

# 7. Symmetric Positive Definite Matrices

---

- Test for Positive Definite Matrix
  - All eigenvalues of  $S$  are positive ( $S$ : real symmetric matrix)
  - Energy  $x^T S x$  is positive for  $x \neq 0$  (best!)
  - $S = A^T A$  with independent columns of  $A$
  - All the pivots are positive
  - All the leading determinants are positive
- Second derivative matrix is positive definite @ a minimum point
- Semidefinite allows zero eigenvalues/energy/pivots/determinants

$\mathbf{S} = \mathbf{S}^T \rightarrow \begin{cases} \text{all } n \text{ eigenvalues } \lambda \text{ are real numbers} \\ \text{all } n \text{ eigenvectors } \mathbf{q} \text{ can be chosen orthogonal (perpendicular each other)} \end{cases}$

$\mathbf{S} = \mathbf{I} \rightarrow \text{all } \lambda = 1, \text{ every nonzero vector } \mathbf{x} \ (\mathbf{Ix} = \mathbf{1x})$

$$\mathbf{q}(\text{unit length}) \rightarrow \mathbf{Q} = \begin{bmatrix} q_1 & \cdots & q_n \end{bmatrix} \xrightarrow{\text{orthonormal}} \begin{cases} \mathbf{Q}^T \mathbf{Q} = \mathbf{I} \\ \mathbf{Q}^T = \mathbf{Q}^{-1} \end{cases}$$

(Spectral Theorem) Every real symmetric matrix has the form  $\mathbf{S} = \mathbf{Q}\Lambda\mathbf{Q}^T$  ( $\leftarrow \mathbf{A} = \mathbf{X}\Lambda\mathbf{X}^{-1}$ ) and every matrix of that form is symmetric.

(quick proofs for orthogonal eigenvectors and real eigenvalues)

$\begin{cases} \text{nonzero / zero eigenvalues: } \mathbf{Sx} = \lambda\mathbf{x} \text{ and } \mathbf{Sy} = 0\mathbf{y} \\ \text{nonzero eigenvalues: } \mathbf{Sx} = \lambda\mathbf{x} \text{ and } \mathbf{Sy} = \alpha\mathbf{y} \end{cases} \rightarrow \text{orthogonal eigenvector?}$   
 complex ?: multiply complex conjugate vector  $\bar{\mathbf{x}}^T \rightarrow \bar{\mathbf{x}}^T \mathbf{Sx} = \lambda \bar{\mathbf{x}}^T \mathbf{x} \rightarrow \text{real eigenvalue?}$

$$\mathbf{S} = \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix} = \bar{\mathbf{S}}^T \rightarrow \lambda = 8, -1 \text{ and } \mathbf{x} = \begin{bmatrix} 1 \\ 1+i \end{bmatrix}, \begin{bmatrix} 1-i \\ -1 \end{bmatrix}$$

$\begin{cases} \text{Test 1: A positive definite matrix has all positive eigenvalues.} \\ \text{Test 2: } \mathbf{S} \text{ is positive definite if the energy } \mathbf{x}^T \mathbf{S} \mathbf{x} \text{ is positive for all vectors } \mathbf{x} \neq \mathbf{0} \end{cases}$

$$\mathbf{x}^T \mathbf{S} \mathbf{x} > 0 \longleftrightarrow \lambda > 0 : \mathbf{S} \mathbf{x} = \lambda \mathbf{x} \rightarrow \mathbf{x}^T \mathbf{S} \mathbf{x} = \lambda \mathbf{x}^T \mathbf{x}, \lambda > 0 \rightarrow \mathbf{x}^T \mathbf{S} \mathbf{x} > 0$$

If  $\mathbf{x}^T \mathbf{S} \mathbf{x} > 0$  for the eigenvectors of  $\mathbf{S}$ , then  $\mathbf{x}^T \mathbf{S} \mathbf{x} > 0$  for every nonzero vector  $\mathbf{x}$ .

