

# Methods

## 1 General formulation

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} f(x) \\ \text{s.t. } & g_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_j(x) = 0, \quad j = 1, \dots, k \end{aligned}$$

Some necessary or/and sufficient conditions are known (See [Optimality conditions. KKT](#) and [Convex optimization problem](#)).

- In fact, there might be very challenging to recognize the convenient form of optimization problem.
- Analytical solution of KKT could be inviable.

### 1.1 Iterative methods

Typically, the methods generate an infinite sequence of approximate solutions

$$\{x_t\},$$

which for a finite number of steps (or better - time) converges to an optimal (at least one of the optimal) solution  $x_*$ .

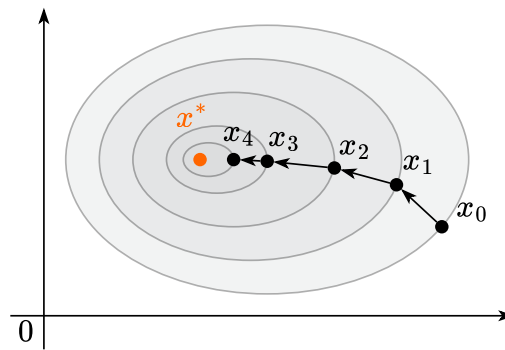
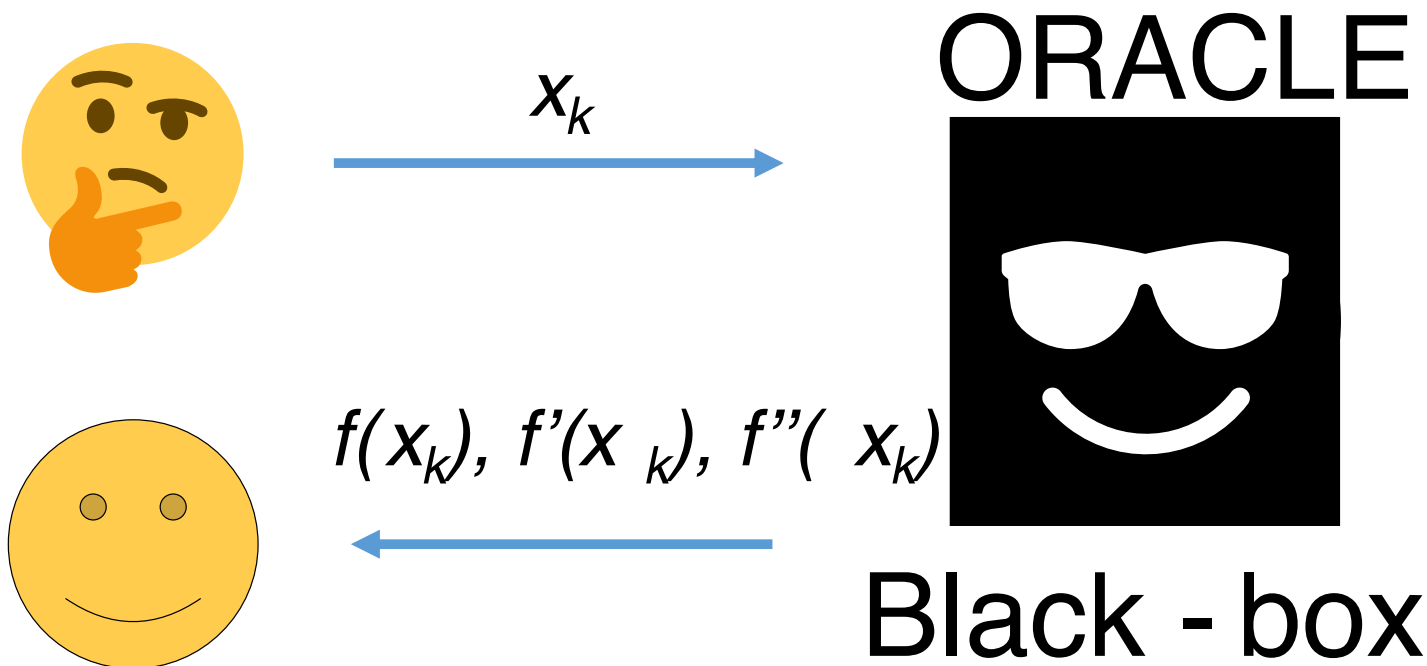


