

# Gradient descent based optimization

## Description and setting guidelines

The goal of this package is to optimize the control parameters (e.g.  $R_E$ ,  $P_{amb}$ ,  $P_A$ , freq...) in order to minimize to-optimize. The argument to-optimize can be any key in the post processing data dictionary, by default it is 'energy-demand'.

### • Setting ranges

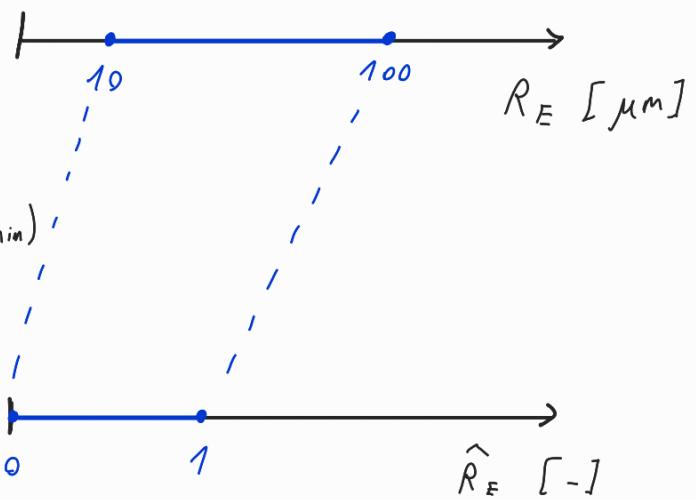
The argument ranges is a regular dictionary, with 1 key corresponding to each control parameters. The value is a list with 1 or 2 elements: 1, if the parameter is not to be optimized and is kept a constant; 2 : if the parameter should be optimized between a lower and upper limit.

Example:

```
ranges = dict(  
    R_E = [10e-6, 100e-6], R_E is optimized between  
10 μm and 100 μm  
    ratio = [1.0, 5.0], ratio ( $R_o/R_E$ ) is optimized  
between 1 and 5  
    P_amb = [1e5], P_amb is kept constant  $10^5$  Pa  
    ...  
)
```

Note, that internally, the control parameters are transformed to be between 0 and 1. In the transformed system, ranges is a multi dimensional unit cube.

$$R_E \in [R_{\min}, R_{\max}]$$



$$\hat{R}_E = \frac{R_E - R_{\min}}{R_{\max} - R_{\min}}$$

$\hat{R}_E \in [0, 1]$

This way a small change, e.g.  $\delta \hat{R}_E = 0.001$  is dependent on ranges (e.g.  $R_E \in [10 \mu\text{m}, 100 \mu\text{m}] \rightarrow \delta R_E = 0.09 \mu\text{m}$   
 $R_E \in [10 \mu\text{m}, 1000 \mu\text{m}] \rightarrow \delta R_E = 0.99 \mu\text{m}$ )

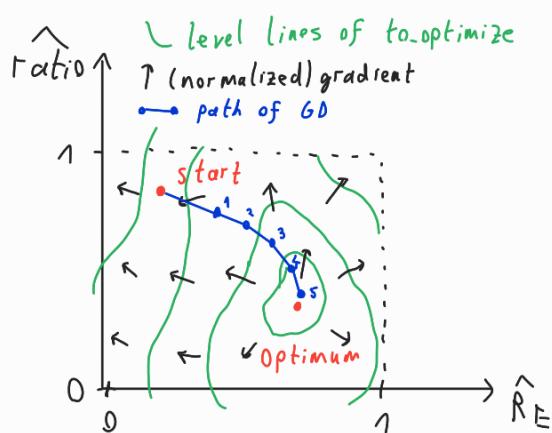
The step size (first\_step, min\_step, and delta) is also defined with these transformed coordinates:

$$0 < \text{step size} = \sqrt{\Delta \hat{R}_E^2 + \Delta \hat{\text{ratio}}_0^2 + \Delta \hat{P}_{amb}^2 + \dots} < 1,$$

where  $\Delta$  denotes the step along one control parameter.

