

## Gradient Descent

© Malay K. Das, 210 Southern Lab, ph-7359, [mkdas@iitk.ac.in](mailto:mkdas@iitk.ac.in)

**Office hours:** W 1030-1130, SL-210

**Previously:**

Steepest descent: theory

**Today:**

Gradient descent, Newton's method: example

**HW due:**

September 03, 2024

**Computing quiz:**

September 04, 2024

1

Malay K. Das, [mkdas@iitk.ac.in](mailto:mkdas@iitk.ac.in)

## Unconstrained Optimization

We wish to minimize  $y = f(\mathbf{x})$       $\mathbf{x} \in \mathbb{R}^n; y \in \mathbb{R}$

search for minimum ——— search in calculated direction with optimum step

**(Gradient Descent Methods)**

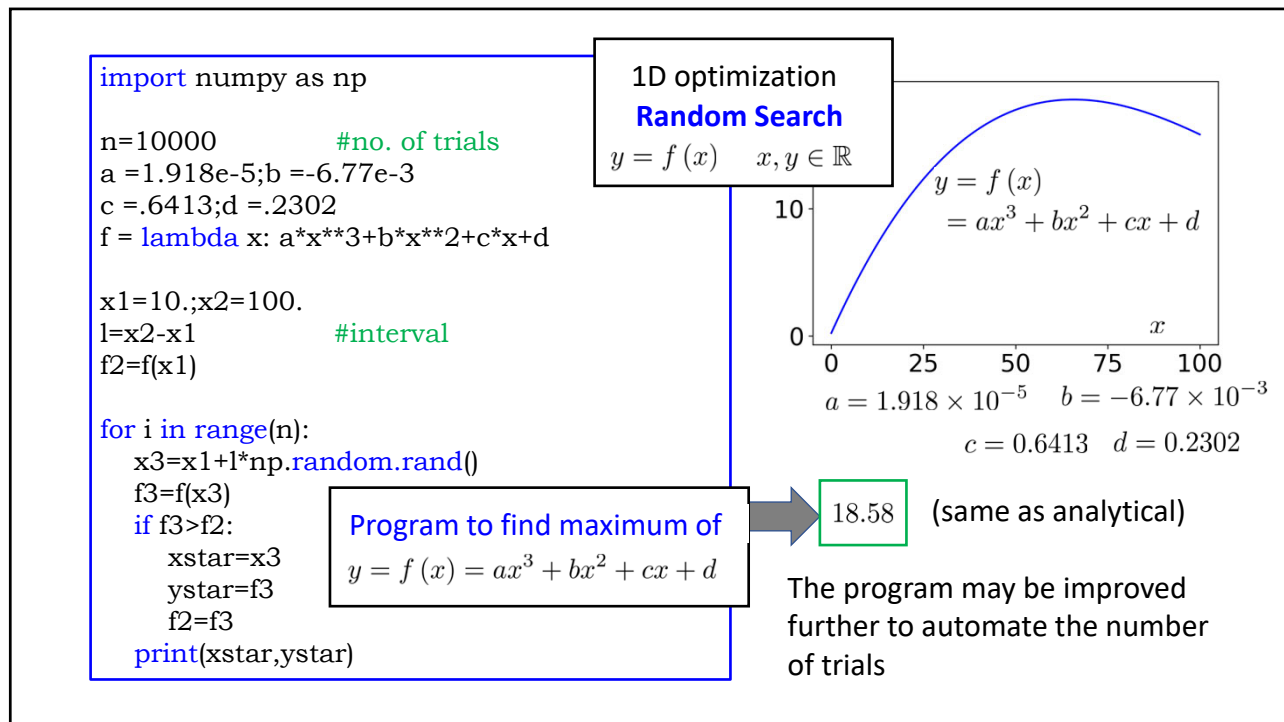
- steepest descent, Newton's method etc.

random search **(Metaheuristic Methods)**

- simulated annealing, genetic algorithm etc.

