $m_k$  is usually derived from the Taylor's theorem.

$$m_k(x_k+p) = f_k + p^T \nabla f_k + \frac{1}{2} p^T B_k p$$

Where  $B_k$  approximates the Hessian of f at  $x_k$ .



# Line Search Methods

# Line Search

For setting the step size, we consider

- Armijo condition and backtracking algorithm
- strong Wolfe conditions and bracketing & zooming

# Line Search

For setting the step size, we consider

- Armijo condition and backtracking algorithm
- strong Wolfe conditions and bracketing & zooming

For setting the direction, we consider

- Gradient descent
- Newton's method
- quasi-Newton methods (BFGS)
- (Conjugate gradients)

We start with the step size.

# Step Size

#### Assume

 $x_{k+1} = x_k + \alpha_k p_k$ 

Where  $p_k$  is a descent direction

 $p_k^\top \nabla f_k < 0$ 



# Step Size

#### Assume

 $x_{k+1} = x_k + \alpha_k p_k$ 

Where  $p_k$  is a descent direction  $p_k^\top 
abla f_k < 0$ 

Define

$$\phi(\alpha) = f(x_k + \alpha p_k)$$



# Step Size

#### Assume

 $x_{k+1} = x_k + \alpha_k p_k$ 

Where  $p_k$  is a descent direction  $p_k^\top \nabla f_k < 0$ 



Define

$$\phi(\alpha) = f(x_k + \alpha p_k)$$

We know that

 $\phi'(\alpha) = \nabla f(x_k + \alpha p_k)^\top p_k$  which means  $\phi'(0) = \nabla f_k^\top p_k$ 

Note that  $\phi'(0)$  must be negative as  $p_k$  is a descent direction.

#### Armijo Condition

The sufficient decrease condition (aka Armijo condition)

$$\phi(\alpha) \le \phi(\mathbf{0}) + \alpha \left(\mu_1 \phi'(\mathbf{0})\right)$$

where  $\mu_1$  is a constant such that  $0 < \mu_1 \leq 1$ 



In practice,  $\mu_1$  is several orders smaller than 1, typically  $\mu_1 = 10^{-4}$ .

# Backtracking Line Search Algorithm

Algorithm 1 Backtracking Line Search

**Input:**  $\alpha_{init} > 0$ ,  $0 < \mu_1 < 1$ ,  $0 < \rho < 1$ 

**Output:**  $\alpha^*$  satisfying sufficient decrease condition

- 1:  $\alpha \leftarrow \alpha_{\mathsf{init}}$
- 2: while  $\phi(\alpha) > \phi(0) + \alpha \mu_1 \phi'(0)$  do
- 3:  $\alpha \leftarrow \rho \alpha$
- 4: end while

# Backtracking Line Search Algorithm

Algorithm 2 Backtracking Line Search Input:  $\alpha_{init} > 0, 0 < \mu_1 < 1, 0 < \rho < 1$ Output:  $\alpha^*$  satisfying sufficient decrease condition 1:  $\alpha \leftarrow \alpha_{init}$ 2: while  $\phi(\alpha) > \phi(0) + \alpha \mu_1 \phi'(0)$  do 3:  $\alpha \leftarrow \rho \alpha$ 

4: end while

The parameter  $\rho$  is typically set to 0.5. It can also be a variable set by a more sophisticated method (interpolation).

The  $\alpha_{init}$  depends on the method for setting the descent direction  $p_k$ . For Newton and quasi-Newton, it is 1.0, but for other methods, it might be different.

There are two scenarios where the method does not perform well:

There are two scenarios where the method does not perform well:

The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.

There are two scenarios where the method does not perform well:

- The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.
- The guess for the initial step immediately satisfies sufficient decrease. However, the function's slope is still highly negative, and we could have decreased the function value by much more if we had taken a more significant step. In this case, our guess for the initial step is far too small.

There are two scenarios where the method does not perform well:

- The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.
- The guess for the initial step immediately satisfies sufficient decrease. However, the function's slope is still highly negative, and we could have decreased the function value by much more if we had taken a more significant step. In this case, our guess for the initial step is far too small.

Even if our original step size is not too far from an acceptable one, the basic backtracking algorithm ignores any information we have about the function values and gradients. It blindly takes a reduced step based on a preselected ratio  $\rho$ .

#### Backtracking Example

$$f(x_1, x_2) = 0.1x_1^6 - 1.5x_1^4 + 5x_1^2 + 0.1x_2^4 + 3x_2^2 - 9x_2 + 0.5x_1x_2$$

2.5 2 1.5 1 0.5 -3 -2 -1 0 1 2 3  $x_1$  $x_1$ 

 $x_2$ 

 $\mu_1 = 10^{-4}$  and  $\rho = 0.7$ .



We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition

 $ig|\phi'(lpha)ig|\leq \mu_2ig|\phi'(0)ig|$  where  $\mu_1<\mu_2<1$  is a

constant.



We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition



Typical values of  $\mu_2$  range from 0.1 to 0.9, depending on the direction setting method.

We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition



Typical values of  $\mu_2$  range from 0.1 to 0.9, depending on the direction setting method.

Note that moving  $\mu_2$  close to 0, the condition enforces  $\phi'(\alpha) \approx 0$ , which would yield an (almost) exact line search.

### Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions* 

#### Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions* 

Sufficient decrease condition

 $\phi(\alpha) \leq \phi(0) + \mu_1 \alpha \phi'(0)$ 

### Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions* 

Sufficient decrease condition

 $\phi(\alpha) \le \phi(0) + \mu_1 \alpha \phi'(0)$ 

Sufficient curvature condition  $\left|\phi'(\alpha)\right| \leq \mu_2 \left|\phi'(0)\right|$  $\phi'(0$  $\phi(0)$  $\mu_1 \phi'(0)$ Sufficient decrease line  $\phi(\alpha)$  $\mu_2 \phi'(0)$  $\alpha = 0$ α Acceptable range Acceptable range

# Satisfiability of Strong Wolfe Conditions

#### Theorem 9

Suppose  $f : \mathbb{R}^n \to \mathbb{R}$  is continuously differentiable. Let  $p_k$  be a descent direction at  $x_k$ , and assume that f is bounded below along the ray  $\{x_k + \alpha p_k \mid \alpha > 0\}$ . Then, if  $0 < \mu_1 < \mu_2 < 1$ , step length intervals exist that satisfy the strong Wolfe conditions.



#### Convergence of Line Search

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

#### Convergence of Line Search

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T \rho_k}{\|\nabla f_k\| \, \|\rho_k\|}$$

Recall that f is L-smooth for some L > 0 if

$$\|
abla f(x) - 
abla f( ilde{x})\| \le L \|x - ilde{x}\|, \quad \text{ for all } x, ilde{x} \in \mathbb{R}^n$$

#### Convergence of Line Search

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T \rho_k}{\|\nabla f_k\| \, \|\rho_k\|}$$

Recall that f is L-smooth for some L > 0 if

$$\|
abla f(x) - 
abla f( ilde{x})\| \le L \|x - ilde{x}\|, \quad \text{ for all } x, ilde{x} \in \mathbb{R}^n$$

#### Theorem 10 (Zoutendijk)

Consider  $x_{k+1} = x_k + \alpha_k p_k$ , where  $p_k$  is a descent direction and  $\alpha_k$  satisfies the strong Wolfe conditions. Suppose that f is bounded below, continuously differentiable, and L-smooth. Then

$$\sum_{k\geq 0}\cos^2\theta_k \|\nabla f_k\|^2 < \infty.$$

### Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

How can we find a step size that satisfies *strong* Wolfe conditions? Use a bracketing and zoom algorithm, which proceeds in the following two phases:

### Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

Use a bracketing and zoom algorithm, which proceeds in the following two phases:

1. The bracketing phase finds an interval within which we are certain to find a point that satisfies the strong Wolfe conditions.

# Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

Use a bracketing and zoom algorithm, which proceeds in the following two phases:

- 1. The bracketing phase finds an interval within which we are certain to find a point that satisfies the strong Wolfe conditions.
- 2. The zooming phase finds a point that satisfies the strong Wolfe conditions within the interval provided by the bracketing phase.

#### Algorithm 3 Bracketing

**Input:**  $\alpha_1 > 0$  and  $\alpha_{max}$ 1: Set  $\alpha_0 \leftarrow 0$ 2:  $i \leftarrow 1$ 3: repeat Evaluate  $\phi(\alpha_i)$ 4: if  $\phi(\alpha_i) > \phi(0) + \alpha_i \mu_1 \phi'(0)$  or  $[\phi(\alpha_i) \ge \phi(\alpha_{i-1})$  and i > 15: then  $\alpha^* \leftarrow \mathbf{zoom}(\alpha_{i-1}, \alpha_i)$  and stop 6: end if 7: Evaluate  $\phi'(\alpha_i)$ 8: if  $|\phi'(\alpha_i)| < \mu_2 |\phi'(0)|$  then 9: set  $\alpha^* \leftarrow \alpha_i$  and stop 10: else if  $\phi'(\alpha_i) > 0$  then 11: set  $\alpha^* \leftarrow \mathbf{zoom}(\alpha_i, \alpha_{i-1})$  and stop 12: end if 13: Choose  $\alpha_{i+1} \in (\alpha_i, \alpha_{\max})$ 14:  $i \leftarrow i + 1$ 15: 16: **until** a condition is met

# Explanation of Bracketing

Note that the sequence of trial steps  $\alpha_i$  is monotonically increasing.

Note that the sequence of trial steps  $\alpha_i$  is monotonically increasing.

Note that **zoom** is called when one of the following conditions is satisfied:

- $\alpha_i$  violates the sufficient decrease condition (lines 5 and 6)
- $\phi(\alpha_i) \ge \phi(\alpha_{i-1})$  (also lines 5 and 6)
- $\phi'(\alpha_i) \ge 0$  (lines 11 and 12)

The last step increases the  $\alpha_i$ . May use, e.g., a constant multiple.
The following algorithm keeps two step lengths:  $\alpha_{\rm lo}$  and  $\alpha_{\rm hi}$ 

The following algorithm keeps two step lengths:  $\alpha_{\rm lo}$  and  $\alpha_{\rm hi}$ 

The following invariants are being preserved:

The interval bounded by α<sub>lo</sub> and α<sub>hi</sub> always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α<sub>lo</sub> ≤ α<sub>hi</sub>

The following algorithm keeps two step lengths:  $\alpha_{\rm lo}$  and  $\alpha_{\rm hi}$ 

The following invariants are being preserved:

- The interval bounded by α<sub>lo</sub> and α<sub>hi</sub> always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α<sub>lo</sub> < α<sub>hi</sub>
- α<sub>lo</sub> is, among all step lengths generated so far and satisfying the sufficient decrease condition, the one giving the smallest value of φ,

The following algorithm keeps two step lengths:  $\alpha_{\rm lo}$  and  $\alpha_{\rm hi}$ 

The following invariants are being preserved:

- The interval bounded by α<sub>lo</sub> and α<sub>hi</sub> always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α<sub>lo</sub> ≤ α<sub>hi</sub>
- α<sub>lo</sub> is, among all step lengths generated so far and satisfying the sufficient decrease condition, the one giving the smallest value of φ,

1: function  $ZOOM(\alpha_{lo}, \alpha_{hi})$ 

2: repeat

- 3: Set  $\alpha$  between  $\alpha_{lo}$  and  $\alpha_{hi}$  using interpolation (bisection, quadratic, etc.)
- 4: Evaluate  $\phi(\alpha)$
- 5: if  $\phi(\alpha) > \phi(0) + \alpha \mu_1 \phi'(0)$  or  $\phi(\alpha) \ge \phi(\alpha_{\mathsf{lo}})$  then
- 6:  $\alpha_{hi} \leftarrow \alpha$

#### else

7:

- 8: Evaluate  $\phi'(\alpha)$
- 9: **if**  $|\phi'(\alpha)| \le \mu_2 |\phi'(0)|$  then 10: Set  $\alpha^* \leftarrow \alpha$  and stop
  - J: Set  $\alpha^* \leftarrow \alpha$
- 11: end if
- 12: **if**  $\phi'(\alpha)(\alpha_{hi} \alpha_{lo}) \ge 0$  **then**

13:  $\alpha_{hi} \leftarrow \alpha_{lo}$ 

- 14: **end if**
- 15:  $\alpha_{\mathsf{lo}} \leftarrow \alpha$
- 16: **end if**
- 17: **until** a condition is met
- 18: end function

#### Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses  $\alpha_{i+1} = 2\alpha_i$ , and the sufficient curvature factor is  $\mu_2 = 0.9$ .



#### Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses  $\alpha_{i+1} = 2\alpha_i$ , and the sufficient curvature factor is  $\mu_2 = 0.9$ .



Bracketing is achieved in the first iteration by using a significant initial step of  $\alpha_{init} = 1.2$  (left). Then, zooming finds an improved point through interpolation.

### Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses  $\alpha_{i+1} = 2\alpha_i$ , and the sufficient curvature factor is  $\mu_2 = 0.9$ .



Bracketing is achieved in the first iteration by using a significant initial step of  $\alpha_{init} = 1.2$  (left). Then, zooming finds an improved point through interpolation.

The small initial step of  $\alpha_{init} = 0.05$  (right) does not satisfy the strong Wolfe conditions, and the bracketing phase moves forward toward a flatter part of the function.

The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x<sub>k</sub>) and f (x<sub>k-1</sub>) may be indistinguishable in finite-precision arithmetic.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x<sub>k</sub>) and f (x<sub>k-1</sub>) may be indistinguishable in finite-precision arithmetic.
- Some procedures also stop if the relative change in x is close to machine accuracy or some user-specified threshold.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x<sub>k</sub>) and f (x<sub>k-1</sub>) may be indistinguishable in finite-precision arithmetic.
- Some procedures also stop if the relative change in x is close to machine accuracy or some user-specified threshold.
- The presented algorithm is implemented in https://docs.scipy.org/doc/scipy/reference/ generated/scipy.optimize.line\_search.html

# Unconstrained Optimization Algorithms

**Descent Direction** 

First-Order Methods

#### Gradient Descent

Consider the *gradient descent* method where

$$x_{k+1} = x_k + \alpha_k p_k$$
  $p_k = -\nabla f(x_k)$ 



#### Gradient Descent

Consider the *gradient descent* method where

$$x_{k+1} = x_k + \alpha_k p_k$$
  $p_k = -\nabla f(x_k)$ 



Unfortunately, the gradient does not possess much information about the step size.

So usually, a normalized gradient is used to obtain the direction, and then a line search is performed:

$$x_{k+1} = x_k + \alpha_k p_k$$
  $p_k = -\frac{\nabla f(x_k)}{||\nabla f(x_k)||}$ 

The line search is *exact* if  $\alpha_k$  minimizes  $f(x_k + \alpha_k p_k)$ . Not practical, we usually find  $\alpha_k$  satisfying the strong Wolfe conditions.

