

# Convex Optimization: Part II

CMPT 419/983

Mo Chen

SFU Computing Science

23/09/2018



• S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2008.

# Outline

- Optimization program
  - Examples and classes
- Convex optimization
  - Convex functions
  - Optimality conditions
- Numerical solutions

• Unconstrained case: minimize f(x)

•  $\nabla f(x) = 0$ 





• Inequality constraints only:

minimize f(x)subject to  $g_i(x) \le 0, i = 1, ..., n$ 

- Penalty view point: penalize constraint violation
  - Lagrangian:  $L(x, \lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x)$ ,  $\lambda_i \ge 0$
- Optimality conditions
  - Stationarity:  $\nabla_{\chi} L(x^*, \lambda^*) = 0$
  - Primal feasibility:  $g_i(x^*) \leq 0$
  - Dual feasibility:  $\lambda^* \ge 0$
  - Complementary slackness:  $\lambda_i^* g_i(x^*) = 0$ , i = 1, ..., n

- Stationarity:  $\nabla_{x} L(x^*, \lambda^*) = 0$ 
  - Lagrangian:

$$L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x), \qquad \lambda_i \ge 0$$

• Take gradient and set to zero:  $0 = \nabla f(x) + \sum_{i=1}^{n} \lambda_i \nabla g_i(x)$ 

$$\nabla f(x) = -\sum_{i=1}^{n} \lambda_i \nabla g_i(x)$$

• Since  $\lambda_i \ge 0$ , gradient of f(x) must point "away" from gradients of active constraint functions



- Stationarity:  $\nabla_{x} L(x^*, \lambda^*) = 0$ 
  - Lagrangian:

$$L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x), \qquad \lambda_i \ge 0$$

• Take gradient and set to zero:  $0 = \nabla f(x) + \sum_{i=1}^{n} \lambda_i \nabla g_i(x)$ 

$$\nabla f(x) = -\sum_{i=1}^{n} \lambda_i \nabla g_i(x)$$

• Since  $\lambda_i \ge 0$ , gradient of f(x) must point "away" from gradients of active constraint functions



- Stationarity:  $\nabla_{x} L(x^*, \lambda^*) = 0$ 
  - Lagrangian:

$$L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x), \qquad \lambda_i \ge 0$$

• Take gradient and set to zero:  $0 = \nabla f(x) + \sum_{i=1}^{n} \lambda_i \nabla g_i(x)$ 

$$\nabla f(x) = -\sum_{i=1}^{n} \lambda_i \nabla g_i(x)$$

• Since  $\lambda_i \ge 0$ , gradient of f(x) must point "away" from gradients of active constraint functions



- Primal feasibility:  $g_i(x^*) \leq 0$ 
  - Constraints must be satisfied
- Dual feasibility:  $\lambda^* \ge 0$ 
  - Penalty view point



- Complementary slackness:  $\lambda_i^* g_i(x^*) = 0, i = 1, ..., n$ 
  - Lagrangian:

$$L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x), \qquad \lambda_i \ge 0$$

- If  $g_i(x^*) < 0$ , then the constraint is not active, so  $\lambda_i^*$  is set to 0 to not decrease the Lagrangian
- If  $g_i(x^*) = 0$ , then the constraint is active, so  $\lambda_i^*$  is free to be positive



• Full optimization problem: minimize f(x)

subject to 
$$g_i(x) \le 0, i = 1, ..., n$$
  
 $a_j^{\mathsf{T}} x = b_j, j = 1, ..., m$ 

- Penalty view point:
  - Lagrangian:  $L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x) + \sum_{j=1}^{m} \mu_j \left( a_j^\top x b_j \right), \ \lambda_i \ge 0$
- Karush-Kuhn-Tucker (KKT) Conditions:
  - Stationarity  $\nabla_{x}L(x^{*},\lambda^{*},\mu^{*})=0$
  - Primal feasibility:  $g_i(x^*) \leq 0$ ,  $a_i^{\top} x^* b_i = 0$
  - Dual feasibility:  $\lambda^* \ge 0$
  - Complementary slackness:  $\lambda_i^* g_i(x^*) = 0$ , i = 1, ..., n
- Solve above systems of equations to obtain optimum