When defining first\_step and min\_step, keep in mind, that if  $\Delta = \text{min\_step}$ , then the largest possible step is  $\Delta R_E \leq \Delta \cdot (R_{\max} - R_{\min})$ .

## • How gradient descent (GD) works?



The basic idea behind GD is that the gradient of a function is always perpendicular to level lines, and points in the direction of the greatest growth.

Thus, by always moving opposite to the local gradient, eventually a local minimum is reached. In practise, the gradient is calculated numerically, and discrete steps are taken.

### - gradient with finite differences (delta)

$g(\underline{y}) : \mathbb{R}^n \rightarrow \mathbb{R}$  function that returns the value of to optimize  
(e.g. 'energy demand' in MJ/kg of  $NH_3$ )

$\underline{y} = [\hat{R}_E, \hat{\text{ratio}}, \hat{P}_{\text{ambi}}, \dots]^T$  vector of transformed control parameters

$y_i$   $i^{th}$  element of  $\underline{y}$  (e.g. here  $y_1 = \hat{R}_E$ ),  $i \in \{1, 2, \dots, n\}$

$\delta$  small number (e.g.  $\delta = 10^{-6}$ )

alias difference

$$\underline{\text{grad}}(g(\underline{y})) = \begin{bmatrix} \frac{\partial g(\underline{y})}{\partial y_1} \\ \vdots \\ \frac{\partial g(\underline{y})}{\partial y_n} \end{bmatrix}$$

forward difference:  $\frac{\partial g(\underline{y})}{\partial y_i} \approx \frac{g(y_1, \dots, y_i + \delta, \dots, y_n) - g(\underline{y})}{\delta}$

- $n+1$  evaluations of  $g()$  for gradient

- bigger numeric error

- default

central difference:  $\frac{\partial g(\underline{y})}{\partial y_i} \approx \frac{g(y_1, \dots, y_i + \delta, \dots, y_n) - g(y_1, \dots, y_i - \delta, \dots, y_n)}{2\delta}$

- $2 \cdot n$  evaluation of  $g()$  for gradient

- smaller numeric error

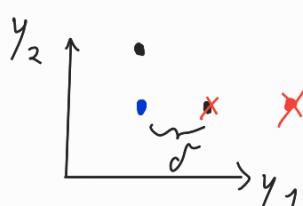
## error management:

(sometimes  $g(\underline{y})$  fails due to timeout or convergence failure)

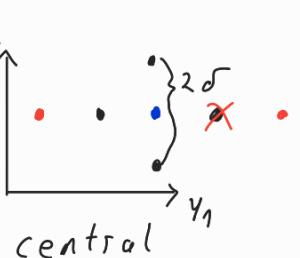
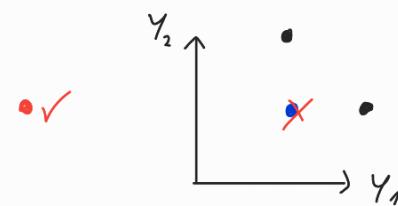
- If a forward point fails, then  $\delta$  is doubled for the current finite difference. (maximum 2 times)
- If the central point fails, then the central difference is used. ( $\delta$  doubling max 2 times)
- If everything fails, the last step is repeated.

(maximum 2 times, otherwise GD terminates)

- $\underline{y}$       • evaluated for gradient (always)      • evaluated in case of errors



forward



central

## - step size control (first-step, min-step)

A step is always opposite to the local gradient:

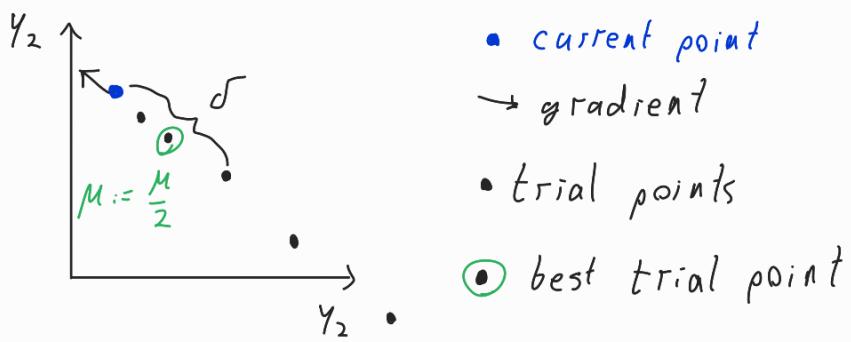
$\mu$ : step size,  $\mu \in (0, 1)$

$\underline{y}[k]$ :  $k^{\text{th}}$  step

$$\underline{y}[k+1] = \underline{y}[k] - \mu \cdot \frac{\underline{\text{grad}}(g(\underline{y}[k]))}{|\underline{\text{grad}}(g(\underline{y}[k]))|}$$

gradient is normalized

Initially,  $\mu = \text{first-step}$ . During each time, 5 trial steps are taken:  $\frac{1}{4}\mu, \frac{1}{2}\mu, \mu, 2\mu, 4\mu$ . Function  $g()$  is evaluated in each of the trial points, and the step size corresponding to the smallest value is used:

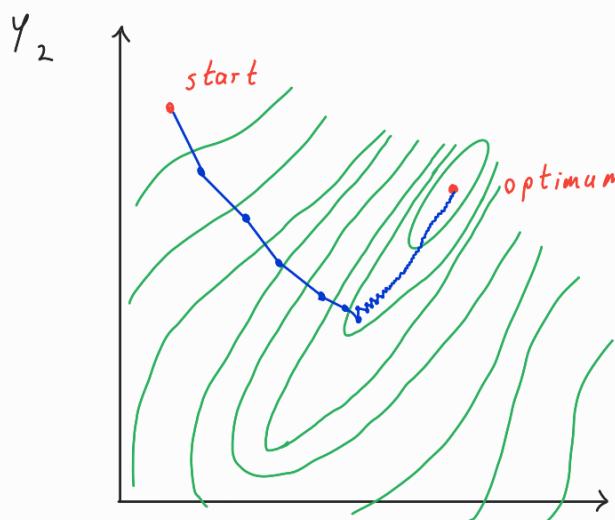


This way, the step size can both decrease for finer tuning (e.g. near the optimum) and increase for faster movement.

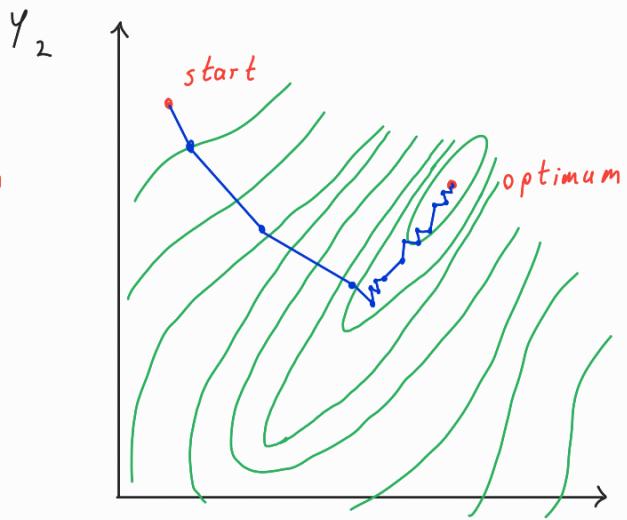
In narrow elongated valleys, the GD quickly moves into the valley, but doing so, the step size is greatly reduced.

Thus, a simpler reducing step size control would barely move forward, and even this adaptive step size control requires several steps:

decreasing step size control



adaptive step size control

- why norm the gradient?

In the original GD (e.g. stochastic gradient descent in machine learning) there is no step size control, the size of the gradient is used instead.

