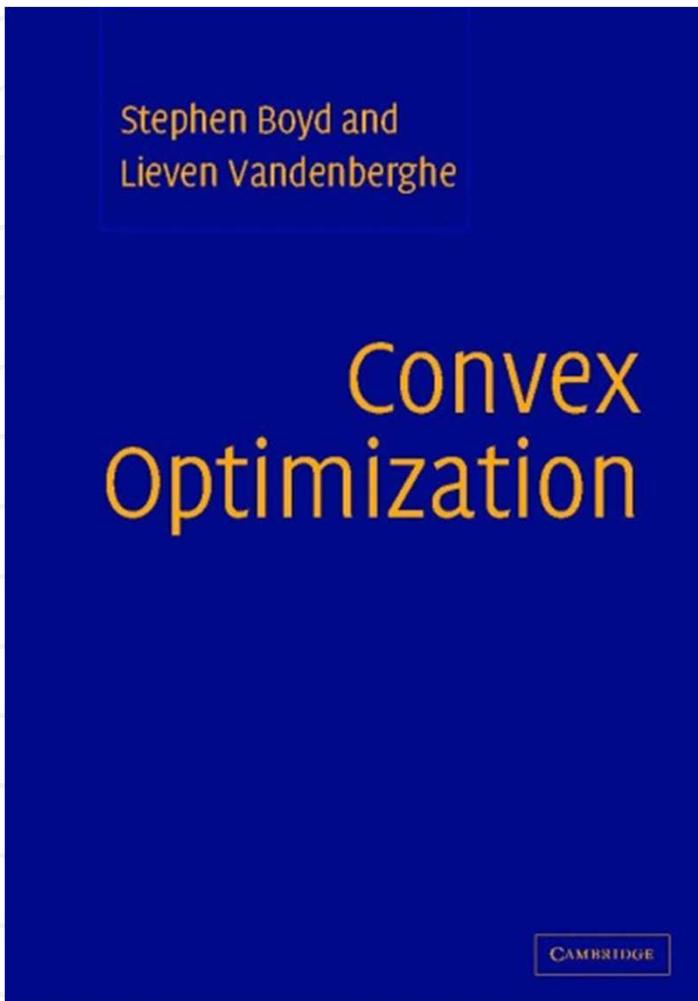


---

# Lecture 15: Subgradient optimization. KKT conditions.

# Reading

---



## Section 5.5

# Recap (again)

The Lagrangian:

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

The dual function:

$$g(\lambda, \nu) = \min_x L(x, \lambda, \nu) \quad (\inf)$$

Let  $x$  to be any feasible primal set of variables, and  $(\lambda, \nu)$  be any feasible dual set of variables (i.e.  $\lambda \geq 0$ ):

$$g(\lambda, \nu) \leq L(x, \lambda, \nu) \leq f_0(x)$$

In particular, this holds for the optimal values:

$$g(\lambda^*, \nu^*) \leq L(x^*, \lambda^*, \nu^*) \leq f_0(x^*)$$

# Partial dualization

# The (partial) Lagrangian:

The (partial) Lagrangian:

**Observation:** suppose  $x \in D$  and  $\lambda \geq o$ . Then

$$L(x, \lambda, v) \leq f_0(x)$$

# Consequence:

$$\min_{x \in D'} L(x, \lambda, v) \leq \min_{x \in D \cap D'} L(x, \lambda, v) \leq \min_{x \in D \cap D'} f_0(x)$$

# The Lagrange dual (partial dualization)

The (partial) dual function:

$$g(\lambda, \nu) = \min_{x \in D'} L(x, \lambda, \nu)$$

Theorem (weak duality):

$$\max_{\lambda \geq 0, \nu} g(\lambda, \nu) \leq \min_{\substack{f_0(x) \leq 0, \\ f'_i(x) \leq 0, \\ h_i(x) = 0}} f_0(x)$$

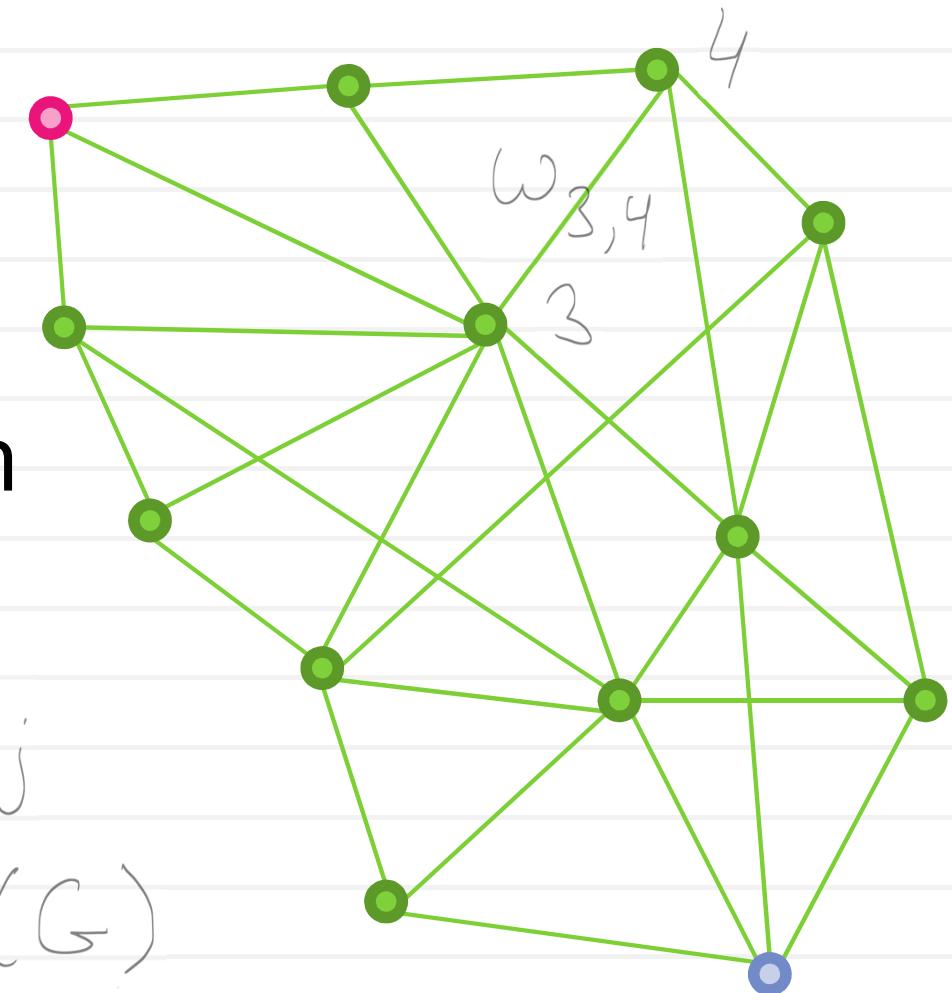
We have a freedom to pick which constraints to dualize and which to leave as they are

# Constrained shortest path

## Shortest path task

- Dijkstra
  - Bellman-Ford
  - Breadth-first search
- .....

$$\begin{aligned} \min_x \quad & \sum w_{ij} \cdot x_{ij} \\ \text{s.t.} \quad & x \in \text{Path}(G) \\ & x_{ij} \in \{0;1\} \\ & \sum x_{ij} = k \end{aligned}$$



# Constrained shortest path

$$\min_x \sum w_{ij} \cdot x_{ij}$$

$$\text{s.t. } x \in \text{Path}(G)$$

$$\sum x_{ij} = k$$

$$g(\lambda) = \min_{x \in \text{Path}(G)} \left( \sum w_{ij} \cdot x_{ij} + \lambda (k - \sum x_{ij}) \right)$$

$$= \lambda k + \min_{x \in \text{Path}(G)} \left( \sum (w_{ij} - \lambda) \cdot x_{ij} \right)$$

Dual can be evaluated using shortest path!

# Constrained shortest path

$$g(\lambda) = \lambda K + \min_{x \in \text{Path}(G)} \left( \sum (\omega_{ij} - \lambda) \cdot x_{ij} \right)$$