# Solving Convex Optimization Problems

- Solve the optimality conditions
- Gradient methods for approximating solutions to convex optimization problems

• Full optimization problem: minimize f(x)

subject to 
$$g_i(x) \le 0, i = 1, ..., n$$
  
 $a_j^{\mathsf{T}} x = b_j, j = 1, ..., m$ 

- Penalty view point:
  - Lagrangian:  $L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x) + \sum_{j=1}^{m} \mu_j (a_j^{\mathsf{T}} x b_j), \ \lambda_i \ge 0$
- Karush-Kuhn-Tucker (KKT) Conditions:
  - Stationarity  $\nabla_{\chi} L(x^*, \lambda^*, \mu^*) = 0$
  - Primal feasibility:  $g_i(x^*) \leq 0$ ,  $a_i^{\mathsf{T}} x^* b_i = 0$
  - Dual feasibility:  $\lambda^* \ge 0$
  - Complementary slackness:  $\lambda_i^* g_i(x^*) = 0$ , i = 1, ..., n
- Solve above systems of equations to obtain optimum



 $\underset{\theta}{\text{minimize}} \| X\theta - Y \|_2^2$ 

- Scalar example:
  - Data:  $\{x_i, y_i\}_{i=1}^n, x_i, y_i \in \mathbb{R}$
  - Model:  $y = mx + b, m, b \in \mathbb{R}$
  - Sum of error of model:  $\sum_{i=1}^{n} (y_i mx_i b)^2$
  - No constraints: allow any m, b
- Error in matrix form:  $e_i = y_i \begin{bmatrix} x_i & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix}$

• Stacking the data points: 
$$E_{i} = \begin{bmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{n} \end{bmatrix} - \begin{bmatrix} x_{1} & 1 \\ x_{2} & 1 \\ \vdots & \vdots \\ x_{n} & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix}$$
$$Y = \begin{bmatrix} x_{1} & 1 \\ x_{2} & 1 \\ \vdots & \vdots \\ x_{n} & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix}$$



• Full optimization problem: minimize f(x)

subject to 
$$g_i(x) \le 0, i = 1, ..., n$$
  
 $a_j^{\mathsf{T}} x = b_j, j = 1, ..., m$ 

- Penalty view point:
  - Lagrangian:  $L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x) + \sum_{j=1}^{m} \mu_j (a_j^\top x b_j), \ \lambda_i \ge 0$
- Karush-Kuhn-Tucker (KKT) Conditions:
  - Stationarity  $\nabla_{x} L(x^*, \lambda^*, \mu^*) = 0$   $\frown$   $\nabla f(x) = 0$
  - Primal feasibility:  $g_i(x^*) \leq 0$ ,  $a_i^{\mathsf{T}} x^* b_i = 0$
  - Dual feasibility:  $\lambda^* \ge 0$
  - Complementary slackness:  $\lambda_i^* g_i(x^*) = 0, \ i = 1, ..., n$
- Solve above systems of equations to obtain optimum

 $\underset{\theta}{\text{minimize}} \| X\theta - Y \|_2^2$ 

- Analytic solution available!
  - Objective:  $f(\theta) = ||X\theta Y||_2^2$ , set derivative to zero
  - $f(\theta) = (X\theta Y)^{\mathsf{T}}(X\theta Y)$
  - $f(\theta) = \theta^{\mathsf{T}} X^{\mathsf{T}} X \theta 2Y^{\mathsf{T}} X \theta + Y^{\mathsf{T}} Y$

$$\frac{\partial f}{\partial \theta} = 2X^{\mathsf{T}}X\theta - 2X^{\mathsf{T}}Y$$
$$0 = 2X^{\mathsf{T}}X\theta - 2X^{\mathsf{T}}Y$$
$$X^{\mathsf{T}}Y = X^{\mathsf{T}}X\theta$$
$$\theta = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}Y$$