On the other hand, energy-demand can change several orders of magnitude during one step, while it can improve less than a percent during other steps:

- **wrong places**: energy-demand  $\sim 10^{10} - \infty \text{ MJ/kg}$ , the resultant target-specie is zero or numeric noise. The gradient is unreliable if it doesn't fail. Shouldn't start or venture in here.
- **steep places**: energy-demand  $\sim 10^3 - 10^{10} \text{ MJ/kg}$  and it can change an order of magnitude per step. Gradient is huge.
- **walley**: energy-demand  $\sim 10^2 - 10^5 \text{ MJ/kg}$ . Close the optimum, the step size is small, and a step means less than 1% improvement. Gradient is small.

#### - STOP conditions (min-step, step-limit)

The GD search will terminate if

- the step size is reduced below min-step  
(normal operation, the GD is near an optimum or a very fine detail)
- step count is greater than step\_limit  
(search terminates before reaching any optimum)
- calculation of the gradient fails repeatedly  
(fatal error, GD is possibly not even remotely close to an optimum — this actually never happened before)

#### - Steps are leaving ranges

When GD would leave the boundary, the step is squized to remain within ranges.

$y_i \in [0+10\delta, 1-10\delta]$  is guaranteed.

## • How to initiate start-point?

A GD search should not start in wrong places, where the gradient cannot be approximated well. Also, GD will move towards the nearest local optimum. To find the global optimum, you should start from several spread out points. These points don't have to be near the optimum,  $10^5$ - $10^{10}$  MJ/kg is also sufficient. You have 2 approaches:

- Do a coarse brute force parameter study along the more important control parameters. Use the best 5-100 results as start points.
- (As in example) Evaluate in several random trial point, and choose the best 5-100. I recommend 200-500 trial points in 2-3 dimensions, 1000-10000 with more dimensions.

## • Settings of gradient-descent()

gradient-descent()

ranges, - according to 1<sup>st</sup> chapter

path, - location to save resulting CSV files. If None, a list of datas is also returned.

to\_optimize, - key in data dict to minimize, usually 'energy-demand'

start\_point, - according to last chapter

step-limit, ↗ between 50-200, if you want results fast

300-500, if you want most searches to converge even higher, if the majority of searches still don't fully converge

first\_step, - you shouldn't step over the optimum. As the step size can grow, you may start small, e.g. 0.001

`min_step`, — determines the final precision of the optimum (see 1<sup>st</sup> chapter). A small `min_step` leads to much longer searches, a too large leads to premature termination.  
 (Near the optimum, in the "valley", very small steps are taken.)  
 I recommend  $10^{-4}$  for fast results,  $10^{-5}$  for good convergence.

`delta`, — should be an order of magnitude smaller, than `min_step`.  
 should be greater than  $10^{-6}$  to avoid numeric noise.

`t-int`, — simulation time interval, default is  $[0, 1]$  sec.  
 the bubble should settle, after that the computational cost is negligible due to adaptive time steps.

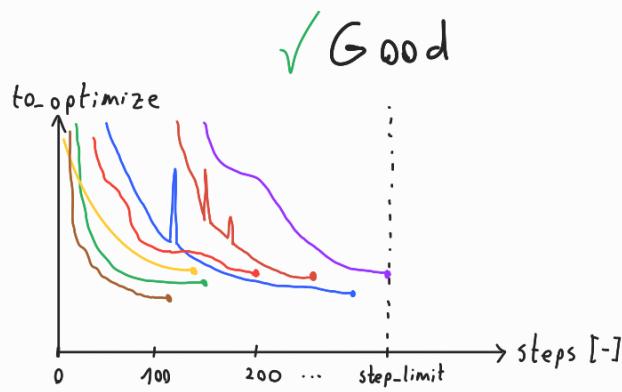
`LSODA-timeout`, — You should have a rough idea about the runtimes.  
 If you observe too much timeouts in `Read-CSV-files.ipynb`  
`Radau-timeout`, — then increase them. Use `Radau-timeout=0` for faster results.