Illustration of iterative method approaches to the solution  $x^*$

```
def GeneralScheme(x, epsilon):
    while not StopCriterion(x, epsilon):
        OracleResponse = RequestOracle(x)
        x = NextPoint(x, OracleResponse)
    return x
```



### 1.2 Oracle conception



Depending on the maximum order of derivative available from the oracle we call the oracles as zero order, first order, second order oracle and etc.

## 2 Unsolvability of numerical optimization problem

In general, **optimization problems are unsolvable.** 😞

Consider the following simple optimization problem of a function over unit cube:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} f(x) \\ \text{s.t. } x \in \mathbb{C}^n \end{aligned}$$

We assume, that the objective function  $f(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$  is Lipschitz continuous on  $\mathbb{B}^n$ :

$$|f(x) - f(y)| \leq L \|x - y\|_{\infty} \forall x, y \in \mathbb{C}^n,$$

with some constant  $L$  (Lipschitz constant). Here  $\mathbb{C}^n$  - the  $n$ -dimensional unit cube

$$\mathbb{C}^n = \{x \in \mathbb{R}^n \mid 0 \leq x_i \leq 1, i = 1, \dots, n\}$$

Our goal is to find such  $\tilde{x} : |f(\tilde{x}) - f^*| \leq \varepsilon$  for some positive  $\varepsilon$ . Here  $f^*$  is the global minima of the problem. Uniform grid with  $p$  points on each dimension guarantees at least this quality:

$$\|\tilde{x} - x_*\|_{\infty} \leq \frac{1}{2p},$$

which means, that

$$|f(\tilde{x}) - f(x_*)| \leq \frac{L}{2p}$$

Our goal is to find the  $p$  for some  $\varepsilon$ . So, we need to sample  $(\frac{L}{2\varepsilon})^n$  points, since we need to measure function in  $p^n$  points. Doesn't look scary, but if we'll take  $L = 2, n = 11, \varepsilon = 0.01$ , computations on the modern personal computers will take 31,250,000 years.

### 2.1 Stopping rules

- Argument closeness:

$$\|x_k - x_*\|_2 < \varepsilon$$

- Function value closeness:

$$\|f_k - f^*\|_2 < \varepsilon$$

- Closeness to a critical point

$$\|f'(x_k)\|_2 < \varepsilon$$

But  $x_*$  and  $f^* = f(x_*)$  are unknown!

Sometimes, we can use the trick:

$$\|x_{k+1} - x_k\| = \|x_{k+1} - x_k + x_* - x_*\| \leq \|x_{k+1} - x_*\| + \|x_k - x_*\| \leq 2\varepsilon$$

**Note:** it's better to use relative changing of these values, i.e.  $\frac{\|x_{k+1} - x_k\|_2}{\|x_k\|_2}$ .

### Example

Suppose, you are trying to estimate the vector  $x_{true}$  with some approximation  $x_{approx}$ . One can choose between two relative errors:

$$\frac{\|x_{approx} - x_{true}\|}{\|x_{approx}\|} \quad \frac{\|x_{approx} - x_{true}\|}{\|x_{true}\|}$$

If both  $x_{approx}$  and  $x_{true}$  are close to each other, then the difference between them is small, while if your approximation is far from the truth (say,  $x_{approx} = 10x_{true}$  or  $x_{approx} = 0.01x_{true}$  they differ drastically).

