

This week: We prove Theorem 16.1 in the notes.

OP : $\min_{\underline{x} \in \mathbb{R}^n} f(\underline{x})$ subject to: $\left\{ \begin{array}{ll} c_i(\underline{x}) = 0 & i \in \Sigma \\ c_i(\underline{x}) \geq 0 & i \in I \end{array} \right\} \quad (1)$

Lagrangian:

$$\mathcal{L}(\underline{x}, \underline{\lambda}) = f(\underline{x}) - \sum_{i \in \Sigma \cup I} \lambda_i c_i(\underline{x})$$

Suppose \underline{x}_* is a minimizer of the OP (1). Furthermore suppose that f and the c_i 's are continuously differentiable and that the LICQ holds at \underline{x}_* . Then there exists a vector $\underline{\lambda}_*$ with n components λ_i^* ($i \in \Sigma \cup I$) such that the following conditions hold:

kKT conditions

$$\left\{ \begin{array}{ll} \nabla_{\underline{x}} \mathcal{L}(\underline{x}_*, \underline{\lambda}_*) = 0 & \\ c_i(\underline{x}_*) = 0 & i \in \Sigma \\ c_i(\underline{x}_*) \geq 0 & i \in I \\ \lambda_i^* \geq 0 & i \in I \\ \lambda_i^* c_i(\underline{x}_*) = 0 & i \in \Sigma \cup I \end{array} \right.$$

Remark: Since $c_i(\underline{x}_*) > 0$ for $i \in I \setminus A(\underline{x}_*)$,

by kKTs we have $\lambda_i^* = 0$ for $i \in I \setminus A(\underline{x}_*)$.

Hence, kKT1 can be re-written as:

$$\dots = \sum \lambda_i^* \nabla c_i(\underline{x}_*)$$

$$\nabla f(x_*) = \sum_{i \in \mathcal{A}(x_*)} \lambda_i^* \nabla c_i(x_*)$$

Some final intermediate results.

Lemma 16.1 (Fundamental Necessary Condition)

If x_* is a minimizer of the OP (1), then

$$\underline{d} \cdot \nabla f(x_*) \geq 0 \quad \forall \underline{d} \in T_{\mathcal{R}}(x_*)$$

Proof by contradiction: Suppose that there exists a $\underline{d} \in T_{\mathcal{R}}(x_*)$

such that $\underline{d} \cdot \nabla f(x_*) < 0$. (*) Let \underline{z}_k be a feasible sequence tending to x_* such that:

$$\underline{d} = \lim_{k \rightarrow \infty} \frac{\underline{z}_k - x_*}{t_k}$$

$$\lim_{t_k \rightarrow 0} \frac{\|\underline{\epsilon}_k\|_2}{t_k} = 0 \Rightarrow \|\underline{\epsilon}_k\|_2 \sim t_k^{1+\delta}$$

Re-arrange:

$$\underline{z}_k = x_* + t_k \underline{d} + \underline{\epsilon}_k, \quad \underline{\epsilon}_k = o(t_k)$$

$$\begin{aligned} f(\underline{z}_k) &= f(x_* + t_k \underline{d} + \underline{\epsilon}_k) \\ &= f(x_*) + t_k \underbrace{\underline{d} \cdot \nabla f(x_*)}_{\text{negative}} + \underbrace{o(t_k)}_{\text{negligible}} \end{aligned}$$

$\Rightarrow f(\underline{z}_k) < f(x_*)$ for t_k sufficiently small.

Contradiction, since x_* is a local minimizer.