## Gradient Descent Algorithm with Line Search

Algorithm 4 Gradient Descent with Line Search

**Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point

1: 
$$k \leftarrow 0$$

2: while 
$$\|\nabla f\|_{\infty} > \varepsilon$$
 do

3: 
$$p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$$

4: Set 
$$\alpha_{init}$$
 for line search

5: 
$$\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$$

6: 
$$x_{k+1} \leftarrow x_k + \alpha_k p_k$$

7: 
$$k \leftarrow k+1$$

8: end while

#### Gradient Descent Algorithm with Line Search

Algorithm 5 Gradient Descent with Line SearchInput:  $x_0$  starting point,  $\varepsilon > 0$ Output:  $x^*$  approximation to a stationary point1:  $k \leftarrow 0$ 2: while  $\|\nabla f\|_{\infty} > \varepsilon$  do3:  $p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$ 4: Set  $\alpha_{init}$  for line search5:  $\alpha_k \leftarrow$  linesearch $(p_k, \alpha_{init})$ 6:  $x_{k+1} \leftarrow x_k + \alpha_k p_k$ 7:  $k \leftarrow k+1$ 

8: end while

Here  $\alpha_{init}$  can be estimated from the previous step size  $\alpha_{k-1}$  by demanding similar decrease in the objective:

$$\alpha_{\text{init}} \boldsymbol{p}_{k}^{\top} \nabla f_{k} \approx \alpha_{k-1} \boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1} \quad \Rightarrow \quad \alpha_{\text{init}} = \alpha_{k-1} \frac{\boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1}}{\boldsymbol{p}_{k}^{\top} \nabla f_{k}}$$



Note that  $p_{k+1}$  and  $p_k$  are always orthogonal.

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2}(2x_2 - x_1^2)^2$$

 $\text{Stopping: } \left|\left|\nabla f\right|\right|_{\infty} \leq 10^{-6}.$ 



#### Global Convergence with Line Search

Recall the Zoutendijk's theorem.

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth for some L > 0 if

$$\|
abla f(x) - 
abla f( ilde{x})\| \le L \|x - ilde{x}\|, \quad ext{ for all } x, ilde{x} \in \mathbb{R}^n$$

#### Theorem 11 (Zoutendijk)

Consider  $x_{k+1} = x_k + \alpha_k p_k$ , where  $p_k$  is a descent direction and  $\alpha_k$  satisfies the strong Wolfe conditions. Suppose that f is bounded below, continuously differentiable, and L-smooth. Then

$$\sum_{k\geq 0}\cos^2\theta_k\,\|\nabla f_k\|^2<\infty.$$

#### Global Convergence of Gradient Descent

Assume that each  $\alpha_k$  satisfies strong Wolfe conditions.

#### Global Convergence of Gradient Descent

Assume that each  $\alpha_k$  satisfies strong Wolfe conditions.

Note that the angle  $\theta_k$  between  $p_k = -\nabla f_k$  and the negative gradient  $-\nabla f_k$  equals 0. Hence,  $\cos \theta_k = 1$ .

#### Global Convergence of Gradient Descent

Assume that each  $\alpha_k$  satisfies strong Wolfe conditions.

Note that the angle  $\theta_k$  between  $p_k = -\nabla f_k$  and the negative gradient  $-\nabla f_k$  equals 0. Hence,  $\cos \theta_k = 1$ .

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\sum_{k\geq 0}\cos^2\theta_k \|\nabla f_k\|^2 = \sum_{k\geq 0} \|\nabla f_k\|^2 < \infty$$

which implies that  $\lim_{k\to\infty} ||\nabla f_k|| = 0$ .

## Local Linear Convergence of Gradient Descent (Optional)

#### Theorem 12

Suppose that  $f : \mathbb{R}^n \to \mathbb{R}$  is twice continuously differentiable, that the line search is exact, and that the descent converges to  $x^*$  where  $\nabla f(x^*) = 0$  and the Hessian matrix  $\nabla^2 f(x^*)$  is positive definite. Then

$$f(x_{k+1}) - f(x^*) \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^2 \left[f(x_k) - f(x^*)\right],$$

where  $\lambda_1 \leq \cdots \leq \lambda_n$  are the eigenvalues of  $\nabla^2 f(x^*)$ .





Here  $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$ 

Two Spring Problem - Gradient Descent



Gradient descent, line search, stop. cond.  $||\nabla f||_{\infty} \leq 10^{-6}$ .

Rosenbrock Function - Gradient Descent Rosenbrock:  $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$ 



Gradient descent, line search, stop. cond.  $||\nabla f||_{\infty} \leq 10^{-6}$ .

The method needs evaluation of ∇f at each x<sub>k</sub>. If f is not differentiable at x<sub>k</sub>, subgradients can be considered (out of the scope of this course).

- The method needs evaluation of ∇f at each x<sub>k</sub>. If f is not differentiable at x<sub>k</sub>, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.

- The method needs evaluation of ∇f at each x<sub>k</sub>. If f is not differentiable at x<sub>k</sub>, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.
- Susceptible to scaling of variables (see the paraboloid example).

- The method needs evaluation of ∇f at each x<sub>k</sub>. If f is not differentiable at x<sub>k</sub>, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.
- Susceptible to scaling of variables (see the paraboloid example).
- THE basis for algorithms training neural networks a huge amount of specific adjustments are developed for working with huge numbers of variables in neural networks (trillions of weights).

# Unconstrained Optimization Algorithms

**Descent Direction** 

Second-Order Methods

#### Newton's Method

Consider an objective  $f : \mathbb{R}^n \to \mathbb{R}$ .

Assume that f is twice differentiable.

#### Newton's Method

Consider an objective  $f : \mathbb{R}^n \to \mathbb{R}$ .

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(x_k+p) \approx f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

Here we denote the gradient  $\nabla f(x_k)$  of f at  $x_k$  by  $\nabla f_k$  and the Hessian  $\nabla^2 f(x_k)$  by  $H_k$ .
## Newton's Method

Consider an objective  $f : \mathbb{R}^n \to \mathbb{R}$ .

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(x_k + p) \approx f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

Here we denote the gradient  $\nabla f(x_k)$  of f at  $x_k$  by  $\nabla f_k$  and the Hessian  $\nabla^2 f(x_k)$  by  $H_k$ .

Define

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

and minimize q w.r.t. p by setting  $\nabla q(p) = 0$ .

### Newton's Method

Consider an objective  $f : \mathbb{R}^n \to \mathbb{R}$ .

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(x_k + p) \approx f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

Here we denote the gradient  $\nabla f(x_k)$  of f at  $x_k$  by  $\nabla f_k$  and the Hessian  $\nabla^2 f(x_k)$  by  $H_k$ .

Define

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

and minimize q w.r.t. p by setting  $\nabla q(p) = 0$ . We obtain:

$$H_k p = -\nabla f_k$$

Denote by  $p_k$  the solution, and set  $x_{k+1} = x_k + p_k$ .

Algorithm 6 Newton's Method

**Input:**  $x_0$  starting point,  $\tau > 0$ 

**Output:**  $x^*$  approximation to a stationary point

- 1:  $k \leftarrow 0$
- 2: while  $\| 
  abla f_k \|_{\infty} > \tau$  do
- 3: Compute  $\nabla f_k = \nabla f(x_k)$
- 4: Solve  $H_k p_k = -\nabla f_k$  for  $p_k$
- 5:  $x_{k+1} \leftarrow x_k + p_k$
- $6: \quad k \leftarrow k+1$
- 7: end while

## Newton's Method - Example

Newton's method finds the minimum of a quadratic function in a single step.



Note that the Newton's method is scale-invariant!

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2} (2x_2 - x_1^2)^2$$

Stopping:  $||\nabla f||_{\infty} \leq 10^{-6}$ .



k = 0

k = 2

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2}(2x_2 - x_1^2)^2$$

Stopping:  $||\nabla f||_{\infty} \leq 10^{-6}$ .



k = 2

k = 8

#### **Convergence** Issues



Also, the computation of the Hessian is costly.

#### Theorem 13

Assume f is twice differentiable and assume that  $\nabla f$  is L-smooth. Let  $x_*$  be a minimizer of f(x) and assume that  $\nabla^2 f(x_*)$  is positive definite.

If  $||x_0 - x_*||$  is sufficiently small, then  $\{x_k\}$  converges quadratically to  $x_*$ .

#### Theorem 13

Assume f is twice differentiable and assume that  $\nabla f$  is L-smooth. Let  $x_*$  be a minimizer of f(x) and assume that  $\nabla^2 f(x_*)$  is positive definite.

If  $||x_0 - x_*||$  is sufficiently small, then  $\{x_k\}$  converges quadratically to  $x_*$ .

Note that the theorem implicitly assumes that  $\nabla^2 f(x_k)$  is nonsingular for every k.

#### Theorem 13

Assume f is twice differentiable and assume that  $\nabla f$  is L-smooth. Let  $x_*$  be a minimizer of f(x) and assume that  $\nabla^2 f(x_*)$  is positive definite.

If  $||x_0 - x_*||$  is sufficiently small, then  $\{x_k\}$  converges quadratically to  $x_*$ .

Note that the theorem implicitly assumes that  $\nabla^2 f(x_k)$  is nonsingular for every k.

As the theorem is concerned only with  $x_k$  approaching  $x^*$ , the continuity of  $\nabla^2 f(x_k)$  and positive definiteness of  $\nabla^2 f(x^*)$  imply that  $\nabla^2 f(x_k)$  is positive definite for all sufficiently large k.

#### Theorem 13

Assume f is twice differentiable and assume that  $\nabla f$  is L-smooth. Let  $x_*$  be a minimizer of f(x) and assume that  $\nabla^2 f(x_*)$  is positive definite.

If  $||x_0 - x_*||$  is sufficiently small, then  $\{x_k\}$  converges quadratically to  $x_*$ .

Note that the theorem implicitly assumes that  $\nabla^2 f(x_k)$  is nonsingular for every k.

As the theorem is concerned only with  $x_k$  approaching  $x^*$ , the continuity of  $\nabla^2 f(x_k)$  and positive definiteness of  $\nabla^2 f(x^*)$  imply that  $\nabla^2 f(x_k)$  is positive definite for all sufficiently large k.

However, what happens if we start far away from a minimizer?

# Newton's Method with Line Search

Algorithm 7 Newton's Method with Line Search **Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point 1  $k \leftarrow 0$ 2:  $\alpha_{\text{init}} \leftarrow 1$ 3: while  $\|\nabla f_k\|_{\infty} > \varepsilon$  do Compute  $\nabla f_k = \nabla f(x_k)$ 4: 5: Solve  $H_k p_k = -\nabla f_k$  for  $p_k$ 6:  $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ 7:  $x_{k+1} \leftarrow x_k + \alpha p_k$  $k \leftarrow k+1$ 8. 9: end while



Here  $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$ 

## Two Spring Problem - Newton's Method



Newton's method, line search and stop. cond.  $||\nabla f||_{\infty} \leq 10^{-6}$ . Compare this with 32 iterations of gradient descent. Rosenbrock Function - Newton's Method Rosenbrock:  $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$ 



Newton's method, line search, stop. cond.  $||\nabla f||_{\infty} \le 10^{-6}$ . Compare this with 10,662 iterations of gradient descent.

## Global Convergence of Line Search

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth for some L > 0 if

$$\|
abla f(x) - 
abla f( ilde{x})\| \le L \|x - ilde{x}\|, \quad \text{ for all } x, ilde{x} \in \mathbb{R}^n$$

#### Theorem 14 (Zoutendijk)

Consider  $x_{k+1} = x_k + \alpha_k p_k$ , where  $p_k$  is a descent direction and  $\alpha_k$  satisfies the strong Wolfe conditions. Suppose that f is bounded below, continuously differentiable, and L-smooth. Then

$$\sum_{k\geq 0}\cos^2\theta_k\,\|\nabla f_k\|^2<\infty.$$

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the Hessians  $H_k$  are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$  for all k

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the Hessians  $H_k$  are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$  for all k

Then  $\theta_k$  between  $p_k = -H_k^{-1} \nabla f_k$  and  $-\nabla f_k$  satisfies

 $\cos \theta_k \ge 1/M$ 

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the Hessians  $H_k$  are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$  for all k

Then  $\theta_k$  between  $p_k = -H_k^{-1} \nabla f_k$  and  $-\nabla f_k$  satisfies

 $\cos \theta_k \geq 1/M$ 

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\frac{1}{M^2}\sum_{k\geq 0} \|\nabla f_k\|^2 \leq \sum_{k\geq 0} \cos^2 \theta_k \|\nabla f_k\|^2 < \infty$$

which implies that  $\lim_{k\to\infty} ||\nabla f_k|| = 0$ .

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the Hessians  $H_k$  are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$  for all k

Then  $\theta_k$  between  $p_k = -H_k^{-1} \nabla f_k$  and  $-\nabla f_k$  satisfies

 $\cos \theta_k \geq 1/M$ 

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\frac{1}{M^2}\sum_{k\geq 0} \|\nabla f_k\|^2 \leq \sum_{k\geq 0} \cos^2 \theta_k \|\nabla f_k\|^2 < \infty$$

which implies that  $\lim_{k\to\infty} ||\nabla f_k|| = 0$ .

What if  $H_k$  is not positive definite or is (nearly) singular?

### **Eigenvalue Modification**

Consider  $H_k = \nabla^2 f(x_k)$  and consider its diagonal form:

$$H_k = QDQ^T$$

Where D contains the eigenvalues of  $H_k$  on the diagonal, i.e.,  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$  and Q is an orthogonal matrix.

## **Eigenvalue Modification**

Consider  $H_k = \nabla^2 f(x_k)$  and consider its diagonal form:

$$H_k = Q D Q^T$$

Where D contains the eigenvalues of  $H_k$  on the diagonal, i.e.,  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$  and Q is an orthogonal matrix.

Observe that

- $H_k$  is not positive definite iff  $\lambda_i \leq 0$  for some *i*
- ▶  $||H_k||$  grows with max{ $\lambda_1, \ldots, \lambda_n$ } going to infinity.
- ||H<sub>k</sub><sup>-1</sup>|| grows with min{λ<sub>1</sub>,...,λ<sub>n</sub>} going to 0 (i.e., the matrix becomes close to a singular matrix)

We want to prevent all three cases, i.e., make sure that for some reasonably large  $\delta > 0$  we have  $\lambda_i \ge \delta$  but not too large.

## **Eigenvalue Modification**

Consider  $H_k = \nabla^2 f(x_k)$  and consider its diagonal form:

$$H_k = QDQ^T$$

Where D contains the eigenvalues of  $H_k$  on the diagonal, i.e.,  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$  and Q is an orthogonal matrix.

Observe that

- $H_k$  is not positive definite iff  $\lambda_i \leq 0$  for some *i*
- ▶  $||H_k||$  grows with max $\{\lambda_1, \ldots, \lambda_n\}$  going to infinity.
- ||H<sub>k</sub><sup>-1</sup>|| grows with min{λ<sub>1</sub>,...,λ<sub>n</sub>} going to 0 (i.e., the matrix becomes close to a singular matrix)

We want to prevent all three cases, i.e., make sure that for some reasonably large  $\delta > 0$  we have  $\lambda_i \ge \delta$  but not too large.