$$\frac{dg}{d\lambda} = (K - \sum \hat{x}_{ij})$$

$$\lambda_{t+1} = \lambda_t + s_t \cdot \frac{dg}{d\lambda}$$

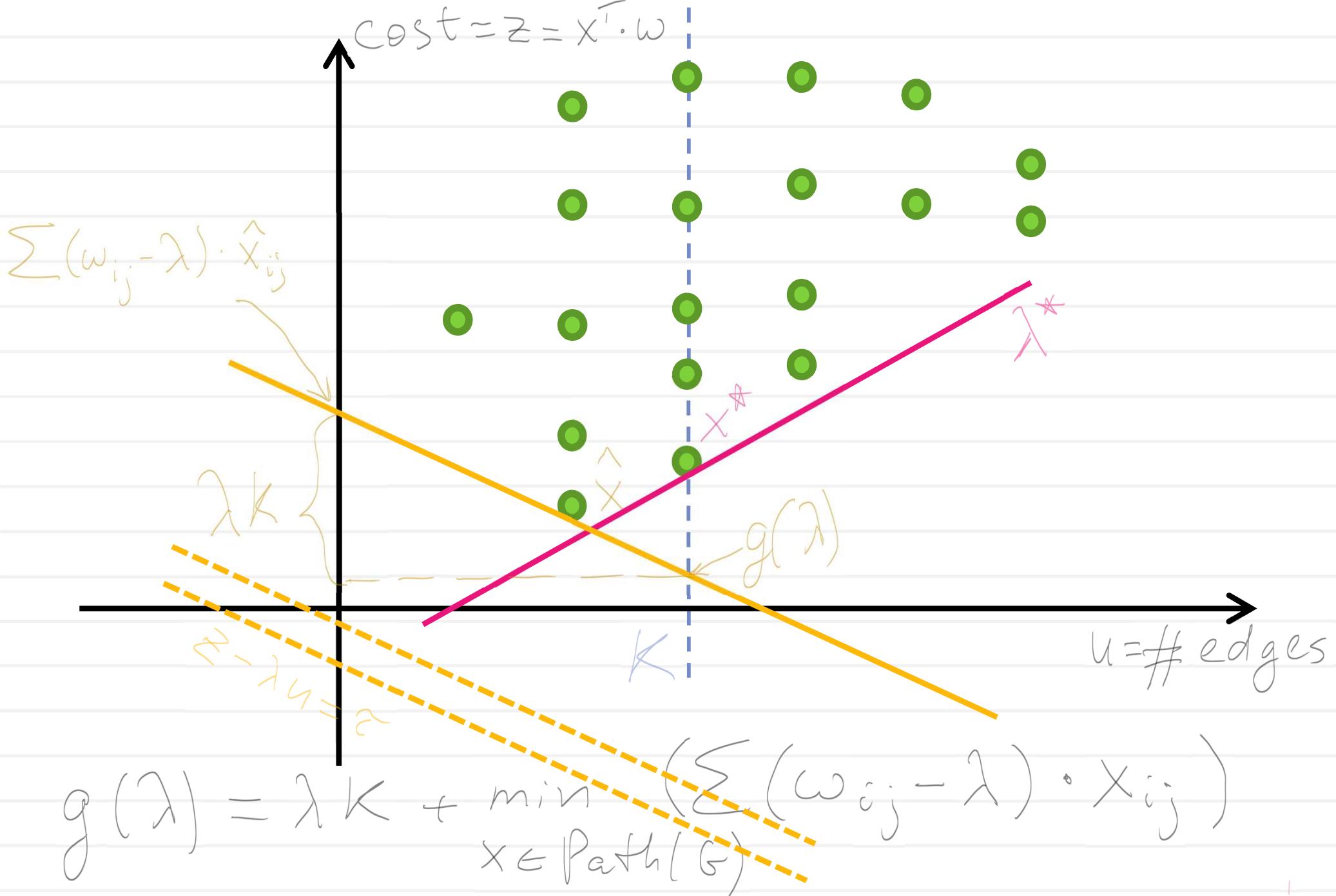
$$s_t \rightarrow 0$$

$$\sum s_t = +\infty$$

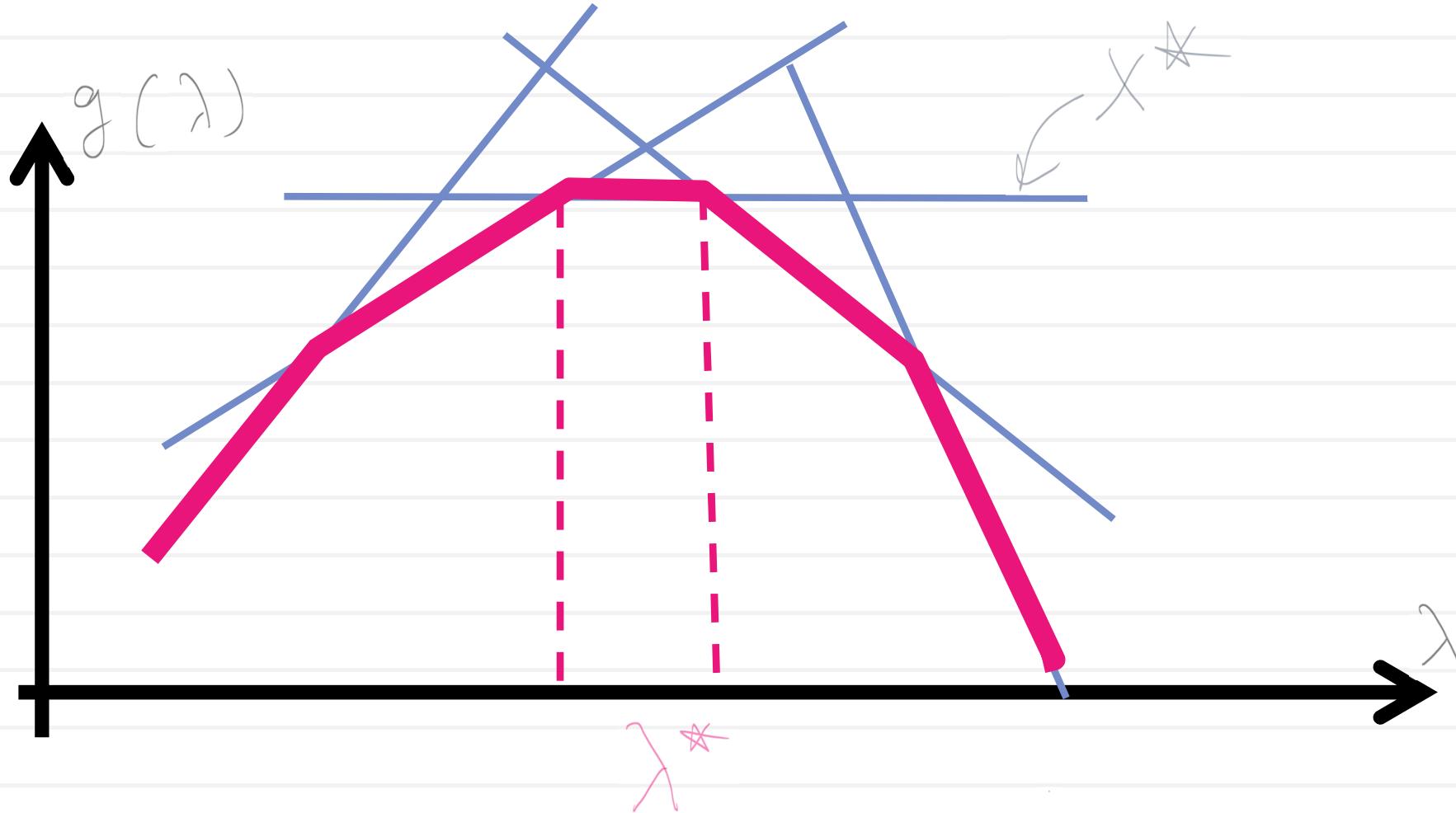
*Subgradient ascent step in this case:*

- If the optimal path is longer than  $K$ , decrease  $\lambda$
- If the optimal path is shorter than  $K$ , increase  $\lambda$