Hence, statement (*) is false, so

Since, statement (*) is false, so

$$\underline{d} \cdot \nabla f(\underline{x}_*) \geq 0 \quad \forall \underline{d} \in T_{\Omega}(\underline{x}_*)$$

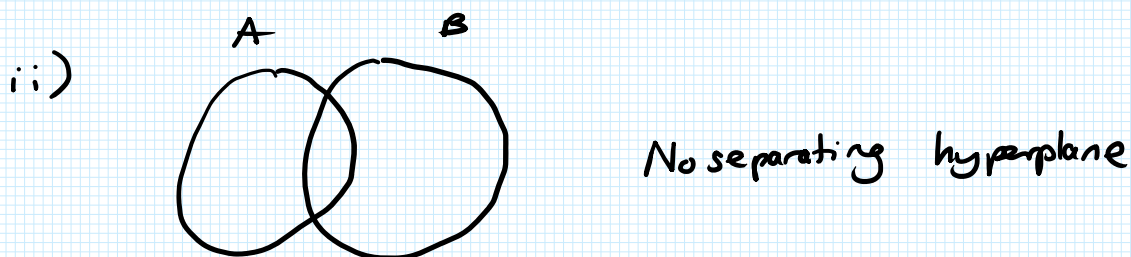
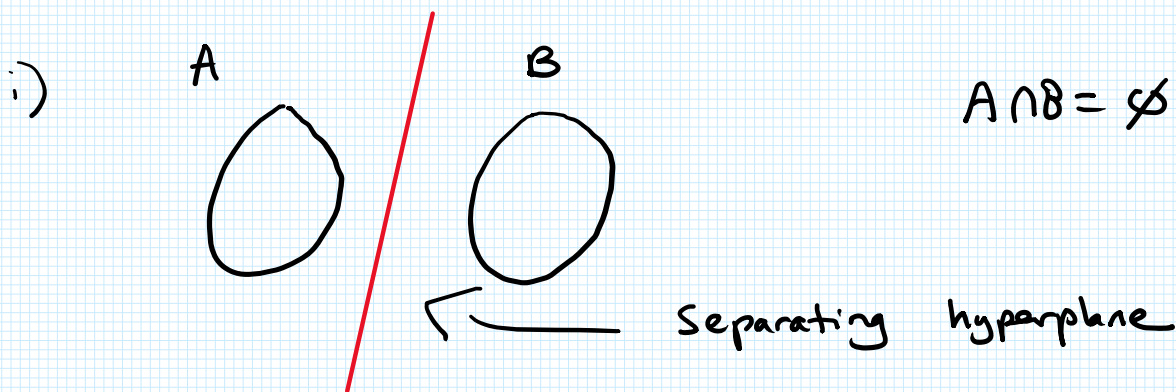
Another key result: Hyperplane Separation Theorem (Theorem 16.2)

Let A and B be two disjoint **non-empty convex subsets** of \mathbb{R}^n . Then there exists a plane which separates A and B . Mathematically, there exists a non-zero vector \underline{n} and a constant c such that:

$$\underline{x} \cdot \underline{n} \geq c \qquad \underline{y} \cdot \underline{n} \leq c$$

for all $\underline{x} \in A$ and all $\underline{y} \in B$.

No proof in this module, but see the pictures below for the idea.



Extension: Let A and B be disjoint non-empty convex subsets of \mathbb{R}^n such that:

- i) A and B are closed \rightarrow contains its own boundary
- ii) At least one of A and B is bounded.

Then there exists a non-zero vector \underline{v} and a constant c such that:

$$\underline{x} \cdot \underline{v} > c \qquad \underline{y} \cdot \underline{v} < c$$

for all $\underline{x} \in A$, $\underline{y} \in B$.

The H.S.T. enables us to prove the final intermediate result called Farkas's Lemma.

Lemma 16.2 (Farkas's Lemma)

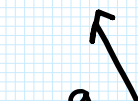
$$\text{Let } \mathcal{K} = \{ B \underline{y} \mid y_i \geq 0, i = 1, \dots, m \}$$

be a cone, where $B \in \mathbb{R}^{n \times m}$. Given any vector $\underline{g} \in \mathbb{R}^n$, either

- $\underline{g} \in \mathcal{K}$
- Or there exists a vector $\underline{d} \in \mathbb{R}^n$ such that $\underline{g} \cdot \underline{d} < 0$ and $(B^T \underline{d})_i \geq 0, i = 1, \dots, m$

Idea:

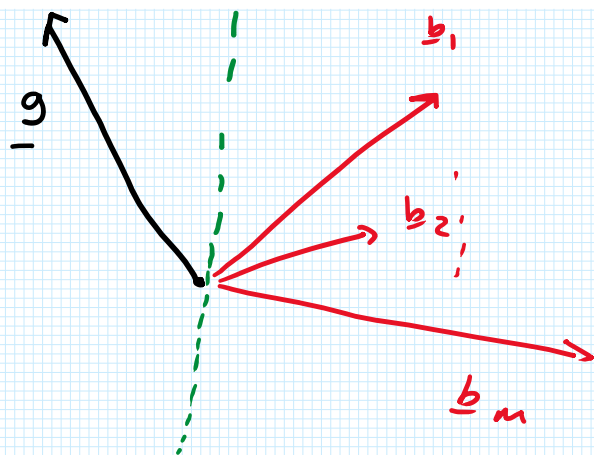
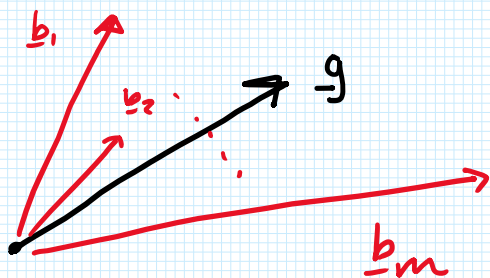
Case 2



Separating hyperplane

\underline{b}_1

Case 1



Proof: We show first that the alternatives are mutually exclusive, i.e. that either

1. $g \in K$
2. $\exists \underline{d} \in \mathbb{R}^n$ s.t. $g \cdot \underline{d} < 0$ AND $(B^T \underline{d})_i \geq 0$,
 $i = 1, \dots, m$

but not both at once.

Therefore, take $g \in K$:
 $g = B \underline{y}$, $y_i \geq 0$. } #1 is true.

Assume for contradiction that #2 is also true.

Hence, $\exists \underline{d} \in \mathbb{R}^n$ s.t.

$$g \cdot \underline{d} < 0$$

$$\Rightarrow \langle g, \underline{d} \rangle < 0$$

$$\stackrel{\#1}{\Rightarrow} \langle B \underline{y}, \underline{d} \rangle < 0$$

$$\Rightarrow \langle \underline{y}, B^T \underline{d} \rangle < 0$$

$$\text{Since } \underline{y} \geq 0 \text{ and } (B^T \underline{d})_i \geq 0$$

$$\text{Components } \sum_{i=1}^m y_i (B^T \underline{d})_i < 0 \quad \star$$

But the y_i 's are non-neg., and the $(B^T \underline{d})_i$'s are non-negative (#2 is true), hence:

$$\sum_{i=1}^m y_i (B^T \underline{d})_i \geq 0 \quad \star \star$$

Contradiction between \star and $\star \star$. Reasoning back up, if #1 is true, then #2 can't be true at the same time, so we have that #1 and #2 are mutually exclusive.

We now show that one option or the other must necessarily hold. If $g \in \mathcal{K}$ then we are done. \square

So we look at $g \notin \mathcal{K}$. We have:

$$A = \{g\}$$

$$B = \mathcal{K} = \{By \mid y_i \geq 0, i=1, \dots, m\}$$

A is bounded, B is closed and convex (cone is closed and convex). $A \cap B = \emptyset$ ($g \notin \mathcal{K}$).

Hence, the H.S.T. applies (strict form). There

exists a constant c and a vector $\underline{d} \in \mathbb{R}^n$ such that:

$$A: \quad \underline{d} \cdot \underline{g} < c$$

$$B: \quad \underline{d} \cdot \underline{s} > c \quad \forall \underline{s} \in \mathcal{R}$$

Notice: $\underline{0} \in \mathcal{R}$, $0 > c$ (c is negative).