## 2.2 Local nature of the methods

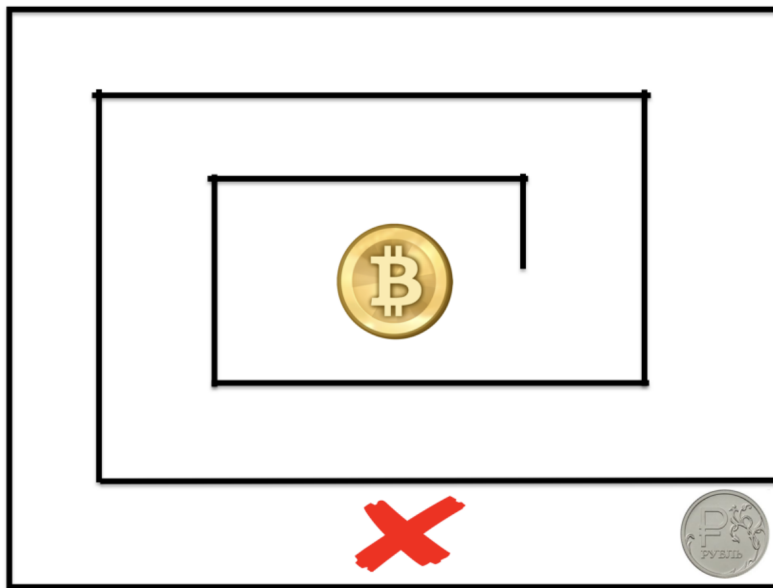
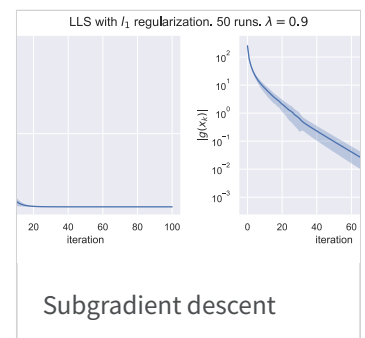
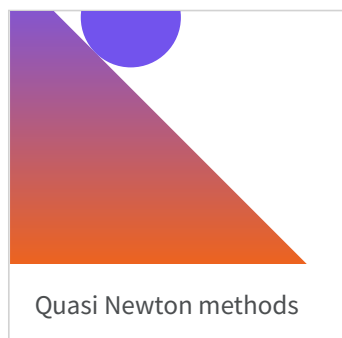
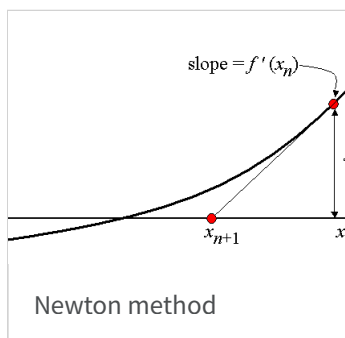
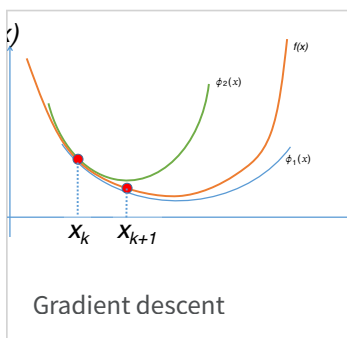


Illustration of the idea of locality in black-box optimization

## 3 Contents of the chapter



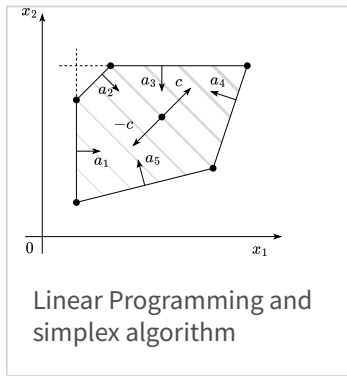
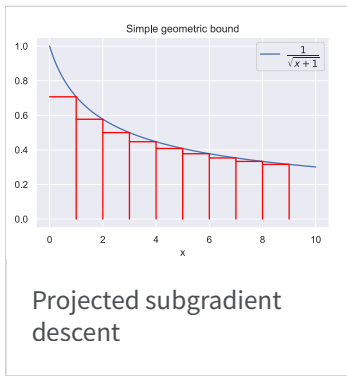