Both the above approaches include variety of techniques

2



3

Malay K. Das, mkdas@iitk.ac.in

**Gradient Descent Methods:** methods using gradient information for optimization

To find  $\mathbf{x}_* = \arg \min_{\mathbf{x}} f(\mathbf{x})$  starting from an initial guess  $\mathbf{x}^{(0)}$

we iteratively update  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k$        $\mathbf{d}_k$  : descent (search) direction

such that  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$        $h_k$  : step size (+ve scalar)

To select a good  $\mathbf{d}$ , we expand in (multivariable) Taylor series

$$f(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)}) + h_k \mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) + \mathcal{O}(h_k^2)$$

neglecting higher-order terms, change in the function  $f$      $\delta f \approx h_k \mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)})$

for minimization of  $f$ , we enforce  $\delta f < 0 \Rightarrow \mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) < 0$     fundamental theory gradient descent

Different choice of  $\mathbf{d}_k, h_k$  leads to different methods    directional derivative in the direction  $\mathbf{d}_k$

4

consider  $f(x, y) = x^2 + 2y^2 - 2xy - 2x$   $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

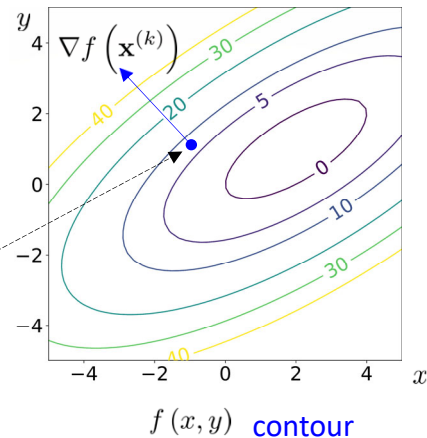
$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x - 2y - 2 \\ 4y - 2x \end{bmatrix}$$

$$\nabla f(\mathbf{x}^{(k)}) = 6 \begin{bmatrix} -1 & 1 \end{bmatrix}^T$$

at  $k$ -th iteration

$$\mathbf{x}^{(k)} = \begin{bmatrix} -1 & 1 \end{bmatrix}^T$$

$$f(-1, 1) = 7$$



for descent direction  $\mathbf{d}_k$   $\mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) < 0$

let us assume  $\mathbf{d}_k = [d_1 \ d_2]^T \Rightarrow \mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) = 6(d_2 - d_1)$

For  $d_2 < d_1$   $\mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) < 0$  and, therefore,  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$

Therefore, all  $\mathbf{d}_k = [d_1 \ d_2]^T$  with  $d_2 < d_1$  are valid descent directions

5

Malay K. Das, mkdas@iitk.ac.in

### Steepest Descent method

$$\mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) < 0 \Rightarrow \|\mathbf{d}_k\| \|\nabla f(\mathbf{x}^{(k)})\| \cos \theta < 0 \quad \theta \text{ is the angle between } \mathbf{d}_k, \nabla f(\mathbf{x}^{(k)})$$

$$\Rightarrow \cos \theta < 0 \quad \text{minimum} \quad \cos \theta = \cos(\pi) = -1 \quad (\text{easy to visualize in 2D})$$

in steepest descent method, the angle between

$\mathbf{d}_k$  and  $\nabla f(\mathbf{x}^{(k)})$  is  $\pi$

Thus  $\mathbf{d}_k = -\nabla f(\mathbf{x}^{(k)})$

$$\mathbf{x}_* = \arg \min_{\mathbf{x}} f(\mathbf{x})$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k$$

$$\mathbf{d}_k^T \nabla f(\mathbf{x}^{(k)}) < 0 \quad h_k > 0$$

The 'rate of minimization' is at its highest at this direction

Rate of minimization refers to the reduction in the value of the function over unit distance