Hence:

$$\underline{d} \cdot \underline{s} > c \quad (c \text{ negative})$$

$$\underline{s} \in \mathcal{R} \Rightarrow \underline{d} \cdot (B\underline{y}) > c$$

$$\Rightarrow \langle \underline{d}, B\underline{y} \rangle > c$$

$$\Rightarrow \langle B^T \underline{d}, \underline{y} \rangle > c$$

$$\Rightarrow \langle \underline{y}, B^T \underline{d} \rangle > c$$

$$\text{Components} \Rightarrow \sum_{i=1}^m y_i (B^T \underline{d})_i > c$$

True for all $\underline{s} \in \mathcal{R}$, therefore true for all $y_i \geq 0$.

So we can rescale the y_i 's, taking $y_i \rightarrow \lambda y_i$,

where λ is positive:

$$\sum_{i=1}^m \lambda y_i (B^T \underline{d})_i > c$$

$$\Rightarrow \sum_{i=1}^m y_i (B^T \underline{d})_i > \frac{c}{\lambda}$$

~~True~~ True for all $\lambda > 0$, so take $\lambda \rightarrow \infty$:

$$\sum_{i=1}^m y_i (B^T \underline{d})_i \geq 0$$

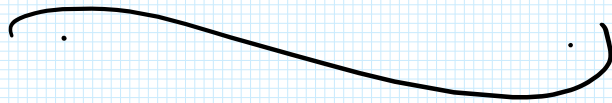
This is exactly Option # 2. H.S.T. gives us \underline{d} ($g_{\underline{d}} < 0$)

and we have shown :

$$\sum_{i=1}^m y_i (B^T \underline{d})_i \geq 0 \quad \forall y_i \geq 0$$

$$\therefore (B^T \underline{d})_i \geq 0 \quad i = 1, \dots, m.$$

Hence, Option # 2 holds. ■



Start with Farkas's Lemma (from Tuesday)

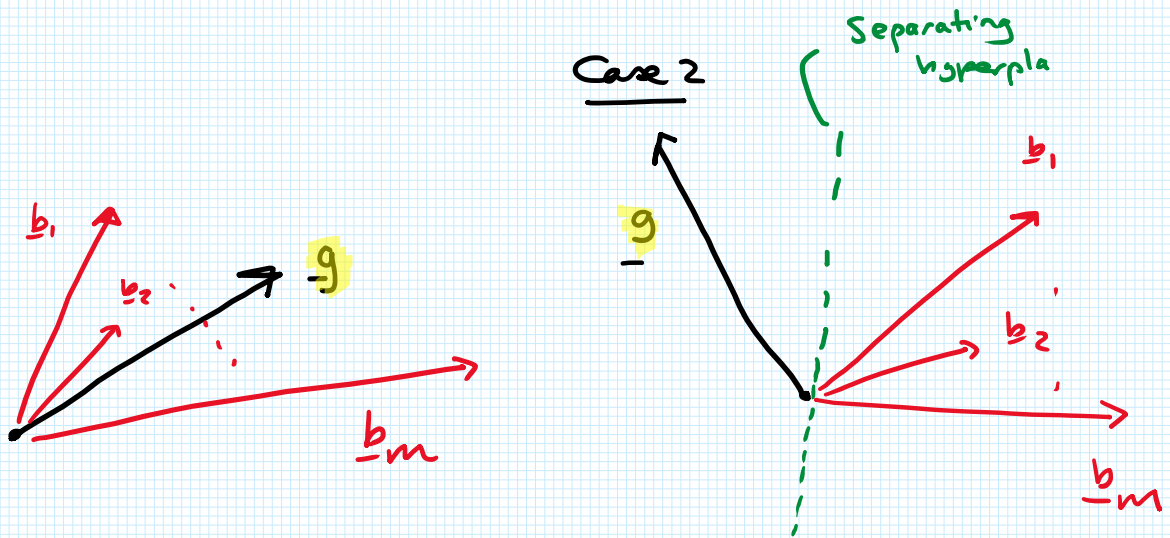
Let

$$K = \{By \mid y_i \geq 0, i=1, \dots, m\}$$

be a cone in \mathbb{R}^n , where $B \in \mathbb{R}^{n \times m}$. Let g be any vector in \mathbb{R}^n . Then, either:

- $g \in K$
- There exists a vector \underline{d} in \mathbb{R}^m such that:
 - $g \cdot \underline{d} < 0$
 - $(B^T \underline{d})_i \geq 0, i=1, \dots, m$

Idea:



Extension of Farkas's Lemma:

Let

$$K = \left\{ By + Cw \mid \begin{array}{l} y_i \geq 0, i=1, \dots, m \\ w_i \in \mathbb{R}, i=1, \dots, p \end{array} \right\} \quad (1)$$

where B and C are matrices of the appropriate dimension.