Two questions are in order:

- What is a reasonably large  $\delta$ ?
- How to modify  $H_k$  so the minimum is large enough?

Consider an example:

$$abla f(x_k) = (1, -3, 2)$$
 and

$$\nabla^2 f(x_k) = \mathsf{diag}(10, 3, -1)$$

Consider an example:

$$abla f(x_k) = (1, -3, 2)$$
 and  $abla^2 f(x_k) = diag(10, 3, -1)$ 

Now, the diagonalization is trivial:

$$abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I$$
 is the identity matrix

Consider an example:

 $abla f(x_k) = (1, -3, 2)$  and  $abla^2 f(x_k) = diag(10, 3, -1)$ 

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I$  is the identity matrix

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say  $\delta = 10^{-8}$ ?

Consider an example:

 $abla f(x_k) = (1, -3, 2)$  and  $abla^2 f(x_k) = diag(10, 3, -1)$ 

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10, 3, -1) Q^\top \quad Q = I$  is the identity matrix

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say  $\delta = 10^{-8}$ ? Obtain

$$B_k = Q \operatorname{diag}(10, 3, 10^{-8}) \ Q^{\top} = \operatorname{diag}(10, 3, 10^{-8})$$

Consider an example:

 $abla f(x_k) = (1, -3, 2)$  and  $abla^2 f(x_k) = diag(10, 3, -1)$ 

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10, 3, -1) Q^\top \quad Q = I \text{ is the identity matrix}$ 

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say  $\delta = 10^{-8}$ ? Obtain

$$B_k = Q \operatorname{diag}(10, 3, 10^{-8}) \ Q^{\top} = \operatorname{diag}(10, 3, 10^{-8})$$

If used in Newton's method, we obtain the following direction:

$$p_k = -B_k^{-1} \nabla f(x_k) = (-1/10, 1, -(2 \cdot 10^8))$$

Thus, a very long vector almost parallel to the third dimension.

Consider an example:

 $abla f(x_k) = (1, -3, 2)$  and  $abla^2 f(x_k) = diag(10, 3, -1)$ 

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10, 3, -1) Q^\top \quad Q = I \text{ is the identity matrix}$ 

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say  $\delta = 10^{-8}$ ? Obtain

$$B_k = Q \operatorname{diag}(10, 3, 10^{-8}) \ Q^{\top} = \operatorname{diag}(10, 3, 10^{-8})$$

If used in Newton's method, we obtain the following direction:

$$p_k = -B_k^{-1} \nabla f(x_k) = (-1/10, 1, -(2 \cdot 10^8))$$

Thus, a very long vector almost parallel to the third dimension. Note that the original Newton's direction is  $-\text{diag}(1/10, 1/3, -1)(1, -3, 2)^{\top} = (-1/10, 1, 2)$  which is completely different.

Other strategies for eigenvalue modification can be devised.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix  $B_k$  should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix  $B_k$  should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix  $B_k$  should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix  $B_k$  should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

The implementation is based on computing  $B_k = H_k + \Delta H_k$  for an appropriate modification matrix  $\Delta H_k$ .

What is  $\Delta H_k$  in our example?

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix  $B_k$  should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

The implementation is based on computing  $B_k = H_k + \Delta H_k$  for an appropriate modification matrix  $\Delta H_k$ .

What is  $\Delta H_k$  in our example?

Various methods for computing  $\Delta H_k$  have been devised in literature. Typically, it is based on some computationally cheaper decomposition than spectral decomposition (e.g., Cholesky).
# Modified Newton's Method

**Algorithm 8** Newton's Method with Line Search **Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point 1:  $k \leftarrow 0$ 2: while  $\|\nabla f_k\|_{\infty} > \varepsilon$  do  $H_{\mu} \leftarrow \nabla^2 f(x_{\mu})$ 3: **4**· if  $H_k$  is **not** sufficiently positive definite **then**  $H_k \leftarrow H_k + \Delta H_k$  so that  $H_k$  is sufficiently pos. definite 5: end if 6: Compute  $\nabla f_k = \nabla f(x_k)$ 7: Solve  $H_k p_k = -\nabla f_k$  for  $p_k$ 8: 9: Set  $x_{k+1} = x_k + \alpha_k p_k$ , here  $\alpha_k$  sat. the Wolfe cond.  $k \leftarrow k + 1$ 10: 11: end while

Convergence theorems are complicated in this case and out of the scope of this course. See Chapter 6 of Numerical Optimization by Nocedal & Wright.

Newton's method is scale invariant.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
  - \$\mathcal{O}(n^2)\$ second derivatives (the Hessian), each may be hard to compute.

Automated derivation methods help but still need store  $O(n^2)$  results.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
  - \$\mathcal{O}(n^2)\$ second derivatives (the Hessian), each may be hard to compute.

Automated derivation methods help but still need store  $O(n^2)$  results.

\$\mathcal{O}(n^3)\$ arithmetic operations to solve the linear system for the direction \$p\_k\$.

May be mitigated by more efficient methods in case of sparse Hessians.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
  - \$\mathcal{O}(n^2)\$ second derivatives (the Hessian), each may be hard to compute.

Automated derivation methods help but still need store  $\mathcal{O}(n^2)$  results.

\$\mathcal{O}(n^3)\$ arithmetic operations to solve the linear system for the direction \$p\_k\$.

May be mitigated by more efficient methods in case of sparse Hessians.

In a sense, Newton's method is an impractical "ideal" with which other methods are compared.

The efficiency issues (and the necessity of second-order derivatives) will be mitigated by using quasi-Newton methods.

Recall that Newton's method step  $p_k$  in  $x_{k+1} = x_k + p_k$  comes from minimization of

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

w.r.t. p by setting  $\nabla q(p) = 0$  and solving

$$H_k p = -\nabla f_k$$

So Newton's method needs the second derivative (Hessian), which is computationally hard to obtain.

Recall that Newton's method step  $p_k$  in  $x_{k+1} = x_k + p_k$  comes from minimization of

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

w.r.t. p by setting  $\nabla q(p) = 0$  and solving

$$H_k p = -\nabla f_k$$

So Newton's method needs the second derivative (Hessian), which is computationally hard to obtain.

Gradient descent needs only the first derivatives but converges slowly.

Recall that Newton's method step  $p_k$  in  $x_{k+1} = x_k + p_k$  comes from minimization of

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

w.r.t. p by setting  $\nabla q(p) = 0$  and solving

$$H_k p = -\nabla f_k$$

So Newton's method needs the second derivative (Hessian), which is computationally hard to obtain.

Gradient descent needs only the first derivatives but converges slowly.

Can we find a compromise?

Recall that Newton's method step  $p_k$  in  $x_{k+1} = x_k + p_k$  comes from minimization of

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

w.r.t. p by setting  $\nabla q(p) = 0$  and solving

$$H_k p = -\nabla f_k$$

So Newton's method needs the second derivative (Hessian), which is computationally hard to obtain.

Gradient descent needs only the first derivatives but converges slowly.

Can we find a compromise?

Quasi-Newton methods use first derivatives to approximate the Hessian  $H_k$  in Newton's method with a matrix  $\tilde{H}_k$ .

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Consider the Hessian  $H_{k+1} = \nabla^2 f(x_{k+1})$  and its approximation denoted by  $\tilde{H}_{k+1}$ .

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Consider the Hessian  $H_{k+1} = \nabla^2 f(x_{k+1})$  and its approximation denoted by  $\tilde{H}_{k+1}$ .

We aim to use  $\tilde{H}_{k+1}$  in the next step, that is, in the equation  $\tilde{H}_{k+1}p = -\nabla f_{k+1}$  yielding  $p_{k+1}$ .

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Consider the Hessian  $H_{k+1} = \nabla^2 f(x_{k+1})$  and its approximation denoted by  $\tilde{H}_{k+1}$ .

We aim to use  $\tilde{H}_{k+1}$  in the next step, that is, in the equation  $\tilde{H}_{k+1}p = -\nabla f_{k+1}$  yielding  $p_{k+1}$ .

What conditions should  $\tilde{H}_{k+1}$  satisfy so that it functions as the "true" Hessian  $H_{k+1}$ ?

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Consider the Hessian  $H_{k+1} = \nabla^2 f(x_{k+1})$  and its approximation denoted by  $\tilde{H}_{k+1}$ .

We aim to use  $\tilde{H}_{k+1}$  in the next step, that is, in the equation  $\tilde{H}_{k+1}p = -\nabla f_{k+1}$  yielding  $p_{k+1}$ .

What conditions should  $\tilde{H}_{k+1}$  satisfy so that it functions as the "true" Hessian  $H_{k+1}$ ?

First, it should be *symmetric positive definite*. To always yield decrease direction.

Suppose we have just obtained the new point  $x_{k+1}$  after a line search starting from  $x_k$  in the direction  $p_k$ .

Consider the Hessian  $H_{k+1} = \nabla^2 f(x_{k+1})$  and its approximation denoted by  $\tilde{H}_{k+1}$ .

We aim to use  $\tilde{H}_{k+1}$  in the next step, that is, in the equation  $\tilde{H}_{k+1}p = -\nabla f_{k+1}$  yielding  $p_{k+1}$ .

What conditions should  $\tilde{H}_{k+1}$  satisfy so that it functions as the "true" Hessian  $H_{k+1}$ ?

First, it should be *symmetric positive definite*. To always yield decrease direction.

Second, extrapolating from the single variable secant method, we demand the *secant condition*:

$$ilde{H}_{k+1}(x_{k+1}-x_k) = 
abla f_{k+1} - 
abla f_k$$

# Secant Condition

Consider the secant condition:

$$\tilde{H}_{k+1}(x_{k+1}-x_k)=\nabla f_{k+1}-\nabla f_k$$

# Secant Condition

Consider the secant condition:

$$\tilde{H}_{k+1}(x_{k+1}-x_k)=\nabla f_{k+1}-\nabla f_k$$

The notation is usually simplified by

$$s_k = x_{k+1} - x_k$$
  $y_k = \nabla f_{k+1} - \nabla f_k$ 

So that the secant condition becomes

$$\tilde{H}_{k+1}s_k = y_k$$

# Secant Condition

Consider the secant condition:

$$ilde{H}_{k+1}(x_{k+1}-x_k) = 
abla f_{k+1} - 
abla f_k$$

The notation is usually simplified by

$$s_k = x_{k+1} - x_k$$
  $y_k = \nabla f_{k+1} - \nabla f_k$ 

So that the secant condition becomes

$$\tilde{H}_{k+1}s_k = y_k$$

Note that even if we demand symmetric positive definite solutions to the secant condition, there are potentially infinitely many. Indeed, there are n(n+1)/2 degrees of freedom in a symmetric matrix, and the secant conditions represent only *n* conditions (Sylvester's criterion).

Moreover, we want to obtain  $\tilde{H}_{k+1}$  from  $\tilde{H}_k$  by

$$ilde{H}_{k+1} = ilde{H}_k + ext{something}$$

To have a nice iterative algorithm.

Note that the information about the solution is present in  $s_k$  and  $y_k$ , so it is natural to compose the solution using these vectors.

Note that the information about the solution is present in  $s_k$  and  $y_k$ , so it is natural to compose the solution using these vectors.

Consider  $u = \left(y_k - \tilde{H}_k s_k\right)$ 

$$\tilde{H}_{k+1} = \tilde{H}_k + \frac{uu^\top}{u^\top s_k}$$

Note that the information about the solution is present in  $s_k$  and  $y_k$ , so it is natural to compose the solution using these vectors.

Consider 
$$u = \left(y_k - \tilde{H}_k s_k\right)$$

$$\tilde{H}_{k+1} = \tilde{H}_k + \frac{uu^\top}{u^\top s_k}$$

Now, the secant condition is satisfied:

$$\tilde{H}_{k+1}s_k = \tilde{H}_k s_k + \frac{uu^{\top}s_k}{u^{\top}s_k} = \tilde{H}_k s_k + u = \tilde{H}_k s_k + \left(y_k - \tilde{H}_k s_k\right) = y_k$$

By the way, the matrix  $\frac{uu^{\top}}{u^{\top}s_k}$  is of rank one and is a unique symmetric rank one matrix which makes a symmetric  $\tilde{H}_{k+1}$  satisfy the secant condition.

Note that the information about the solution is present in  $s_k$  and  $y_k$ , so it is natural to compose the solution using these vectors.

Consider 
$$u = \left(y_k - \tilde{H}_k s_k\right)$$

$$\tilde{H}_{k+1} = \tilde{H}_k + \frac{uu^\top}{u^\top s_k}$$

Now, the secant condition is satisfied:

$$\tilde{H}_{k+1}s_k = \tilde{H}_k s_k + \frac{uu^\top s_k}{u^\top s_k} = \tilde{H}_k s_k + u = \tilde{H}_k s_k + \left(y_k - \tilde{H}_k s_k\right) = y_k$$

By the way, the matrix  $\frac{uu^{\top}}{u^{\top}s_k}$  is of rank one and is a unique symmetric rank one matrix which makes a symmetric  $\tilde{H}_{k+1}$  satisfy the secant condition.

To obtain a quasi-Newton method, it suffices to initialize  $\tilde{H}_0$ , typically to the identity *I*, and use  $\tilde{H}_k$  instead of the Hessian  $H_k = \nabla^2 f_k$  in Newton's method.

# Symmetric Rank One Update

#### Algorithm 9 SR1

**Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point  $k \leftarrow 0, \alpha_{\text{init}} \leftarrow 1, \ddot{H}_0 \leftarrow I$ while  $\|\nabla f_k\|_{\infty} > \varepsilon$  do Compute  $\nabla f_k = \nabla f(x_k)$ Solve for  $p_k$  in  $\tilde{H}_k p_k = -\nabla f_k$  $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$  $x_{k+1} \leftarrow x_k + \alpha p_k$  $s \leftarrow x_{k+1} - x_k$  $y \leftarrow \nabla f_{k+1} - \nabla f_k$  $u \leftarrow v - H_k s$  $\tilde{H}_{k+1} \leftarrow \tilde{H}_k + \frac{uu^{\top}}{u^{\top}}$  $k \leftarrow k + 1$ end while

Note that the denominator  $u^{\top}s_k$  can be 0, in which case the update is impossible. The usual strategy is to skip the update and set  $\tilde{H}_{k+1} = \tilde{H}_k$ .

We will look at a three-dimensional quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = egin{pmatrix} 2 & 0 & 0 \ 0 & 3 & 0 \ 0 & 0 & 4 \end{pmatrix} \quad ext{ and } \quad c = egin{pmatrix} -8 \ -9 \ -8 \end{pmatrix},$$

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

The initial guesses are  $\tilde{H}_0 = I$  and  $x_0 = (0, 0, 0)^{\top}$ .