why?  $\mathbf{x} = c_1 \mathbf{x}_1 + \dots + c_n \mathbf{x}_n \rightarrow \mathbf{x}^T \mathbf{S} \mathbf{x} = c_1^2 \lambda_1 \mathbf{x}_1^T \mathbf{x}_1 + \dots + c_n^2 \lambda_n \mathbf{x}_n^T \mathbf{x}_n > 0$  if every  $\lambda_i > 0$

If  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are symmetric positive definite, so is  $\mathbf{S}_1 + \mathbf{S}_2$

**Test 3:**  $\mathbf{S} = \mathbf{A}^T \mathbf{A}$  for a matrix  $\mathbf{A}$  with independent columns

$$\mathbf{x}^T \mathbf{S} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^T \mathbf{A} \mathbf{x} = \|\mathbf{A} \mathbf{x}\|^2 \geq 0, \mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \rightarrow \mathbf{S} ?$$

$\begin{cases} \text{Test 4: All the leading determinants } D_1, D_2, \dots, D_n \text{ of } \mathbf{S} \text{ are positive} \\ \text{Test 5: All the pivots of } \mathbf{S} \text{ are positive (in elimination)} \end{cases}$

The  $k$ th pivot equals the ratio  $\frac{D_k}{D_{k-1}}$  of the leading determinants (size  $k$  and  $k-1$ )

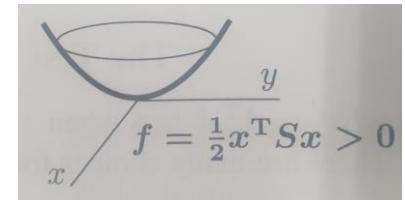
$$\mathbf{S} = \mathbf{L} \mathbf{U} = \mathbf{L} \mathbf{D} \mathbf{L}^T = \mathbf{A}^T \mathbf{A} \rightarrow \mathbf{A} = \sqrt{\mathbf{D}} \mathbf{L}^T \text{ (Cholesky factorization)}$$

$$\mathbf{S} = \mathbf{Q} \Lambda \mathbf{Q}^T \rightarrow \mathbf{A} = \mathbf{Q} \sqrt{\Lambda} \mathbf{Q}^T = \mathbf{A}^T$$

$$\mathbf{S} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \Rightarrow \left\{ \begin{array}{l} D_1 = 2, D_2 = 3, D_4 = 4 \\ \xrightarrow{l_{21} = \frac{-1}{2}, l_{31} = 0} \begin{bmatrix} 2 & -1 & 0 \\ 0 & 3/2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \xrightarrow{l_{32} = \frac{-2}{3}} \begin{bmatrix} 2 & -1 & 0 \\ 0 & 3/2 & -1 \\ 0 & 0 & 4/3 \end{bmatrix} \end{array} \right.$$

$$\mathbf{S} = \begin{bmatrix} 1 & & \\ -1/2 & 1 & \\ 0 & -2/3 & 1 \end{bmatrix} \begin{bmatrix} 2 & -1 & 0 \\ 3/2 & -1 & \\ 4/3 & & \end{bmatrix} = \mathbf{L}\mathbf{U} = \begin{bmatrix} 1 & & \\ -1/2 & 1 & \\ 0 & -2/3 & 1 \end{bmatrix} \begin{bmatrix} 2 & & \\ 3/2 & & \\ 4/3 & & \end{bmatrix} \begin{bmatrix} 1 & -1/2 & 0 \\ 1 & & -2/3 \\ 1 & & \end{bmatrix} = \mathbf{LDL}^T$$

$$= \begin{bmatrix} \sqrt{2} & & \\ -\sqrt{1/2} & \sqrt{3/2} & \\ 0 & -\sqrt{2/3} & \sqrt{4/3} \end{bmatrix} \begin{bmatrix} \sqrt{2} & -\sqrt{1/2} & 0 \\ \sqrt{3/2} & -\sqrt{2/3} & \\ \sqrt{4/3} & & \end{bmatrix} = \mathbf{A}^T \mathbf{A}$$



