

# Ch 7. Continuous Optimization

Find  $\min_x f(x)$ , minimum value of  $f(x)$

- Unconstrained Optimization(Optimization using Gradient Descent)
  - 1. Gradient Descent
  - 2. Gradient Descent with Momentum
  - 3. Stochastic Gradient Descent SGD
- Constrained Optimization
  - 1. Linear Programming
  - 2. Quadratic Programming
  - 3. Convex set - Convex Optimization

# Optimization using Gradient Descent

## Unconstrained Optimization

- 1. Gradient Descent & Step-size/Learning rate
- 2. Gradient Descent with Momentum
- 3. Stochastic Gradient Descent SGD & Standard Gradient Descent

# Optimization using Gradient Descent

- 1. Gradient Descent

- $x_1 = x_0 - \gamma((\nabla f)(x_0))^T$  for small step size  $\gamma$

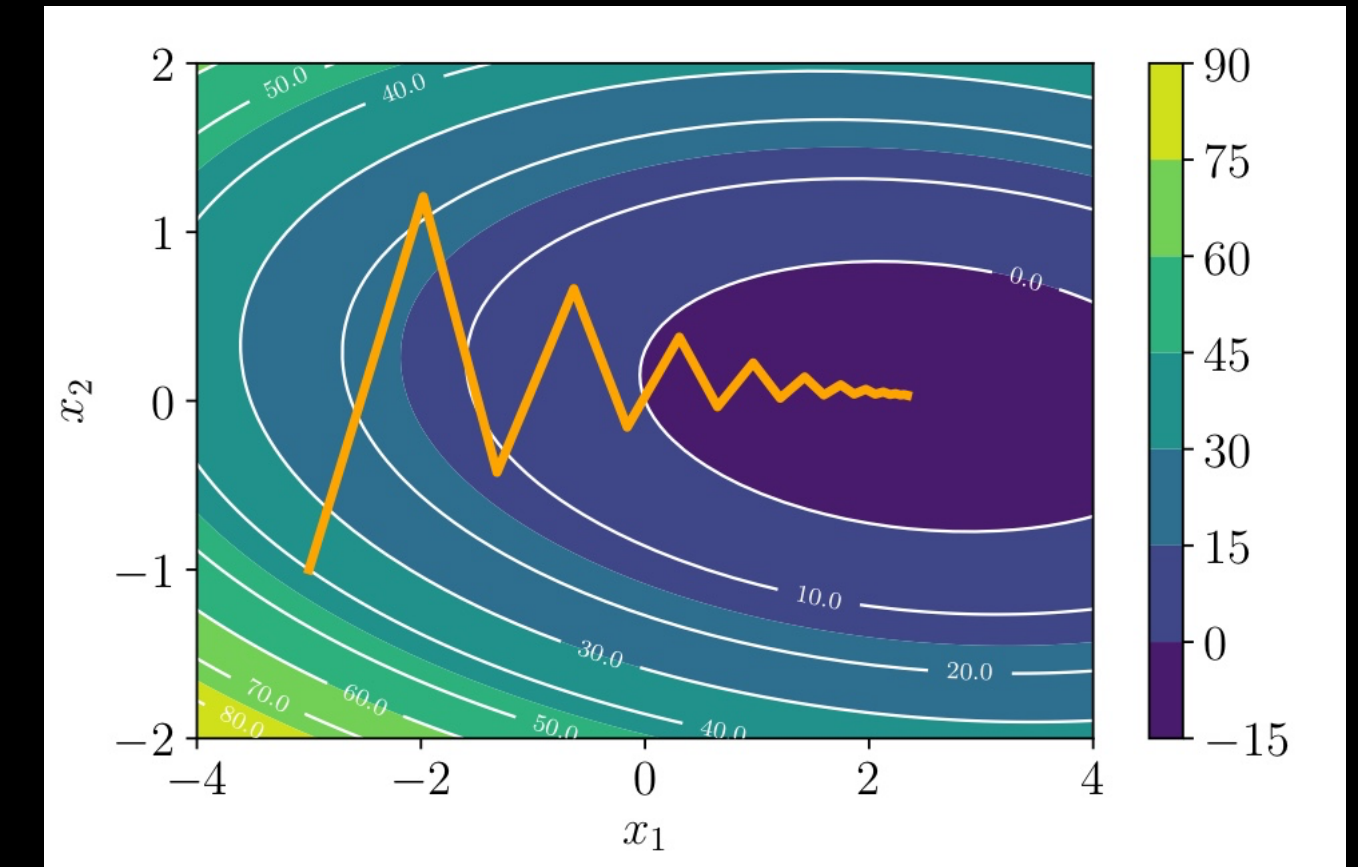
- $\implies x_{i+1} = x_i - \gamma_i((\nabla f)(x_i))^T$