We will look at a three-dimensional quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix}$$
 and  $c = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix}$ ,

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

The initial guesses are  $\tilde{H}_0 = I$  and  $x_0 = (0, 0, 0)^{\top}$ .

At the initial point,  $\|\nabla f(x_0)\|_{\infty} = \|-c\|_{\infty} = 9$ , so this point is not optimal.

We will look at a three-dimensional quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = egin{pmatrix} 2 & 0 & 0 \ 0 & 3 & 0 \ 0 & 0 & 4 \end{pmatrix}$$
 and  $c = egin{pmatrix} -8 \ -9 \ -8 \end{pmatrix}$ ,

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

The initial guesses are  $\tilde{H}_0 = I$  and  $x_0 = (0, 0, 0)^{\top}$ .

At the initial point,  $\|\nabla f(x_0)\|_{\infty} = \|-c\|_{\infty} = 9$ , so this point is not optimal. The first search direction is

$$p_0 = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix}.$$

The exact line search gives  $\alpha_0 = 0.3333$ .

The new estimate of the solution, the update vectors, and the new Hessian approximation are:

$$x_1 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, \nabla f_1 = \begin{pmatrix} 2.66 \\ 0 \\ -2.66 \end{pmatrix}, s_0 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, y_0 = \begin{pmatrix} -5.33 \\ -9.00 \\ -10.66 \end{pmatrix}$$

The new estimate of the solution, the update vectors, and the new Hessian approximation are:

$$x_1 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, \nabla f_1 = \begin{pmatrix} 2.66 \\ 0 \\ -2.66 \end{pmatrix}, s_0 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, y_0 = \begin{pmatrix} -5.33 \\ -9.00 \\ -10.66 \end{pmatrix}$$

and

$$\tilde{H}_1 = I + \frac{(y_0 - Is_0)(y_0 - Is_0)^\top}{(y_0 - Is_0)^\top s_0} = \begin{pmatrix} 1.1531 & 0.3445 & 0.4593 \\ 0.3445 & 1.7751 & 1.0335 \\ 0.4593 & 1.0335 & 2.3780 \end{pmatrix}$$

The new estimate of the solution, the update vectors, and the new Hessian approximation are:

$$x_1 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, \nabla f_1 = \begin{pmatrix} 2.66 \\ 0 \\ -2.66 \end{pmatrix}, s_0 = \begin{pmatrix} -2.66 \\ -3.00 \\ -2.66 \end{pmatrix}, y_0 = \begin{pmatrix} -5.33 \\ -9.00 \\ -10.66 \end{pmatrix}$$

and

$$\tilde{H}_1 = I + \frac{(y_0 - Is_0)(y_0 - Is_0)^{\top}}{(y_0 - Is_0)^{\top}s_0} = \begin{pmatrix} 1.1531 & 0.3445 & 0.4593 \\ 0.3445 & 1.7751 & 1.0335 \\ 0.4593 & 1.0335 & 2.3780 \end{pmatrix}$$

At this new point  $\|\nabla f(x_1)\|_{\infty} = 2.66$  so we keep going, obtaining the search direction

$$p_1 = \begin{pmatrix} -2.9137 \\ -0.5557 \\ 1.9257 \end{pmatrix},$$

and the step length  $\alpha_1=$  0.3942.

This gives the new estimates:

$$x_2 = \begin{pmatrix} -3.81 \\ -3.21 \\ -1.90 \end{pmatrix}, \quad \nabla f_2 = \begin{pmatrix} 0.36 \\ -0.65 \\ 0.36 \end{pmatrix}, \quad s_1 = \begin{pmatrix} -1.14 \\ -0.21 \\ 0.75 \end{pmatrix}, \quad y_1 = \begin{pmatrix} -2.29 \\ -0.65 \\ 3.03 \end{pmatrix}$$

and

	/ 1.6568	0.6102	-0.3432	
$\tilde{H}_2 =$	0.6102	1.9153	0.6102	
	\_0.3432	0.6102	3.6568 /	

At the point x\_2,  $\|
abla f(x_2)\|_{\infty} = 0.65$  so we keep going, with

$$p_2 = \begin{pmatrix} -0.4851\\ 0.5749\\ -0.2426 \end{pmatrix},$$

and  $\alpha = 0.3810$ .

#### This gives

$$x_3 = \begin{pmatrix} -4 \\ -3 \\ -2 \end{pmatrix}, \quad \nabla f_3 = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \quad s_2 = \begin{pmatrix} -0.18 \\ 0.21 \\ -0.09 \end{pmatrix}, \quad y_2 = \begin{pmatrix} -0.36 \\ 0.65 \\ -0.36 \end{pmatrix},$$

and  $\tilde{H}_3 = Q$ . Now  $\|\nabla f(x_3)\|_{\infty} = 0$ , so we stop.

# Properties of SR1

Does symmetric rank one update satisfy our demands? We want every  $\tilde{H}_k$  to be a symmetric positive definite solution to the secant condition.

# Properties of SR1

Does symmetric rank one update satisfy our demands? We want every  $\tilde{H}_k$  to be a symmetric positive definite solution to the secant condition.

Unfortunately, though  $\tilde{H}_k$  is a symmetric positive definite, the updated matrix  $\tilde{H}_{k+1}$  does not have to be positive definite.
# Properties of SR1

Does symmetric rank one update satisfy our demands? We want every  $\tilde{H}_k$  to be a symmetric positive definite solution to the secant condition.

Unfortunately, though  $\tilde{H}_k$  is a symmetric positive definite, the updated matrix  $\tilde{H}_{k+1}$  does not have to be positive definite.

Still, the symmetric rank one approximation is used in practice, especially in trust region methods.

# Properties of SR1

Does symmetric rank one update satisfy our demands? We want every  $\tilde{H}_k$  to be a symmetric positive definite solution to the secant condition.

Unfortunately, though  $\tilde{H}_k$  is a symmetric positive definite, the updated matrix  $\tilde{H}_{k+1}$  does not have to be positive definite.

Still, the symmetric rank one approximation is used in practice, especially in trust region methods.

However, for line search, let us try a bit "richer" solution to the secant condition.

#### Symmetric Rank Two Update Consider

$$\tilde{H}_{k+1} = \tilde{H}_k - \frac{\left(\tilde{H}_k s_k\right) \left(\tilde{H}_k s_k\right)^\top}{s_k^\top \tilde{H}_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k}$$

Once again, verifying  $\tilde{H}_{k+1}s_k = y_k$  is not difficult.

#### Symmetric Rank Two Update Consider

$$\tilde{H}_{k+1} = \tilde{H}_k - \frac{\left(\tilde{H}_k s_k\right) \left(\tilde{H}_k s_k\right)^\top}{s_k^\top \tilde{H}_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k}$$

Once again, verifying  $\tilde{H}_{k+1}s_k = y_k$  is not difficult.

#### Crucial observation:

Assume that  $\tilde{H}_k$  is symmetric positive definite.

If the next approximation  $x_{k+1}$  is computed using a line search satisfying the strong Wolfe conditions, then  $\tilde{H}_{k+1}$  is also symmetric positive definite.

For proof see Lemma 12.10 and Exercise 3.9 of "Linear and Nonlinear Optimization" by Griva et al.

Thus, starting with a symmetric positive definite  $\tilde{H}_0$  and doing line search satisfying the strong Wolfe conditions, every  $\tilde{H}_k$  is symmetric positive definite and satisfies the secant condition.

# BFGS

#### Algorithm 10 BFGS v1

**Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point  $k \leftarrow 0, \alpha_{\text{init}} \leftarrow 1, \tilde{H}_0 \leftarrow I$ while  $\|\nabla f_k\|_{\infty} > \tau$  do Compute  $\nabla f_k = \nabla f(x_k)$ Solve for  $p_k$  in  $\tilde{H}_k p_k = -\nabla f_k$  $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$  $x_{k+1} \leftarrow x_k + \alpha p_k$  $s \leftarrow x_{k+1} - x_k$  $y \leftarrow \nabla f_{k+1} - \nabla f_k$  $ilde{H}_{k+1} \leftarrow ilde{H}_k - rac{\left( ilde{H}_k s
ight) \left( ilde{H}_k s
ight)^{ op}}{s^{ op} ilde{H}_k s} + rac{yy^{ op}}{y^{ op}}$  $k \leftarrow k + 1$ end while

Note that we still have to solve a linear system for  $p_k$ .

Consider the quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix}$$
 and  $c = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix}$ ,

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

Consider the quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix} \quad \text{and} \quad c = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix},$$

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

Choose  $\tilde{H}_0 = I$  and  $x_0 = (0, 0, 0)^T$ .

Consider the quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix} \quad \text{and} \quad c = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix},$$

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

Choose 
$$\tilde{H}_0 = I$$
 and  $x_0 = (0, 0, 0)^T$ .

At iteration  $0, \|\nabla f(x_0)\|_{\infty} = 9$ , so this point is not optimal.

Consider the quadratic problem  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x$  with

$$Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix} \quad \text{and} \quad c = \begin{pmatrix} -8 \\ -9 \\ -8 \end{pmatrix},$$

whose solution is  $x_* = (-4, -3, -2)^{\top}$ . Use the exact line search.

Choose 
$$\tilde{H}_0 = I$$
 and  $x_0 = (0, 0, 0)^T$ .

At iteration  $0, \|\nabla f(x_0)\|_{\infty} = 9$ , so this point is not optimal.

The search direction is

$$p_0 = \left(\begin{array}{c} -8\\ -9\\ -8 \end{array}\right)$$

and  $\alpha_0 = 0.3333$ .

The new estimate of the solution and the new Hessian approximation are

$$x_1 = \begin{pmatrix} -2.6667 \\ -3.0000 \\ -2.6667 \end{pmatrix} \text{ and } \tilde{H}_1 = \begin{pmatrix} 1.1021 & 0.3445 & 0.5104 \\ 0.3445 & 1.7751 & 1.0335 \\ 0.5104 & 1.0335 & 2.3270 \end{pmatrix}$$

.

The new estimate of the solution and the new Hessian approximation are

$$x_1 = \begin{pmatrix} -2.6667 \\ -3.0000 \\ -2.6667 \end{pmatrix} \quad \text{and} \quad \tilde{H}_1 = \begin{pmatrix} 1.1021 & 0.3445 & 0.5104 \\ 0.3445 & 1.7751 & 1.0335 \\ 0.5104 & 1.0335 & 2.3270 \end{pmatrix}$$

At iteration 1,  $\left\| \nabla f(x_1) \right\|_{\infty} = 2.6667$ , so we continue. The next search direction is

$$p_1 = \left(\begin{array}{c} -3.2111 \\ -0.6124 \\ 2.1223 \end{array}\right)$$

and  $\alpha_1 = 0.3577$ .

This gives the estimates.

$$x_2 = \begin{pmatrix} -3.8152 \\ -3.2191 \\ -1.9076 \end{pmatrix} \quad \text{and} \quad \tilde{H}_2 = \begin{pmatrix} 1.6393 & 0.6412 & -0.3607 \\ 0.6412 & 1.8600 & 0.6412 \\ -0.3607 & 0.6412 & 3.6393 \end{pmatrix}$$

At iteration 2,  $\left\| 
abla f(x_2) \right\|_\infty = 0.6572$ , so we continue, computing

$$p_2 = \left(\begin{array}{c} -0.5289\\ 0.6268\\ -0.2644\end{array}\right)$$

and  $\alpha_2 = 0.3495$ . This gives

$$x_3 = \begin{pmatrix} -4 \\ -3 \\ -2 \end{pmatrix}$$
 and  $\tilde{H}_3 = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix}$ .

Now  $\left\|\nabla f(x_3)\right\|_{\infty} = 0$ , so we stop.

Notice that we got the same  $x_1, x_2, x_3$  as for SR1. This follows from using the exact line search and the quadratic problem. It does not hold in general.



Here  $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$ 

## Two Spring Problem - BFGS



BFGS, line search, stop. cond.  $||\nabla f||_{\infty} \leq 10^{-6}$ . Compare this with 32 iterations of gradient descent and 12 iterations of Newton's method.

## Rosenbrock Function - BFGS *Rosenbrock:* $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$



BFGS, line search, stop. cond.  $||\nabla f||_{\infty} \leq 10^{-6}$ . Compare with 10,662 iterations of gradient descent and 24 iterations of Newton's method.

**Problem:** SR1 and BFGS solve  $\tilde{H}_k p = -\nabla f_k$  repeatedly. What if we could iteratively update  $H_k^{-1}$ ?

**Problem:** SR1 and BFGS solve  $\tilde{H}_k p = -\nabla f_k$  repeatedly. What if we could iteratively update  $H_k^{-1}$ ?

The equation would be solved by  $p_k = -H_k^{-1} \nabla f_k$ .

**Problem:** SR1 and BFGS solve  $\tilde{H}_k p = -\nabla f_k$  repeatedly. What if we could iteratively update  $H_k^{-1}$ ?

The equation would be solved by  $p_k = -H_k^{-1} \nabla f_k$ .

Ideally, we would like to compute  $\tilde{H}_k^{-1}$  iteratively along the optimization, i.e.,

$$ilde{H}_{k+1}^{-1} = ilde{H}_k^{-1} + ext{something}$$

**Problem:** SR1 and BFGS solve  $\tilde{H}_k p = -\nabla f_k$  repeatedly. What if we could iteratively update  $H_k^{-1}$ ?

The equation would be solved by  $p_k = -H_k^{-1} \nabla f_k$ .

Ideally, we would like to compute  $\tilde{H}_k^{-1}$  iteratively along the optimization, i.e.,

 $ilde{H}_{k+1}^{-1} = ilde{H}_k^{-1} + ext{something}$ 

To get such a "something" we use the following Sherman–Morrison–Woodbury (SMW) formula:

$$(A + UV^{T})^{-1} = A^{-1} - A^{-1}U(I + V^{T}A^{-1}U)^{-1}V^{T}A^{-1}$$

Here A is a  $(n \times n)$ -matrix, U, V are  $(n \times m)$ -matrices with  $m \le n$ .

## Rank 1 – Iterative Inverse Hessian Approximation

Applying SMW to the rank one update

$$ilde{H}_{k+1} = ilde{H}_k + rac{\left(y_k - ilde{H}_k s_k
ight) \left(y_k - ilde{H}_k s_k
ight)^ op}{\left(y_k - ilde{H}_k s_k
ight)^ op s_k}$$

## Rank 1 – Iterative Inverse Hessian Approximation

Applying SMW to the rank one update

$$ilde{H}_{k+1} = ilde{H}_k + rac{\left(y_k - ilde{H}_k s_k
ight) \left(y_k - ilde{H}_k s_k
ight)^ op}{\left(y_k - ilde{H}_k s_k
ight)^ op s_k}$$

yields

$$\tilde{H}_{k+1}^{-1} = \tilde{H}_{k}^{-1} + \frac{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}}{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}y_{k}}$$

Yes, only y and s swapped places.

