

PV027 Optimization

Tomáš Brázdil

Resources & Prerequisites

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 & Prerequisites

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 \mathbb{R}^n
Operations with vectors and matrices, bases, diagonalization.
- ▶ Multi-variable calculus (i.e., in \mathbb{R}^n)
Partial derivatives, gradients, Hessians, Taylor's theorem.

We will refresh our memories during lectures and tutorials.

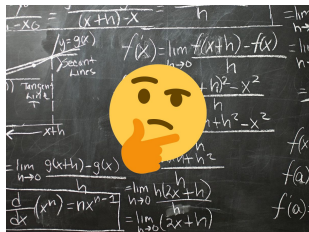
Evaluation

Oral exam - You will get a manual describing the knowledge necessary for **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*.

Optimization

People optimize in

- ▶ scheduling
 - ▶ transportation,
 - ▶ education,
 - ▶ ...

Optimization

People optimize in

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

Optimization

People optimize in

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

Optimization

People optimize in

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

Optimization

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, x0, 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)

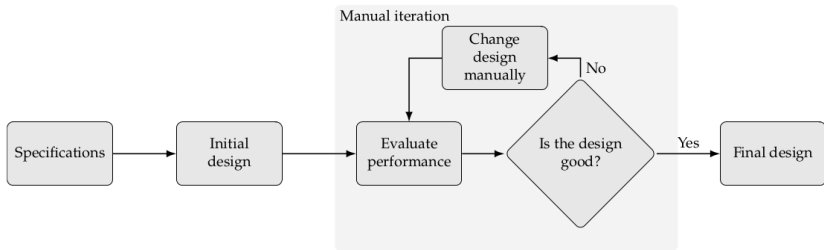
Optimization Algorithms

`sklearn.linear_model.LogisticRegression`

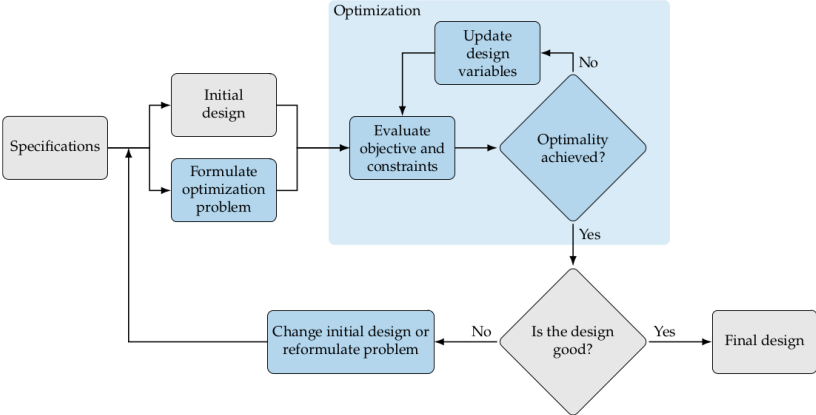
```
class sklearn.linear_model.LogisticRegression(penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None)
```

`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



A bit naive example:

- ▶ 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.

A bit naive example:

- ▶ 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.

A bit naive example:

- ▶ 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 \Rightarrow , introducing constraints on the maximum production of the plants.

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

A bit naive example:

- ▶ 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 \Rightarrow , 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.
- ▶ ...

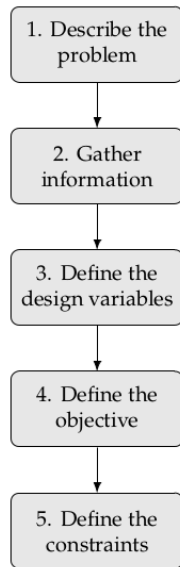
Optimization Problem Formulation

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.



Optimization Problem Formulation

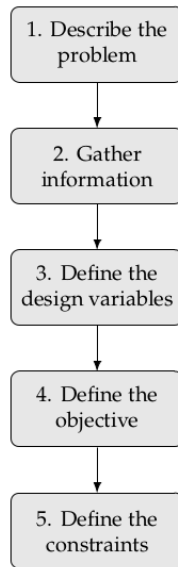
3. Define the design **variables**

- ▶ 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).



Optimization Problem Formulation

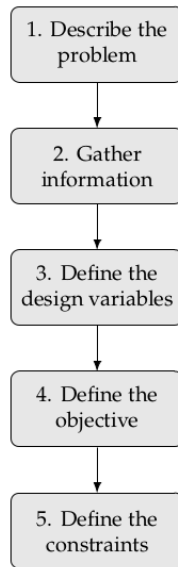
4. Define the **objective**

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

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

(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.



Optimization Problem Formulation

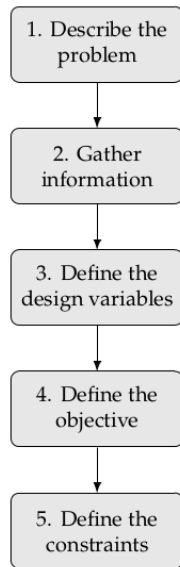
5. Define the **constraints**

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

$$c(x) \leq 0 \text{ or } c(x) \geq 0 \text{ or } c(x) = 0$$

(e.g., time cannot be negative, bounded amount of money to invest)

Where $c : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function depending on the variables.



Modelling and Optimization

The **Optimization Problem** consists of

- ▶ **variables**
- ▶ **objective**
- ▶ **constraints**

The above components constitute a **model**.