- Relatively slow close to the minimum, increasingly zigzags as the near point-oscillation problem

- Step-size/Learning rate

- Too small - slow, too large - fail to converge, diverge

- Speed of convergence of gradient-descent  $>$  Condition number  $\kappa = \frac{\sigma(A)_{max}}{\sigma(B)_{max}}$



# Optimization using Gradient Descent

- 2. Gradient descent with Momentum
  - An additional term to remember what happened in the previous iteration
    - Dampen oscillation, smoothes out the gradient updates
  - $\Rightarrow x_{i+1} = x_i - \gamma_i((\nabla f)(x_0))^{\top} + \alpha \Delta x_i$ 
    - $\Delta x_i = x_i - x_{i-1} = \alpha \Delta x_{i-1} - \gamma_{i-1}((\nabla f)(x_{i-1}))^{\top}$
  - Momentum term averages out different noisy estimates of the gradient.

# Optimization using Gradient Descent

- 3. Stochastic Gradient Descent SGD

- By constraining the probability distribution of the approximate gradients

- $\theta_{i+1} = \theta_i - \gamma_i(\Delta L(\theta_i))^{\top}$

- Loss  $L(\theta) = \sum_{n=1}^N L_n(\theta) = - \sum_{n=1}^N \log p(y_n | x_n, \theta)$

- Standard Gradient Descent/Batch Optimization

- When training set is enormous and/or no simple formulas exist -> evaluation is expensive
- Mini-batch GD: Randomly choose a single  $L_n$  to estimate the gradient

# Constrained Optimization

- 1. Linear Programming
- 2. Quadratic Programming
- 3. Legendre-Fenchel Transform and Convex Conjugate + Convex Set

# Constrained Optimization

- 1. Linear Programming
  - $\min_{x \in R^d} c^\top x$  subject to  $Ax \leq b$
  - $\mathcal{L}(x, \lambda) = f(x) + \lambda^\top g(x) = c^\top x + \lambda^\top (Ax - b)$
  - $\mathcal{L}(x, \lambda) = (c + A^\top \lambda)^\top x - \lambda^\top b$
  - Derivative  $\Rightarrow c + A^\top \lambda = 0$
  - $\Rightarrow \mathcal{D}(\lambda) = -\lambda^\top b$
- Dual Optimization Problem
  - $\min_{\lambda \in R^m} -b^\top \lambda$  (subject to  $c + A^\top \lambda = 0, \lambda \geq 0$ )