 $L(x,\lambda) = f(x) + \sum_{i=1}^{n} \lambda_i g_i(x)$ 

$$\begin{array}{ll} \text{minimize} & \|X\theta - Y\|_2^2\\ \text{subject to} & \theta_1^2 + \theta_2^2 \le 1 \end{array}$$

• Lagrangian:

• Stationarity 
$$\nabla_x L(x^*, \lambda^*, \mu^*) = 0$$

$$\begin{array}{ll} \underset{\theta}{\text{minimize}} & \|X\theta - Y\|_2^2 \\ \text{subject to} & \|\theta\|_2^2 - 1 \le 0 \end{array}$$

 $L(\theta, \lambda) = \|X\theta - Y\|_2^2 + \lambda(\|\theta\|_2^2 - 1)$ 

 $\nabla_{\theta} L(\theta, \lambda) = 2X^{\mathsf{T}} X \theta - 2X^{\mathsf{T}} Y + 2\lambda \theta$  $0 = X^{\mathsf{T}} X \theta - X^{\mathsf{T}} Y + \lambda \theta$  $X^{\mathsf{T}} Y = (X^{\mathsf{T}} X + \lambda I) \theta$  $\|\theta\|_{2}^{2} - 1 \leq 0$ 

 $\lambda > 0$ 

- Primal feasibility:  $g_i(x^*) \le 0$ ,  $a_i^{\mathsf{T}} x^* b_i = 0$
- Dual feasibility:  $\lambda^* \ge 0$
- Complementary slackness:  $\lambda_i^* g_i(x^*) = 0$ , i = 1, ..., n

 $\lambda(\|\theta\|_2^2 - 1) = 0$  $\lambda = 0 \text{ or } \|\theta\|_2^2 = 1$ 

- Case 1: If  $\lambda = 0$ , then
  - $\lambda \ge 0$  is satisfied automatically
  - $X^{\top}Y = (X^{\top}X)\theta \Rightarrow \theta = (X^{\top}X)^{-1}X^{\top}Y$
  - If  $\|\theta\|_2^2 1 \le 0$  happens to be true, we are done
  - Otherwise, try case 2
- Case 2: If  $\|\theta\|_2^2 = 1$ , then
  - $\|\theta\|_2^2 1 \le 0$  is satisfied automatically
  - $X^{\top}Y = (X^{\top}X + \lambda I)\theta \Rightarrow \theta = (X^{\top}X + \lambda I)^{-1}X^{\top}Y$
  - Solve  $\|\theta\|_2^2 = 1$  and  $\theta = (X^T X + \lambda I)^{-1} X^T Y$  for  $\theta$  and  $\lambda$
  - If  $\lambda \geq 0$ , we are done

#### **KKT conditions:**

- $X^{\top}Y = (X^{\top}X + \lambda I)\theta$
- $\|\theta\|_2^2 1 \le 0$
- $\lambda \ge 0$
- $\lambda = 0$  or  $\|\theta\|_2^2 = 1$

# Solving the Optimality Conditions

minimize f(x)

- Equations to solve: KKT conditions
  - Stationarity  $\nabla_{x} L(x^*, \lambda^*, \mu^*) = 0$
  - Primal feasibility:  $g_i(x^*) \leq 0$ ,  $a_i^{\mathsf{T}} x^* b_i = 0$
  - Dual feasibility:  $\lambda^* \ge 0$
  - Complementary slackness:  $\lambda_i^* g_i(x^*) = 0$ , i = 1, ..., n
- Use numerical equation solvers, or do it by hand (as much as possible)
- For convex problems, KKT conditions are necessary and sufficient
- For non-convex problems, KKT conditions are just necessary

subject to  $g_i(x) \le 0, i = 1, ..., n$  $a_j^{\mathsf{T}} x = b_j, j = 1, ..., m$ 