For any vector $g \in \mathbb{R}^n$, either:

For any vector $g \in \mathbb{R}^n$, either:

- $g \in \mathcal{R}$ or
- There exists a vector \underline{d} in \mathbb{R}^n such that
 - $g \cdot \underline{d} < 0$
 - $(A^T \underline{d})_i \geq 0, \quad i = 1, \dots, m$
 - $(C^T \underline{d})_i = 0, \quad i = 1, \dots, p$

We take this extension of Farkas's Lemma as given (Nocedal & Wright).

We apply the extended Farkas's Lemma.

Let

$$N = \left\{ \begin{array}{l} \sum_{i \in A(x_*) \cap I} \lambda_i \nabla c_i(x_*) + \sum_{i \in A(x_*) \cap E} \lambda_i \nabla c_i(x_*) \\ \left. \begin{array}{l} \lambda_i \geq 0, i \in A(x_*) \cap I \\ \lambda_i \in \mathbb{R}, i \in A(x_*) \cap E \end{array} \right\} \end{array} \right\}$$

where x_* is a generic feasible point.

Re-write:

$$N = \left\{ A^T \underline{\lambda} \mid \lambda_i \geq 0, i \in A(x_*) \cap I \right\} \quad (2)$$

where

$$A = \begin{pmatrix} \frac{\partial c_1}{\partial x_1}(x_*) & \dots & \frac{\partial c_1}{\partial x_n}(x_*) \\ \vdots & & \\ \frac{\partial c_m}{\partial x_1}(x_*) & \dots & \frac{\partial c_m}{\partial x_n}(x_*) \end{pmatrix} \in \mathbb{R}^{m \times n}$$

where we have relabelled constraints such that

$$A(x_*) = \{1, \dots, m\}.$$

The set (2) is of the same form as (1). So extended Farkas's Lemma applies.

Take $g = \nabla f(x_*)$. Then, either:

- $\nabla f(x_*) \in N$ or
- there exists a \underline{d} in \mathbb{R}^n such that:
 - $\underline{d} \cdot \nabla f(x_*) < 0$
 - $(A\underline{d})_i \geq 0, \quad i \in A(x_*) \cap I$
 - $(A\underline{d})_i = 0, \quad i \in A(x_*) \cap E.$

Hence, $\underline{d} \in F_{\Omega}(x_*)$.

Proof of KKT conditions

Let x_* be a local minimizer of the OP.

Then, by the fundamental necessary condition (lemma 16.1)

$$\underline{d} \cdot \nabla f(x_*) \geq 0, \quad \forall \underline{d} \in T_{\Omega}(x_*).$$

By the hypothesis of the KKT theorem, the LICQ holds at x_* , hence $T_{\Omega}(x_*) = F_{\Omega}(x_*)$, hence

holds at x_* , hence $T_{\Omega}(x_*) = F_{\Omega}(x_*)$, hence

$$\underline{d} \cdot \nabla f(x_*) \geq 0 \quad \forall \underline{d} \in F_{\Omega}(x_*).$$

So $\underline{d} \cdot \nabla f(x_*)$ can't be negative. This rules out option #2 in Farkas's Lemma (for $g = \nabla f(x_*)$).