Modelling and Optimization

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.

Modelling and Optimization

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.

Modelling and Optimization

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

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ an *objective function*,

x a vector of real *variables*,

g_1, \dots, g_{n_g} *inequality constraint functions* $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$.

h_1, \dots, h_{n_h} *equality constraint functions* $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$.

Optimization Problem Formally

Denote by

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ an *objective function*,

x a vector of real *variables*,

g_1, \dots, g_{n_g} *inequality constraint functions* $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$.

h_1, \dots, h_{n_h} *equality constraint functions* $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$.

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

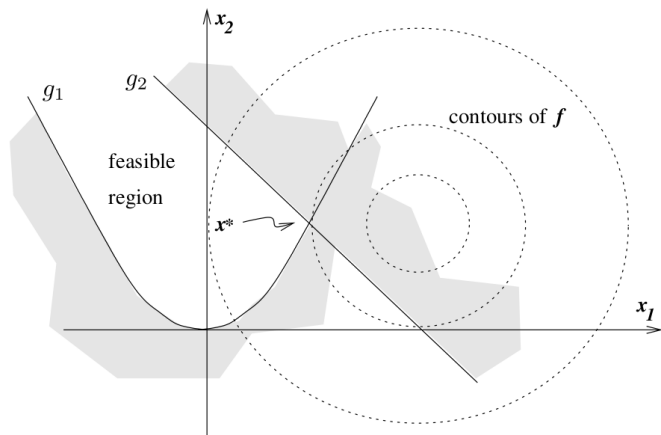
$$\text{minimize } (x_1 - 2)^2 + (x_2 - 1)^2 \quad \text{subject to } \begin{cases} x_1^2 - x_2 \leq 0, \\ x_1 + x_2 - 2 \leq 0. \end{cases}$$

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 $\{x \in \mathbb{R}^n \mid f(x) = c\}$

Constraints

Consider the constraints

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

$$h_j(x) = 0 \quad j = 1, \dots, n_h$$

Constraints

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*.

Constraints

Consider the constraints

$$\begin{aligned}g_i(x) &\leq 0 & i = 1, \dots, n_g \\h_j(x) &= 0 & j = 1, \dots, n_h\end{aligned}$$

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 $g_i(x) \geq 0$ can be easily transformed to the inequality constraints $-g_i(x) \leq 0$

Constraints

Consider the constraints

$$\begin{aligned}g_i(x) &\leq 0 & i = 1, \dots, n_g \\h_j(x) &= 0 & j = 1, \dots, n_h\end{aligned}$$

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 $g_i(x) \geq 0$ can be easily transformed to the inequality constraints $-g_i(x) \leq 0$

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

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

Constraints

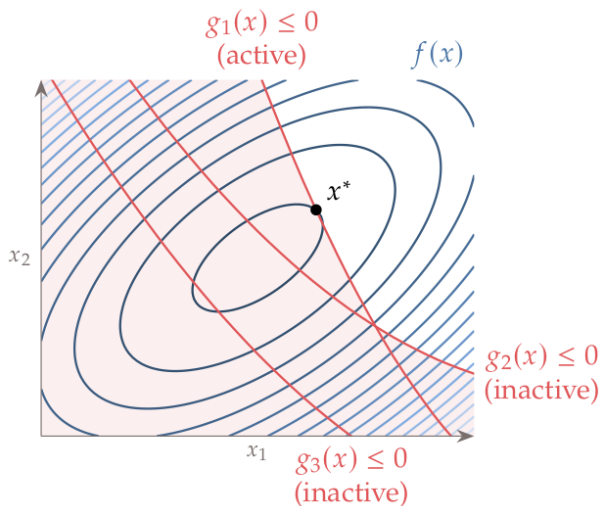
Inequality constraints $g_i(x) \leq 0$ can be *active* or *inactive*.

active at x

$$g_i(x) = 0$$

inactive at x

$$g_i(x) < 0$$



More Practical Example

The problem formulation:

- ▶ A company has two chemical factories F_1 and F_2 , and a dozen retail outlets R_1, \dots, R_{12} .
- ▶ Each F_i can produce (maximum of) a_i 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} .

More Practical Example

The problem formulation:

- ▶ A company has two chemical factories F_1 and F_2 , and a dozen retail outlets R_1, \dots, R_{12} .
- ▶ Each F_i can produce (maximum of) a_i 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.

More Practical Example

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

The objective:

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

More Practical Example

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

The objective:

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

subject to

$$\sum_{j=1}^{12} x_{ij} \leq a_i, \quad i = 1, 2$$

$$\sum_{i=1}^2 x_{ij} \geq b_j, \quad j = 1, \dots, 12,$$

$$x_{ij} \geq 0, \quad i = 1, 2, \quad j = 1, \dots, 12.$$

More Practical Example

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

The objective:

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

subject to

$$\sum_{j=1}^{12} x_{ij} \leq a_i, \quad i = 1, 2$$

$$\sum_{i=1}^2 x_{ij} \geq b_j, \quad j = 1, \dots, 12,$$

$$x_{ij} \geq 0, \quad i = 1, 2, \quad j = 1, \dots, 12.$$

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

$$x_i \in \mathbb{Z}$$

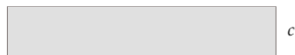
It constrains x_i only to integer values. This leads to so-called *integer programming*.