# Case 1: strong duality holds



# The shape of the dual function

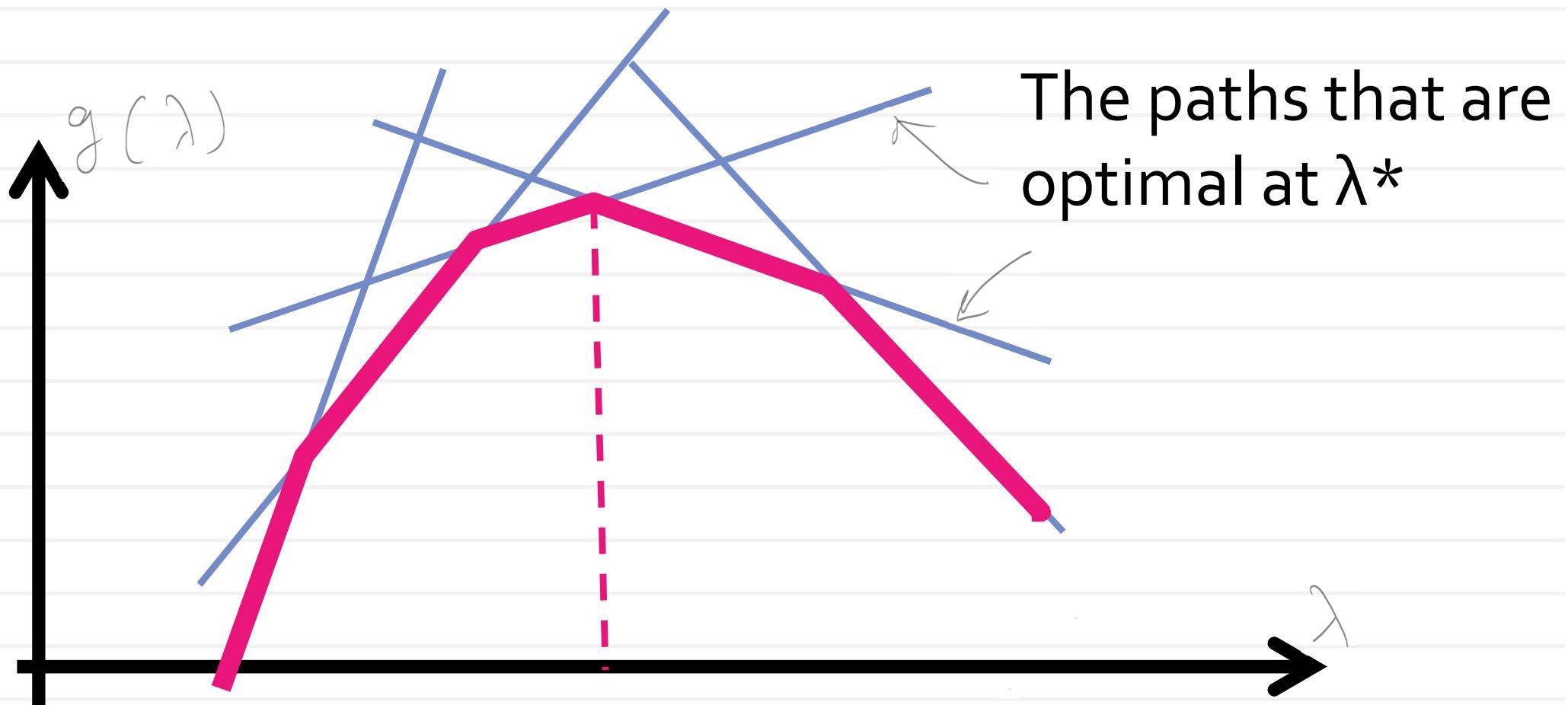


## Case 2: strong duality does not hold



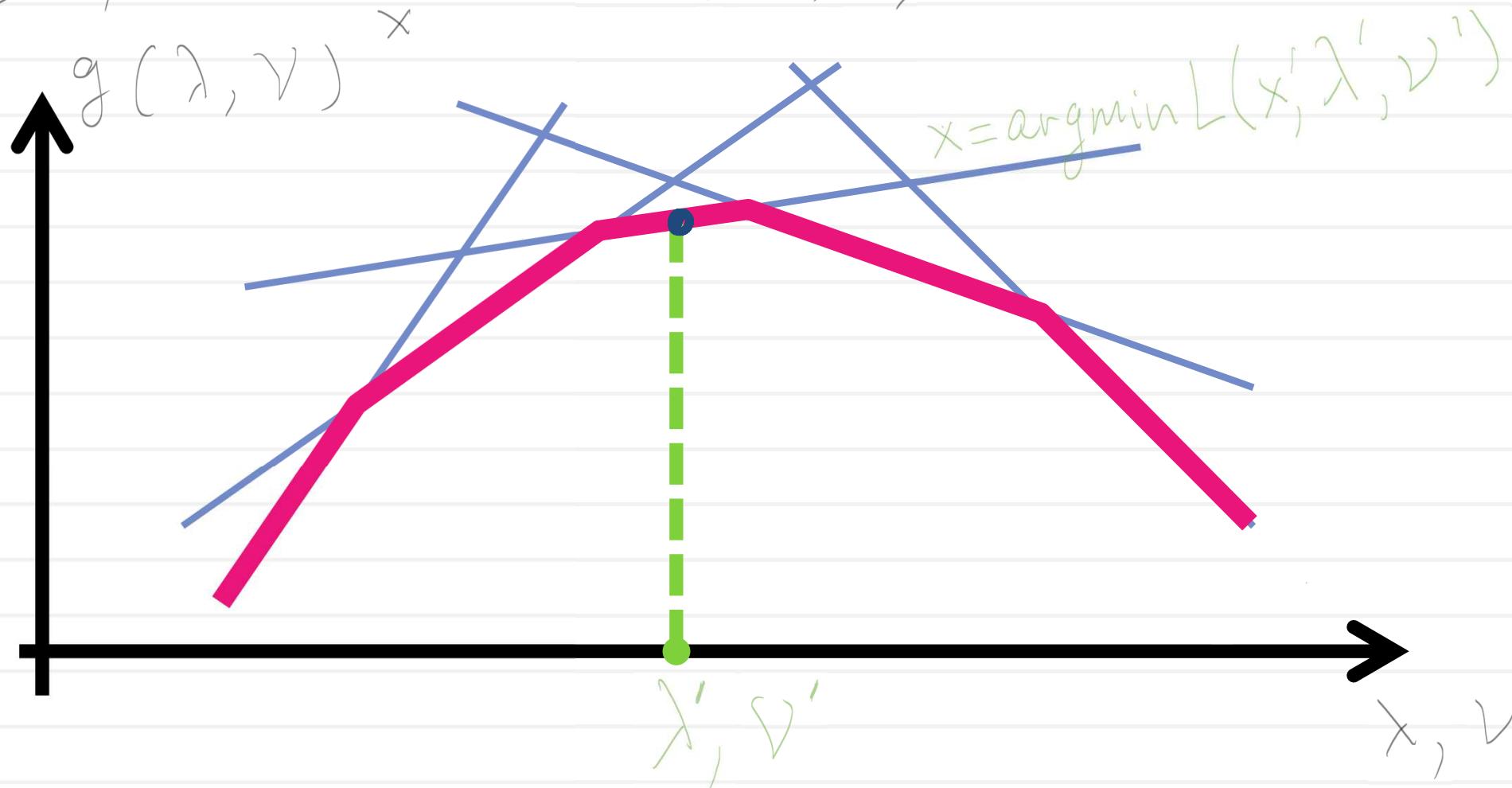
$$g(\lambda) = \lambda K + \min_{x \in \text{Path}(G)} \left( \sum (\omega_{ij} - \lambda) \cdot x_{ij} \right)$$

# The shape of the dual function



# Recap

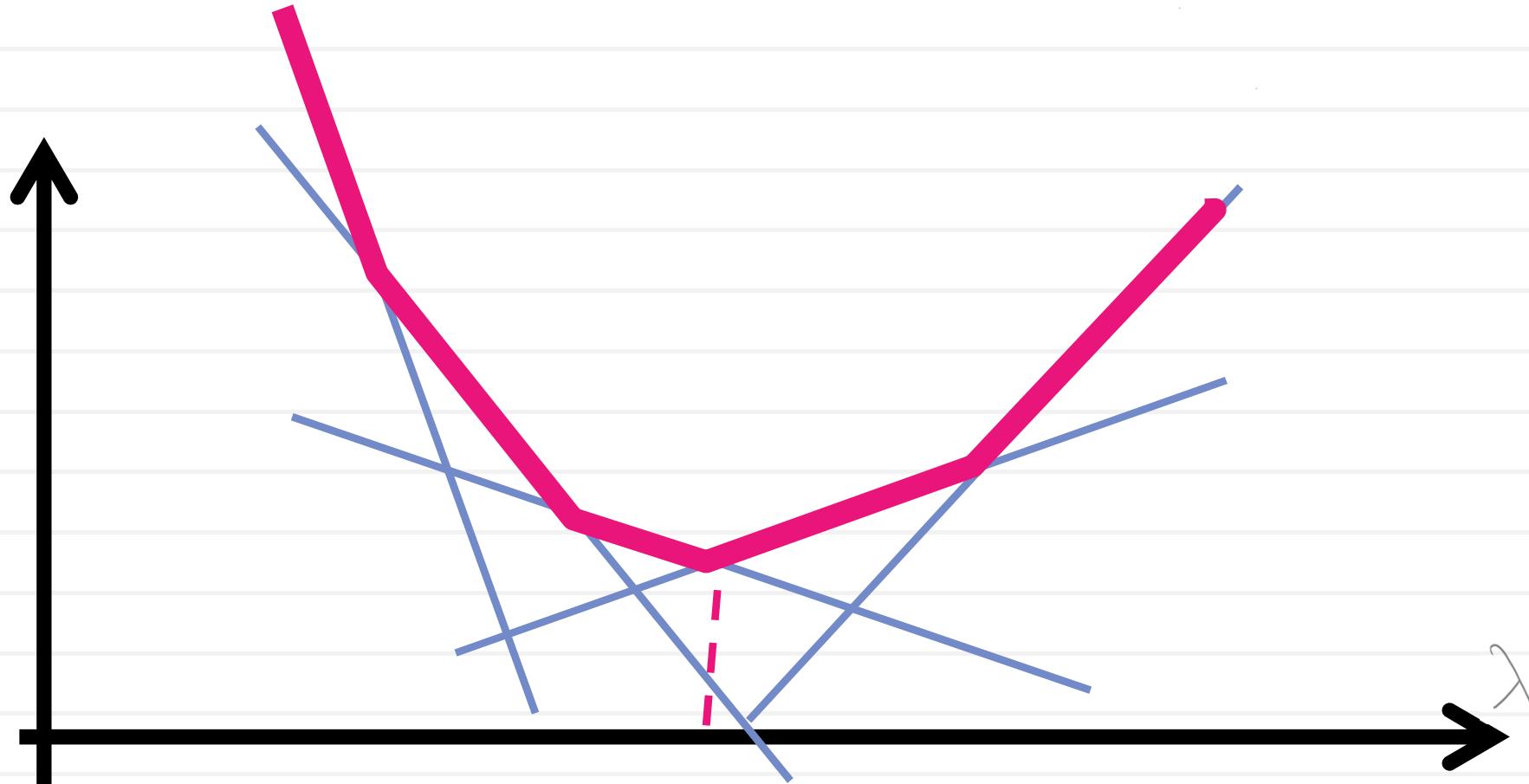
$$g(\lambda, \nu) = \min_x L(x, \lambda, \nu)$$



$$g(\lambda, \nu) \leq g(\lambda', \nu') + f(x')^\top (\lambda - \lambda') + h(x')^\top (\nu - \nu')$$

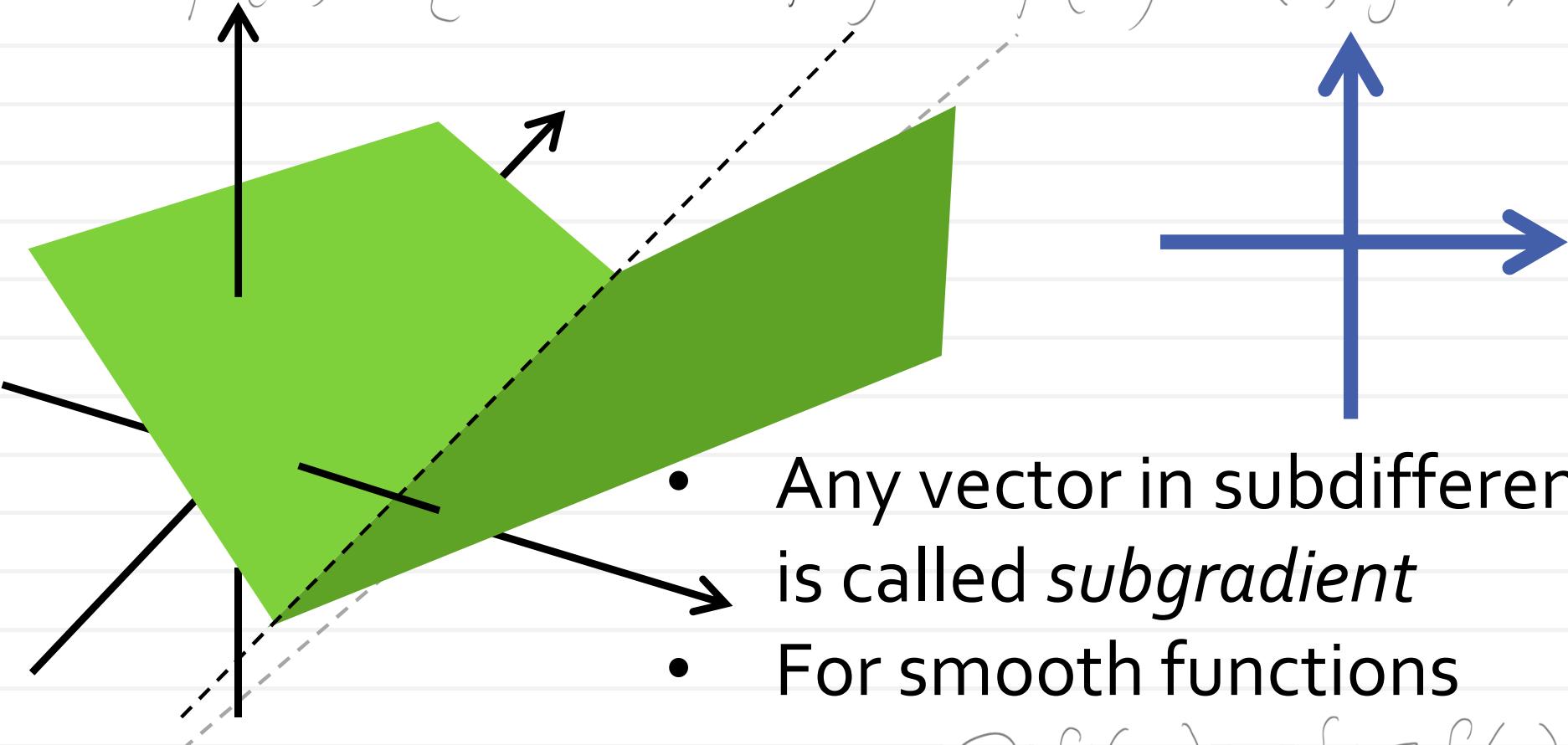
$\nabla g = (f_i(x), h_i(x))$  - (sub)gradient of  $g$

# Minimizing non-differentiable convex functions



# Subdifferential

$$\partial f(x) = \{ s \in \mathbb{R}^n : f(y) \geq f(x) + \langle s, y-x \rangle \}$$



- Any vector in subdifferential is called *subgradient*
  - For smooth functions
- $$\partial f(x) = \{ \nabla f(x) \}$$
- Convex functions have non-empty subdifferential (why?)