$$\mathbf{S} = \begin{bmatrix} a & b \\ b & c \end{bmatrix} \rightarrow \left\{ \begin{array}{l} \text{determinants: } a > 0, ac - b^2 > 0 \\ \text{pivots: } a > 0, c - b \frac{b}{a} = \frac{ac - b^2}{a} > 0 \\ \text{eigenvalues: } \lambda_1 > 0, \lambda_2 > 0 \\ \text{energy: } ax^2 + 2bxy + cy^2 > 0 \end{array} \right. \xrightarrow{\frac{a=c=5}{b=4}} \left\{ \begin{array}{l} E = \mathbf{x}^T \mathbf{S} \mathbf{x} = [x \ y] \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 5x^2 + 8xy + 5y^2 > 0 \\ \rightarrow \text{bowl opening upwards, strictly convex} \\ \text{Test: matrix of second derivatives is positive definite at all points} \end{array} \right.$$

Test for a minimum (maximum, saddle point)

$$f(x): \frac{df}{dx} = 0 \text{ and } \frac{d^2f}{dx^2} > 0 \text{ at } x = x_0$$

$$f(x, y): \left( \frac{df}{dx} = 0, \frac{df}{dy} = 0 \right) \text{ and } \begin{bmatrix} \partial^2 f / \partial x^2 & \partial^2 f / \partial x \partial y \\ \partial^2 f / \partial x \partial y & \partial^2 f / \partial y^2 \end{bmatrix} \text{ is positive definite at } x_0, y_0$$

Optimization and Machine Learning

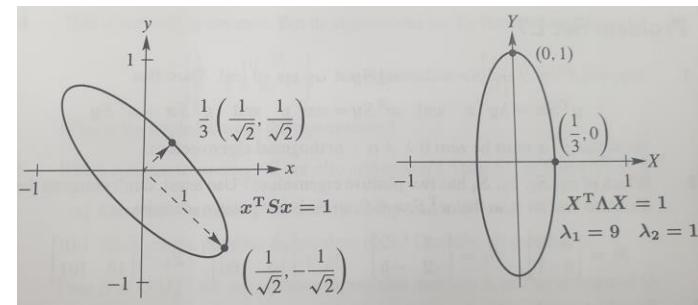
$$\begin{cases} \text{Calculus: partial derivatives of } f \text{ are all zero at } \mathbf{x}^* \rightarrow \frac{\partial f}{\partial x_i} = 0 \rightarrow \text{direction to move} \\ \text{Linear Algebra: matrix } \mathbf{S} \text{ of second derivatives is positive definite} \end{cases}$$

The graph (cut through the bowl) of  $\mathbf{x}^T \mathbf{S} \mathbf{x} = 1$  is an ellipse , with its axes pointing along the eigenvectors of  $\mathbf{S}$

$$\mathbf{x}^T \mathbf{S} \mathbf{x} = \mathbf{x}^T (\mathbf{Q} \Lambda \mathbf{Q}^T) \mathbf{x} = (\mathbf{x}^T \mathbf{Q}) \Lambda (\mathbf{Q}^T \mathbf{x}) = [x \ y] \mathbf{Q} \Lambda \mathbf{Q}^T \begin{bmatrix} x \\ y \end{bmatrix} = [X \ Y] \Lambda \begin{bmatrix} X \\ Y \end{bmatrix}$$

$$\mathbf{S} = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 9 & 1 \\ 1 & 1 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$5x^2 + 8xy + 5y^2 = 1 \rightarrow 9 \left( \frac{x+y}{\sqrt{2}} \right)^2 + 1 \left( \frac{x-y}{\sqrt{2}} \right)^2 = 1 \rightarrow 9X^2 + Y^2 = 1$$