Discrete optimization problems have discrete and finite variables.

Wing Design Example

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

Assume a rectangular wing.

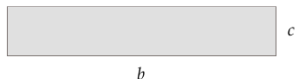


The parameters are called *span* b and *chord* c .

Wing Design Example

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

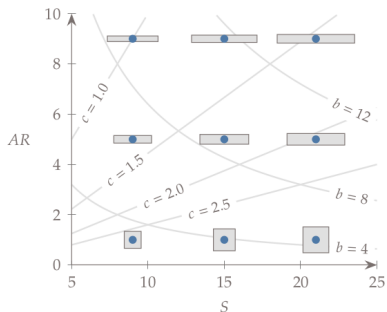
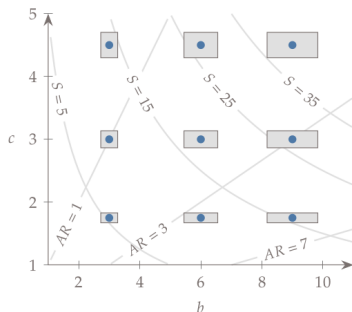
Assume a rectangular wing.



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

$$S = bc \quad AR = b^2/S$$



Wing Design Example

What exactly are the objectives and constraints?

Wing Design Example

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.
- ▶ η 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.

Wing Design Example

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

Wing Design Example

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 $L = W$.

Wing Design Example

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.

Wing Design Example

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 ρ is the air density, C_L is a lift coefficient (depending on the wing shape).

Wing Design Example

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 ρ is the air density, C_L is a lift coefficient (depending on the wing shape).

Thus, we may obtain the lift velocity as

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

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

Wing Design Example

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

Wing Design Example

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.

Wing Design Example

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

$$D_f = k C_f q S$$

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:

$$Re = \rho v c / \mu$$

Here μ is the air dynamic viscosity.

Wing Design Example

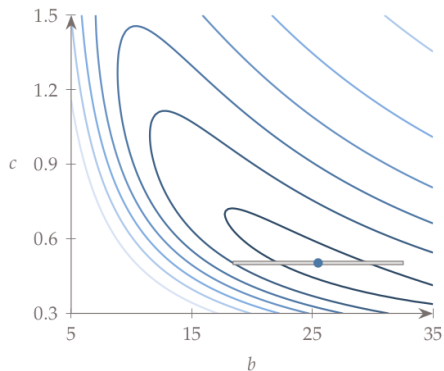
The propulsion efficiency η can be roughly approximated by the Gaussian efficiency curve.

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

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

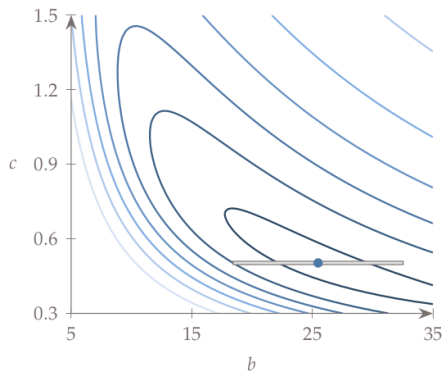
Wing Design Example

The objective function contours:



Wing Design Example

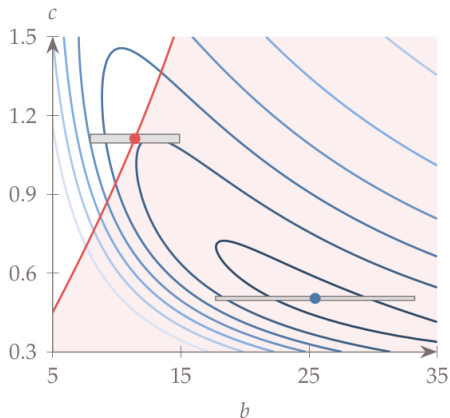
The objective function contours:



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!

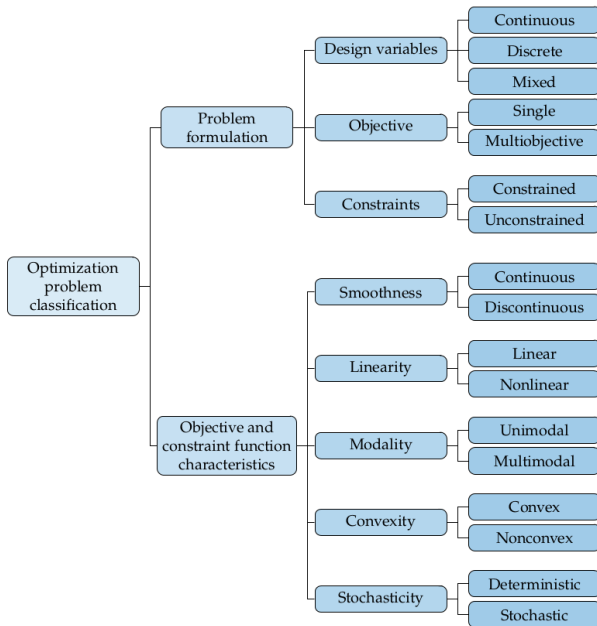
Wing Design Example

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

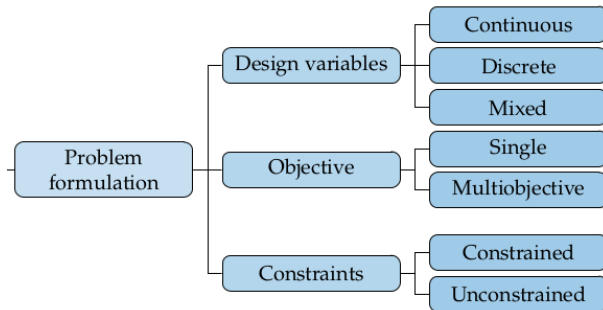


It looks like a reasonable wing ...