# Subgradient descent

- A subgradient is not necessarily a descent direction
- However, persistence pays off:

for  $k = 1..K$

$g = \text{GetSubgradient}(f, x(k))$

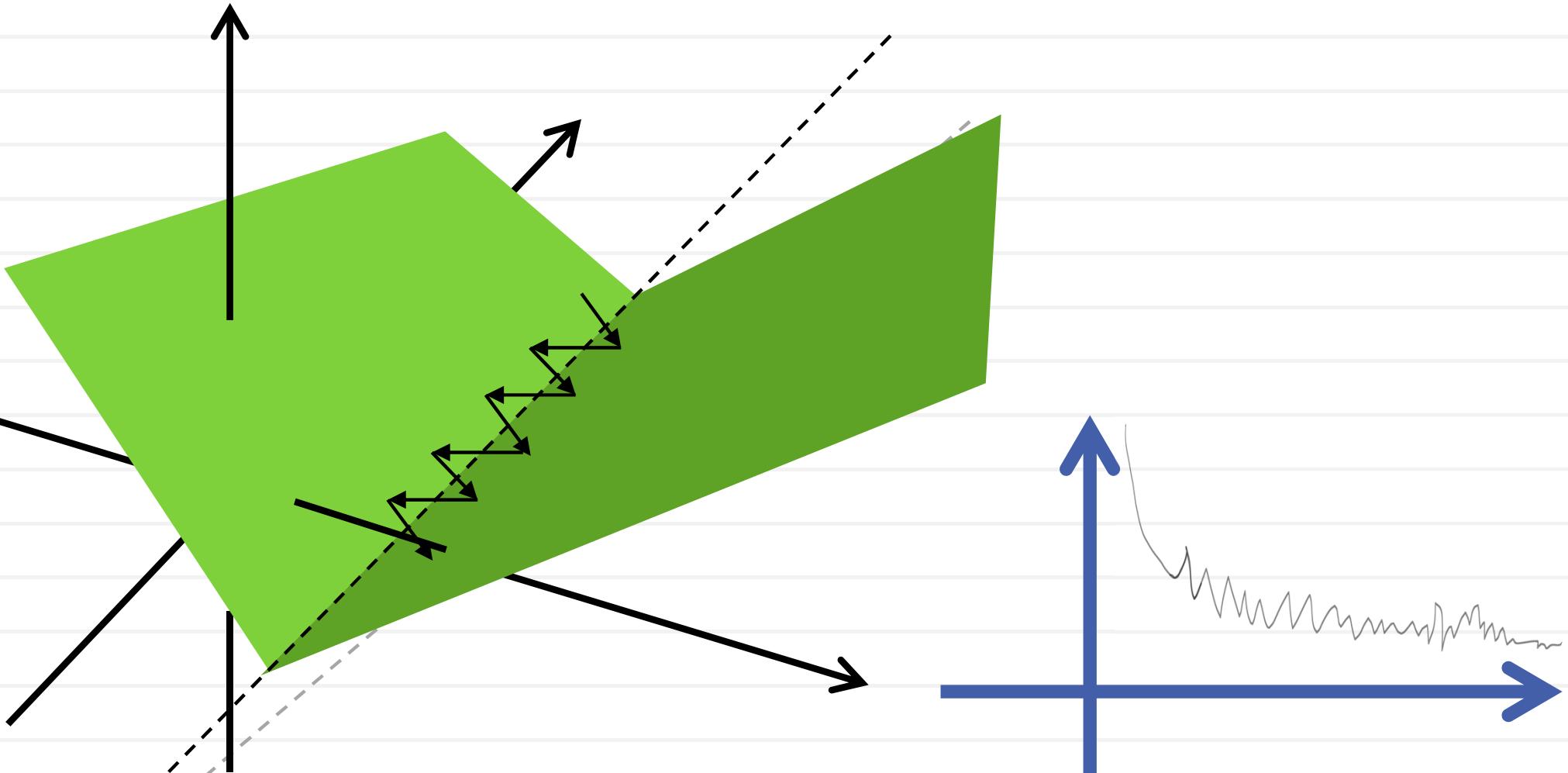
$x(k+1) = x(k) - \alpha_k g / \|g\|$

end

**Corollary:** if  $\alpha_k \geq 0, \alpha_k \rightarrow 0, \sum \alpha_k \rightarrow \infty$

then  $f(x(k)) \rightarrow \min_x f(x)$

# Subgradient descent



- A subgradient is not necessarily a descent direction
- However, persistence pays off

# Subgradient descent

```
xbest = x(0), fbest = f(x(0))
for k = 0..K
    g = GetSubgradient(f,x(k))
    x(k+1) = x(k)-1/(sqrt(k)+βk)*g
    if f(x(k+1)) < fbest
        fbest = f(x(k+1))
        xbest = x(k+1)
    end
end
return xbest
```

In practical implementation:

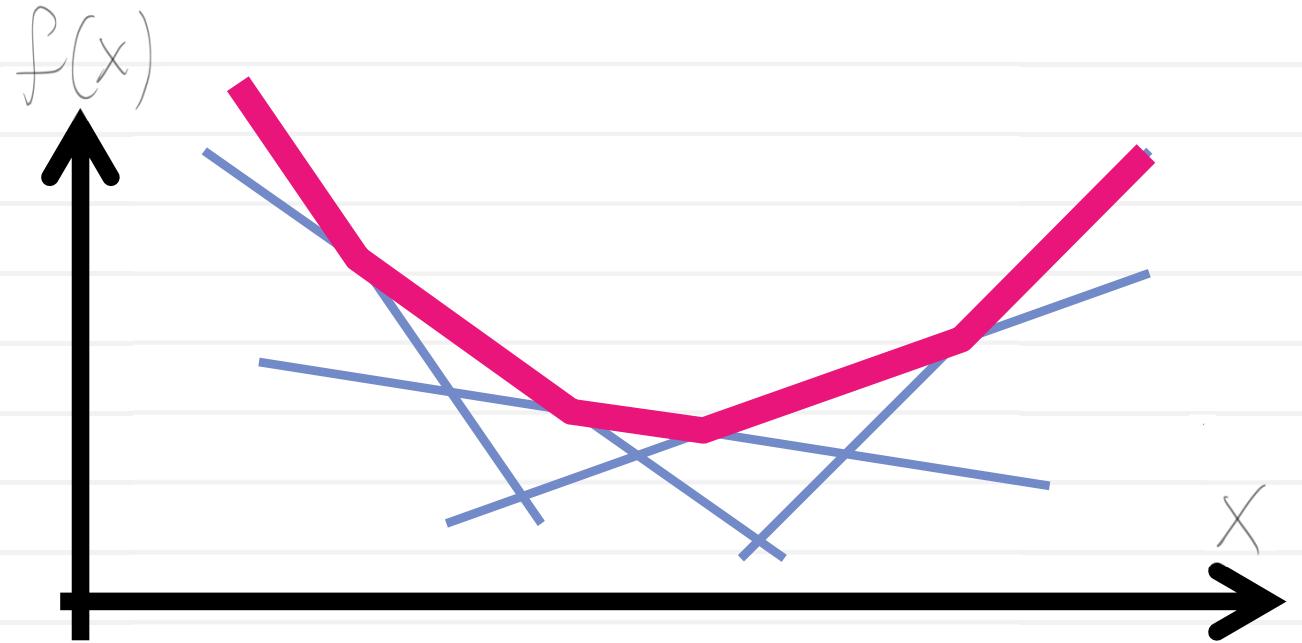
- “stashing” good solution is important
- Tuning schedule is important
- Even for a good schedule can be very slow

# Cutting plane method

$$f(x) = \max_y f(y) + \tilde{\nabla}f(y)(x-y)$$

Minimizing  $f$  is equivalent to solving the LP:

$$\min_{x, \{ \}} \{$$



$$\forall y: \{ \geq f(y) + \tilde{\nabla}f(y)(x-y)$$

Often converges much faster!

# Cutting plane method

$$f(x) = \max_y f(y) + \tilde{\nabla}f(y)(x-y)$$

$$\begin{aligned} & \min_{x, \xi} \quad \{ \\ & \quad y \in Y \quad \} \geq f(y) + \tilde{\nabla}f(y)(x-y) \end{aligned} \quad (*)$$

Initialize  $Y = \{y_1, y_2, \dots, y_n\}$

**for**  $k = 1..K$

$(y, \xi) = \operatorname{argmin} (*)$

Add  $y$  to  $Y$

**end**

# Cutting plane method

$$\begin{aligned} \min_{x} \quad & \{ \\ x, \{ \quad & (*) \\ y \in Y \quad \{ \geq f(y) + \tilde{\nabla}f(y)(x-y) \end{aligned}$$

Initialize  $Y = \{y_1, y_2, \dots, y_n\}$

for  $k = 1..K$

$(y, \xi) = \operatorname{argmin} (*)$

Add  $y$  to  $Y$

end

Problems:

- LPs grows and become too slow and too big
- The optimal points oscillate wildly (“zigzagging”)

# Bundle method

$$\min_{x, \{ \}} \{ + \lambda (x - x(k))^2 \quad (*)$$

$$y \in Y \} \geq f(y) + \tilde{\nabla}f(y)(x-y)$$

- Penalize deviation from the last solution
- The sequence of solutions changes slowly
- As with trust-region methods, can tweak  $\lambda$  on the fly
- Keep only the last T points in Y (limits time/memory)
- The devil is in details (but “black box” implementations exist)

# Plan for the second part

---

- What happens when we consider the primal and dual problems together?
- What do dual variables tell us about the primal program?
- What does duality mean for “smooth” problems?

