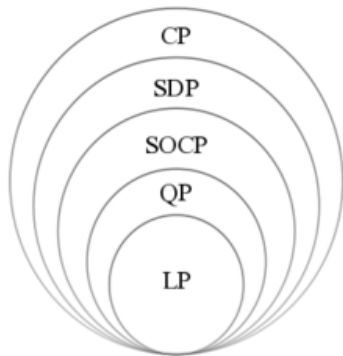


Optimization in Machine Learning

Linear Programming



Learning goals

- Instances of LPs underlying statistical estimation
- Definition of an LP
- Geometric intuition of LPs

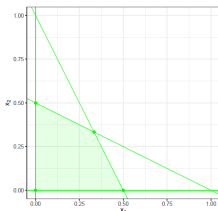
LINEAR PROGRAMMING

Special case: Linear programming (LP)

objective function and constraints are **linear functions**.

Example:

$$\begin{array}{ll}\min & -x_1 - x_2 \\ \text{s.t.} & x_1 + 2x_2 \leq 1 \\ & 2x_1 + x_2 \leq 1 \\ & x_1, x_2 \geq 0\end{array}$$



LINEAR PROGRAMMING

- (Sparse) Quantile regression:

$$\begin{aligned} \min_{\beta \in \mathbb{R}^p} \quad & \frac{1}{n} \sum_{i=1}^n \rho_{\tau} \left(y^{(i)} - \beta_0 - \beta^T \mathbf{x}^{(i)} \right) \\ \text{s.t.} \quad & \|\beta\|_1 \leq t \end{aligned}$$

where scalar β_0 and $\beta \in \mathbb{R}^p$ are the quantile regression coefficients for $\tau \in [0, 1]$, and $\rho_{\tau}(\cdot)$ is the check function defined as

$$\rho_{\tau}(s) = \begin{cases} \tau \cdot s & \text{for } s > 0 \\ -1(1 - \tau) \cdot s & \text{for } s \leq 0 \end{cases}$$

When parameter $\tau = 1/2$, quantile regression amounts to median regression, least absolute error (LAE), or least absolute deviation (LAD). As in Lasso, $t \geq 0$ is a tuning parameter.

LINEAR PROGRAMMING

- Dantzig selector:

$$\begin{aligned} \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad & \|\boldsymbol{\beta}\|_1 \\ \text{s.t.} \quad & \|\mathbf{X}^T(\mathbf{X}\boldsymbol{\beta} - \mathbf{y})\|_\infty \leq \lambda \end{aligned}$$

where $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times p}$, and $\lambda > 0$ is a tuning parameter. The infinity norm is defined as $\|x\|_\infty = \max\{|x_1|, \dots, |x_i|, \dots, |x_n|\}$ is

The Dantzig selector is similar (and behaves similar) to the Lasso and was introduced for variable selection in the seminal paper by Terence Tao and Emmanuel Candès (see moodle page for reference).

Details about LPs in statistical estimation can be found, e.g., in the PhD thesis of [Yonggong Gao](#)).

LINEAR PROGRAMMING

W.l.o.g. Linear programming is specified using the so-called **standard form**.

$$\begin{array}{ll}\max_{\mathbf{x} \in \mathbb{R}^n} & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & \mathbf{Ax} \leq \mathbf{b} \\ & \mathbf{x} \geq 0\end{array}$$

The inequality constraints $\mathbf{Ax} \leq \mathbf{b}$ ($\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$) and $\mathbf{x} \geq 0$ are to be understood componentwise.

The condition $\mathbf{x} \geq 0$ is known as the **non-negativity constraint** and the vector \mathbf{c} is known as the **cost vector**.

LINEAR PROGRAMMING

Linear optimization problems can be converted to the standard form by the following operations:

- Maximization instead of minimization: multiplication of the cost vector \mathbf{c} by -1
- Less than or equal instead of greater than or equal: multiply the inequality by -1
- Equality instead of inequality: replace $\mathbf{a}_i\mathbf{x} = b_i$ with two conditions of inequality $\mathbf{a}_i\mathbf{x} \geq b_i$ and $\mathbf{a}_i\mathbf{x} \leq b_i$
- Variable without non-negativity constraint: replace x_i with $x_i^+ - x_i^-$ with $x_i^+, x_i^- \geq 0$ (positive or negative part).

In the following we assume that the LP is given in standard form.

LINEAR PROGRAMMING

Example:

The example above

$$\begin{array}{ll}\min & -x_1 - x_2 \\ \text{s.t.} & x_1 + 2x_2 \leq 1 \\ & 2x_1 + x_2 \leq 1 \\ & x_1, x_2 \geq 0\end{array}$$

can also be formulated as

$$\begin{array}{ll}\max & (1, 1) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ \text{s.t.} & \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix} \mathbf{x} \leq \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ & \mathbf{x} \geq 0\end{array}$$

GEOMETRIC INTERPRETATION

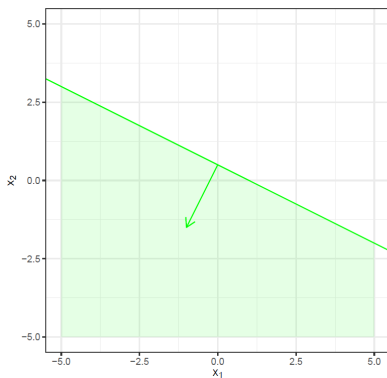
Linear programming can be interpreted geometrically.