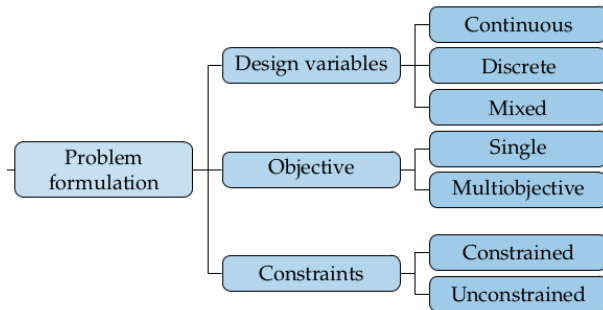
Optimization Problem Classification



Optimization Problem Classification

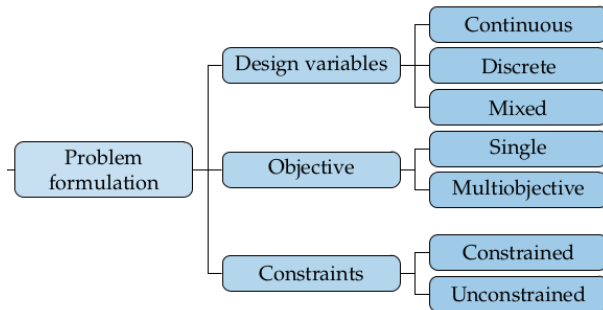


Optimization Problem Classification



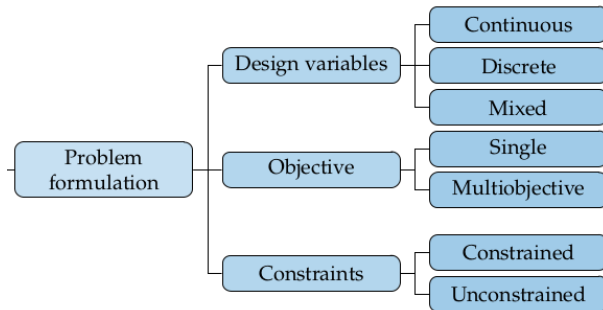
- ▶ *Continuous* allows only $x_i \in \mathbb{R}$, *discrete* allows only $x_i \in \mathbb{Z}$, mixed allows variables of both kinds.

Optimization Problem Classification



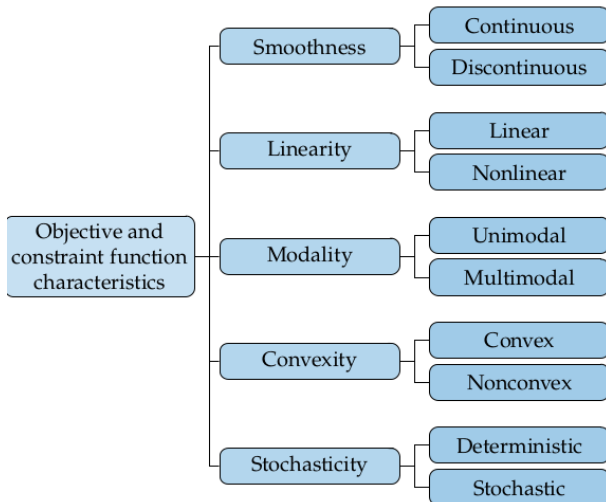
- ▶ *Continuous* allows only $x_i \in \mathbb{R}$, *discrete* allows only $x_i \in \mathbb{Z}$, mixed allows variables of both kinds.
- ▶ *Single-objective*: $f : \mathbb{R}^n \rightarrow \mathbb{R}$, *Multi-objective*: $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Optimization Problem Classification



- ▶ *Continuous* allows only $x_i \in \mathbb{R}$, *discrete* allows only $x_i \in \mathbb{Z}$, mixed allows variables of both kinds.
- ▶ *Single-objective*: $f : \mathbb{R}^n \rightarrow \mathbb{R}$, *Multi-objective*: $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- ▶ *Unconstrained*: No constraints, just the objective function.

Optimization Problem Classification

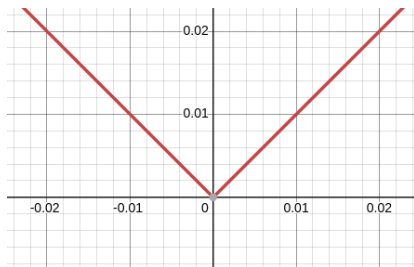


Smoothness

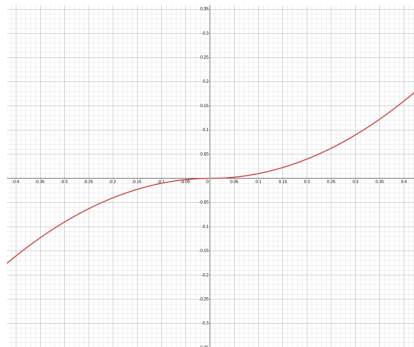
We consider various classes of problems depending on the smoothness properties of the objective/constraint functions:

- ▶ C^0 : Continuous function **Continuity allows us to estimate value in small neighborhoods.**
Discontinuous functions exist.
- ▶ C^1 : 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^2 : Continuous second derivatives **The second derivatives inform about curvature.**
Continuously differentiable functions without second derivatives and twice differentiable functions with discontinuous second derivatives exist.