So option #1 in (extended) Farkas's Lemma applies:

$$g = \nabla f(x_*) \in N.$$

Hence, there exist scalars λ_i such that:

$$- \lambda_i \geq 0 \quad i \in A(x_*) \cap I \quad (3)$$

$$- \lambda_i \in \mathbb{R} \quad i \in A(x_*) \cap E$$

such that

$$\overbrace{g = \nabla f(x_*)} = \sum_{i \in A(x_*)} \lambda_i \nabla C_i(x_*)$$

We may take:

$$\lambda_i^* = \begin{cases} \lambda_i & , \quad i \in A(x_*) \\ 0 & , \quad i \in I \cap A(x_*) \end{cases} \quad (4)$$

Hence:

$$\nabla f(x_*) = \sum_{i \in E \cup I} \lambda_i^* \nabla C_i(x_*) \quad \text{KKT 1}$$

Since x_* is feasible by assumption:

$$c_i(x^*) = 0, \quad i \in \mathcal{E} \quad \text{KKT 2}$$

$$c_i(x^*) \geq 0, \quad i \in \mathcal{I} \quad \text{KKT 3}$$

From (3) and (4), it follows that

$$\lambda_i^* \geq 0, \quad i \in \mathcal{I} \quad \text{KKT 4}$$

Finally:

$$i \in \mathcal{A}(x^*) \cap \mathcal{E} \quad c_i(x^*) = 0 \Rightarrow \lambda_i^* c_i(x^*) = 0$$

$$i \in \mathcal{A}(x^*) \cap \mathcal{I} \quad c_i(x^*) = 0 \Rightarrow \lambda_i^* c_i(x^*) = 0$$

$$i \in \mathcal{I} \setminus \mathcal{A}(x^*) \stackrel{\text{Eq (3)}}{\Rightarrow} \lambda_i^* = 0 \Rightarrow \lambda_i^* c_i(x^*) = 0$$

This covers all $i \in \mathcal{I} \cup \mathcal{E}$, therefore:

$$\lambda_i^* c_i(x^*) = 0, \quad i \in \mathcal{I} \cup \mathcal{E} \quad \text{KKT 5}$$

This establishes the KKT conditions. \square

Plan:

- Today
 - Second-derivative test for constrained optimization
 - Dual formulation
- Next Tuesday
 - Recorded lecture (website) Exercises #6

- Next Thursday
 - Exercises #6 wrap-up (duality)
 - Structure of final exam

Second-derivative test (Ch. 17)

Critical cone:

$$\mathcal{C}(x_*, \lambda_*) = \left\{ \underline{w} \in \mathcal{F}_\Omega(x_*) \left\{ \begin{array}{l} \underline{w} \cdot \nabla c_i(x_*) = 0, \quad i \in E \\ \underline{w} \cdot \nabla c_i(x_*) = 0, \quad i \in \mathcal{A}(x_*) \cap I \\ \text{and } \lambda_i^* > 0 \\ \underline{w} \cdot \nabla c_i(x_*) \geq 0 \quad i \in \mathcal{A}(x_*) \cap I \\ \text{and } \lambda_i^* = 0 \end{array} \right. \right\}$$

Since $\lambda_i^* = 0$ for $i \in I \setminus \mathcal{A}(x_*)$

$$\lambda_i^* \underline{w} \cdot \nabla c_i(x_*) = 0$$

for all $i \in I \cup E$ and all $\underline{w} \in \mathcal{C}(x_*, \lambda_*)$.

Suppose that the KKT conditions are satisfied at x_* :

$$\nabla f(x_*) = \sum_{i \in \mathcal{A}(x_*)} \lambda_i^* \nabla c_i(x_*)$$

Dot both sides with $\underline{w} \in \mathcal{C}(x_*, \lambda_*)$:

$$\underline{w} \cdot \nabla f(x_*) = \sum_{i \in \mathcal{A}(x_*)} \lambda_i^* \underline{w} \cdot \nabla c_i(x_*)$$

KKT conditions can be re-written as a first-derivative test:

$$\underline{w} \cdot \nabla f(x_*) = 0, \quad \forall \underline{w} \in \mathcal{C}(x_*, \lambda_*)$$

$$\underline{w} \cdot \nabla f(x_*) = 0, \quad \forall \underline{w} \in \mathcal{C}(x_*, \lambda_*)$$

Theorem 17.1 Suppose that x_* is a local s.t.p. of the OP and that the LICQ holds at x_* . Let λ_* be the Lagrange multiplier for which the KKT conditions hold.