# Multivariate Calculus

one function / one variable:  $F(x + \Delta x) \approx F(x) + (\Delta x) \frac{dF}{dx} + \frac{1}{2}(\Delta x)^2 \frac{d^2 F}{dx^2}$

one function /  $n$  variables:  $F(\mathbf{x} + \Delta \mathbf{x}) \approx F(\mathbf{x}) + (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) + \frac{1}{2}(\Delta \mathbf{x})^T \mathbf{H}(\mathbf{x})(\Delta \mathbf{x})$

$$\mathbf{x} = (x_1, \dots, x_n), \quad \nabla F(\mathbf{x}) = \underbrace{\begin{bmatrix} \frac{\partial F}{\partial x_1} \\ \vdots \\ \frac{\partial F}{\partial x_n} \end{bmatrix}}_{n \text{ derivatives}}, \quad H_{jk} = \underbrace{\frac{\partial^2 F}{\partial x_j \partial x_k}}_{\text{symmetric: } n^2 \rightarrow \left(\frac{1}{2}n^2 + \frac{1}{2}n\right) \text{ derivatives}} = H_{kj}$$

$m$  functions /  $n$  variables:  $\mathbf{f} = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))$

$\mathbf{f}(\mathbf{x} + \Delta \mathbf{x}) \approx \mathbf{f}(\mathbf{x}) + \mathbf{J}(\mathbf{x})(\Delta \mathbf{x})$

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} (\nabla f_1)^T \\ \vdots \\ (\nabla f_m)^T \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}, \quad J_{jk} = \underbrace{\frac{\partial f_j}{\partial x_k}}_{\text{Jacobian matrix}} \leftarrow \text{multivariable calculus}$$

# Optimization (1)

	unknowns	constraints	condition
Applied Mathematics	vectors	matrices	first derivative = zero
Calculus of Variation	functions	integrals	first variation = zero

$\arg \min F(\mathbf{x})$  = value(s) of  $\mathbf{x}$  where  $F$  reaches its minimum =  $\mathbf{x}^*$

minimize  $F(\mathbf{x}) \rightarrow \begin{cases} \text{solving } \nabla F(\mathbf{x}) = 0 \\ \text{solving } \mathbf{f}(\mathbf{x}) [= \nabla F(\mathbf{x})] = 0 \rightarrow n \text{ equations and } n \text{ unknowns} \end{cases}$

1) steepest descent:  $\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla F(\mathbf{x}_k)$

$s_k$  = step size (learning rate)  $\rightarrow$  exact line search for best  $s_k$

2) Newon's method:  $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}(\mathbf{x}_k)^{-1} \nabla F(\mathbf{x}_k)$

$\mathbf{H}(\mathbf{x}_k)$  = Hessian, Jacobian of the gradient

$$\nabla F(\mathbf{x}_{k+1}) \approx \nabla F(\mathbf{x}_k) + \mathbf{H}(\mathbf{x}_k)(\mathbf{x}_{k+1} - \mathbf{x}_k) = 0$$

# Optimization (2)

convergence  $\begin{cases} 1) \text{ steepest descent: linear, } \|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq C \|\mathbf{x}_k - \mathbf{x}^*\| \quad (C < 1) \\ 2) \text{ Newton's method: quadratic, } \|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq C \|\mathbf{x}_k - \mathbf{x}^*\|^2 \end{cases}$

$$\mathbf{f}(\mathbf{x}_{k+1}) = 0 \rightarrow \mathbf{f}(\mathbf{x}_k) + J(\mathbf{x}_k)(\mathbf{x}_{k+1} - \mathbf{x}_k) = 0 \rightarrow \mathbf{x}_{k+1} = \mathbf{x}_k - J(\mathbf{x}_k)^{-1} \mathbf{f}(\mathbf{x}_k)$$

$$\text{example: } f(x) = x^2 - 9 = 0$$

# Gradient Descent: Downhill to a Minimum

"gradient descent" uses the derivatives  $\frac{\partial f}{\partial x_i}$  to find a direction that reduces  $f(\mathbf{x})$