Now  $h_k$  may be obtained from any suitable Line Search technique

6

### Steepest Descent method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k \quad \mathbf{d}_k = -\nabla f(\mathbf{x}^{(k)})$$

$$\mathbf{x}_* = \arg \min_{\mathbf{x}} f(\mathbf{x})$$

$\mathbf{d}_k$  : descent (search) direction

**Line Search** to find  $h_k$

$h_k$  : step size (+ve scalar)

since we now know  $\mathbf{x}^{(k)}, \mathbf{d}_k$

$$f(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)} + h_k \mathbf{d}_k) = \phi(h_k) \quad \text{is a function of } h_k \text{ only}$$

We define  $h_k = \arg \min_h \phi(h)$  and minimize  $\phi(h)$  using 1D optimization techniques

When minimization of  $\phi(h)$  becomes computationally expensive, alternatives methods are used

alternatives methods partially minimize  $\phi(h)$

7

Malay K. Das, mkdas@iitk.ac.in

**Example:** minimize  $f(x, y) = x^2 + 2y^2 - 2xy - 2x$

initial guess  $\mathbf{x}^{(0)} = [-1 \ 1]^T$

$$\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

**steepest descent**

$$\mathbf{d}_0 = -\nabla f(\mathbf{x}^{(0)}) = -6 \begin{bmatrix} -1 & 1 \end{bmatrix}^T$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k$$

$$\begin{aligned} \phi(h) &= f(\mathbf{x}^{(0)} + h \mathbf{d}_0) \quad \mathbf{x}^{(0)} + h \mathbf{d}_0 = (1 - 6h) \begin{bmatrix} -1 \\ 1 \end{bmatrix} \\ &= 5(1 - 6h)^2 + 2(1 - 6h) \end{aligned}$$

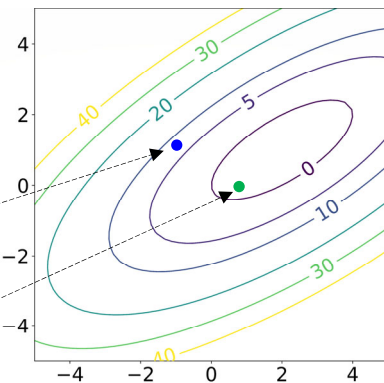
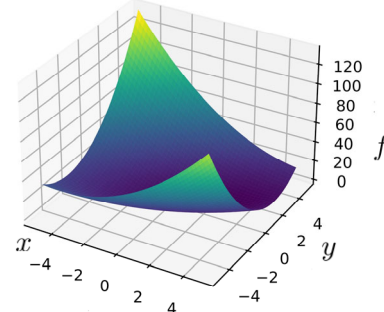
$$\phi'(h) = -60(1 - 6h) - 2 \quad \phi'(h = h_0) = 0 \Rightarrow h_0 = 0.2$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + h_0 \mathbf{d}_0 = [0.2 \ -0.2]^T$$

$$f(\mathbf{x}^{(0)}) = 7 \quad f(\mathbf{x}^{(1)}) = -0.36$$

**solution**

**converges when**  $\|\mathbf{d}_k\| < \epsilon$  **tolerance**



8

### Steepest Descent Method: Summary

To find  $\mathbf{x}_* = \arg \min_{\mathbf{x}} f(\mathbf{x})$  starting from an initial guess  $\mathbf{x}^{(0)}$

we iteratively update  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k$        $\mathbf{d}_k$  : descent (search) direction  
in steepest descent method  $\mathbf{d}_k = -\nabla f(\mathbf{x}^{(k)})$        $h_k$  : step size (+ve scalar)

$h_k$  may be obtained by various **Line Search** techniques, such as

$$h_k = \arg \min_h \phi(h) \quad \text{where} \quad \phi(h) = f(\mathbf{x}^{(k)} + h \mathbf{d}_k)$$

Such step size calculations are difficult to implement, sometimes **constant step size** is used in computer program