## Rank 1 – Iterative Inverse Hessian Approximation

Applying SMW to the rank one update

$$ilde{H}_{k+1} = ilde{H}_k + rac{\left(y_k - ilde{H}_k s_k
ight) \left(y_k - ilde{H}_k s_k
ight)^ op}{\left(y_k - ilde{H}_k s_k
ight)^ op s_k}$$

yields

$$\tilde{H}_{k+1}^{-1} = \tilde{H}_{k}^{-1} + \frac{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}}{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}y_{k}}$$

Yes, only y and s swapped places.

This allows us to avoid solving  $\tilde{H}_k p_k = -\nabla f_k$  for  $p_k$  in every iteration.

Masochists may study details of the proof, e.g., in "On the derivation of quasi-Newton formulas for optimization in function spaces" by Vuchkov et al. Journal of Numerical Functional Analysis and Optimization, 2021

## Rank One Update V2

Algorithm 11 Rank 1 update v1 **Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point 1:  $k \leftarrow 0$ ,  $\alpha_{init} \leftarrow 1$ .  $H_0 \leftarrow I$ 2: while  $\|\nabla f_k\|_{\infty} > \varepsilon$  do 3: Compute  $\nabla f_k = \nabla f(x_k)$ 4:  $p_k \leftarrow -\tilde{H}_{l_k}^{-1} \nabla f_k$ 5:  $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ 6:  $x_{k+1} \leftarrow x_k + \alpha p_k$ 7:  $s \leftarrow x_{k+1} - x_k$ 8:  $\mathbf{v} \leftarrow \nabla f_{k+1} - \nabla f_k$  $ilde{H}_{k+1}^{-1} \leftarrow ilde{H}_{k}^{-1} + rac{\left(s - ilde{H}_{k}^{-1} y
ight) \left(s - ilde{H}_{k}^{-1} y
ight)^{ op}}{\left(s - ilde{H}_{k}^{-1} y
ight)^{ op} y}$ 9:  $k \leftarrow k + 1$ 10:

11: end while

## BFGS

Applying SMW to the BFGS Hessian update

$$\tilde{H}_{k+1} = \tilde{H}_k - \frac{\left(\tilde{H}_k s_k\right) \left(\tilde{H}_k s_k\right)^\top}{s_k^\top \tilde{H}_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k}$$

## BFGS

Applying SMW to the BFGS Hessian update

$$\tilde{H}_{k+1} = \tilde{H}_k - \frac{\left(\tilde{H}_k s_k\right) \left(\tilde{H}_k s_k\right)^\top}{s_k^\top \tilde{H}_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k}$$

yields

$$\tilde{H}_{k+1}^{-1} = \left(I - \frac{s_k y_k^\top}{s_k^\top y_k}\right) \tilde{H}_k^{-1} \left(I - \frac{y_k s_k^\top}{s_k^\top y_k}\right) + \frac{s_k s_k^\top}{s_k^\top y_k}$$

We avoid solving the linear system for  $p_k$ .

For a proof see, e.g., in "On the derivation of quasi-Newton formulas for optimization in function spaces" by Vuchkov et al. Journal of Numerical Functional Analysis and Optimization, 2021

## BFGS V2

#### Algorithm 12 BFGS v2

**Input:**  $x_0$  starting point,  $\varepsilon > 0$ **Output:**  $x^*$  approximation to a stationary point 1:  $k \leftarrow 0$ .  $\alpha_{\text{init}} \leftarrow 1$ .  $\tilde{H}_0 \leftarrow I$ 2: while  $\|\nabla f_k\|_{\infty} > \varepsilon$  do Compute  $\nabla f_k = \nabla f(x_k)$ 3: 4:  $p_k \leftarrow -\tilde{H}_k^{-1} \nabla f_k$ 5:  $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ 6:  $x_{k+1} \leftarrow x_k + \alpha p_k$ 7:  $s \leftarrow x_{k+1} - x_k$ 8:  $\mathbf{v} \leftarrow \nabla f_{k+1} - \nabla f_k$  $\tilde{H}_{k+1}^{-1} \leftarrow \left(I - \frac{sy^{\top}}{s^{\top}y}\right) \tilde{H}_{k}^{-1} \left(I - \frac{ys^{\top}}{s^{\top}y}\right) + \frac{ss^{\top}}{s^{\top}y}$ 9:  $k \leftarrow k + 1$ 10: 11: end while

Still, we must drag the quadratic size matrix  $\tilde{H}_{k+1}^{-1}$  along.

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS.

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS. Observe that  $\tilde{H}_k^{-1}$  is determined completely by  $H_0^{-1}$  and the two sequences  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS. Observe that  $\tilde{H}_k^{-1}$  is determined completely by  $H_0^{-1}$  and the two sequences  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

So, the matrix  $\tilde{H}_k^{-1}$  does not have to be stored if the algorithm remembers the values  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

Note that this would be more space efficient for k < n/2.

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS. Observe that  $\tilde{H}_k^{-1}$  is determined completely by  $H_0^{-1}$  and the two sequences  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

So, the matrix  $\tilde{H}_k^{-1}$  does not have to be stored if the algorithm remembers the values  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

Note that this would be more space efficient for k < n/2.

However, we may go further and observe that typically only a few, say *m*, past values of *s* and *y* are sufficient for a good approximation of  $\tilde{H}_k^{-1}$  when we set  $\tilde{H}_{k-m-1}^{-1} = I$ .

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS. Observe that  $\tilde{H}_k^{-1}$  is determined completely by  $H_0^{-1}$  and the two sequences  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

So, the matrix  $\tilde{H}_k^{-1}$  does not have to be stored if the algorithm remembers the values  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

Note that this would be more space efficient for k < n/2.

However, we may go further and observe that typically only a few, say *m*, past values of *s* and *y* are sufficient for a good approximation of  $\tilde{H}_k^{-1}$  when we set  $\tilde{H}_{k-m-1}^{-1} = I$ .

This is the basic idea behind limited-memory BFGS, which stores only the running window  $s_{k-m}, \ldots, s_k$  and  $y_{k-m}, \ldots, y_k$  and computes  $\tilde{H}_k^{-1} \nabla f_k$  using these values.

Let us denote by  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$  the values of the variables s and y, resp., during the iterations  $1, \ldots, k$  of BFGS. Observe that  $\tilde{H}_k^{-1}$  is determined completely by  $H_0^{-1}$  and the two sequences  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

So, the matrix  $\tilde{H}_k^{-1}$  does not have to be stored if the algorithm remembers the values  $s_0, \ldots, s_k$  and  $y_0, \ldots, y_k$ .

Note that this would be more space efficient for k < n/2.

However, we may go further and observe that typically only a few, say *m*, past values of *s* and *y* are sufficient for a good approximation of  $\tilde{H}_k^{-1}$  when we set  $\tilde{H}_{k-m-1}^{-1} = I$ .

This is the basic idea behind limited-memory BFGS, which stores only the running window  $s_{k-m}, \ldots, s_k$  and  $y_{k-m}, \ldots, y_k$  and computes  $\tilde{H}_k^{-1} \nabla f_k$  using these values.

The space complexity becomes nm, which is beneficial when n is large.

# Another View on BFGS (Optional)

We search for  $\tilde{H}_{k+1}^{-1}$  where  $\tilde{H}_{k+1}$  satisfies  $\tilde{H}_{k+1}s_k = y_k$ . Search for a solution  $\tilde{V}$  for  $\tilde{V}y_k = s_k$ .

The idea is to use  $\tilde{V}$  close to  $\tilde{H}_k^{-1}$  (in some sense):

$$\min_{ ilde{H}} \left\| ilde{V} - ilde{H}_k^{-1} 
ight\|$$
  
subject to  $ilde{V} = ilde{V}^ op, \quad ilde{V} y_k = s_k$ 

Here the norm is weighted Frobenius norm:

$$\|A\| \equiv \left\| W^{1/2} A W^{1/2} \right\|_F,$$

where  $\|\cdot\|_F$  is defined by  $\|C\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n c_{ij}^2$ . The weight W can be chosen as any matrix satisfying the relation  $Wy_k = s_k$ .

BFGS is obtained with  $W = \overline{G}_k^{-1}$  where  $\overline{G}_k$  is the average Hessian defined by  $\overline{G}_k = \left[\int_0^1 \nabla^2 f(x_k + \tau \alpha_k p_k) d\tau\right]$ 

Solving this gives precisely the BFGS formula for  $\tilde{H}_{k+1}^{-1}$ .

## Global Convergence of Line Search

Denote by  $\theta_k$  the angle between  $p_k$  and  $-\nabla f_k$ , i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth for some L > 0 if

$$\|
abla f(x) - 
abla f( ilde{x})\| \le L \|x - ilde{x}\|, \quad \text{ for all } x, ilde{x} \in \mathbb{R}^n$$

#### Theorem 15 (Zoutendijk)

Consider  $x_{k+1} = x_k + \alpha_k p_k$ , where  $p_k$  is a descent direction and  $\alpha_k$  satisfies the strong Wolfe conditions. Suppose that f is bounded below, continuously differentiable, and L-smooth. Then

$$\sum_{k\geq 0}\cos^2\theta_k\,\|\nabla f_k\|^2<\infty.$$

## Global Convergence of Quasi-Newton's Method

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the approximations to the Hessians  $\tilde{H}_k$  are positive definite with a uniformly bounded condition number:

$$\left|\left| ilde{H}_k
ight|\right|\left|\left| ilde{H}_k^{-1}
ight|
ight|\leq M$$
 for all  $k$ 

## Global Convergence of Quasi-Newton's Method

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the approximations to the Hessians  $\tilde{H}_k$  are positive definite with a uniformly bounded condition number:

$$\left|\left|\tilde{H}_{k}\right|\right|\left|\left|\tilde{H}_{k}^{-1}\right|\right| \leq M$$
 for all  $k$ 

Then  $\theta_k$  between  $p_k = -\tilde{H}_k^{-1} \nabla f_k$  and  $-\nabla f_k$  and satisfies

 $\cos \theta_k \ge 1/M$
## Global Convergence of Quasi-Newton's Method

Assume that all  $\alpha_k$  satisfy strong Wolfe conditions.

Assume that the approximations to the Hessians  $\tilde{H}_k$  are positive definite with a uniformly bounded condition number:

$$\left|\left|\widetilde{H}_{k}\right|\right|\left|\left|\widetilde{H}_{k}^{-1}\right|\right| \leq M$$
 for all  $k$ 

Then  $\theta_k$  between  $p_k = -\tilde{H}_k^{-1} \nabla f_k$  and  $-\nabla f_k$  and satisfies

 $\cos \theta_k \ge 1/M$ 

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\frac{1}{M^{2}} \sum_{k \ge 0} \|\nabla f_{k}\|^{2} \le \sum_{k \ge 0} \cos^{2} \theta_{k} \|\nabla f_{k}\|^{2} < \infty$$

which implies that  $\lim_{k\to\infty} ||\nabla f_k|| = 0$ .

# Behavior of BFGS

▶ It may happen that  $\tilde{H}_k$  becomes a poor approximation of the Hessian  $H_k$ . If, e.g.,  $y_k^{\top}$  is tiny, then  $\tilde{H}_{k+1}$  will be huge.

However, it has been proven experimentally that if  $\tilde{H}_k$  wrongly estimates the curvature of f and this estimate slows down the iteration, then the approximation will tend to correct the bad Hessian approximations.

The above self-correction works only if an appropriate line search is performed (strong Wolfe conditions).

# Behavior of BFGS

▶ It may happen that  $\tilde{H}_k$  becomes a poor approximation of the Hessian  $H_k$ . If, e.g.,  $y_k^{\top}$  is tiny, then  $\tilde{H}_{k+1}$  will be huge.

However, it has been proven experimentally that if  $\tilde{H}_k$  wrongly estimates the curvature of f and this estimate slows down the iteration, then the approximation will tend to correct the bad Hessian approximations.

The above self-correction works only if an appropriate line search is performed (strong Wolfe conditions).

There are more sophisticated ways of setting the initial Hessian approximation H<sub>0</sub>.

See Numerical Optimization, Nocedal & Wright, page 201.

Each iteration is performed for O(n<sup>2</sup>) operations as opposed to O(n<sup>3</sup>) for methods involving solutions of linear systems.

- Each iteration is performed for O(n<sup>2</sup>) operations as opposed to O(n<sup>3</sup>) for methods involving solutions of linear systems.
- There is even a memory-limited variant (L-BFGS) that uses only information from past *m* steps, and its single iteration complexity is O(*mn*).

- Each iteration is performed for O(n<sup>2</sup>) operations as opposed to O(n<sup>3</sup>) for methods involving solutions of linear systems.
- There is even a memory-limited variant (L-BFGS) that uses only information from past *m* steps, and its single iteration complexity is O(*mn*).
- Compared with Newton's method, no second derivatives are computed.

- Each iteration is performed for O(n<sup>2</sup>) operations as opposed to O(n<sup>3</sup>) for methods involving solutions of linear systems.
- There is even a memory-limited variant (L-BFGS) that uses only information from past *m* steps, and its single iteration complexity is O(*mn*).
- Compared with Newton's method, no second derivatives are computed.
- Local superlinear convergence can be proved under specific conditions.

Compare with local quadratic convergence of Newton's method and linear convergence of gradient descent.



Here  $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$ 



Rosenbrock: 
$$f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$$



Rosenbrock:

$$f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$$



# Computational Complexity

Algorithm	Computational Complexity
Steepest Descent	O(n) per iteration
Newton's Method	$O(n^3)$ to compute Hessian and solve system
BFGS	$O(n^2)$ to update Hessian approximation

Table: Summary of the computational complexity for each optimization algorithm.

- Steepest Descent: Simple but often slow, requiring many iterations.
- Newton's Method: Fast convergence but expensive per iteration.
- BFGS: Quasi-Newton, no Hessian needed, good speed and iteration count balance.

This slide was a test of ChatGPT 4.0. The prompt was something like "Give me a beamer slide with complexity analysis of gradient descent, Newton's method, BFGS." The algorithm preferred steepest descent though.

# **Constrained Optimization**

#### Constrained Optimization Problem

Recall that the constrained optimization problem is

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i = 1, \dots, n_g \\ & h_j(x) = 0 \quad j = 1, \dots, n_h \end{array}$$

 $x^*$  is now a constrained minimizer if

$$f(x^*) \leq f(x)$$
 for all  $x \in \mathcal{F}$ 

where  ${\cal F}$  is the feasibility region

$$\mathcal{F} = \{x \mid g_i(x) \leq 0, h_j(x) = 0, i = 1, \dots, n_g, j = 1, \dots, n_h\}$$

Thus, to find a constrained minimizer, we have to inspect unconstrained minima of f inside of  $\mathcal{F}$  and points along the boundary of  $\mathcal{F}$ .

#### COP - Example



# Equality Constraints

Let us restrict our problem only to the equality constraints:

minimize f(x)by varying xsubject to  $h_j(x) = 0$   $j = 1, ..., n_h$ 

Assume that f and  $h_i$  have continuous second derivatives.