$$z = f(x, y) \rightarrow \nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right), \|\nabla f\| = \sqrt{\left( \frac{\partial f}{\partial x} \right)^2 + \left( \frac{\partial f}{\partial y} \right)^2} = (\text{steepest slope of } f)$$

$$\begin{cases} f = \text{const} \text{ (level set, surface)} \\ \nabla f = 0 \text{ on the surface} \\ \mathbf{H} = \nabla^2 f \rightarrow \begin{cases} \text{positive semi-definite: convexity} \\ \text{positive definite: strict convexity} \end{cases} \end{cases}$$

pure quadratic, unknown  $x, y$

$$f = \frac{1}{2} \mathbf{x}^T \mathbf{S} \mathbf{x} = \frac{1}{2} (x^2 + b y^2) \rightarrow \mathbf{S} = \begin{bmatrix} 1 & 0 \\ 0 & b \end{bmatrix} \xrightarrow{b < 1} \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{1}{b} = (\text{condition number})$$

$$f(x, y) = 2x + 5y \rightarrow \nabla f = \begin{bmatrix} 2 \\ 5 \end{bmatrix} \rightarrow \mathbf{H} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

# Example

$$\left. \begin{array}{l} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{S} \mathbf{x} - \mathbf{a}^T \mathbf{x} (-b) \\ \nabla f(\mathbf{x}) = \mathbf{S} \mathbf{x} - \mathbf{a} \\ \mathbf{H} = \mathbf{S} \end{array} \right\} \rightarrow \begin{cases} \text{minimum of } f(\mathbf{x}) \text{ is at } \mathbf{x}^* = \mathbf{S}^{-1} \mathbf{a} = \arg \min f(\mathbf{x}) \\ f^* = \frac{1}{2} (\mathbf{S}^{-1} \mathbf{a})^T \mathbf{S} (\mathbf{S}^{-1} \mathbf{a}) - \mathbf{a}^T (\mathbf{S}^{-1} \mathbf{a}) \\ = \frac{1}{2} \mathbf{a}^T (\mathbf{S}^{-1} \mathbf{S}) \mathbf{S}^{-1} \mathbf{a} - \mathbf{a}^T \mathbf{S}^{-1} \mathbf{a} = -\frac{1}{2} \mathbf{a}^T \mathbf{S}^{-1} \mathbf{a} \end{cases}$$

$$f(\mathbf{x}) = \log(\det \mathbf{X}) \rightarrow \frac{\partial f}{\partial (\det \mathbf{X})} \frac{\partial (\det \mathbf{X})}{\partial x_{ij}} = \frac{C_{ij}}{\det \mathbf{X}} \rightarrow \nabla f(\mathbf{x}) = \text{entries of } \mathbf{X}^{-1} \rightarrow \nabla^2 f(\mathbf{x})?$$

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \cdots & x_{mn} \end{bmatrix} \rightarrow \begin{cases} \det \mathbf{X} = x_{11} |\text{minor}| - x_{12} |\text{minor}| + \cdots + x_{1n} |\text{minor}| \\ = x_{11} C_{11} - x_{12} C_{12} + \cdots + x_{1n} C_{1n} \end{cases}$$

$$\frac{\partial (\det \mathbf{X})}{\partial x_{11}} = C_{11} \rightarrow x_{11}^{-1} = \frac{C_{11}}{\det \mathbf{X}}$$

# Gradient Descent: Convergence

gradient (steepest) descent:  $\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla f(\mathbf{x}_k)$

$s_k$ : step size  $\rightarrow$    
  $\begin{cases} \text{too big: oscillate all over the place} \\ \text{too small: take too long} \\ \text{exact line search: choose } s_k \text{ to make } f(\mathbf{x}_{k+1}) \text{ a minimum in search direction} \\ \text{bracketing: } s_0, \alpha s_0, \alpha^2 s_0 \end{cases}$

convergence analysis

convexity:  $\mathbf{H} = \nabla^2 f$  has eigenvalues between  $0 < m \leq \lambda \leq M$  at all  $\mathbf{x}$