9

Malay K. Das, mkdas@iitk.ac.in

### Exercise

Minimize  $f(x, y) = x^2 + 2y^2 - 2xy - 2x$

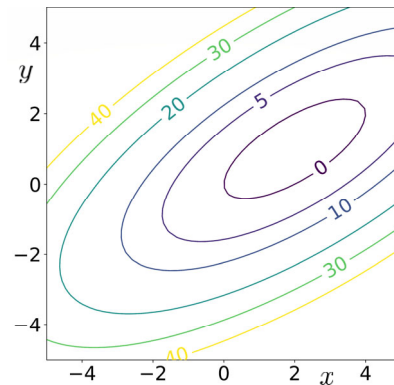
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k \quad \mathbf{d}_k = -\nabla f(\mathbf{x}^{(k)}) \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x - 2y - 2 \\ 4y - 2x \end{bmatrix}$$

Let's first try using a constant step size  $h_k = 0.01$

we know the analytical solution

$$\min f(x, y) = -2 \quad \text{at} \quad \begin{bmatrix} x_* \\ y_* \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



$f(x, y)$  contour

10

```
import numpy as np
```

```
maxit=10000;epsilon=1.e-6;stepsize=.01
```

```
objective_function= lambda x: x[0]**2 + 2*x[1]**2-2*x[0]*x[1]-2*x[0]
```

```
gradient_function= lambda x: np.array([2 * x[0]-2*x[1]-2, 4 * x[1]-2*x[0]])
```

```
x=np.array([-1.,1.]) #initial guess
```

```
for i in range(maxit):
```

```
    gradient=gradient_function(x);b=np.linalg.norm(gradient)
```

```
    if b < epsilon:
```

```
        break
```

```
    x -= stepsize * gradient
```

```
    print(i,x,b)
```

```
minimum_value = objective_function(x)
```

```
print("Minimum value:", minimum_value)
```

```
print("Minimum location:", x)
```

constant step size

objective function

gradient of  
objective function

norm of gradient;  
used as stopping  
criterion

**Program output:** Minimum value: -1.999999, Minimum location: [1.999999,.999999]

**Analytical:** Minimum value: -2, Minimum location: [2,1]

11

Malay K. Das, mkdas@iitk.ac.in

### Backtracking Line Search: Armijo condition

we enforce  $f(x^{(k)} + h\mathbf{d}_k) < f(x^{(k)})$  but do not minimize  $f(x^{(k)} + h\mathbf{d}_k)$

Armijo condition suggests

$$f(x^{(k)} + h\mathbf{d}_k) \leq f(x^{(k)}) + \beta h \mathbf{d}_k^T \nabla f(x^{(k)}) \quad \text{for some constant } \beta \in (0, 1)$$

With guess small values of  $\beta \in (0, 1)$   $\tau \in (0, 1)$  initial guess of  $h = 1$  (usually)

1. check the Armijo condition is satisfied or not

2. if not satisfied  $h \leftarrow \tau h$   $\tau \in (0, 1)$

typical value of  $\tau = 0.5$

3. continue until the Armijo condition is satisfied

12

### Recall 1D Newton-Raphson Method

To minimize  $f(x)$

we begin with an initial guess  $x^{(0)}$

$$x^{(1)} = x^{(0)} + h_0 d_0 \quad d_0 = -f'(x^{(0)})$$

The step size  $h_0$  hopefully minimizes  $f(x)$

$$\Rightarrow f'(x = x^{(1)}) = 0$$

$d_0$  : descent direction

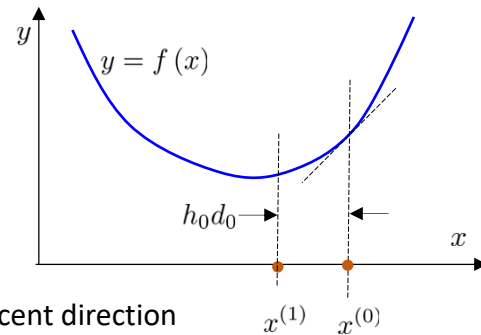
$h_0$  : step size

Expanding in Taylor series

$$f'(x = x^{(1)}) = 0 = f'(x^{(0)} + h_0 d_0) = f'(x^{(0)}) + h_0 d_0 f''(x^{(0)}) + \dots$$