Then:

$$\langle \underline{w}, \nabla_{xx} \mathcal{L}(x_*, \lambda_*) \underline{w} \rangle \geq 0 \quad \forall \underline{w} \in \mathcal{C}(x_*, \lambda_*)$$

where $\nabla_{xx} \mathcal{L}$ is the Hessian matrix.

Theorem 17.2: Suppose that for some feasible point x_* there exists a vector of Lagrange multipliers satisfying the KKT conditions. Suppose also that

$$\langle \underline{w}, \nabla_{xx} \mathcal{L}(x_*, \lambda_*) \underline{w} \rangle > 0 \quad \forall \underline{w} \in \mathcal{C}(x_*, \lambda_*).$$

Then x_* is a strict local minimizer of the OP.

Duality (§ 17.2)

To fix ideas, assume:

- Inequality constraints only:

$$c_i(x) \geq 0, \quad i \in \{1, \dots, m\}$$

- f and $-c_i$ are convex functions (essential).

- f and $-c_i$ are convex functions (essential).

OP:

$$\min_{\underline{x} \in \mathbb{R}^n} f(\underline{x}) \quad \text{subject to} \quad c_i(\underline{x}) \geq 0, \quad i \in \{1, \dots, m\} \quad \underline{(1)}$$

Introduce a general (no constraints) Lagrangian:

$$\mathcal{L}(\underline{x}, \underline{\lambda}) = f(\underline{x}) - \sum_{i=1}^m \lambda_i c_i(\underline{x})$$

$$\underline{\lambda} = (\lambda_1, \dots, \lambda_m)^T.$$

Dual objective function:

$$q: \mathbb{R}^m \rightarrow \mathbb{R}$$

$$\underline{\lambda} \mapsto q(\underline{\lambda})$$

where

$$q(\underline{\lambda}) = \inf_{\underline{x} \in \mathbb{R}^n} \mathcal{L}(\underline{x}, \underline{\lambda}) \quad \text{"Legendre transformation"}$$

We require q to be bounded below, so we restrict to $\underline{\lambda} \in \mathcal{D}$ where

$$\mathcal{D} = \{ \underline{\lambda} \in \mathbb{R}^m \mid q(\underline{\lambda}) > -\infty \}$$

Under certain conditions, solving the OP (1) is equivalent to solving the dual problem:

$$\max_{\underline{\lambda} \in \mathcal{D}} q(\underline{\lambda}) \quad \text{subject to} \quad \lambda_i \geq 0 \quad i \in \{1, \dots, m\}$$

$$\max_{\underline{\lambda} \in \mathbb{R}^m} q(\underline{x}) \quad \text{subject to} \quad \lambda_i \geq 0 \quad i \in \{1, \dots, m\}$$

\nwarrow
Eq. (2)

Terminology:

- OP (1) is called the primal problem.
- OP (2) is called the dual problem.

Application: Support-Vector Machines, in Machine Learning.

Theorem 17.4: Suppose that \underline{x}^* solves the OP (1).

Then any $\underline{\lambda}$ for which $(\underline{x}^*, \underline{\lambda})$ solve the KKT conditions solves the dual problem (2).

The other way around is true as well: A solution $\underline{\lambda}^*$ to the dual problem gives rise to a solⁿ of the primal problem (Theorem 17.5).

Why? High-dimensional parameter space (\mathbb{R}^n). Solution of KKT conditions requires solving for $n+m$ unknowns (O.P. in n dimensions). Solution of dual problem is an O.P. in m dimensions. So it makes sense to solve the dual problem when $m < n$.