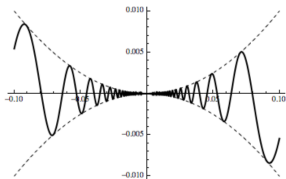
$f(x) = |x|$ is continuous, f is not differentiable at 0



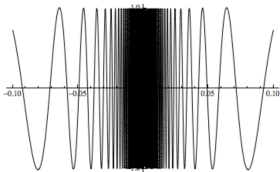
$f(x) = x|x|$ is differentiable on \mathbb{R} , f' has no second derivative at 0



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



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



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} ,

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

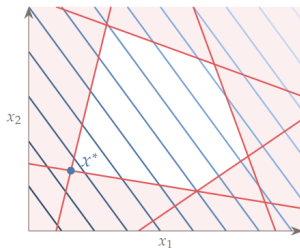
f' is differentiable on \mathbb{R} ,

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

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

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

Multimodality

Denote by \mathcal{F} the feasibility set.

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

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

x^* is a (weak) *global minimiser* if

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

Multimodality

Denote by \mathcal{F} the feasibility set.

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

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

x^* is a (weak) *global minimiser* if

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

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

Multimodality

Denote by \mathcal{F} the feasibility set.

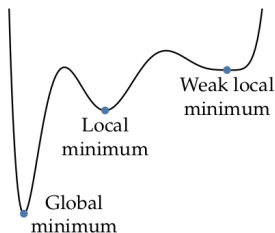
x^* is a (weak) *local minimiser* if there is $\varepsilon > 0$ such that

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

x^* is a (weak) *global minimiser* if

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

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

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

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

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

A *standard form convex optimization* assumes

- ▶ convex objective f and convex inequality constraint functions g_i
- ▶ affine equality constraint functions h_j

Implications:

- ▶ Every local minimum is a global minimum.
- ▶ If the above inequality is strict for all $x \neq y$, 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, \dots 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, \dots 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.

Optimization Algorithm Properties

Typically, we are concerned with the following issues:

Optimization Algorithm Properties

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.

Optimization Algorithm Properties

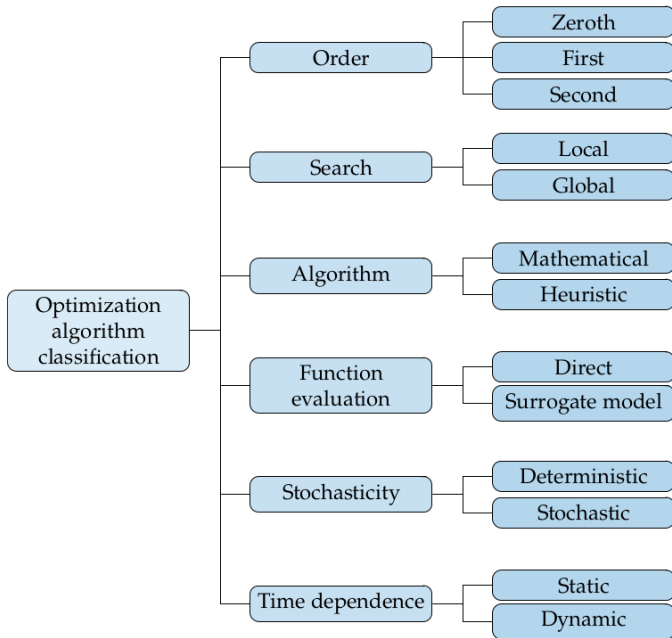
Typically, we are concerned with the following issues:

- ▶ *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.

Optimization Algorithm Properties

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)

Mathematical vs Heuristic

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.

Mathematical vs Heuristic

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.

Mathematical vs Heuristic

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.

Mathematical vs Heuristic

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.

Mathematical vs Heuristic

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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

Find x^* such that

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

Assume that

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad \text{for } x \in \mathbb{R}$$

is continuous on \mathbb{R} .

Denote by \mathcal{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 α_k to α_{k+1} based on a certain strategy
- 6: Increment k
- 7: **end while**
- 8: **return** x_k

Convergence of Single Variable Gradient Descent

Theorem 1

Assume that f is

- ▶ differentiable, i.e., that f' exists,
- ▶ bounded below, i.e., there is $B \in \mathbb{R}$ such that $f(x) \geq B$ for all $x \in \mathbb{R}$,
- ▶ L -smooth, i.e., there is $L > 0$ such that $|f'(x) - f'(x')| \leq L|x - x'|$ for all $x, x' \in \mathbb{R}$.

Consider a sequence x_0, x_1, \dots computed by the gradient descent algorithm for f . Assume a constant step length $\alpha \leq \frac{1}{L}$.