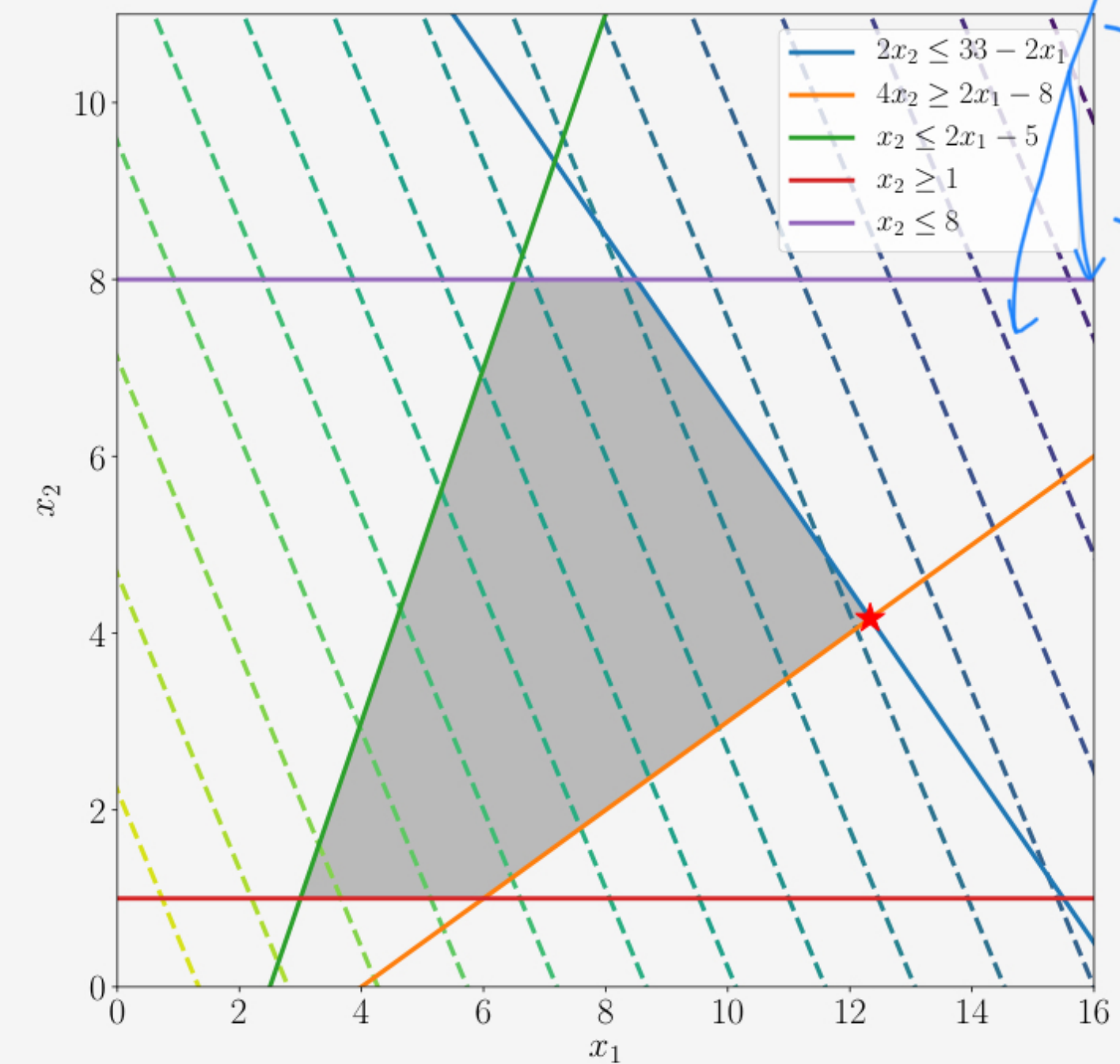
# Linear Programming

## Example 7.5 (Linear Program)

Consider the linear program

$$\begin{aligned} \min_{x \in \mathbb{R}^2} \quad & - \begin{bmatrix} 5 \\ 3 \end{bmatrix}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ \text{subject to} \quad & \begin{bmatrix} 2 & 2 \\ 2 & -4 \\ -2 & 1 \\ 0 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 33 \\ 8 \\ 5 \\ -1 \\ 8 \end{bmatrix} \end{aligned} \quad (7.44)$$

with two variables. This program is also shown in Figure 7.9. The objective function is linear, resulting in linear contour lines. The constraint set in standard form is translated into the legend. The optimal value must lie in the shaded (feasible) region, and is indicated by the star.