# Duality gap and certificates of optimality

Let  $x$  to be any feasible primal set of variables, and  $(\lambda, \nu)$  be any feasible dual set of variables (i.e.  $\lambda \geq 0$ ). Then:

$$g(\lambda, \nu) \leq g(\lambda^*, \nu^*) \leq f_o(x^*) \leq f_o(x)$$



$$f_o(x) - f_o(x^*) \leq f_o(x) - g(\lambda, \nu) = \text{gap}(x, \lambda, \nu)$$

Any feasible  $(x, \lambda, \nu)$  give us a *certificate* on how far is  $f_o(x)$  from optimality.

# Strong duality and saddle point

Assume the strong duality holds (and optimal values are attained):

$$g(\lambda^*, \nu^*) = L(x^*, \lambda^*, \nu^*) = f_0(x^*)$$

$$\forall x \quad L(x, \lambda^*, \nu^*) \geq g(\lambda^*, \nu^*) = L(x^*, \lambda^*, \nu^*)$$

$$x^* = \underset{x}{\operatorname{argmin}} \ L(x, \lambda^*, \nu^*)$$

$$\forall \lambda \geq 0, \nu \quad L(x^*, \lambda, \nu) \leq f_0(x^*) = L(x^*, \lambda^*, \nu^*)$$

$$(\lambda^*, \nu^*) = \underset{\lambda \geq 0, \nu}{\operatorname{argmax}} \ L(x^*, \lambda, \nu)$$

Thus  $(x^*, \lambda^*, \nu^*)$  correspond to a saddle point of the Lagrangian.

# Duality for smooth functions

Let us assume that  $f_i, h_i$  are differentiable and strong duality holds.

We know that

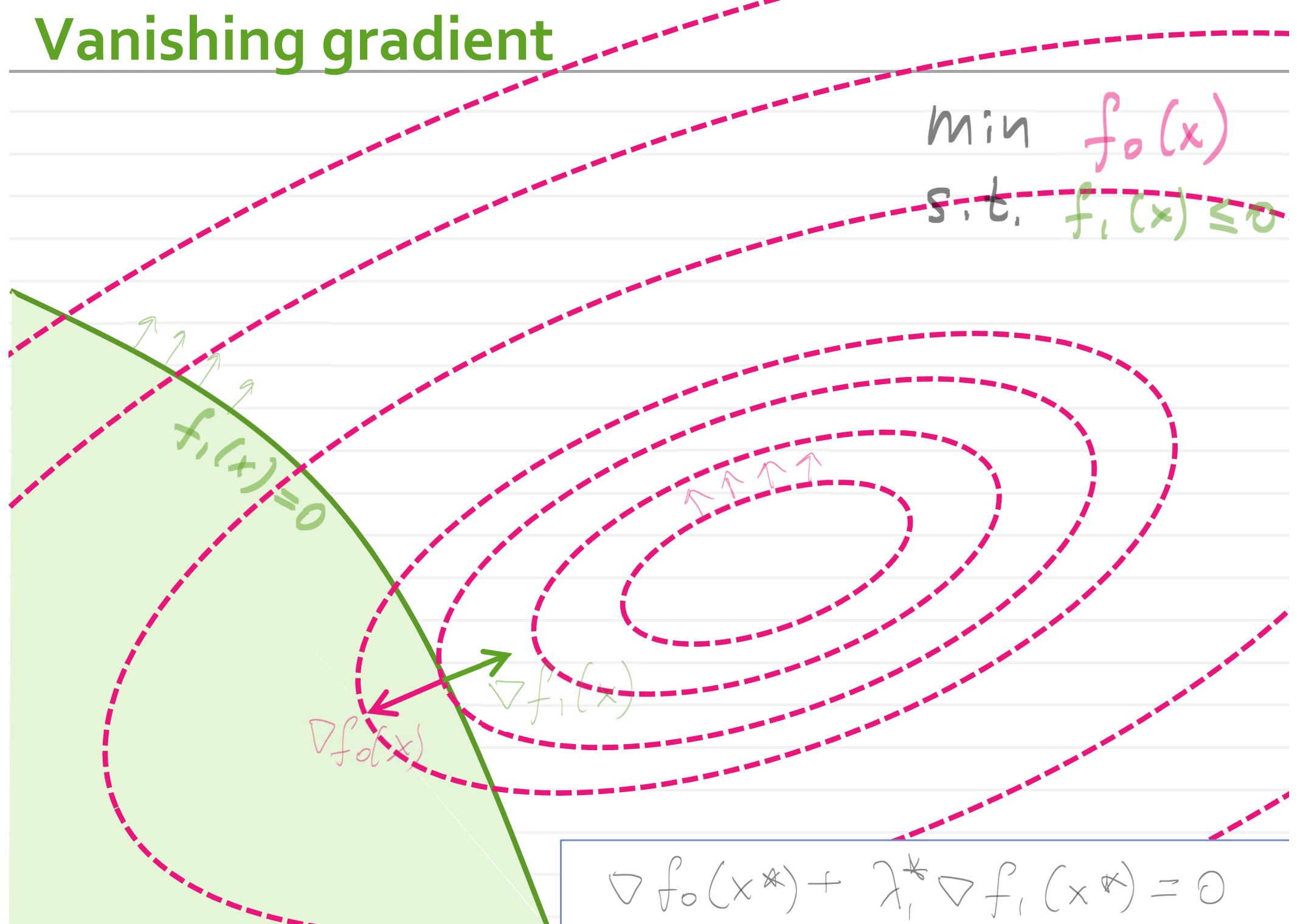
$$x^* = \underset{x}{\operatorname{argmin}} L(x, \lambda^*, v^*)$$

Then the gradient of the Lagrangian w.r.t.  $x$  must vanish:

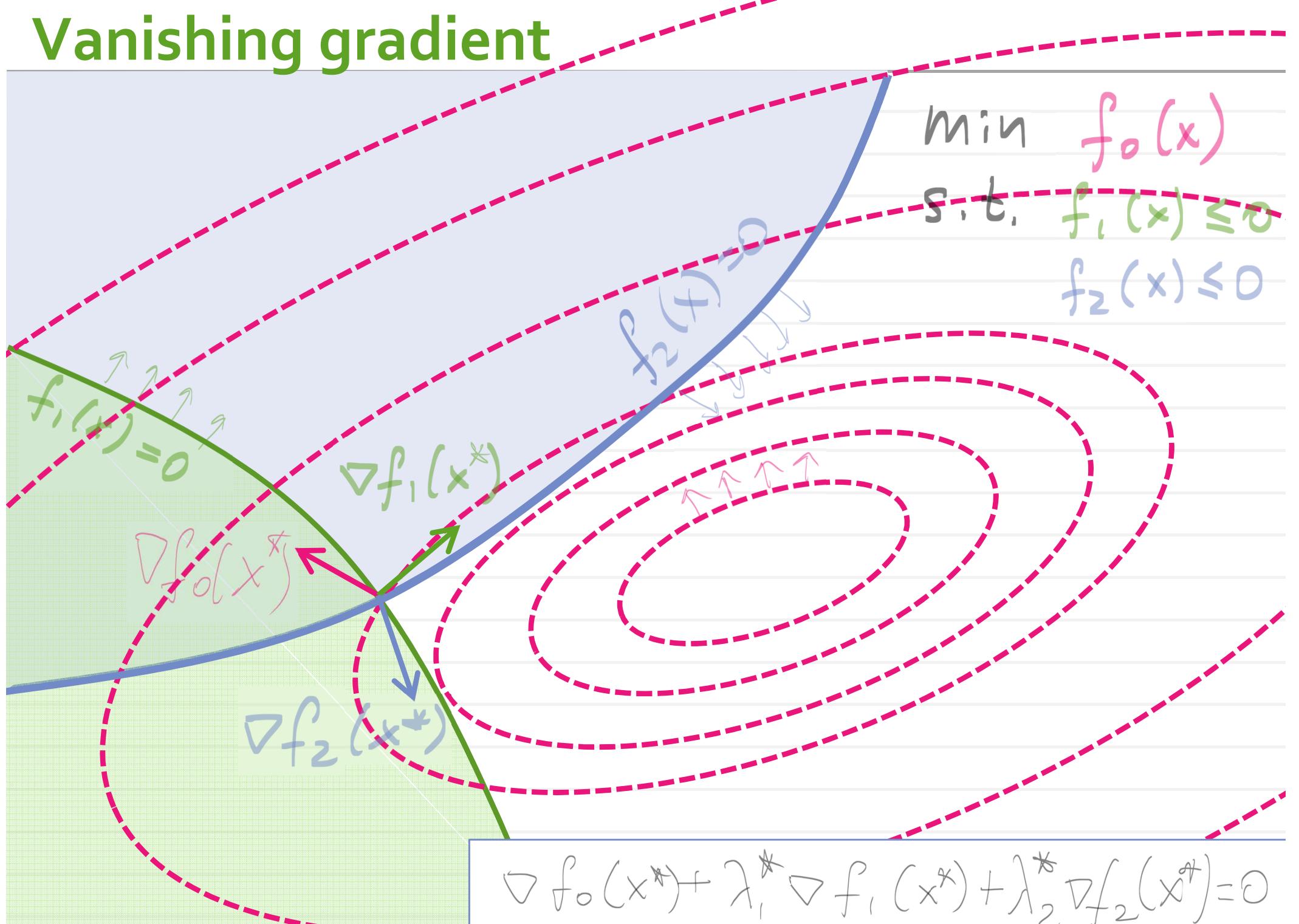
$$\frac{\partial L}{\partial x}(x^*, \lambda^*, v^*) = 0$$

$$\nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^p v_i^* \nabla h_i(x^*) = 0$$

# Vanishing gradient



# Vanishing gradient



# Complementary slackness

Assume the strong duality holds (and optimal values are attained):

$$L(x^*, \lambda^*, \nu^*) = f_0(x^*)$$

$$f_0(x^*) = f_0(x^*) + \underbrace{\sum_{i=1}^m \lambda_i^* f_i(x^*)}_{\parallel 0} + \underbrace{\sum_{i=1}^p \nu_i^* h_i(x^*)}_{\parallel 0}$$