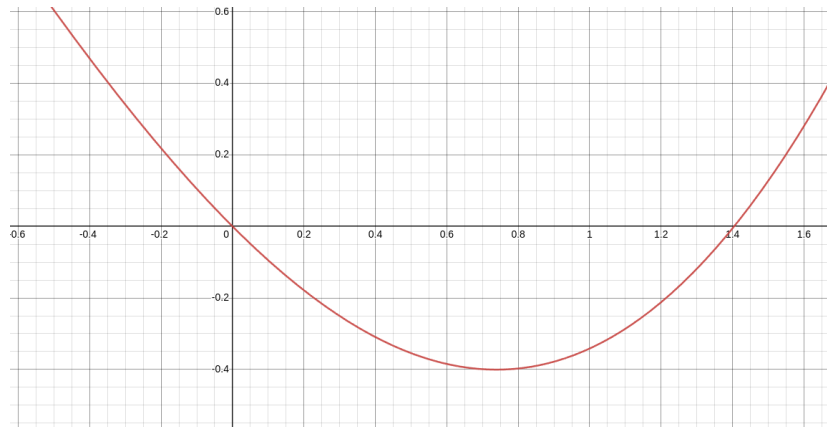
Then $\lim_{k \rightarrow \infty} |f'(x_k)| = 0$ and, moreover,

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

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,

$$\begin{aligned}x_1 &= 0.5 - (0.5 - \cos 0.5) \\ &= 0.5 - (-0.37758) \\ &= 0.87758\end{aligned}$$

Example

Continuing in the same way:

$$x_1 = 0.87758$$

$$x_2 = 0.63901$$

$$x_3 = 0.80269$$

$$x_4 = 0.69478$$

$$x_5 = 0.76820$$

$$x_6 = 0.71917$$

$$x_7 = 0.75236$$

$$x_8 = 0.73008$$

$$x_9 = 0.74512$$

$$x_{10} = 0.73501$$

$$x_{11} = 0.74183$$

$$x_{12} = 0.73724$$

$$x_{13} = 0.74033$$

$$x_{14} = 0.73825$$

$$x_{15} = 0.73965$$

$$x_{16} = 0.73870$$

$$x_{17} = 0.73934$$

$$x_{18} = 0.73891$$

$$x_{19} = 0.73920$$

$$x_{20} = 0.73901$$

$$x_{21} = 0.73914$$

$$x_{22} = 0.73905$$

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

Example

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....

Example

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

$$x_1 = 0.50000$$

$$x_2 = 0.87758$$

$$x_3 = 0.75830$$

$$x_4 = 0.74753$$

$$x_5 = 0.74399$$

$$x_6 = 0.74235$$

$$x_7 = 0.74144$$

$$x_8 = 0.74087$$

$$x_9 = 0.74050$$

$$x_{10} = 0.74024$$

$$x_{11} = 0.74004$$

$$x_{12} = 0.73990$$

$$x_{13} = 0.73978$$

$$x_{14} = 0.73969$$

...

$$x_{115} = 0.739100605$$

$$x_{116} = 0.739100379$$

$$x_{117} = 0.739100159$$

$$x_{118} = 0.739099944$$

$$x_{119} = 0.739099734$$

$$x_{120} = 0.739099529$$

$$x_{121} = 0.739099328$$

$$x_{122} = 0.739099132$$

$$x_{123} = 0.739098940$$

$$x_{124} = 0.739098752$$

$$x_{125} = 0.739098568$$

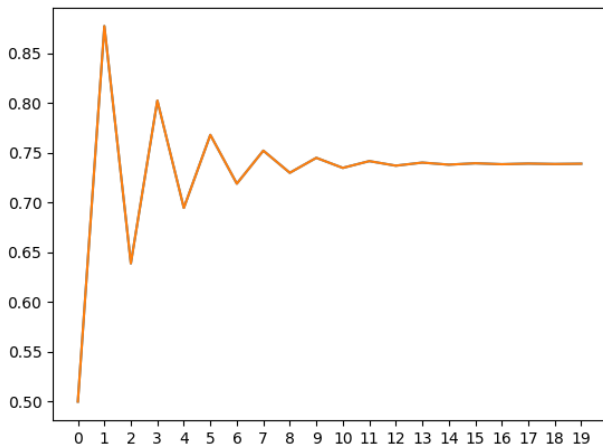
$$x_{126} = 0.739098388$$

$$x_{127} = 0.739098212$$

$$x_{128} = 0.739098040$$

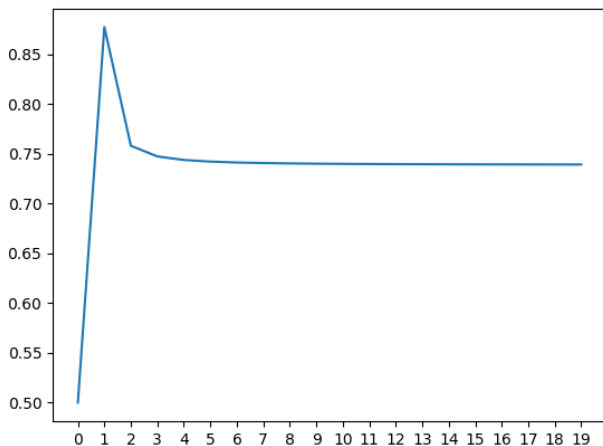
Example

Gradient descent with the step length = 1.0:



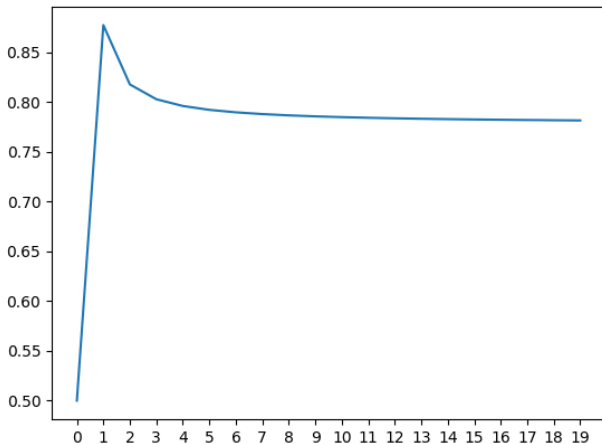
Example

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



Example

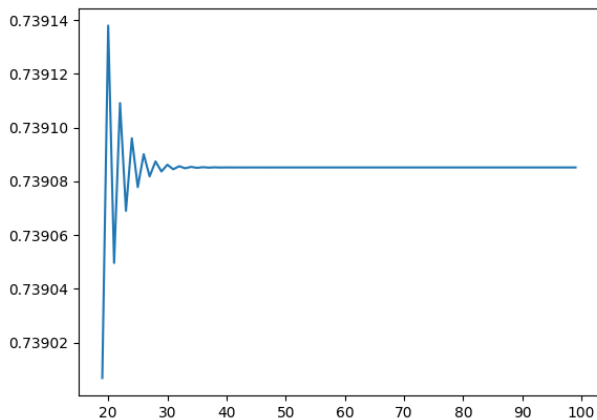
Gradient descent with the step length = $1/k^2$:



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

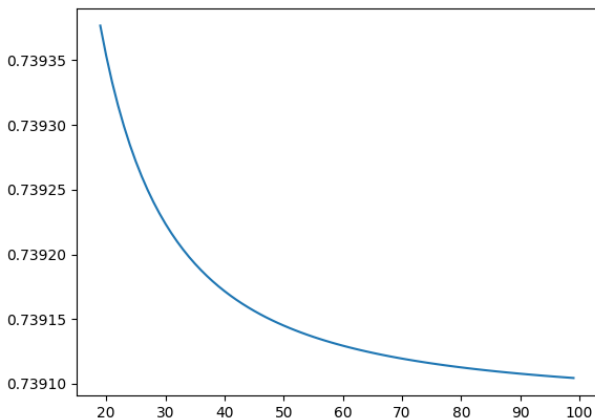
Example

Gradient descent with the step length = 1.0:



Example

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



Properties of Gradient Descent

- ▶ 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.

Properties of Gradient Descent

- ▶ 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.

Properties of Gradient Descent

- ▶ 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).

Properties of Gradient Descent

- ▶ 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).

Properties of Gradient Descent

- ▶ 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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

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

Find x^* such that

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

Assume that

$$f''(x) = \lim_{h \rightarrow 0} \frac{f'(x+h) - f'(x)}{h} \quad \text{for } x \in \mathbb{R}$$

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 $f, g : \mathbb{R} \rightarrow \mathbb{R}$ we write $f = o(g)$ if

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

Taylor Series Approximation

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

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

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$ 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 $f, g : \mathbb{R} \rightarrow \mathbb{R}$ we write $f = o(g)$ if

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

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$ 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 $x_0, x_1, \dots, x_k, \dots$ as the GD.

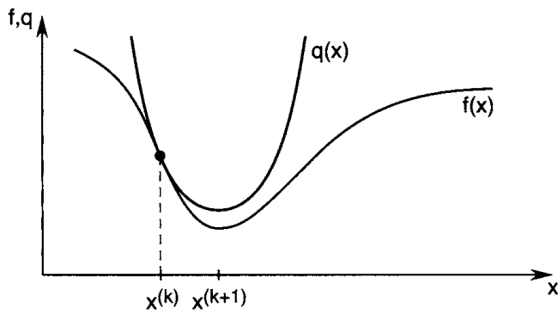
Newton's Method Idea

The method computes successive approximations $x_0, x_1, \dots, x_k, \dots$ 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 .



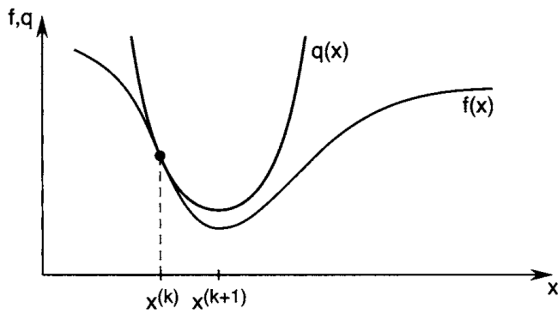
Newton's Method Idea

The method computes successive approximations $x_0, x_1, \dots, x_k, \dots$ 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., $q'(x_{k+1}) = 0$).

Newton's Method Algorithm

Now note that for

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

Newton's Method Algorithm

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

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

Newton's Method Algorithm

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

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

Newton's Method Algorithm

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

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

Newton's method then sets

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

Newton's Method Algorithm

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_k

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

Example

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 $|x_{k+1} - x_k| \leq \epsilon$.

Example

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 $|x_{k+1} - x_k| \leq \epsilon$.

We compute

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

Example

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 $|x_{k+1} - x_k| \leq \epsilon$.