`verbose` — prints debug data, if True.  
 use only with one GA at a time

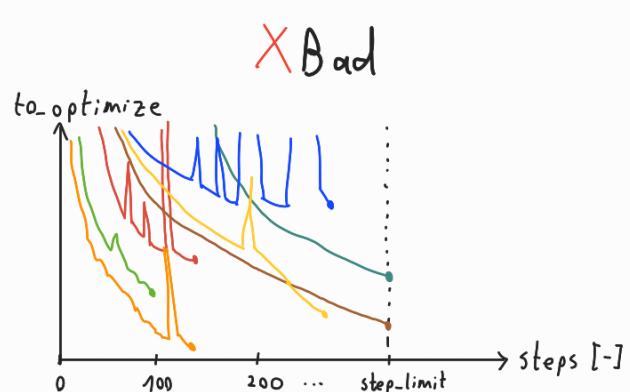
- Verify convergence

- Convergence of `to-optimize`

Included in `Gradient-descent.ipynb` example. This plot shows, how the value of `to-optimize` decreases as function of `steps`:



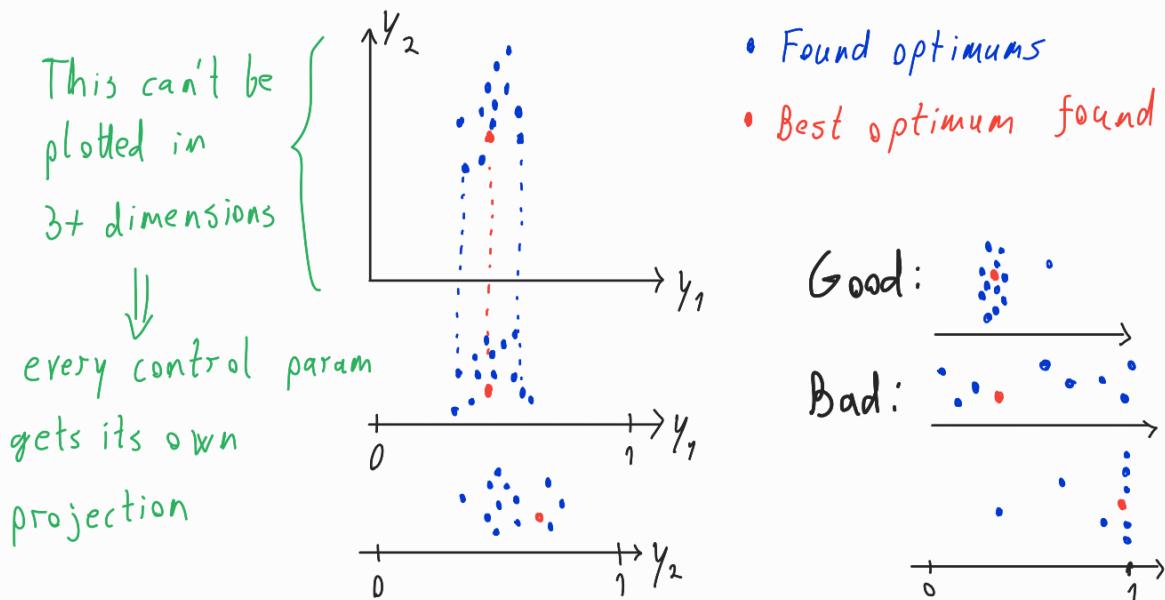
- most searches terminate before reaching `step_limit`
- searches seemingly converge (flatten)
- the final values of `to-optimize` are close to each other



- searches reach `step_limit` (`step_limit` is too small)
- searches terminate, while `to-optimize` is still rapidly decreasing (too large `min-step` or too small timeouts)
- searches diverge (too large/small `min-step`)

## - Locations of optimums

Included in Gradient\_descent.ipynb example. This plots show the deviation of the found optimums (last steps)



Most points should be close to each other, and not at the edge of ranges.

## - Quality of an optimum

Included in Create\_plots.ipynb example. On these plots, all control parameters are fixed (at the optimum), and only one of them is changed at a time.