Now, we try to imitate the theory from the unconstrained case and characterize minima using gradients.

This time, we must consider the gradients of f and  $h_i$ .

### Unconstrained Minimizer

Consider the first-order Taylor approximation of f at x

 $f(x+p) \approx f(x) + \nabla f(x)^{\top} p$ 

#### Unconstrained Minimizer

Consider the first-order Taylor approximation of f at x

$$f(x+p) \approx f(x) + \nabla f(x)^{\top} p$$

Note that if  $x^*$  is an unconstrained minimizer of f, then

 $f(x^* + p) \ge f(x^*)$ 

for all p small enough.

#### Unconstrained Minimizer

Consider the first-order Taylor approximation of f at x

$$f(x+p) \approx f(x) + \nabla f(x)^{\top} p$$

Note that if  $x^*$  is an unconstrained minimizer of f, then

 $f(x^* + p) \ge f(x^*)$ 

for all p small enough.

Together with the Taylor approximation, we obtain

$$f(x^*) + 
abla f(x^*)^{ op} p \geq f(x^*)$$

and hence

$$\nabla f(x^*)^\top p \geq 0$$



The hyperplane defined by  $\nabla f^{\top} p = 0$  contains directions p of zero variation in f.

In the unconstrained case,  $x^*$  is minimizer only if  $\nabla f(x^*) = 0$  because otherwise there would be a direction p satisfying  $\nabla f(x^*)p < 0$ , a *decrease direction*.

In COP, p is a decrease direction in  $x \in \mathcal{F}$  if  $\nabla f(x)^{\top} p < 0$  and if p is a *feasible direction*!

That is, points into the feasible region.

In COP, p is a decrease direction in  $x \in \mathcal{F}$  if  $\nabla f(x)^{\top} p < 0$  and if p is a *feasible direction*!

That is, points into the feasible region.

How do we characterize feasible directions?

In COP, p is a decrease direction in  $x \in \mathcal{F}$  if  $\nabla f(x)^{\top} p < 0$  and if p is a *feasible direction*!

That is, points into the feasible region.

How do we characterize feasible directions?

Consider Taylor approximation of  $h_j$  for all j:

$$h_j(x+p) \approx h_j(x) + \nabla h_j(x)^\top p$$

In COP, p is a decrease direction in  $x \in \mathcal{F}$  if  $\nabla f(x)^{\top} p < 0$  and if p is a *feasible direction*!

That is, points into the feasible region.

How do we characterize feasible directions?

Consider Taylor approximation of  $h_j$  for all j:

$$h_j(x+p) \approx h_j(x) + \nabla h_j(x)^\top p$$

Assuming  $x \in \mathcal{F}$ , we have  $h_j(x) = 0$  for all j and thus

 $h_j(x+p) \approx \nabla h_j(x)^\top p$ 

In COP, p is a decrease direction in  $x \in \mathcal{F}$  if  $\nabla f(x)^{\top} p < 0$  and if p is a *feasible direction*!

That is, points into the feasible region.

How do we characterize feasible directions?

Consider Taylor approximation of  $h_j$  for all j:

$$h_j(x+p) \approx h_j(x) + \nabla h_j(x)^\top p$$

Assuming  $x \in \mathcal{F}$ , we have  $h_j(x) = 0$  for all j and thus

$$h_j(x+p) \approx \nabla h_j(x)^\top p$$

As p is a feasible direction iff  $h_j(x + p) = 0$ , we obtain that

p is a *feasible direction* iff  $\nabla h_j(x)^\top p = 0$  for all j

Feasible Points and Directions



Here, the only feasible direction at x is p = 0.

## Feasible Points and Directions



Here the feasible directions at  $x^*$  point along the red line, i.e.,

$$abla h_1(x^*) p = 0$$
  $abla h_2(x^*) p = 0$ 

Consider a direction p. Observe that

If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.

Consider a direction p. Observe that

- If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.
- If  $\nabla h_j(x)^\top p = 0$  for all j and
  - ▶  $\nabla f(x)^\top p > 0$ , then moving a short step in the direction p increases f and stays in  $\mathcal{F}$ .

Consider a direction p. Observe that

- If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.
- If  $\nabla h_j(x)^\top p = 0$  for all j and
  - ▶  $\nabla f(x)^\top p > 0$ , then moving a short step in the direction p increases f and stays in  $\mathcal{F}$ .
  - ∇f(x)<sup>T</sup>p < 0, then moving a short step in the direction p decreases f and stays in F.</p>

Consider a direction p. Observe that

- If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.
- If  $\nabla h_j(x)^\top p = 0$  for all j and
  - $\nabla f(x)^{\top} p > 0$ , then moving a short step in the direction p increases f and stays in  $\mathcal{F}$ .
  - ∇f(x)<sup>T</sup>p < 0, then moving a short step in the direction p decreases f and stays in F.</p>
  - ∇f(x)<sup>T</sup>p = 0, then moving a short step in the direction p does not change f and stays F.

Consider a direction p. Observe that

If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.

• If 
$$\nabla h_j(x)^\top p = 0$$
 for all  $j$  and

- ∇f(x)<sup>T</sup>p > 0, then moving a short step in the direction p increases f and stays in F.
- ∇f(x)<sup>T</sup>p < 0, then moving a short step in the direction p decreases f and stays in F.</p>
- ∇f(x)<sup>T</sup>p = 0, then moving a short step in the direction p does not change f and stays F.

To be a minimizer,  $x^*$  must be feasible and every direction satisfying  $\nabla h_j(x^*)^\top p = 0$  for all j must also satisfy  $\nabla f(x^*)^\top p \ge 0$ .

Consider a direction p. Observe that

If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.

• If 
$$\nabla h_j(x)^\top p = 0$$
 for all  $j$  and

- ∇f(x)<sup>T</sup>p > 0, then moving a short step in the direction p increases f and stays in F.
- ∇f(x)<sup>T</sup>p < 0, then moving a short step in the direction p decreases f and stays in F.</p>
- ∇f(x)<sup>T</sup>p = 0, then moving a short step in the direction p does not change f and stays F.

To be a minimizer,  $x^*$  must be feasible and every direction satisfying  $\nabla h_j(x^*)^\top p = 0$  for all j must also satisfy  $\nabla f(x^*)^\top p \ge 0$ .

Note that if p is a feasible direction, then -p is also, and thus  $\nabla f(x^*)^{\top}(-p) \ge 0$ . So finally,

Consider a direction p. Observe that

If ∇h<sub>j</sub>(x)<sup>T</sup>p ≠ 0, then moving a short step in the direction p violates the constraint h<sub>j</sub>(x) = 0.

• If 
$$\nabla h_j(x)^\top p = 0$$
 for all  $j$  and

- ∇f(x)<sup>T</sup>p > 0, then moving a short step in the direction p increases f and stays in F.
- ∇f(x)<sup>T</sup>p < 0, then moving a short step in the direction p decreases f and stays in F.</p>
- ∇f(x)<sup>T</sup>p = 0, then moving a short step in the direction p does not change f and stays F.

To be a minimizer,  $x^*$  must be feasible and every direction satisfying  $\nabla h_j(x^*)^\top p = 0$  for all j must also satisfy  $\nabla f(x^*)^\top p \ge 0$ .

Note that if p is a feasible direction, then -p is also, and thus  $\nabla f(x^*)^{\top}(-p) \ge 0$ . So finally,

If  $x^*$  is a *constrained minimizer*, then  $\nabla f(x^*)^\top p = 0$  for all p satisfying  $(\forall j : \nabla h_j(x^*)^\top p = 0)$ 

# Lagrange Multipliers



Left: f increases along p. Right: f does not change along p.

# Lagrange Multipliers



Left: f increases along p. Right: f does not change along p.

Observe that at an optimum,  $\nabla f$  lies in the space spanned by the gradients of constraint functions.
# Lagrange Multipliers



Left: f increases along p. Right: f does not change along p.

Observe that at an optimum,  $\nabla f$  lies in the space spanned by the gradients of constraint functions.

There are Lagrange multipliers  $\lambda_1, \lambda_2$  satisfying

$$\nabla f(x^*) = -(\lambda_1 \nabla h_1 + \lambda_2 \nabla h_2)$$

The minus sign is arbitrary for equality constraints but will be significant when dealing with inequality constraints.

#### Lagrange Multipliers

We know that if  $x^*$  is a constrained minimizer, then.

 $\nabla f(x^*)^\top p = 0$  for all p satisfying  $(\forall j : \nabla h_j(x^*)^\top p = 0)$ 

#### Lagrange Multipliers

We know that if  $x^*$  is a constrained minimizer, then.

$$\nabla f(x^*)^\top p = 0$$
 for all  $p$  satisfying  $(\forall j : \nabla h_j(x^*)^\top p = 0)$ 

But then, from the geometry of the problem, we obtain

#### Theorem 16

Consider the COP with only equality constraints and f and all  $h_j$  twice continuously differentiable.

Assume that  $x^*$  is a constrained minimizer and that  $x^*$  is regular, which means that  $\nabla h_j(x^*)$  are linearly independent. Then there are  $\lambda_1, \ldots, \lambda_{n_h} \in \mathbb{R}$  satisfying

$$abla f(x^*) = -\sum_{j=1}^{n_h} \lambda_j 
abla h_j(x^*)$$

The coefficients  $\lambda_1, \ldots, \lambda_{n_h}$  are called *Lagrange multipliers*.

Try to transform the constrained problem into an unconstrained one by moving the constraints  $h_i(x) = 0$  into the objective.

Try to transform the constrained problem into an unconstrained one by moving the constraints  $h_i(x) = 0$  into the objective.

Consider Lagrangian function  $\mathcal{L}: \mathbb{R}^n \times \mathbb{R}^{n_h} \to \mathbb{R}$  defined by

 $\mathcal{L}(x,\lambda) = f(x) + \lambda^{\top} h(x)$  here  $h(x) = (h_1(x), \dots, h_{n_h}(x))^{\top}$ 

Try to transform the constrained problem into an unconstrained one by moving the constraints  $h_i(x) = 0$  into the objective.

Consider Lagrangian function  $\mathcal{L}: \mathbb{R}^n \times \mathbb{R}^{n_h} \to \mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x)$$
 here  $h(x) = (h_1(x), \dots, h_{n_h}(x))^{ op}$ 

Note that the stationary point of  $\mathcal L$  gives us the Lagrange multipliers:

$$abla_{\mathbf{x}}\mathcal{L} = 
abla f(\mathbf{x}) + \sum_{j=1}^{n_h} \lambda_j 
abla h_j(\mathbf{x})$$
 $abla_{\lambda}\mathcal{L} = h(\mathbf{x})$ 

Try to transform the constrained problem into an unconstrained one by moving the constraints  $h_i(x) = 0$  into the objective.

Consider Lagrangian function  $\mathcal{L}: \mathbb{R}^n \times \mathbb{R}^{n_h} \to \mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x)$$
 here  $h(x) = (h_1(x), \dots, h_{n_h}(x))^{ op}$ 

Note that the stationary point of  $\mathcal{L}$  gives us the Lagrange multipliers:

$$abla_{\mathbf{x}}\mathcal{L} = 
abla f(\mathbf{x}) + \sum_{j=1}^{n_h} \lambda_j 
abla h_j(\mathbf{x})$$
 $abla_{\mathbf{\lambda}}\mathcal{L} = h(\mathbf{x})$ 

Now putting  $\nabla \mathcal{L}(x) = 0$ , we obtain precisely the above properties of the constrained minimizer:

$$h(x) = 0$$
 and  $abla f(x) = -\sum_{j=1}^{n_h} \lambda_j 
abla h_j(x)$ 

So we can now use methods for searching stationary points. This will lead to the Lagrange-Newton method.

$$\begin{array}{ll} \underset{x_{1},x_{2}}{\text{minimize}} & f\left(x_{1},x_{2}\right) = x_{1} + 2x_{2} \\ \text{subject to} & h\left(x_{1},x_{2}\right) = \frac{1}{4}x_{1}^{2} + x_{2}^{2} - 1 = 0 \end{array}$$

The Lagrangian function

$$\mathcal{L}(x_1, x_2, \lambda) = x_1 + 2x_2 + \lambda \left(\frac{1}{4}x_1^2 + x_2^2 - 1\right)$$

$$\begin{array}{ll} \underset{x_{1},x_{2}}{\text{minimize}} & f\left(x_{1},x_{2}\right) = x_{1} + 2x_{2} \\ \text{subject to} & h\left(x_{1},x_{2}\right) = \frac{1}{4}x_{1}^{2} + x_{2}^{2} - 1 = 0 \end{array}$$

The Lagrangian function

$$\mathcal{L}(x_1, x_2, \lambda) = x_1 + 2x_2 + \lambda \left(\frac{1}{4}x_1^2 + x_2^2 - 1\right)$$

Differentiating this to get the first-order optimality conditions,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_1} &= 1 + \frac{1}{2}\lambda x_1 = 0 \qquad \frac{\partial \mathcal{L}}{\partial x_2} = 2 + 2\lambda x_2 = 0\\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \frac{1}{4}x_1^2 + x_2^2 - 1 = 0. \end{aligned}$$

$$\begin{array}{ll} \underset{x_{1},x_{2}}{\text{minimize}} & f\left(x_{1},x_{2}\right) = x_{1} + 2x_{2} \\ \text{subject to} & h\left(x_{1},x_{2}\right) = \frac{1}{4}x_{1}^{2} + x_{2}^{2} - 1 = 0 \end{array}$$

The Lagrangian function

$$\mathcal{L}(x_1, x_2, \lambda) = x_1 + 2x_2 + \lambda \left(\frac{1}{4}x_1^2 + x_2^2 - 1\right)$$

Differentiating this to get the first-order optimality conditions,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_1} &= 1 + \frac{1}{2}\lambda x_1 = 0 \qquad \frac{\partial \mathcal{L}}{\partial x_2} = 2 + 2\lambda x_2 = 0\\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \frac{1}{4}x_1^2 + x_2^2 - 1 = 0. \end{aligned}$$

Solving these three equations for the three unknowns  $(x_1, x_2, \lambda)$ , we obtain two possible solutions:

$$x_A = (x_1, x_2) = (-\sqrt{2}, -\sqrt{2}/2), \quad \lambda_A = \sqrt{2}$$
  
 $x_B = (x_1, x_2) = (\sqrt{2}, \sqrt{2}/2), \quad \lambda_A = -\sqrt{2}$ 



As in the unconstrained case, the first-order conditions characterize any "stable" point (minimum, maximum, saddle).

As in the unconstrained case, the first-order conditions characterize any "stable" point (minimum, maximum, saddle).

Consider Lagrangian Hessian:

$$H(x,\lambda) = H_f(x) + \sum_{j=1}^{n_h} \lambda_j H_{h_j}(x)$$

Here  $H_f$  is the Hessian of f, and each  $H_{h_j}$  is the Hessian of  $h_j$ . Note that Lagrangian Hessian is NOT the Hessian of the Lagrangian!