We compute

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

Hence,

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

Example

Proceeding similarly, we obtain

$$x_2 = x_1 - \frac{f'(x_1)}{f''(x_1)} = x_1 - \frac{0.02710}{1.685} = 0.7391$$

$$x_3 = x_2 - \frac{f'(x_2)}{f''(x_2)} = x_2 - \frac{9.461 \times 10^{-5}}{1.673} = 0.7390851339$$

$$x_4 = x_3 - \frac{f'(x_3)}{f''(x_3)} = x_3 - \frac{1.17 \times 10^{-9}}{1.673} = 0.7390851332$$

...

Example

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

$$|x_4 - x_3| < \epsilon = 10^{-5}$$

$$f'(x_4) = -8.6 \times 10^{-6} \approx 0$$

$$f''(x_4) = 1.673 > 0$$

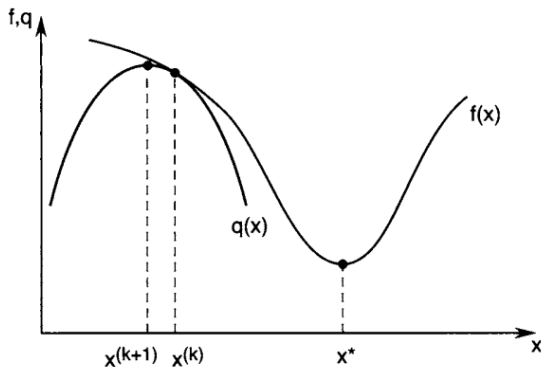
So, we conclude that $x^* \approx x_4$ is a strict minimizer.

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:

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

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:

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

Superlinear Convergence

Convergence is superlinear if:

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

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:

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

Quadratic Convergence of Newton's Method

Quadratic Convergence

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

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

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

Theorem 2 (Quadratic Convergence of Newton's Method)

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ satisfy $f \in \mathcal{C}^2$ and suppose x^ is a minimizer of f such 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

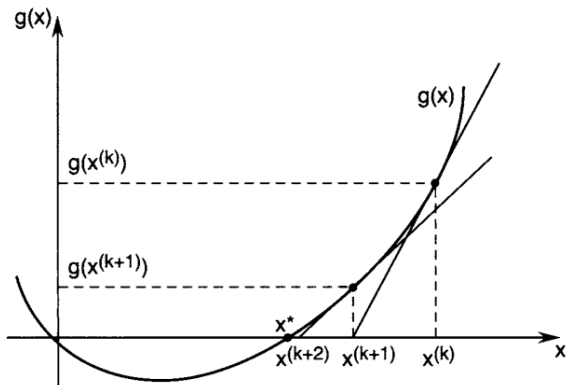
Newton's method is also a technique for finding roots of functions. In our case, this means finding a root of f' .

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:

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



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

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

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

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

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?

Example

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.

Example

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,

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

Example

Continuing, we obtain:

$$x_0 = 0.5$$

$$x_1 = 1.0$$

$$x_2 = 0.72548$$

$$x_3 = 0.73839$$

$$x_4 = 0.739087$$

$$x_5 = 0.739085132$$

$$x_6 = 0.739085133$$

Example

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

$$x_0 = 0.5$$

$$x_1 = 0.7552$$

$$x_2 = 0.7381$$

$$x_3 = 0.739081$$

$$x_5 = 0.7390851339$$

$$x_6 = 0.7390851332$$

...

Compare with Newton's method:

$$x_0 = 0.5$$

$$x_1 = 0.7552$$

$$x_2 = 0.7391$$

$$x_3 = 0.7390851339$$

$$x_4 = 0.73908513321516067229$$

$$x_5 = 0.73908513321516067229$$

...

Superlinear Convergence of Secant Method

Theorem 3 (Superlinear Convergence of Secant Method)

Assume $f : \mathbb{R} \rightarrow \mathbb{R}$ 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,

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

Secant Method for Root Finding

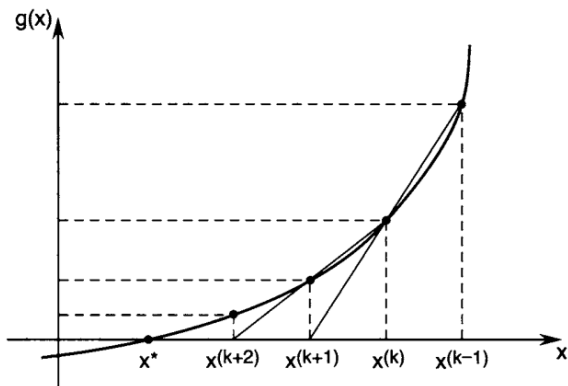
As for Newton's method of tangents, the secant method can be seen as a method for finding a root of f' .

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

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



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:

- ▶ $a_k = 1$ gives the **gradient descent**,
- ▶ $a_k = f''(x_k)$ gives **Newton's method**,
- ▶ $a_k = \frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}}$ gives the **secant method**,
- ▶ $a_k = f''(x_m)$ where $m = \lfloor k/p \rfloor p$ gives **Shamanskii method**.

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 $f : \mathbb{R} \rightarrow \mathbb{R}$

A variable x

A constraint

$$a_0 \leq x \leq b_0$$

Consider the following cases:

- ▶ f continuously differentiable on $[a_0, b_0]$
- ▶ 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).