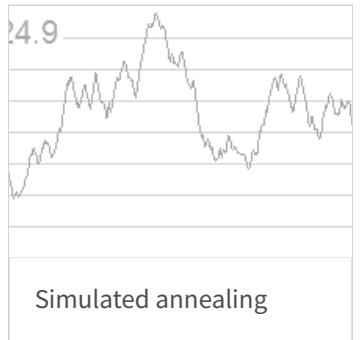
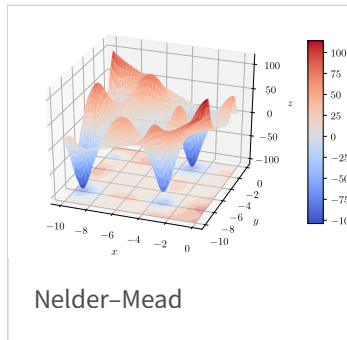
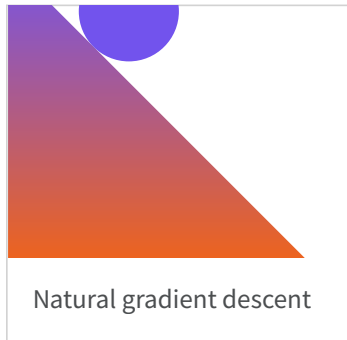
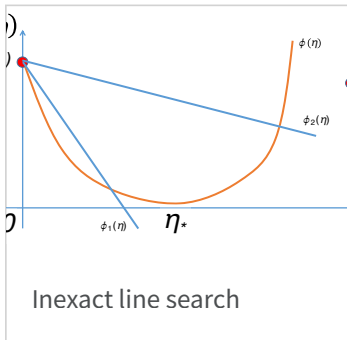
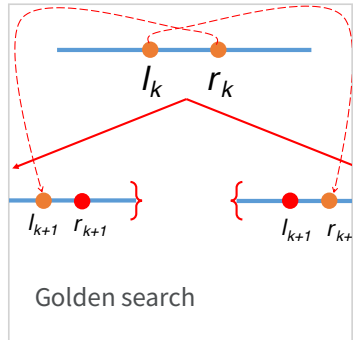
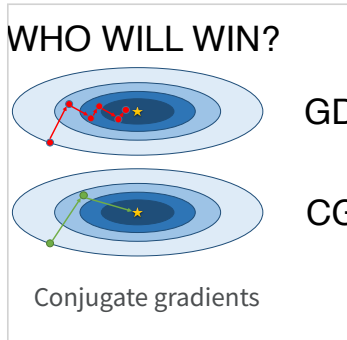
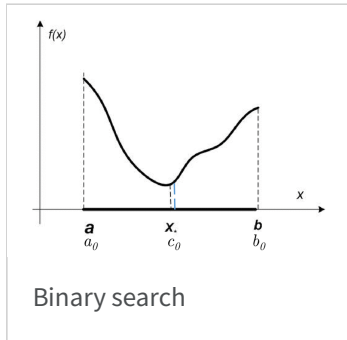
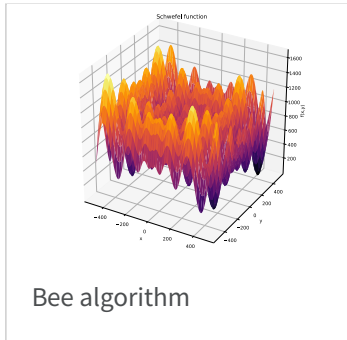
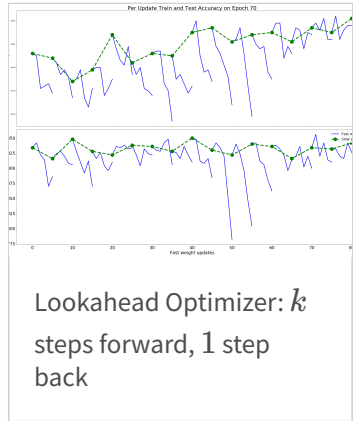
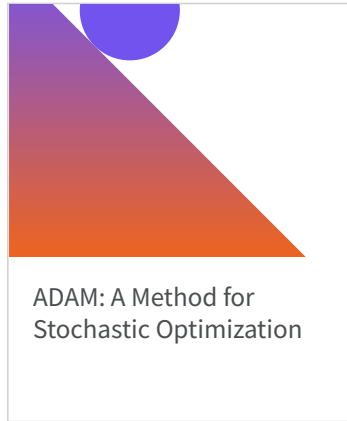
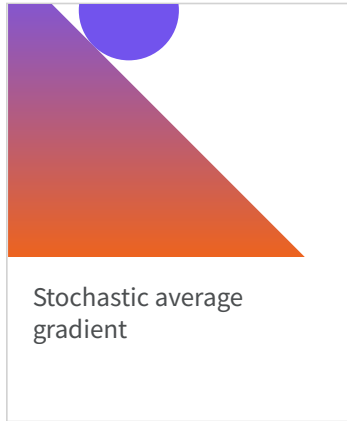
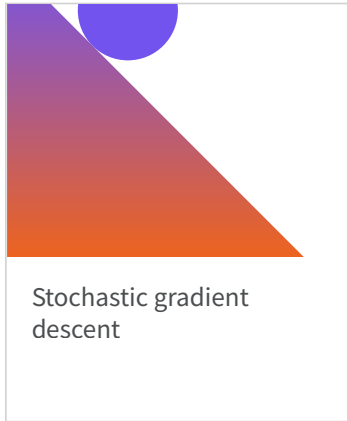
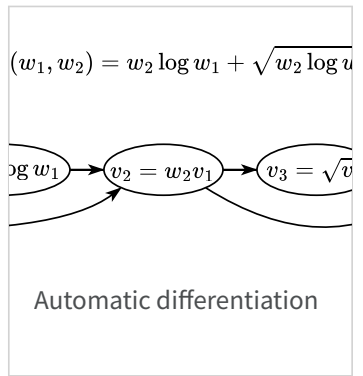
Table 2.1  
Common real functions and the corresponding divergences.