Since  $\lambda_i^* \geq 0$   $f_i(x^*) \leq 0$

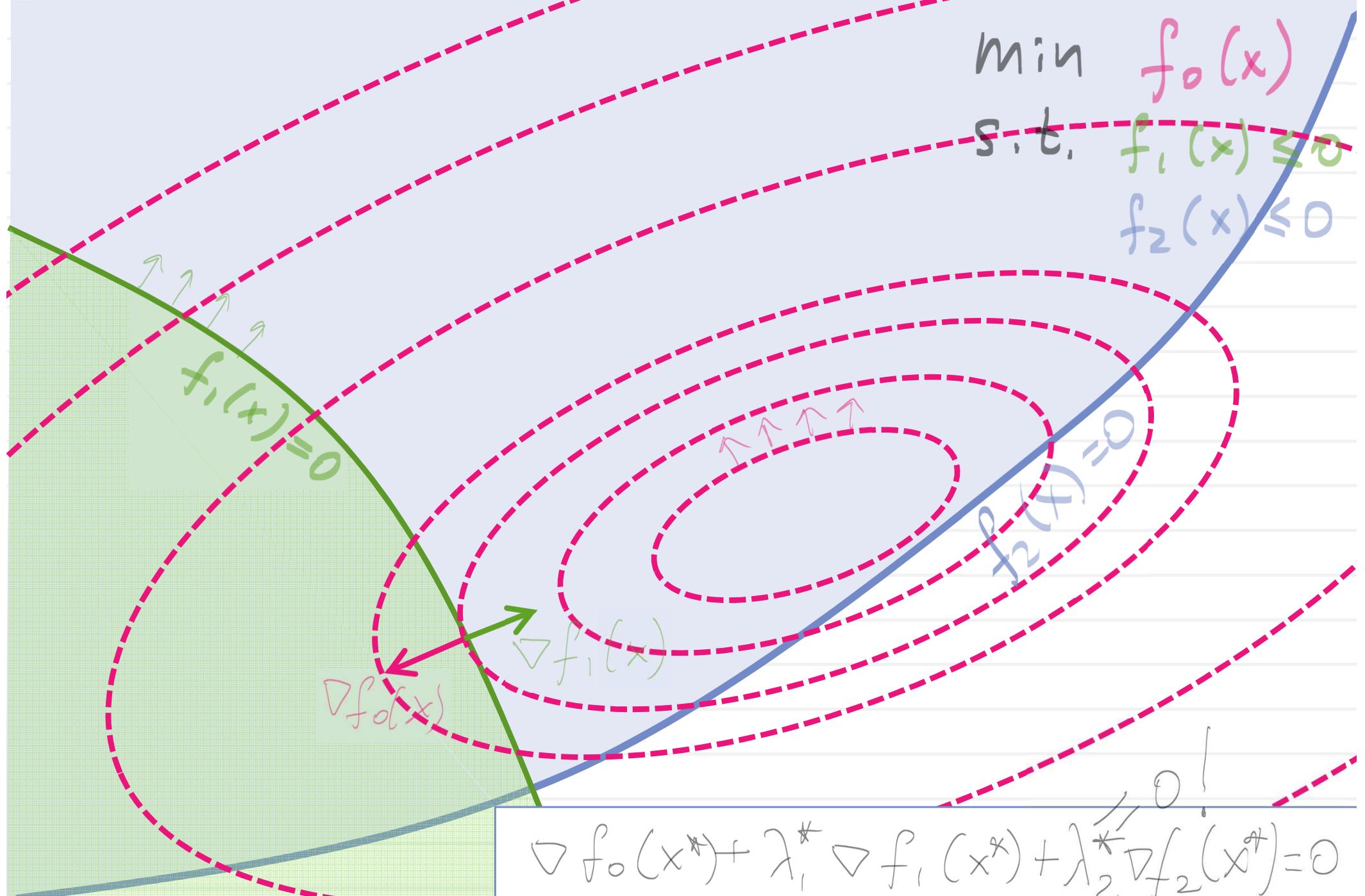
We have:  $\forall i \quad \lambda_i^* f_i(x^*) \leq 0$

Since they all sum to 0, they all should equal 0.

*Complementary slackness*

$$\forall i \quad \lambda_i^* f_i(x^*) = 0$$

# Complementary slackness



# Karush-Kuhn-Tucker conditions

$$\begin{aligned} \min_x \quad & f_0(x) \\ \text{s.t.:} \quad & f_i(x) \leq 0 \\ & h_i(x) = 0 \end{aligned}$$

If  $(x^*, \lambda^*, \nu^*)$  are a set of optimal primal and dual variables and strong duality holds, then the following must be true:

$$\left\{ \begin{array}{l} f_i(x^*) \leq 0 \quad i = 1 \dots m \quad h_i(x^*) = 0 \quad i = 1 \dots p \\ \lambda_i^* \geq 0 \quad i = 1 \dots m \\ \nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^p \nu_i^* \nabla h_i(x^*) = 0 \\ \lambda_i^* f_i(x^*) = 0 \quad i = 1 \dots m \end{array} \right.$$

# Karush-Kuhn-Tucker conditions

$$\left\{ \begin{array}{l} f_i(x^*) \leq 0 \quad i=1 \dots m \quad h_i(x^*) = 0 \quad i=1 \dots p \\ \lambda_i^* \geq 0 \quad i=1 \dots m \\ \nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^p v_i^* \nabla h_i(x^*) = 0 \\ \lambda_i^* f_i(x^*) = 0 \quad i=1 \dots m \end{array} \right.$$

# Sufficiency of KKT

**Theorem:** if the primal problem is convex and  $(x', \lambda', \nu')$  are a set of optimal primal and dual variables satisfying KKT, then they are primal and dual optimal.

**Proof:**

$$\begin{aligned} & x' = \underset{x}{\operatorname{argmin}} L(x, \lambda', \nu') \\ g(\lambda', \nu') &= L(x', \lambda', \nu') = \\ &= f_0(x') + \sum \lambda'_i f_i(x') + \sum \nu'_i h_i(x') = \\ &= f_0(x') \end{aligned}$$

# Example 1: unconstrained optimization

$$\min_x f_0(x)$$

KKT are turned into:

$$\nabla f_0(x) = 0$$

We get a Newton method.

KKT = “primal feasibility” + “dual feasibility” + “vanishing gradient” + “complementary slackness”

# Equality-constrained quadratic program

$$\min \frac{1}{2} x^T Q x + p^T x$$

$$\text{s.t.: } Ax = b$$

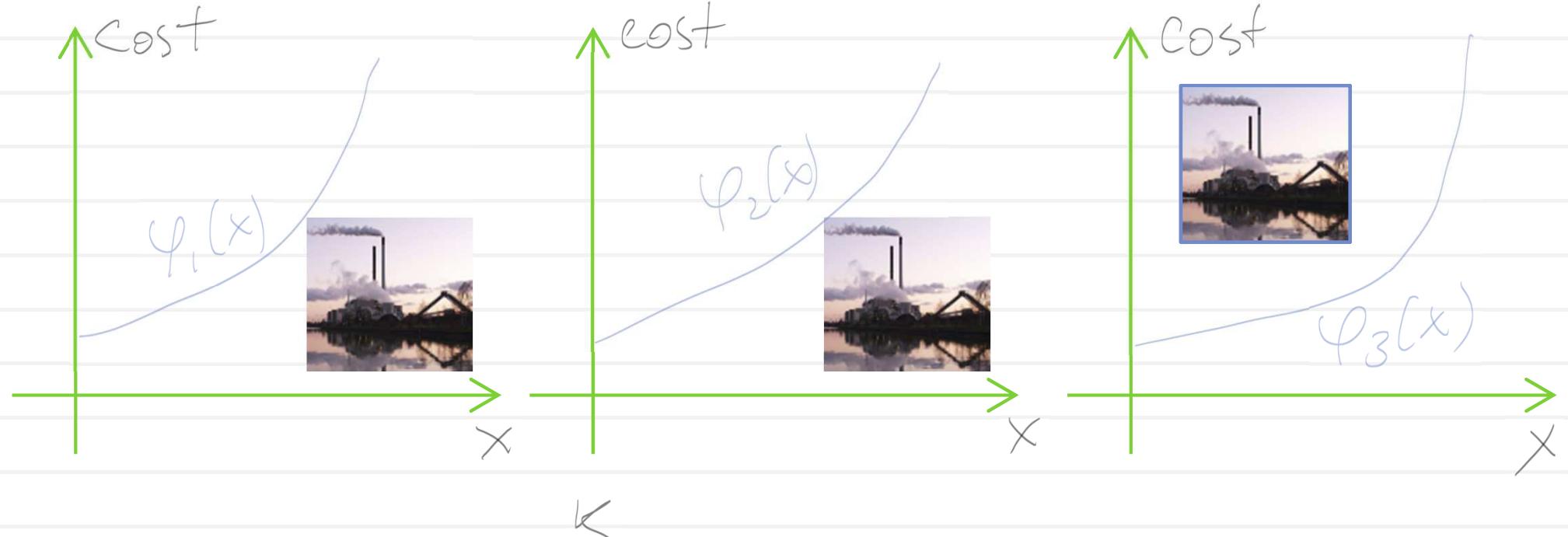
KKT results in:

$$\begin{cases} Qx + p + A^T v = 0 \\ Ax = b \end{cases}$$

*M+N equations on M+N variables*

KKT = “primal feasibility” + “dual feasibility” + “vanishing gradient” + “complementary slackness”

# Economic dispatch



$$\min_x \sum_{i=1}^k \varphi_i(x_i)$$

Total cost

$$\text{s. t.: } \sum x_i = D$$

Meet the demand

# Economic dispatch

$$\min_x \sum_{i=1}^k \varphi_i(x_i)$$

$$\text{s.t. : } \sum x_i = D$$

Note: we denote the dual variable as  $\lambda$  because of the traditional notation for this particular application (should have been  $v$ ).