### Numerical Solution: Gradient Methods

- Start from  $x^0$  and construct a sequence  $x^k$  such that  $x^k \rightarrow x^*$ 
  - Calculate  $x^{k+1}$  from  $x^k$  by "going down the gradient"
  - Unconstrained case:  $x^{k+1} = x^k \alpha^k \nabla f(x)$ ,  $\alpha^k > 0$





### Numerical Solution: Gradient Methods

- Start from  $x^0$  and construct a sequence  $x^k$  such that  $x^k \rightarrow x^*$ 
  - Calculate x<sup>k+1</sup> from x<sup>k</sup> by "going down the gradient"
  - Unconstrained case:  $x^{k+1} = x^k \alpha^k \nabla f(x)$ ,  $\alpha^k > 0$
- More generally,  $x^{k+1} = x^k + \alpha^k d^k$  for some d such that  $\nabla f(x^k) \cdot d^k < 0$
- Tuning parameters: descent direction  $d^k$ , and step size  $\alpha^k$



### **Descent Direction**

- Steepest descent:  $d^k = -\nabla f(x^k)$ 
  - $x^{k+1} = x^k \alpha^k \nabla f(x)$
  - Simple but sometimes leads to slow convergence

### Steepest Descent (Gradient Descent) Example

• Line fitting: 
$$f(\theta) = ||X\theta - Y||_2^2$$
  
•  $\frac{\partial f}{\partial \theta} = 2X^T X \theta - 2X^T Y$ 

```
theta_last = [-2; -2];
dtheta = inf;
maxIter = 500;
```

end

```
alpha = 0.1/k;
theta = theta_last - alpha*(2*X'*X*theta_last - 2*X'*Y);
dtheta = theta_last - theta;
theta_last = theta;
```



end

### Steepest Descent (Gradient Descent) Example





### **Descent Direction**

- Steepest descent:  $d^k = -\nabla f(x^k)$ 
  - $x^{k+1} = x^k \alpha^k \nabla f(x)$
  - Simple but sometimes leads to slow convergence

• Newton's method: 
$$d^k = \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k)$$

Minimize the quadratic approximation:

$$f^{k}(x) = f(x^{k}) + \nabla f(x^{k})^{\mathsf{T}}(x - x^{k}) + \frac{1}{2}(x - x^{k})^{\mathsf{T}} \nabla^{2} f(x^{k})(x - x^{k})$$

0

-1

-2

-2.8

-2.6

-2.4

-2.2

-2

-1.8

-1.6

-1.4

• Set gradient to zero to obtain next iterate

$$\nabla f^{k}(x) = \nabla f(x^{k}) + \nabla^{2} f(x^{k})(x - x^{k}) = 0$$
  
$$\Rightarrow x^{k+1} = x^{k} - \left(\nabla^{2} f(x^{k})\right)^{-1} \nabla f(x^{k})$$

- Fast convergence, but matrix inverse required
- Alternatively, use an algorithm to minimize a quadratic function

# Step Size

- Recall  $x^{k+1} = x^k + \alpha^k d^k$ , with  $\nabla f(x^k)^T d^k < 0$
- Line search: choose  $\alpha^k = \min_{\alpha \ge 0} f(x^k + \alpha^k d^k)$ 
  - Requires minimization
- Constant step size:  $\alpha^k = \alpha$ 
  - May not converge
- Diminishing step size:  $\alpha^k \to 0$ 
  - Still need to explore all regions  $\sum \alpha^k = \infty$
  - For example:  $\alpha^k = \frac{\alpha^0}{k}$



• Steepest descent,  $\alpha^k = \alpha^0/k$ 





• Steepest descent,  $\alpha^k = \alpha^0$  (small steps)



• Steepest descent,  $\alpha^k = \alpha^0$  (large steps)



• Steepest descent,  $\alpha^k = \alpha^0/k^2$  (steps do not sum to  $\infty$ :  $\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6}$ )



# Dealing with Constraints

- Idea 1: Apply descent step, and project point to feasible set
  - Proximal gradient methods
  - Difficulty: Computing the projected point
- Idea 2: Set penalty to  $\infty$  for constraint violation
  - Barrier functions