Function name	$\phi(x)$	$\text{dom}(\phi(x))$	$V_\phi(y)$
Squared norm	$\frac{1}{2}x^2$	$(-\infty, +\infty)$	$\frac{1}{2}(x-y)^2$
Shannon entropy	$x \log x - x$	$(0, +\infty)$	$x \log \frac{x}{y} - x + y$
Bit entropy	$x \log_2 x + (1-x) \log_2(1-x)$	$(0, 1)$	$x \log_2 \frac{x}{y} + (1-x) \log_2 \frac{1-x}{1-y}$
Burg entropy	$-\log x$	$(0, +\infty)$	$1 - \log \frac{x}{y} - 1$
Bellings	$-\sqrt{x-x^2}$	$[0, 1]$	$(1-x)(1-y)^{1/2} - (1-x^2)^{1/2}$
$\phi_p$ quantum	$-x^p$	$(0, c, +\infty)$	$-x^p + y^{p-1} - (y-x)^{p/2}$
$\phi_p$ norm	$ x ^p$	$(-\infty, +\infty)$	$ x ^p - p \log  x ^{p-1} + (y-1) y ^p$
Exponential	$\exp x$	$(-\infty, +\infty)$	$\exp x - (x-y + 1) \exp y$
Inverse	$1/x$	$(0, +\infty)$	$1/x + x/y^2 - 2/y$

Table 2.2  
Common exponential families and the corresponding divergences.