As in the unconstrained case, the first-order conditions characterize any "stable" point (minimum, maximum, saddle).

Consider Lagrangian Hessian:

$$H(x,\lambda) = H_f(x) + \sum_{j=1}^{n_h} \lambda_j H_{h_j}(x)$$

Here  $H_f$  is the Hessian of f, and each  $H_{h_j}$  is the Hessian of  $h_j$ . Note that Lagrangian Hessian is NOT the Hessian of the Lagrangian!

The second-order sufficient conditions are as follows: Assume  $x^*$  is regular and feasible. Also, assume that there is  $\lambda^*$  s.t.

$$\nabla f(x^*) = \sum_{j=1}^{n_h} -\lambda_j^* \nabla h_j(x^*)$$

As in the unconstrained case, the first-order conditions characterize any "stable" point (minimum, maximum, saddle).

Consider Lagrangian Hessian:

$$H(x,\lambda) = H_f(x) + \sum_{j=1}^{n_h} \lambda_j H_{h_j}(x)$$

Here  $H_f$  is the Hessian of f, and each  $H_{h_j}$  is the Hessian of  $h_j$ . Note that Lagrangian Hessian is NOT the Hessian of the Lagrangian!

The second-order sufficient conditions are as follows: Assume  $x^*$  is regular and feasible. Also, assume that there is  $\lambda^*$  s.t.

$$\nabla f(x^*) = \sum_{j=1}^{n_h} -\lambda_j^* \nabla h_j(x^*)$$

and that

 $p^{\top}H(x^*,\lambda^*)p > 0$  for all p satisfying  $(\forall j : \nabla h_j(x^*)^{\top}p = 0)$ Then,  $x^*$  is a constrained minimizer of f.

# Inequality Constraints

Recall that the constrained optimization problem is

 $\begin{array}{ll} \mbox{minimize} & f(x) \\ \mbox{by varying} & x \\ \mbox{subject to} & g_i(x) \leq 0 \quad i=1,\ldots,n_g \\ & h_j(x)=0 \quad j=1,\ldots,n_h \end{array}$ 

# Inequality Constraints

Recall that the constrained optimization problem is

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i = 1, \dots, n_g \\ & h_j(x) = 0 \quad j = 1, \dots, n_h \end{array}$$

Lagrange multipliers and the Lagrangian function can be extended to deal with inequality constraints.

The resulting necessary conditions for constrained minima are called Karush-Tucker-Kuhn (KKT) conditions.

In this course, Lagrange methods are considered only for equality-constrained problems. So, we omit further discussion of KKT.

# $\underset{\text{Penalty Methods}}{\text{Constrained Optimization}}$

**The idea:** Transform a constrained problem into an unconstrained one by adding a penalty to the objective function when constraints are violated or close to being violated.

**The idea:** Transform a constrained problem into an unconstrained one by adding a penalty to the objective function when constraints are violated or close to being violated.

Assuming an objective function f, the penalized objective is of the form

$$\hat{f}(x) = f(x) + \mu \pi(x)$$

Here,  $\mu$  is a fixed constant determining how strong the penalty should be, and  $\pi$  is the penalty function.

**The idea:** Transform a constrained problem into an unconstrained one by adding a penalty to the objective function when constraints are violated or close to being violated.

Assuming an objective function f, the penalized objective is of the form

 $\hat{f}(x) = f(x) + \mu \pi(x)$ 

Here,  $\mu$  is a fixed constant determining how strong the penalty should be, and  $\pi$  is the penalty function.

Now we may apply the unconstrained optimization methods (e.g., L-BFGS) to  $\hat{f}$  and obtain an approximation of a minimizer of f.

**The idea:** Transform a constrained problem into an unconstrained one by adding a penalty to the objective function when constraints are violated or close to being violated.

Assuming an objective function f, the penalized objective is of the form

 $\hat{f}(x) = f(x) + \mu \pi(x)$ 

Here,  $\mu$  is a fixed constant determining how strong the penalty should be, and  $\pi$  is the penalty function.

Now we may apply the unconstrained optimization methods (e.g., L-BFGS) to  $\hat{f}$  and obtain an approximation of a minimizer of f.

There are two kinds of penalty methods:

- exterior penalizing infeasible x
- interior penalizing x close to being infeasible

## Interior vs Exterior Penalty



Exterior Penalty Methods - Quadratic Penalty

Consider equality-constrained problems:

minimize f(x)by varying xsubject to  $h_j(x) = 0$   $j = 1, ..., n_h$  Exterior Penalty Methods - Quadratic Penalty

Consider equality-constrained problems:

minimize f(x)by varying xsubject to  $h_j(x) = 0$   $j = 1, ..., n_h$ 

Consider *quadratic penalty*:

$$\hat{f}(x;\mu) = f(x) + rac{\mu}{2} \sum_{j=1}^{n_h} h_j(x)^2$$

If f is continuously differentiable,  $\hat{f}$  is as well (w.r.t. x).

#### Quadratic Penalty



The true solution would be recovered for  $\mu = \infty$ .

## Quadratic Penalty



The true solution would be recovered for  $\mu = \infty$ .

However, large  $\mu$  means large condition number of the Hessian of  $\hat{f}$ Intuitively, large curvature of  $\hat{f}$ , not good for optimization.

Need to choose  $\mu$  carefully, possibly iteratively.

#### Algorithm 13 Exterior Penalty Method

- 1: Choose starting point  $x_0$
- 2: Choose an initial penalty parameter  $\mu_0$
- 3: Choose a penalty increase factor ho>1
- 4:  $k \leftarrow 0$
- 5: repeat
- 6:  $x_{k+1} \leftarrow x$  minimizing  $\hat{f}(x; \mu_k)$
- 7:  $\mu_{k+1} \leftarrow \rho \mu_k$
- 8:  $k \leftarrow k+1$
- 9: until convergence

## Convergence of Quadratic Penalty Method

#### Theorem 17

Assume that f and all  $h_j$  have continuous second derivatives. Suppose that each  $x_k$  is the exact global minimizer of  $\hat{f}(x; \mu_k)$  and that  $\lim_{k\to\infty} \mu_k = \infty$ . Then, every limit point  $x^*$  of the sequence  $\{x_k\}$  solves the constrained optimization problem.

In practice, inexact methods are used to minimize  $\hat{f}(x; \mu_k)$ 

## Convergence of Quadratic Penalty Method

#### Theorem 17

Assume that f and all  $h_j$  have continuous second derivatives. Suppose that each  $x_k$  is the exact global minimizer of  $\hat{f}(x; \mu_k)$  and that  $\lim_{k\to\infty} \mu_k = \infty$ . Then, every limit point  $x^*$  of the sequence  $\{x_k\}$  solves the constrained optimization problem.

In practice, inexact methods are used to minimize  $\hat{f}(x; \mu_k)$ 

Let  $x^*$  be a limit point of  $x_k$  and let  $\lambda^*$  be such that  $(x^*, \lambda^*)$  satisfy the Lagrange conditions for the constrained problem.

## Convergence of Quadratic Penalty Method

#### Theorem 17

Assume that f and all  $h_j$  have continuous second derivatives. Suppose that each  $x_k$  is the exact global minimizer of  $\hat{f}(x; \mu_k)$  and that  $\lim_{k\to\infty} \mu_k = \infty$ . Then, every limit point  $x^*$  of the sequence  $\{x_k\}$  solves the constrained optimization problem.

In practice, inexact methods are used to minimize  $\hat{f}(x; \mu_k)$ 

Let  $x^*$  be a limit point of  $x_k$  and let  $\lambda^*$  be such that  $(x^*, \lambda^*)$  satisfy the Lagrange conditions for the constrained problem.

Then, for a subsequence of points  $x_k$ , which converges to  $x^*$ , we have that

$$\lim_{k\to\infty} -\mu_k h_j(x_k) = \lambda_j^*$$

## **Practical Problems**

- Small µ may result in so weak penalty that f unbounded below results in f unbounded as well
- As  $\mu = \infty$  is impossible, the solution is always slightly infeasible
- Growing curvature of  $\hat{f}$  as  $\mu$  grows makes the Hessian of  $\hat{f}$  almost singular

$$\hat{f}(x;\mu) = x_1 + 2x_2 + rac{\mu}{2} \left(rac{1}{4}x_1^2 + x_2^2 - 1
ight)^2$$



 $\mu = 0.5$ 

 $\mu = 3.0$ 

$$\hat{f}(x;\mu) = x_1 + 2x_2 + rac{\mu}{2} \left(rac{1}{4}x_1^2 + x_2^2 - 1
ight)^2$$



Quadratic Penalty for Inequality Constraints



Minimizer approached from the infeasible side.

# Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 + \frac{\mu}{2} \max\left(0, \frac{1}{4}x_1^2 + x_2^2 - 1\right)^2$$


Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 + \frac{\mu}{2} \max\left(0, \frac{1}{4}x_1^2 + x_2^2 - 1\right)^2$$



 $\mu = 10.0$ 

## Augmented Lagrangian (Optional)

We may augment the Lagrangian  $\mathcal{L} = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x)$  with penalty and optimize the augmented Lagrangian

$$\hat{f}(x; \lambda, \mu) = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x) + \frac{\mu}{2} \sum_{j=1}^{n_h} h_j(x)^2$$

### Augmented Lagrangian (Optional)

We may augment the Lagrangian  $\mathcal{L} = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x)$  with penalty and optimize the augmented Lagrangian

$$\hat{f}(x;\lambda,\mu) = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x) + \frac{\mu}{2} \sum_{j=1}^{n_h} h_j(x)^2$$

Note the relationship between optimality conditions for  ${\cal L}$  and  $\hat{f}$ 

$$abla_x \hat{f}(x;\lambda,\mu) = 
abla f(x) + \sum_{j=1}^{n_h} \left(\lambda_j + \mu h_j(x)\right) 
abla h_j(x) = 0$$

$$abla_{x}\mathcal{L}\left(x^{*},\lambda^{*}
ight)=
abla f\left(x^{*}
ight)+\sum_{j=1}^{n_{h}}\lambda_{j}^{*}
abla h_{j}\left(x^{*}
ight)=0.$$

### Augmented Lagrangian (Optional)

We may augment the Lagrangian  $\mathcal{L} = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x)$  with penalty and optimize the augmented Lagrangian

$$\hat{f}(x;\lambda,\mu) = f(x) + \sum_{j=1}^{n_h} \lambda_j h_j(x) + \frac{\mu}{2} \sum_{j=1}^{n_h} h_j(x)^2$$

Note the relationship between optimality conditions for  $\mathcal L$  and  $\hat f$ 

$$abla_x \hat{f}(x;\lambda,\mu) = 
abla f(x) + \sum_{j=1}^{n_h} \left(\lambda_j + \mu h_j(x)\right) 
abla h_j(x) = 0$$

$$abla_{x}\mathcal{L}\left(x^{*},\lambda^{*}\right)=
abla f\left(x^{*}\right)+\sum_{j=1}^{n_{h}}\lambda_{j}^{*}
abla h_{j}\left(x^{*}\right)=0.$$

Comparing these two conditions suggests an approximation:

$$\lambda_j^* \approx \lambda_j + \mu h_j.$$

## Augmented Lagrangian Penalty Method (Optional)

#### Inputs:

- x<sub>0</sub>: Starting point
- $\lambda_0 = 0$ : Initial Lagrange multiplier
- $\mu_0 > 0$ : Initial penalty parameter
- $\rho > 1$ : Penalty increase factor

## Outputs:

- x\*: Optimal point
- $f(x^*)$ : Corresponding function value

## Algorithm:

k = 0

#### repeat

 $x_{k+1} \leftarrow x \text{ minimizing } \hat{f}(x; \lambda_k, \mu_k)$  $\lambda_{k+1} = \lambda_k + \mu_k h(x_k)$  $\mu_{k+1} \leftarrow \rho \mu_k$  $k \leftarrow k+1$ **until** convergence

# Comparison of Quadratic and Lagrangian Penalty (Optional)

Compare

$$h_j pprox rac{1}{\mu} \left( \lambda_j^* - \lambda_j 
ight).$$

with the corresponding approximation of  $h_j$  in the quadratic penalty method is

$$h_j pprox rac{\lambda_j^*}{\mu}$$

Thus, the quadratic penalty relies solely on increasing  $\mu$ .

# Comparison of Quadratic and Lagrangian Penalty (Optional)

Compare

$$h_j pprox rac{1}{\mu} \left( \lambda_j^* - \lambda_j 
ight).$$

with the corresponding approximation of  $h_j$  in the quadratic penalty method is

$$h_j pprox rac{\lambda_j^*}{\mu}$$

Thus, the quadratic penalty relies solely on increasing  $\mu$ .

However, the augmented Lagrangian also controls the numerator via estimating  $\lambda_i$ .

If  $\lambda_j$  is close to  $\lambda_j^*$ , we may obtain a close solution for modest values of  $\mu$ .

Several variants of the Lagrangian penalty exist for inequality constraints; see Nocedal & Wright.

## Interior Penalty Methods

Always seek to maintain feasibility as opposed to the exterior methods.

Instead of adding a penalty only when constraints are violated; add a penalty as the constraint is approached from the feasible region.

## Interior Penalty Methods

Always seek to maintain feasibility as opposed to the exterior methods.

Instead of adding a penalty only when constraints are violated; add a penalty as the constraint is approached from the feasible region.

Desirable if the objective function is ill-defined outside the feasible region.

The interior methods are also referred to as *barrier methods* because the penalty function acts as a barrier preventing iterates from leaving the feasible region.

Consider inequality-constrained problems:

```
\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i=1,\ldots,n_g \end{array}
```

Consider inequality-constrained problems:

```
\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i=1,\ldots,n_g \end{array}
```

Minimize the augmented objective function.

 $\hat{f}(x;\mu) = f(x) + \mu \pi(x)$ 

Here  $\pi$  is a penalty function.

Consider inequality-constrained problems:

minimize 
$$f(x)$$
  
by varying  $x$   
subject to  $g_i(x) \le 0$   $i = 1, ..., n_g$ 

Minimize the augmented objective function.

 $\hat{f}(x;\mu) = f(x) + \mu \pi(x)$ 

Here  $\pi$  is a penalty function.

Inverse barrier

Logarithmic barrier

$$\pi(x) = \sum_{i=1}^{n_g} -\frac{1}{g_i(x)} \qquad \qquad \pi(x) = \sum_{i=1}^{n_g} -\ln(-g_i(x))$$

Algorithms based on these penalties must be prevented from evaluating infeasible points.