# Constrained Optimization

- 2. Quadratic Programming

- $\min_{x \in \mathbb{R}^d} \frac{1}{2} x^\top Q x + c^\top x$  subject to  $Ax \leq b$

- $\mathcal{L}(x, \lambda) = \frac{1}{2} x^\top Q x + c^\top x + \lambda^\top (Ax - b)$

- Derivative  $\Rightarrow Qx + (c + A^\top \lambda) = 0$

- $\Rightarrow x = -Q^{-1}(c + A^\top \lambda)$

- $\Rightarrow \mathcal{D}(\lambda) = -\frac{1}{2}(c + A^\top \lambda)^\top Q^{-1}(c + A^\top \lambda) - \lambda^\top b$

- Dual Optimization Problem

- $\min_{\lambda \in \mathbb{R}^m} -\frac{1}{2}(c + A^\top \lambda)^\top Q^{-1}(c + A^\top \lambda) - \lambda^\top b$  (subject to  $\lambda \geq 0$ )

# Quadratic Programming

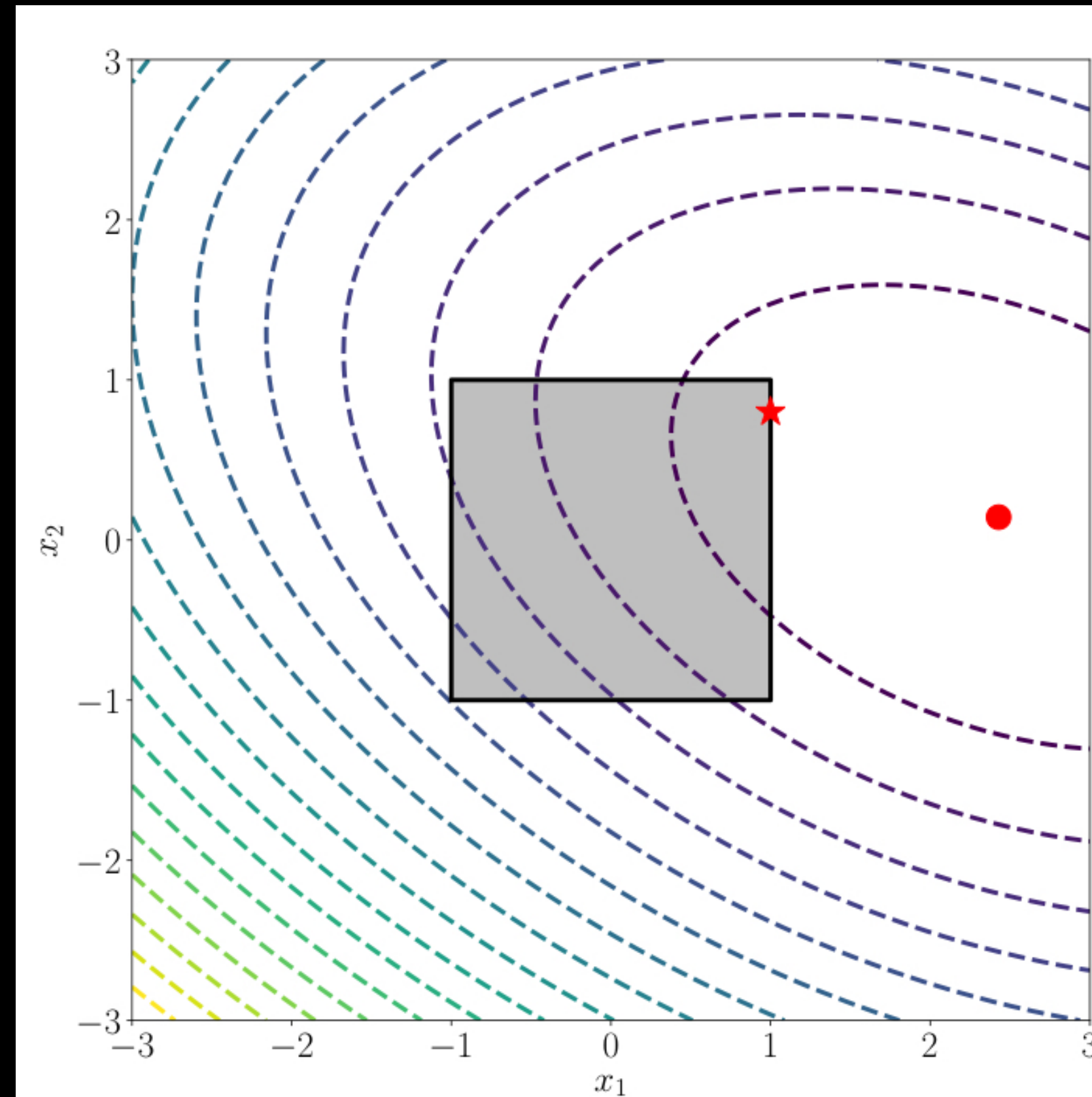
## Example 7.6 (Quadratic Program)