KKT results in:

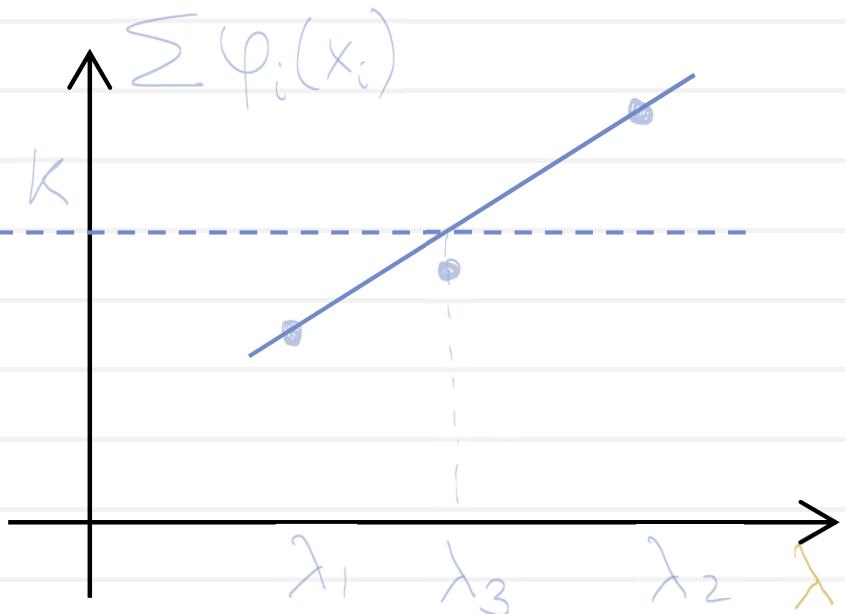
$$\left\{ \begin{array}{l} \sum \nabla \varphi_i(x) - \begin{pmatrix} \lambda \\ \lambda \\ \vdots \\ \lambda \end{pmatrix} = 0 \\ \sum x_i = D \end{array} \right. \xrightarrow{\hspace{1cm}} \left\{ \begin{array}{l} \varphi_i'(x_i) = \lambda \\ \sum x_i = D \end{array} \right.$$

KKT = “primal feasibility” + “dual feasibility” + “vanishing gradient” + “complementary slackness”

# Economic dispatch: lambda search



$$\left\{ \begin{array}{l} \varphi'_i(x_i) = \lambda \\ \sum x_i = D \end{array} \right.$$



# Inequality constrained QP

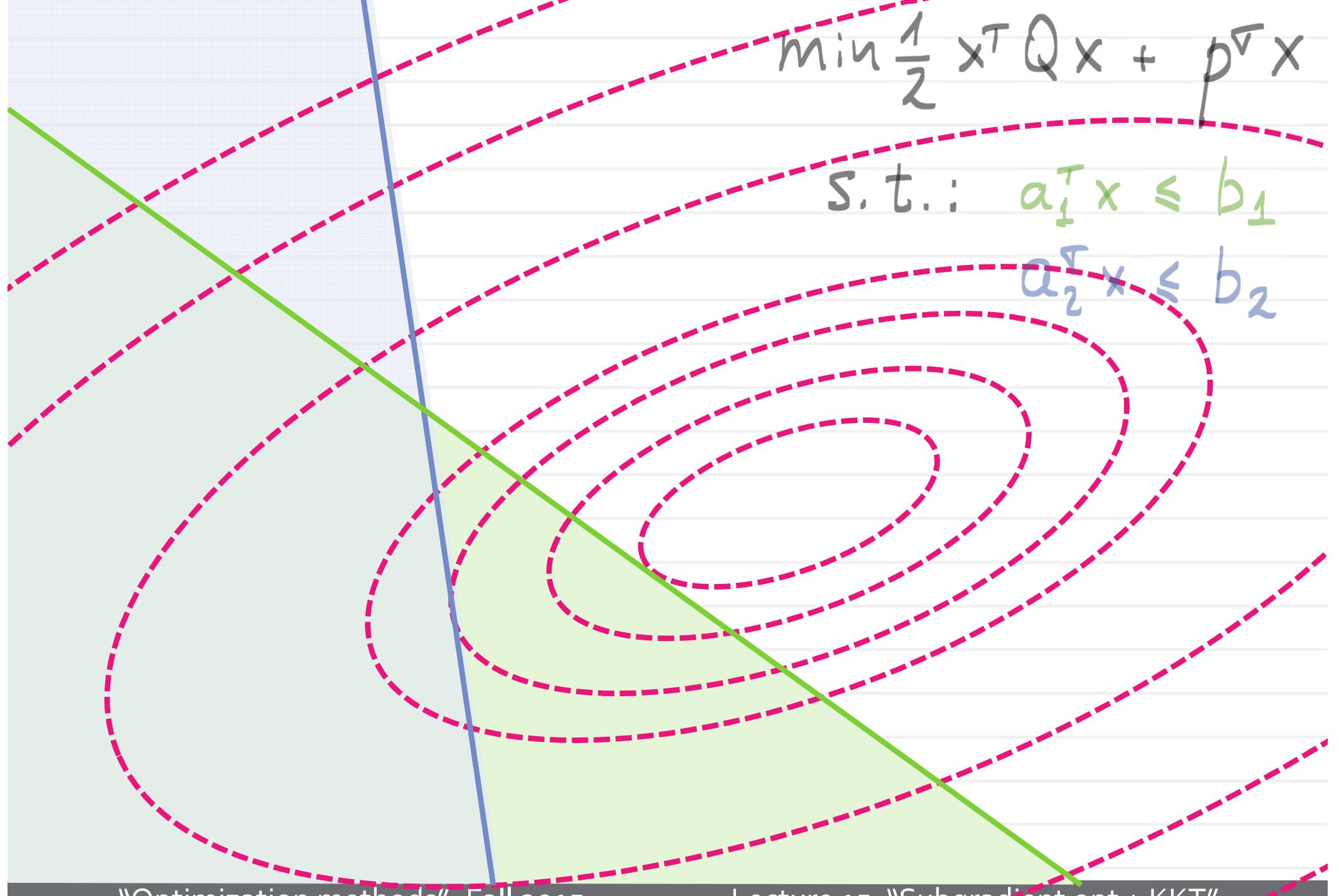
$$\min \frac{1}{2} x^T Q x + p^T x$$

$$\text{s.t.: } a_1^T x \leq b_1$$

$$a_2^T x \leq b_2$$

$$\left\{ \begin{array}{l} Qx + P + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^T x \leq b_1 \quad \lambda_1 \geq 0 \\ a_2^T x \leq b_2 \quad \lambda_2 \geq 0 \\ \lambda_1(a_1^T x - b_1) = 0 \quad \lambda_2(a_2^T x - b_2) = 0 \end{array} \right.$$

# Inequality constrained QP



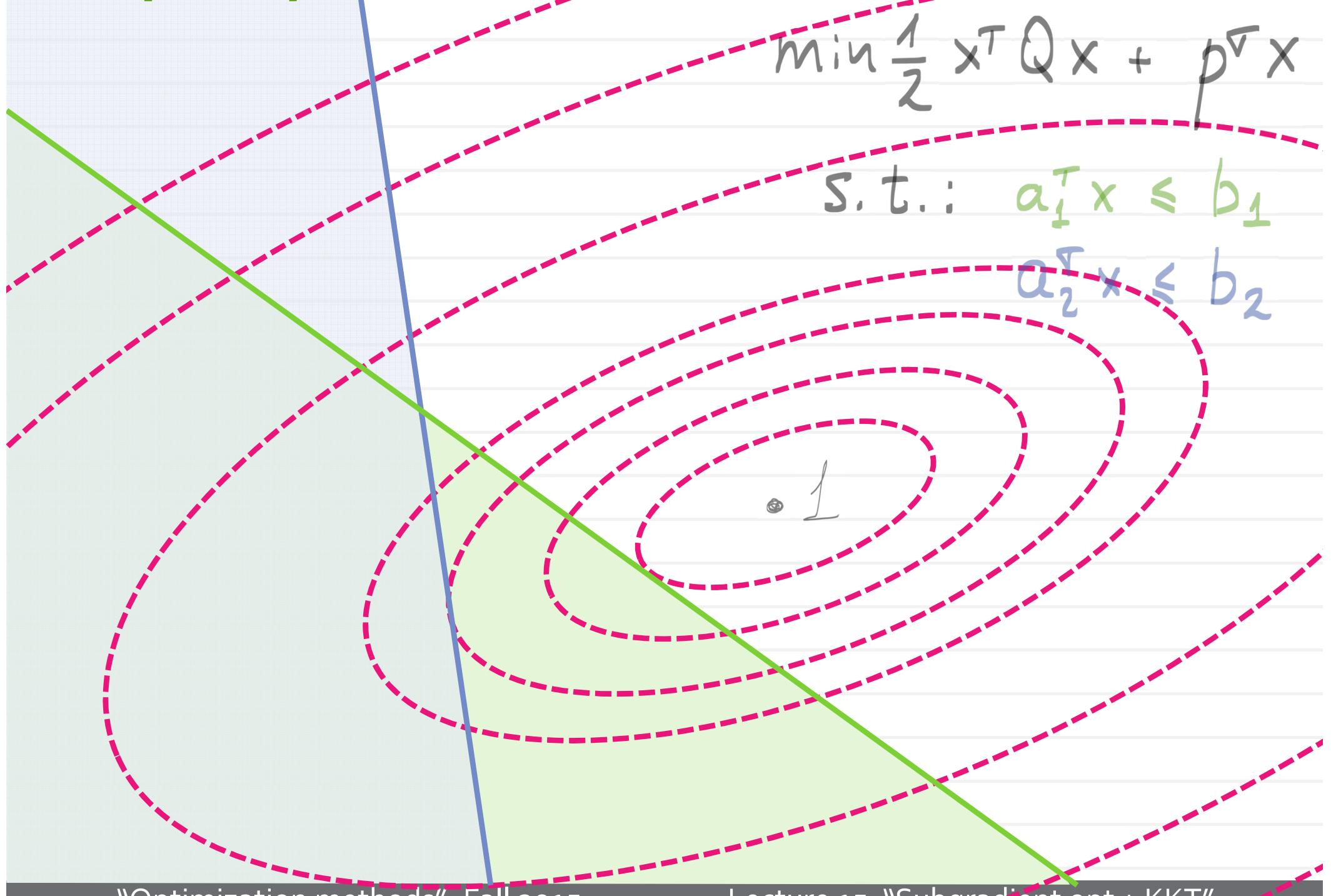
# Inequality constrained QP

$$\left\{ \begin{array}{l} Qx + P + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^T x \leq b_1 \quad \lambda_1 \geq 0 \\ a_2^T x \leq b_2 \quad \lambda_2 \geq 0 \\ \lambda_1(a_1^T x - b_1) = 0 \quad \lambda_2(a_2^T x - b_2) = 0 \end{array} \right.$$