Exponential family	$\psi(\theta)$	$\text{dom}(\psi)$	$\mu(\theta)$	$\phi(x)$	Divergence
Gaussian ( $\theta^T$ level)	$\frac{1}{2}\theta^T \theta$	$(-\infty, +\infty)$	$\theta$	$\frac{1}{2}x^2$	Fuchsian
Poisson	$\exp \theta$	$(-\infty, +\infty)$	$\exp \theta$	$x \log x - x$	Relative entropy
Bernoulli	$\log(1 + \exp \theta)$	$(-\infty, +\infty)$	$\frac{\exp \theta}{1 + \exp \theta}$	$x \log x + (1-x) \log(1-x)$	Logistic loss
Gamma ( $\alpha$ fixed)	$-\alpha \log(-\theta)$	$(-\infty, 0)$	$-\alpha/\theta$	$-\alpha \log x + \alpha \log \alpha - \alpha$	Bakura-Saito

Mirror descent





# Rates of convergence

## 1 Speed of convergence

In order to compare performance of algorithms we need to define a terminology for different types of convergence. Let  $r_k = \{\|x_k - x^*\|_2\}$  be a sequence in  $\mathbb{R}^n$  that converges to zero.

### 1.1 Linear convergence

We can define the *linear* convergence in a two different forms:

$$\|x_{k+1} - x^*\|_2 \leq Cq^k \quad \text{or} \quad \|x_{k+1} - x^*\|_2 \leq q\|x_k - x^*\|_2,$$

for all sufficiently large  $k$ . Here  $q \in (0, 1)$  and  $0 < C < \infty$ . This means that the distance to the solution  $x^*$  decreases at each iteration by at least a constant factor bounded away from 1. Note, that sometimes this type of convergence is also called *exponential* or *geometric*. We call the  $q$  the convergence rate.

#### Question

Suppose, you have two sequences with linear convergence rates  $q_1 = 0.1$  and  $q_2 = 0.7$ , which one is faster?

#### Example

Let us have the following sequence:

$$r_k = \frac{1}{3^k}$$

One can immediately conclude, that we have a linear convergence with parameters  $q = \frac{1}{3}$  and  $C = 0$ .

#### Question

Let us have the following sequence:

$$r_k = \frac{4}{3^k}$$

Will this sequence be convergent? What is the convergence rate?

### 1.2 Sublinear convergence

If the sequence  $r_k$  converges to zero, but does not have linear convergence, the convergence is said to be sublinear. Sometimes we can consider the following class of sublinear convergence:

$$\|x_{k+1} - x^*\|_2 \leq Ck^q,$$

where  $q < 0$  and  $0 < C < \infty$ . Note, that sublinear convergence means, that the sequence is converging slower, than any geometric progression.

### 1.3 Superlinear convergence

The convergence is said to be *superlinear* if:

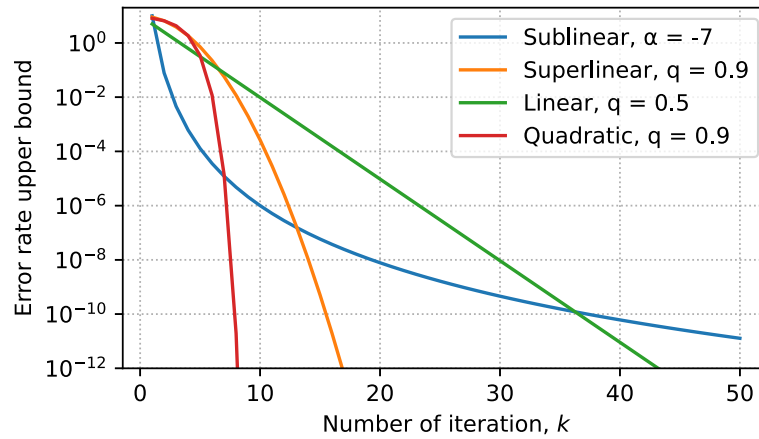
$$\|x_{k+1} - x^*\|_2 \leq Cq^{k^2} \quad \text{or} \quad \|x_{k+1} - x^*\|_2 \leq C_k\|x_k - x^*\|_2,$$