$$\Rightarrow h_0 = -\frac{f'(x^{(0)})}{d_0 f''(x^{(0)})} = \frac{1}{f''(x^{(0)})}$$

$$\Rightarrow x^{(1)} = x^{(0)} - \frac{f'(x^{(0)})}{f''(x^{(0)})}$$



in general

$$x^{(k+1)} = x^{(k)} - \left( \frac{f'}{f''} \right)_{x=x^{(k)}}$$

13

Malay K. Das, mkdas@iitk.ac.in

### Multidimensional Newton's (or Newton-Raphson) Method

Extension of 1D Newton-Raphson method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + h_k \mathbf{d}_k$$

Recall 1D  $x^{(k+1)} = x^{(k)} - \left( \frac{f'}{f''} \right)_{x=x^{(k)}}$

$\mathbf{d}_k$  : descent direction

$h_k$  : step size

Following the same argument

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left[ \mathbf{H}f(\mathbf{x}^{(k)}) \right]^{-1} \nabla f(\mathbf{x}^{(k)})$$

As in 1D technique, the method may not converge for certain initial guesses, and for highly nonlinear functions

As in 1D technique, when converges, shows rapid convergence otherwise

To improve convergence, the method is often modified using Armijo condition

14

## Backtracking Line Search

Armijo condition suggests

for some constant

$$f(x^{(k+1)}) - f(x^{(k)}) \leq \beta h \left[ \mathbf{H}f(\mathbf{x}^{(k)}) \right]^{-1} \nabla f(\mathbf{x}^{(k)}) \quad \beta \in (0, 1)$$

With a guess small values of  $\beta \in (0, 1)$   $\tau \in (0, 1)$  initial guess of  $h = 1$

1. check the Armijo condition is satisfied or not
2. if not satisfied  $h \leftarrow \tau h$   $\tau \in (0, 1)$  typical value of  $\tau = 0.5$
3. continue until the Armijo condition is satisfied

Once Armijo condition is satisfied, update

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - h \left[ \mathbf{H}f(\mathbf{x}^{(k)}) \right]^{-1} \nabla f(\mathbf{x}^{(k)})$$

15

Malay K. Das, mkdas@iitk.ac.in

## Exercise

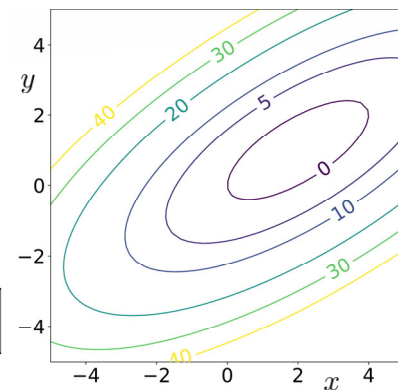
Minimize  $f(x, y) = x^2 + 2y^2 - 2xy - 2x$   $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - h \left[ \mathbf{H}f(\mathbf{x}^{(k)}) \right]^{-1} \nabla f(\mathbf{x}^{(k)})$$

$$\begin{aligned} \nabla f &= \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{bmatrix}^T & \mathbf{H}f &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2 & -2 \\ -2 & 4 \end{bmatrix} \\ &= \begin{bmatrix} 2x - 2y - 2 \\ 4y - 2x \end{bmatrix} \end{aligned}$$

$$(\mathbf{H}f)^{-1} = \frac{1}{\det(\mathbf{H}f)} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(\mathbf{H}f)^{-1} \nabla f = \frac{1}{12} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 2x - 2y - 2 \\ 4y - 2x \end{bmatrix} = \frac{1}{3} \begin{bmatrix} x - 2 \\ y - 1 \end{bmatrix}$$



$f(x, y)$  contour

analytical solution

$$\min f(x, y) = -2$$

$$\text{at } \begin{bmatrix} x_* \\ y_* \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

16