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + (\nabla f)^T (\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{M}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$$

$$= f(\mathbf{x}_k) - s \|\nabla f\|^2 + \frac{Ms^2}{2} \|\nabla f\|^2 \rightarrow \text{best: } s = \frac{1}{M}$$

$$\left. \begin{array}{l} f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2M} \|\nabla f(\mathbf{x}_k)\|^2 \\ f(\mathbf{x}^*) \geq f(\mathbf{x}_k) - \frac{1}{2m} \|\nabla f(\mathbf{x}_k)\|^2 \end{array} \right\} \rightarrow f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{M}\right) [f(\mathbf{x}_k) - f(\mathbf{x}^*)]$$

# Example with Zig-Zag

$$f = \frac{1}{2}(x^2 + by^2) \rightarrow \begin{bmatrix} x \\ y \end{bmatrix}_{k+1} = \begin{bmatrix} x \\ y \end{bmatrix}_k - s_k \begin{bmatrix} x \\ by \end{bmatrix}_k \text{ with } (x_0, y_0) = (b, 1), f_0 = \frac{1}{2}(b^2 + b)$$

$$(x_k, y_k) = \left( b \left( \frac{b-1}{b+1} \right)^k, \left( \frac{1-b}{b+1} \right)^k \right), f_k = \left( \frac{1-b}{b+1} \right)^{2k} f_0$$

$$\text{small } b? (b=0.1) \begin{bmatrix} 0.1 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0.1 \left( \frac{-0.9}{1.1} \right) \\ \frac{0.9}{1.1} \end{bmatrix} \rightarrow \begin{bmatrix} 0.1 \left( \frac{-0.9}{1.1} \right)^2 \\ \left( \frac{0.9}{1.1} \right)^2 \end{bmatrix}, \text{ zig-zag (bouncing back and forth)}$$

$(1-b)$  is the critical quantity for convergence

accelerated gradient descent  
1. momentum add (heavy ball)  
2. Nesterov formula

↔

stochastic gradient descent  
mini-batch of samples each step  
not one nor millions

# Descent with Momentum

\* momentum: means of dampening oscillations and speeding up the iteration leading to faster convergence, it allows a large range of step size

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla f(\mathbf{x}_k) \xrightarrow{\text{descent with momentum}} \begin{cases} \mathbf{x}_{k+1} = \mathbf{x}_k - s \mathbf{z}_k \\ \mathbf{z}_k = \nabla f(\mathbf{x}_k) + \beta \underbrace{\mathbf{z}_{k-1}}_{\substack{\text{memory} \\ \text{momentum}}} \end{cases} \rightarrow \begin{cases} \mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbf{z}_k \\ \mathbf{z}_{k+1} = \nabla f(\mathbf{x}_{k+1}) + \beta \mathbf{z}_k \end{cases}$$

$$\frac{d^2y}{dt^2} + b \frac{dy}{dt} + ky = 0 \rightarrow \begin{cases} \frac{d}{dt}(y) = 0 \\ \frac{d}{dt}\left(\frac{dy}{dt}\right) = -ky - b \frac{dy}{dt} \end{cases} \rightarrow \frac{d}{dt} \begin{bmatrix} y \\ \frac{dy}{dt} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -k & -b \end{bmatrix} \begin{bmatrix} y \\ \frac{dy}{dt} \end{bmatrix}$$