**Case 1:**  $\lambda_1 = 0 \quad \lambda_2 = 0$

$$\left\{ \begin{array}{l} Qx + P = 0 \\ a_1^T x \leq b_1 \\ a_2^T x \leq b_2 \end{array} \right.$$

# Inequality constrained QP



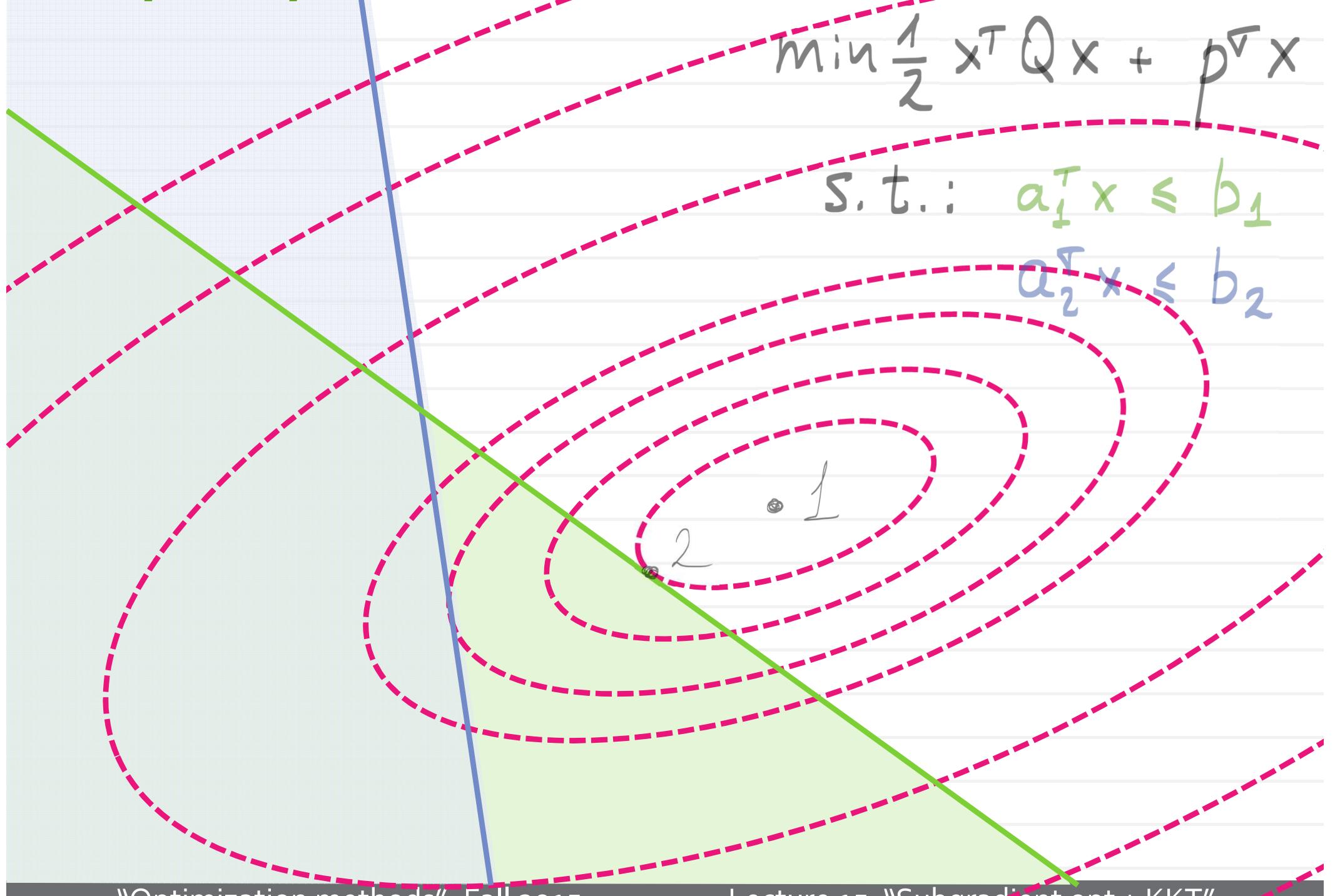
# Inequality constrained QP

$$\left\{ \begin{array}{l} Qx + P + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^\top x \leq b_1 \quad \lambda_1 \geq 0 \\ a_2^\top x \leq b_2 \quad \lambda_2 \geq 0 \\ \lambda_1(a_1^\top x - b_1) = 0 \quad \lambda_2(a_2^\top x - b_2) = 0 \end{array} \right.$$

**Case 2:**  $\lambda_1 > 0, \lambda_2 = 0$

$$\left\{ \begin{array}{l} Qx + P + a_1 \lambda_1 = 0 \\ a_1^\top x = b_1 \\ a_2^\top x \leq b_2 \quad \lambda_1 > 0 \end{array} \right.$$

# Inequality constrained QP



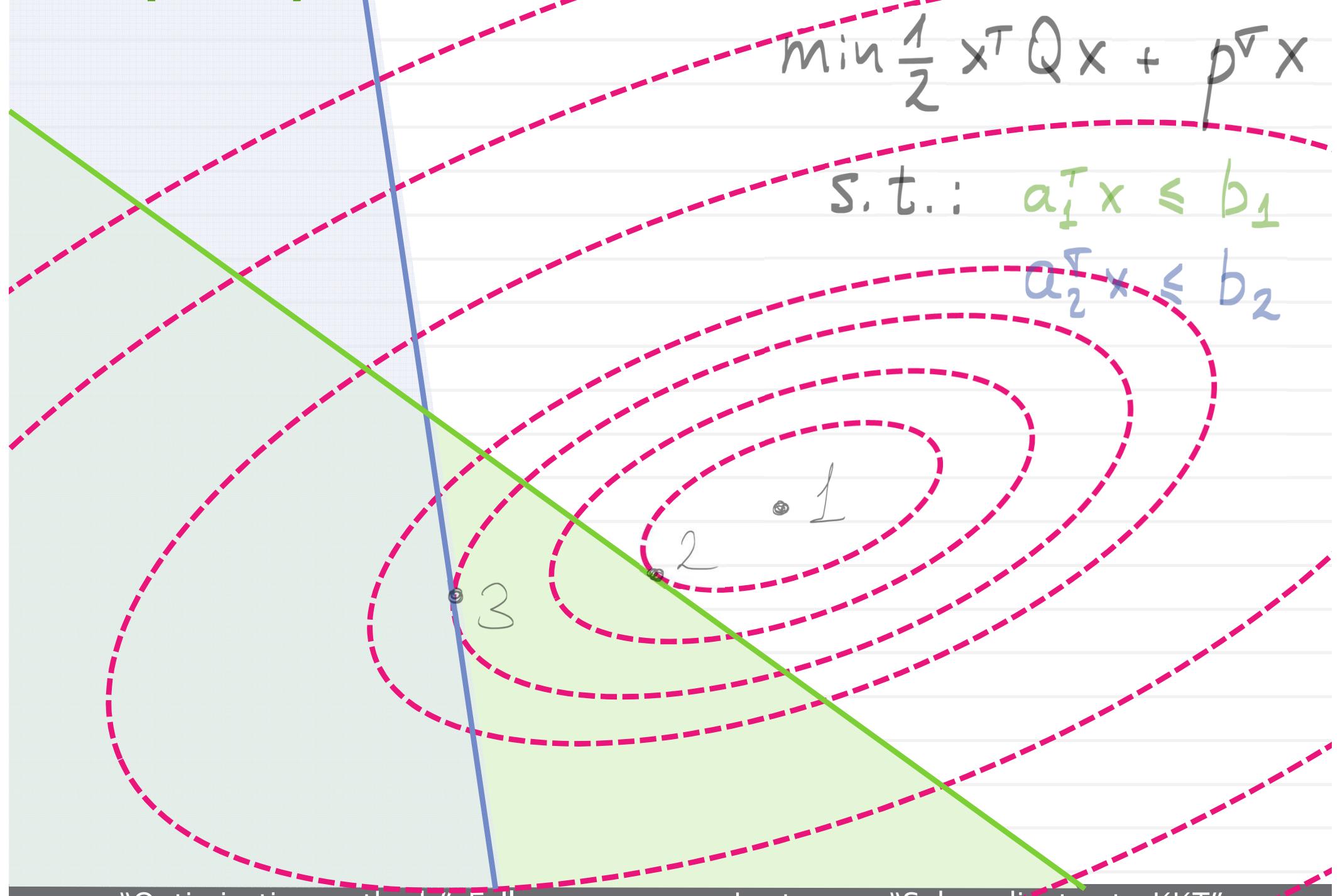
# Inequality constrained QP

$$\left\{ \begin{array}{l} Qx + P + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^\top x \leq b_1 \quad \lambda_1 \geq 0 \\ a_2^\top x \leq b_2 \quad \lambda_2 \geq 0 \\ \lambda_1(a_1^\top x - b_1) = 0 \quad \lambda_2(a_2^\top x - b_2) = 0 \end{array} \right.$$

**Case 3:**  $\lambda_1 = 0, \lambda_2 > 0$

$$\left\{ \begin{array}{l} Qx + P + a_2 \lambda_2 = 0 \\ a_2^\top x = b_2 \\ a_1^\top x \leq b_1 \quad \lambda_2 > 0 \end{array} \right.$$

# Inequality constrained QP



# Inequality constrained QP

$$\left\{ \begin{array}{l} Qx + p + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^T x \leq b_1 \quad \lambda_1 \geq 0 \\ a_2^T x \leq b_2 \quad \lambda_2 \geq 0 \\ \lambda_1(a_1^T x - b_1) = 0 \quad \lambda_2(a_2^T x - b_2) = 0 \end{array} \right.$$

**Case 4:**  $\lambda_1 > 0 \quad \lambda_2 > 0$

$$\left\{ \begin{array}{l} Qx + p + a_1 \lambda_1 + a_2 \lambda_2 = 0 \\ a_1^T x = b_1 \quad \lambda_1 > 0 \quad \lambda_2 > 0 \\ a_2^T x = b_2 \end{array} \right.$$

# Inequality constrained QP