Consider the quadratic program

$$\min_{\mathbf{x} \in \mathbb{R}^2} \quad \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 5 \\ 3 \end{bmatrix}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (7.46)$$

$$\text{subject to} \quad \begin{bmatrix} 1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (7.47)$$

of two variables. The program is also illustrated in Figure 7.4. The objective function is quadratic with a positive semidefinite matrix  $\mathbf{Q}$ , resulting in elliptical contour lines. The optimal value must lie in the shaded (feasible) region, and is indicated by the star.



**Figure 7.4**  
Illustration of constrained optimization. The unconstrained problem (indicated by the contour lines) has a minimum on the right side (indicated by the circle). The box constraints ( $-1 \leq x \leq 1$  and  $-1 \leq y \leq 1$ ) require that the optimal solution is within the box, resulting in an optimal value indicated by the star.

# Constrained Optimization - Convex Set

- Convex set
  - Set  $C$  for any  $x, y \in C, 0 \leq \theta \leq 1 : \theta x + (1 - \theta)y \in C$
  - Convex function:  $f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$
  - -> described by its supporting hyperplanes
  - -> described by a function of their gradients : Legendre transform/Convex conjugate
- Convex Optimization Problem
  - $\min_x f(x)$  subject to  $g_i(x) \leq 0$  for all  $i = 1 \dots m, h_j(x) = 0$  for all  $j = 1 \dots n$ 
    - $f, g$  are convex functions,  $h$  are convex sets

# Constrained Optimization

- 3. Legendre-Fenchel transform / Convex conjugate

- $$f^*(x) = \sup_{x \in \mathbb{R}^D} (\langle s, x \rangle - f(x))$$

## Example 7.7 (Convex Conjugates)

To illustrate the application of convex conjugates, consider the quadratic function

$$f(\mathbf{y}) = \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} \quad (7.59)$$

based on a positive definite matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$ . We denote the primal variable to be  $\mathbf{y} \in \mathbb{R}^n$  and the dual variable to be  $\boldsymbol{\alpha} \in \mathbb{R}^n$ .

Applying Definition 7.4, we obtain the function

$$f^*(\boldsymbol{\alpha}) = \sup_{\mathbf{y} \in \mathbb{R}^n} \langle \mathbf{y}, \boldsymbol{\alpha} \rangle - \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}. \quad (7.60)$$

Since the function is differentiable, we can find the maximum by taking the derivative and with respect to  $\mathbf{y}$  setting it to zero.

$$\frac{\partial [\langle \mathbf{y}, \boldsymbol{\alpha} \rangle - \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}]}{\partial \mathbf{y}} = (\boldsymbol{\alpha} - \lambda \mathbf{K}^{-1} \mathbf{y})^\top \quad (7.61)$$

and hence when the gradient is zero we have  $\mathbf{y} = \frac{1}{\lambda} \mathbf{K} \boldsymbol{\alpha}$ . Substituting into (7.60) yields

$$f^*(\boldsymbol{\alpha}) = \frac{1}{\lambda} \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} - \frac{\lambda}{2} \left( \frac{1}{\lambda} \mathbf{K} \boldsymbol{\alpha} \right)^\top \mathbf{K}^{-1} \left( \frac{1}{\lambda} \mathbf{K} \boldsymbol{\alpha} \right) = \frac{1}{2\lambda} \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}. \quad (7.62)$$