# Descent with Momentum: Quadratic Model

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{S} \mathbf{x} \rightarrow \nabla f = \mathbf{S} \mathbf{x} \rightarrow \begin{cases} \mathbf{S} \mathbf{q} = \lambda \mathbf{q} \\ \mathbf{x}_k = c_k \mathbf{q} \\ \mathbf{z}_k = d_k \mathbf{q} \end{cases} \rightarrow \nabla f(\mathbf{x}_k) = \mathbf{S} \mathbf{x}_k = \lambda \mathbf{x}_k = \lambda c_k \mathbf{q}$$

$$\left. \begin{array}{l} c_{k+1} = c_k - s d_k \\ d_{k+1} = \lambda c_{k+1} + \beta d_k \end{array} \right\} \rightarrow \begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} \rightarrow \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & -s \\ \lambda & \beta - \lambda s \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} = \mathbf{R} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

both eigenvalues  $e_1$  and  $e_2$  of  $\mathbf{R}$  to be as small as possible

choose  $s$  and  $\beta$  to minimize  $\max [|e_1(\lambda)|, |e_2(\lambda)|]$  for  $\lambda_{\min}(\mathbf{S}) \leq \lambda \leq \lambda_{\max}(\mathbf{S})$

$$\rightarrow s = \left( \frac{2}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^2, \beta = \left( \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^2$$

accelerated descent factor:  $\left( \frac{1-b}{1+b} \right)^2 \rightarrow \left( \frac{1-\sqrt{b}}{1+\sqrt{b}} \right)^2$  when  $b = \frac{1}{100}$ ,  $\left( \frac{0.99}{1.01} \right)^2 = 0.96 \rightarrow \left( \frac{0.9}{1.1} \right)^2 = 0.67$

# Nesterov Acceleration

evaluating the gradient  $\nabla f$  at  $\mathbf{x}_k \rightarrow \mathbf{x}_k + \gamma_k (\mathbf{x}_k - \mathbf{x}_{k-1})$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla f(\mathbf{x}_k) \xrightarrow[\text{Nesterov}]{\gamma_k=\beta} \begin{cases} \mathbf{x}_{k+1} = \mathbf{y}_k - s \nabla f(\mathbf{y}_k) \\ \mathbf{y}_k = \mathbf{x}_k + \beta(\mathbf{x}_k - \mathbf{x}_{k-1}) \end{cases} \rightarrow \begin{cases} \mathbf{x}_{k+1} = \mathbf{y}_k - s \nabla f(\mathbf{y}_k) \\ \mathbf{y}_{k+1} = \mathbf{x}_{k+1} + \beta(\mathbf{x}_{k+1} - \mathbf{x}_k) \end{cases}$$

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{S} \mathbf{x} \rightarrow \nabla f = \mathbf{S} \mathbf{x} \rightarrow \begin{cases} \mathbf{S} \mathbf{q} = \lambda \mathbf{q} \\ \mathbf{x}_k = c_k \mathbf{q} \\ \mathbf{y}_k = d_k \mathbf{q} \end{cases} \rightarrow \nabla f(\mathbf{y}_k) = \mathbf{S} \mathbf{y}_k = \lambda \mathbf{y}_k = \lambda d_k \mathbf{q}$$

$$\left. \begin{array}{l} c_{k+1} = d_k - s \lambda d_k \\ d_{k+1} = c_{k+1} + \beta(c_{k+1} - c_k) \end{array} \right\} \rightarrow \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 0 & 1-s\lambda \\ -\beta & (1+\beta)(1-s\lambda) \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} = \mathbf{R} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

both eigenvalues  $e_1$  and  $e_2$  of  $\mathbf{R}$  to be as small as possible

choose  $s$  and  $\beta$  to minimize  $\max(|e_1(\lambda)|, |e_2(\lambda)|)$  for  $\lambda_{\min}(\mathbf{S}) \leq \lambda \leq \lambda_{\max}(\mathbf{S})$

$$\rightarrow s = \frac{1}{\lambda_{\max}}, \beta = \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \rightarrow \max(|e_1(\lambda)|, |e_2(\lambda)|) = \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}}}$$

(example) convergence factor:  $1 - \sqrt{b}$

# Momentum vs. Nesterov