Inverse barrier

$$\pi(x) = \sum_{i=1}^{n_g} -\frac{1}{g_i(x)}$$

Logarithmic barrier

$$\pi(x) = \sum_{i=1}^{n_g} - \ln(-g_i(x))$$



Solve a sequence of unconstrained problems for  $\hat{f}$  with  $\mu \rightarrow 0.$ 

## Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 - \mu \ln \left(-\frac{1}{4}x_1^2 - x_2^2 + 1\right)$$



As for exterior methods, the Hessian becomes increasingly ill-conditioned as  $\mu \rightarrow 0.$ 

Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 - \mu \ln \left(-\frac{1}{4}x_1^2 - x_2^2 + 1\right)$$



As for exterior methods, the Hessian becomes increasingly ill-conditioned as  $\mu \rightarrow 0.$ 

## Comments on Interior Penalty Methods

Interior penalty methods must stay in the feasible region:

- Every unconstrained optimization must start at an initial point feasible for the constrained problem.
- The line search must check for feasibility and backtrack from steps to infeasible points.

## Comments on Interior Penalty Methods

Interior penalty methods must stay in the feasible region:

- Every unconstrained optimization must start at an initial point feasible for the constrained problem.
- The line search must check for feasibility and backtrack from steps to infeasible points.

Convergence issues:

- As  $\mu \to 0$  solutions of  $\hat{f}$  converge to solutions of the constrained problem.
- On the other hand, with  $\mu \rightarrow 0$  the Hessian of  $\hat{f}$  becomes increasingly ill-conditioned.

Various modifications exist to alleviate the problem with ill-conditioned Hessians.

These methods lead to a class of modern interior point methods.

Penalty methods penalize approximations that either leave the feasible region (exterior methods), or are close to the border of the feasible region (interior methods).

Penalty methods are simple and easy to implement.

Both exterior and interior methods lead to ill-conditioned Hessians when approaching the correct solutions to the constrained problem.

## Constrained Optimization Sequential Quadratic Programming

The quadratic optimization problem with equality constraints is to

minimize  $\frac{1}{2}x^{\top}Qx + q^{\top}x$ by varying xsubject to Ax + b = 0

The quadratic optimization problem with equality constraints is to

minimize  $\frac{1}{2}x^{\top}Qx + q^{\top}x$ by varying xsubject to Ax + b = 0

Here

- Q is a n × n symmetric matrix. For simplicity assume positive definite.
- A is a  $m \times n$  matrix. Assume full rank.



How to solve the quadratic program?

How to solve the quadratic program? Consider the Lagrangian function

$$L(x,\lambda) = \frac{1}{2}x^{\top}Qx + q^{\top}x + \lambda^{\top}(Ax + b)$$

How to solve the quadratic program? Consider the Lagrangian function

$$L(x,\lambda) = \frac{1}{2}x^{\top}Qx + q^{\top}x + \lambda^{\top}(Ax + b)$$

and its partial derivatives:

$$abla_x L(x) = Qx + q + A^\top \lambda = 0$$
  
 $abla_\lambda L(x) = Ax + b = 0$ 

How to solve the quadratic program?

Consider the Lagrangian function

$$L(x,\lambda) = \frac{1}{2}x^{\top}Qx + q^{\top}x + \lambda^{\top}(Ax + b)$$

and its partial derivatives:

$$abla_x L(x) = Qx + q + A^\top \lambda = 0$$
  
 $abla_\lambda L(x) = Ax + b = 0$ 

For Q positive definite, we know that a solution to the above system is a minimizer.

So in order to solve the quadratic program, it suffices to solve the system of linear equations.

Now consider an arbitrary  $f : \mathbb{R}^n \to \mathbb{R}$  and arbitrary constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

Consider the Lagrangian function  $\mathcal{L}:\mathbb{R}^n\times\mathbb{R}^{n_h}\to\mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x) \quad ext{here} \quad h(x) = (h_1(x),\ldots,h_{n_h}(x))^{ op}$$

Now consider an arbitrary  $f : \mathbb{R}^n \to \mathbb{R}$  and arbitrary constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

Consider the Lagrangian function  $\mathcal{L}:\mathbb{R}^n imes\mathbb{R}^{n_h} o\mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x) \quad ext{here} \quad h(x) = (h_1(x),\ldots,h_{n_h}(x))^{ op}$$

We search for the stationary point of  $\mathcal{L}$ , that is  $(x^*, \lambda^*)$  satisfying

$$\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}^*,\lambda^*) = \nabla f(\mathbf{x}^*) + \sum_{j=1}^{n_h} \lambda_j^* \nabla h_j(\mathbf{x}^*) = 0$$
$$\nabla_{\lambda}\mathcal{L}(\mathbf{x}^*,\lambda^*) = h(\mathbf{x}^*) = 0$$

These are  $n + n_h$  equations in unknowns  $(x^*, \lambda^*)$ .

Now consider an arbitrary  $f : \mathbb{R}^n \to \mathbb{R}$  and arbitrary constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

Consider the Lagrangian function  $\mathcal{L}:\mathbb{R}^n\times\mathbb{R}^{n_h}\to\mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x) \quad ext{here} \quad h(x) = (h_1(x),\ldots,h_{n_h}(x))^{ op}$$

We search for the stationary point of  $\mathcal{L}$ , that is  $(x^*, \lambda^*)$  satisfying

$$\nabla_{x}\mathcal{L}(x^{*},\lambda^{*}) = \nabla f(x^{*}) + \sum_{j=1}^{n_{h}} \lambda_{j}^{*} \nabla h_{j}(x^{*}) = 0$$
$$\nabla_{\lambda}\mathcal{L}(x^{*},\lambda^{*}) = h(x^{*}) = 0$$

These are  $n + n_h$  equations in unknowns  $(x^*, \lambda^*)$ .

From Lagrange theorem: If  $x^*$  is regular and solves the COP, then there exists  $\lambda^*$  such that  $(x^*, \lambda^*)$  solves the system of equations.

Now consider an arbitrary  $f : \mathbb{R}^n \to \mathbb{R}$  and arbitrary constraint functions  $h_j : \mathbb{R}^n \to \mathbb{R}$ .

Consider the Lagrangian function  $\mathcal{L}:\mathbb{R}^n\times\mathbb{R}^{n_h}\to\mathbb{R}$  defined by

$$\mathcal{L}(x,\lambda) = f(x) + \lambda^{ op} h(x) \quad ext{here} \quad h(x) = (h_1(x),\ldots,h_{n_h}(x))^{ op}$$

We search for the stationary point of  $\mathcal{L}$ , that is  $(x^*, \lambda^*)$  satisfying

$$\nabla_{x}\mathcal{L}(x^{*},\lambda^{*}) = \nabla f(x^{*}) + \sum_{j=1}^{n_{h}} \lambda_{j}^{*} \nabla h_{j}(x^{*}) = 0$$
$$\nabla_{\lambda}\mathcal{L}(x^{*},\lambda^{*}) = h(x^{*}) = 0$$

These are  $n + n_h$  equations in unknowns  $(x^*, \lambda^*)$ .

From Lagrange theorem: If  $x^*$  is regular and solves the COP, then there exists  $\lambda^*$  such that  $(x^*, \lambda^*)$  solves the system of equations.

We use Newton's method to solve the system of equations.

Start with some  $(x_0, \lambda_0)$  and compute  $(x_1, \lambda_1), \ldots, (x_k, \lambda_k), \ldots$ 

Start with some  $(x_0, \lambda_0)$  and compute  $(x_1, \lambda_1), \ldots, (x_k, \lambda_k), \ldots$ 

In every step we compute  $(x_{k+1}, \lambda_{k+1})$  from  $(x_k, \lambda_k)$  using Newton's step.

Start with some  $(x_0, \lambda_0)$  and compute  $(x_1, \lambda_1), \ldots, (x_k, \lambda_k), \ldots$ 

In every step we compute  $(x_{k+1}, \lambda_{k+1})$  from  $(x_k, \lambda_k)$  using Newton's step.

Consider the gradient of the Lagrangian:

$$\nabla \mathcal{L}(x_k, \lambda_k) = (\nabla_x \mathcal{L}(x_k, \lambda_k), \nabla_\lambda \mathcal{L}(x_k, \lambda_k))^\top$$
$$= (\nabla f(x_k) + \sum_{j=1}^{n_h} \lambda_{kj} \nabla h_j(x_k), \quad h(x_k))^\top \in \mathbb{R}^{n+n_h}$$

Start with some  $(x_0, \lambda_0)$  and compute  $(x_1, \lambda_1), \ldots, (x_k, \lambda_k), \ldots$ 

In every step we compute  $(x_{k+1}, \lambda_{k+1})$  from  $(x_k, \lambda_k)$  using Newton's step.

Consider the gradient of the Lagrangian:

$$\nabla \mathcal{L}(x_k, \lambda_k) = (\nabla_x \mathcal{L}(x_k, \lambda_k), \nabla_\lambda \mathcal{L}(x_k, \lambda_k))^\top$$
$$= (\nabla f(x_k) + \sum_{j=1}^{n_h} \lambda_{kj} \nabla h_j(x_k), \quad h(x_k))^\top \in \mathbb{R}^{n+n_h}$$

and the Hessian matrix of the (complete) Lagrangian

$$abla^2 \mathcal{L}(x_k,\lambda_k) \in \mathbb{R}^{n+n_h} imes \mathbb{R}^{n+n_h}$$

We compute this Hessian in the next slide.

Start with some  $(x_0, \lambda_0)$  and compute  $(x_1, \lambda_1), \ldots, (x_k, \lambda_k), \ldots$ 

In every step we compute  $(x_{k+1}, \lambda_{k+1})$  from  $(x_k, \lambda_k)$  using Newton's step.

Consider the gradient of the Lagrangian:

$$\nabla \mathcal{L}(x_k, \lambda_k) = (\nabla_x \mathcal{L}(x_k, \lambda_k), \nabla_\lambda \mathcal{L}(x_k, \lambda_k))^\top$$
$$= (\nabla f(x_k) + \sum_{j=1}^{n_h} \lambda_{kj} \nabla h_j(x_k), \quad h(x_k))^\top \in \mathbb{R}^{n+n_h}$$

and the Hessian matrix of the (complete) Lagrangian

$$abla^2 \mathcal{L}(x_k,\lambda_k) \in \mathbb{R}^{n+n_h} imes \mathbb{R}^{n+n_h}$$

We compute this Hessian in the next slide.

The Newton's step is then computed by

$$\begin{aligned} x_{k+1} &= x_k + p_k \qquad \lambda_{k+1} = \lambda_k + \mu_k \\ (p_k, \mu_k) &= - \left( \nabla^2 \mathcal{L}(x_k, \lambda_k) \right)^{-1} \nabla \mathcal{L}(x_k, \lambda_k) \end{aligned}$$

## Hessian of Lagrangian

Note that

$$\nabla^{2}\mathcal{L}(x_{k},\lambda_{k}) = \begin{pmatrix} \nabla_{xx}\mathcal{L}(x_{k},\lambda_{k}) & \nabla_{x\lambda}\mathcal{L}(x_{k},\lambda_{k}) \\ \nabla_{\lambda x}\mathcal{L}(x_{k},\lambda_{k}) & \nabla_{\lambda\lambda}\mathcal{L}(x_{k},\lambda_{k}) \end{pmatrix}$$
$$= \begin{pmatrix} H(x_{k},\lambda_{k}) & \nabla h(x_{k}) \\ \nabla h(x_{k})^{\top} & 0 \end{pmatrix}$$

Here H is the Lagrangian-Hessian:

$$H(x_k,\lambda_k) = H_f(x_k) + \sum_{j=1}^{n_h} \lambda_{kj} H_{h_j}(x_k)$$

Here  $H_f$  is the Hessian of f, and each  $H_{h_i}$  is the Hessian of  $h_j$ .

 $\nabla h(x_k) = (\nabla h_1(x_k) \cdots \nabla h_{n_h}(x_k))$ 

is the matrix of columns  $\nabla h_j(x_k)$  for  $j = 1, \ldots, n_h$ .
### Lagrange-Newton for Equality Constraints

#### Algorithm 14 Lagrange-Newton

- 1: Choose starting point  $x_0$
- 2:  $k \leftarrow 0$
- 3: repeat

4: Compute 
$$\nabla f(x_k)$$
,  $\nabla h(x_k)$ ,  $h(x_k)$ 

- 5: Compute  $\nabla \mathcal{L}(x_k, \lambda_k)$
- 6: Compute Hessians  $H_f(x_k), H_{h_j}(x_k)$  for  $j = 1, ..., n_h$
- 7: Compute Lagrangian-Hessian  $H(x_k, \lambda_k)$
- 8: Compute  $\nabla^2 \mathcal{L}(x_k, \lambda_k)$
- 9: Compute  $(p_k, \mu_k)^{\top} = -(\nabla^2 \mathcal{L}(x_k, \lambda_k))^{-1} \nabla \mathcal{L}(x_k, \lambda_k)$
- 10:  $x_{k+1} \leftarrow x_k + p_k$
- 11:  $\lambda_{k+1} \leftarrow \lambda_k + \mu_k$
- 12:  $k \leftarrow k+1$

13: **until** convergence

# Sequential Quadratic Programming for Inequality Constraints

Introducing inequality constraints brings serious problems.

The main problem is caused by the fact that active constraints behave differently from inactive ones.

# Sequential Quadratic Programming for Inequality Constraints

Introducing inequality constraints brings serious problems.

The main problem is caused by the fact that active constraints behave differently from inactive ones.

Roughly speaking, algorithms proceed by searching through possible combinations of active/inactive constraints and solve for each combination as if only equality constraints were present. This is very closely related to the support enumeration algorithm from game theory.

# Sequential Quadratic Programming for Inequality Constraints

Introducing inequality constraints brings serious problems.

The main problem is caused by the fact that active constraints behave differently from inactive ones.

Roughly speaking, algorithms proceed by searching through possible combinations of active/inactive constraints and solve for each combination as if only equality constraints were present. This is very closely related to the support enumeration algorithm from game theory.

We will consider this type of algorithm only for linear programming (the simplex algorithm).

We have considered optimization for differentiable f and  $h_j$ 's.

We have considered both constrained and unconstrained optimization problems.

Primarily line-search methods: Local search, in every step set a direction and a step length.

The step length should satisfy the strong Wolfe conditions.

### Summary of Unconstrained Methods

Consider only f without constraints.

For setting direction we used several methods

- Gradient descent
  Go downhill. Only first-order derivatives needed. Zig-zags.
- Newton's method

Always minimize the local quadratic approximation of f. Second-order derivatives needed. Better behavior than GD, computationally heavy.

quasi-Newton (SR1, BFGS, L-BFGS) Approximate the quadratic approximation of *f*. Only first-order derivatives needed. Behaves similarly to Newton's method. Much more computationally efficient.

### Summary of Constrained Optimization

Penalty methods, both exterior and interior. Penalize minimizer approximations out of the feasible region (exterior), or close to the border (interior).

#### Exterior

Penalize minimizer approximations out of the feasible region.

Quadratic penalty, both for equality and inequality constraints.

Interior

Penalize minimizer approximations close to the border (interior). Inverse barrier, logarithmic barrier, only for inequality constraints.

Finally, we have considered the Lagrange-Newton method for equality constraints.