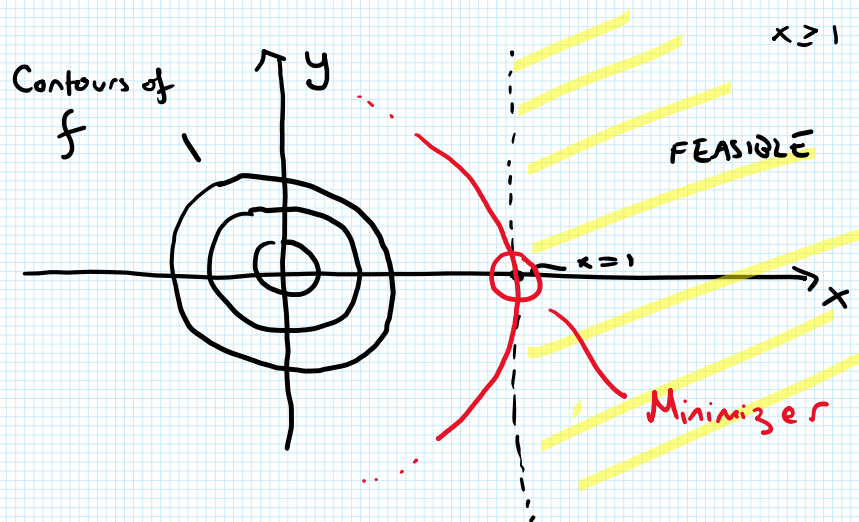
Example (§ 17.2.1)

$$\min_{\underline{x} \in \mathbb{R}^2} \frac{1}{2} (x^2 + y^2) \quad \text{subject to} \quad x \geq 1 \quad (x-1 \geq 0).$$

Primal problem: 2D optimization

Dual problem (λ): 1D optimization

Graphical solution:



$$\underline{x}_* = (1, 0)^T.$$

$$\mathcal{L}(\underline{x}, \lambda) = \frac{1}{2} (x^2 + y^2) - \lambda (x-1)$$

$$q(\lambda) = \inf_{\underline{x} \in \mathbb{R}^2} \mathcal{L}(\underline{x}, \lambda)$$

Compute:

$$\frac{\partial \mathcal{L}}{\partial x} = 0, \quad \frac{\partial \mathcal{L}}{\partial y} = 0.$$

↓

$$x - \lambda = 0 \quad y = 0 \quad \Rightarrow \quad \begin{cases} x = \lambda \\ y = 0 \end{cases}$$

$$\begin{aligned} \therefore q(\lambda) &= \mathcal{L}(x=\lambda, y=0, \lambda) \\ &= \left[\frac{1}{2} (x^2 + y^2) - \lambda (x-1) \right]_{(x=\lambda, y=0)} \end{aligned}$$

$$= \left[\frac{1}{2}(x^2 + y^2) - \lambda(x-1) \right]_{(x=\lambda, y=0)}$$

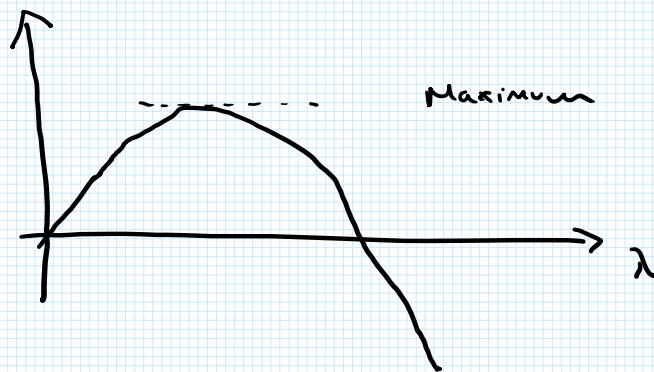
$$= \frac{1}{2}\lambda^2 - \lambda(\lambda-1)$$

$$= \frac{1}{2}\lambda^2 - \lambda^2 + \lambda$$

$$\Rightarrow \boxed{q(\lambda) = -\frac{1}{2}\lambda^2 + \lambda}$$

Dual problem:

$$\max_{\lambda \in \mathbb{R}} q(\lambda), \quad \lambda \geq 0.$$



$$\frac{dq}{d\lambda} = -\lambda + 1, \quad \frac{dq}{d\lambda} = 0 \Rightarrow \lambda = 1$$

Solution: $\lambda_* = 1.$

Back to: $x = \lambda, y = 0.$

Therefore: $x_* = 1, y_* = 0.$ ▣