# Lecture 23 : Accelerating Gradient Descent (Use momentum)

- Gradient descent $\quad x_{k+1} = x_k - s\nabla f(x_k)$

- **Accelerated gradient descent**
  ① **momentum added**
  ② Nesterov formula
- Stochastic gradient descent

**Momentum**

$$\boxed{\begin{array}{l} x_{k+1} = x_k - s z_k \\[6pt] z_k = \nabla f_k + \beta z_{k-1} \end{array}}$$

$f = \frac{1}{2} x^T S x$
$\nabla f = S x$

$\longrightarrow$

$$\boxed{\begin{array}{l} x_{k+1} = x_k - s z_k \\[6pt] z_{k+1} - s x_{k+1} = \beta z_k \end{array}}$$

$$\begin{bmatrix} 1 & 0 \\ -s & 1 \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}_{k+1} = \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}_k$$

$$\begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

$$\begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -s \\ \lambda & \beta - \lambda s \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} = R \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

Simple, beautiful steps come from tracking one eigenvector, which makes the whole problem scalar.

$x_k$ ↓ $\qquad$ $x_{k+1}$ ↓

$$\begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} q = \begin{bmatrix} -s & 1 \\ \beta & 0 \end{bmatrix} \begin{bmatrix} d_k \\ c_k \end{bmatrix} q$$

$c_{k+1} q = c_k q - s d_k q$

$d_{k+1} q - \lambda c_{k+1} q = \beta d_k q$

$Sq = \lambda q$

Suppose $x_k = c_k q$

$x_k$ is tracking eigenvector

$\qquad z_k = d_k q$

$S x_k = c_k \lambda q$

↓

$\nabla f_k$

$\begin{bmatrix} c_k \\ d_k \end{bmatrix}$ is multiplied by $R$ every step

$$\begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 5 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 5x + 2y \\ 2x + 3y \end{bmatrix} = 5x^2 + 4xy + 3y^2$$

How to choose $\lambda$ ?

$$0 < m \leq \lambda \leq M$$

↳ eigenvalue of S

$\dfrac{M}{m} = K =$ condition #of S

$e_1, e_2 =$ eigenvalues of R

depends on eigenvalues of S

※ Choose s and $\beta$ to minimize $\max \left[ |e_1(\lambda)|, |e_2(\lambda)| \right]$

for $\lambda_{min}(S) \leq \lambda \leq \lambda_{max}(S)$

$m \leq \lambda \leq M$

It turned out that:

$$S_{optimal} = \left( \frac{2}{\sqrt{M} + \sqrt{m}} \right)^2 = \left( \frac{2}{1 + \sqrt{b}} \right)^2$$

$$\beta_{optimal} = \left( \frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}} \right)^2 = \left( \frac{1 - \sqrt{b}}{1 + \sqrt{b}} \right)^2$$

for 2x2 example

$$S = \begin{bmatrix} 1 & 0 \\ 0 & b \end{bmatrix}$$

$$\left| \begin{matrix} \text{eigenvalues} \\ \text{of R} \end{matrix} \right| < \left( \frac{1 - \sqrt{b}}{1 + \sqrt{b}} \right)^2$$