where  $q \in (0, 1)$  or  $0 < C_k < \infty, C_k \rightarrow 0$ . Note, that superlinear convergence is also linear convergence (one can even say, that it is linear convergence with  $q = 0$ ).

### 1.4 Quadratic convergence

$$\|x_{k+1} - x^*\|_2 \leq Cq^{2^k} \quad \text{or} \quad \|x_{k+1} - x^*\|_2 \leq C\|x_k - x^*\|_2^q,$$

where  $q \in (0, 1)$  and  $0 < C < \infty$ .



Difference between the convergence speed

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

## 2 How to determine convergence type

### 2.1 Root test

Let  $\{r_k\}_{k=m}^{\infty}$  be a sequence of non-negative numbers, converging to zero, and let

$$q = \limsup_{k \rightarrow \infty} r_k^{1/k}$$

- If  $0 \leq q < 1$ , then  $\{r_k\}_{k=m}^{\infty}$  has linear convergence with constant  $q$ .
- In particular, if  $q = 0$ , then  $\{r_k\}_{k=m}^{\infty}$  has superlinear convergence.
- If  $q = 1$ , then  $\{r_k\}_{k=m}^{\infty}$  has sublinear convergence.
- The case  $q > 1$  is impossible.

### 2.2 Ratio test

Let  $\{r_k\}_{k=m}^{\infty}$  be a sequence of strictly positive numbers converging to zero. Let

$$q = \lim_{k \rightarrow \infty} \frac{r_{k+1}}{r_k}$$

- If there exists  $q$  and  $0 \leq q < 1$ , then  $\{r_k\}_{k=m}^{\infty}$  has linear convergence with constant  $q$ .
- In particular, if  $q = 0$ , then  $\{r_k\}_{k=m}^{\infty}$  has superlinear convergence.
- If  $q$  does not exist, but  $q = \limsup_{k \rightarrow \infty} \frac{r_{k+1}}{r_k} < 1$ , then  $\{r_k\}_{k=m}^{\infty}$  has linear convergence with a constant not exceeding  $q$ .
- If  $\liminf_{k \rightarrow \infty} \frac{r_{k+1}}{r_k} = 1$ , then  $\{r_k\}_{k=m}^{\infty}$  has sublinear convergence.
- The case  $\liminf_{k \rightarrow \infty} \frac{r_{k+1}}{r_k} > 1$  is impossible.
- In all other cases (i.e., when  $\liminf_{k \rightarrow \infty} \frac{r_{k+1}}{r_k} < 1 \leq \limsup_{k \rightarrow \infty} \frac{r_{k+1}}{r_k}$ ) we cannot claim anything concrete about the convergence rate  $\{r_k\}_{k=m}^{\infty}$ .

#### Example

Let us have the following sequence:

$$r_k = \frac{1}{k}$$

Determine the convergence



### Example

Let us have the following sequence:

$$r_k = \frac{1}{k^2}$$

Determine the convergence



### Example

Let us have the following sequence:

$$r_k = \frac{1}{k^q}, q > 1$$

Determine the convergence



### Try to use root test here

Let us have the following sequence:

$$r_k = \frac{1}{k^k}$$

Determine the convergence

## 3 References

- Code for convergence plots - [Open In Colab](#)
- [CMC seminars \(ru\)](#)
- Numerical Optimization by J.Nocedal and S.J.Wright