Feasible set:

- Let $\mathbf{a}_i \mathbf{x} \geq b_i$ be the i -th line of the inequality conditions.
- The points that satisfy the linear system $\mathbf{a}_i \mathbf{x} = b_i$ form a hyperplane in n -dimensional space.
- The vector \mathbf{a}_i is perpendicular to the plane and is called the **normal vector**.

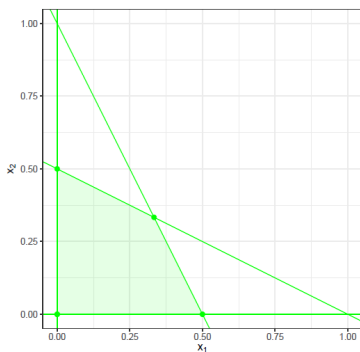
GEOMETRIC INTERPRETATION

- The set of points $\{\mathbf{x} : \mathbf{a}_i \mathbf{x} \geq b_i\}$ consists of points on the side of the hyperplane into which the normal vector points (**half-space**).



GEOMETRIC INTERPRETATION

- Each of the m inequalities divides the n -dimensional space into two halves.
- **Claim:** The points that satisfy **all** inequalities form a **convex polytope**.



GEOMETRIC INTERPRETATION

In geometry, a **polytope** denotes a generalized polygon in an arbitrary dimension. A polytope consists of several subpolytopes.

- A 0 polytope is a single point,
- A 1 polytope is a line,
- A 2 polytope is a polygon, ...

In general, a d polytope is formed from several $(d - 1)$ polytopes (so-called facets) which in turn can have a $(d - 2)$ polytope in common. Thus a 3 polytope (e.g. a cube) has several sides / facets, some of which have common edges, etc.

GEOMETRIC INTERPRETATION

The points for which $\mathbf{a}_i \mathbf{x} = b_i$ applies lie on the **facet** of the polytope.

The polytope which is defined by the inequalities $\mathbf{Ax} \geq \mathbf{b}$ is convex. For two points $\mathbf{x}_1, \mathbf{x}_2$, which lie in the polytope, any point that results from the convex combination of the two points, is again inside the polytope:

$$\begin{aligned}\mathbf{A}(\mathbf{x}_1 + t(\mathbf{x}_2 - \mathbf{x}_1)) &= \mathbf{Ax}_1 + t(\mathbf{Ax}_2 - \mathbf{Ax}_1) \\ &= (1-t) \underbrace{\mathbf{Ax}_1}_{\geq \mathbf{b}} + t \underbrace{\mathbf{Ax}_2}_{\geq \mathbf{b}} \\ &\geq (1-t)\mathbf{b} + t\mathbf{b} = \mathbf{b} \quad \text{for } t \in [0, 1]\end{aligned}$$

A polytope formed by the convex hull of $(n+1)$ affine independent points in \mathbb{R}^n is also called **n -simplex**.

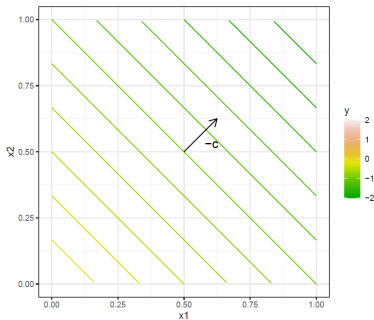
GEOMETRIC INTERPRETATION

Objective function:

- For the objective function, we consider the **contour lines**. Points that lie on the same contour line have the same objective function value.
- In the linear case, the contours $y = \mathbf{c}^T \mathbf{x}$ are also a hyperplane for fixed y .
- The vector \mathbf{c} is again perpendicular to the respective contour lines.
- The vector \mathbf{c} can also be interpreted as a gradient. The **negative** gradient $-\mathbf{c}$ points in the direction of the “steepest” descent.

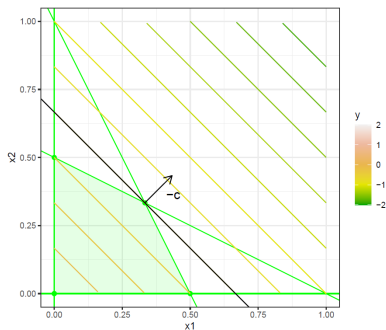
GEOMETRIC INTERPRETATION

- The value of the function becomes smaller when we go in the direction of the **negative gradient** — $-c$.



GEOMETRIC INTERPRETATION

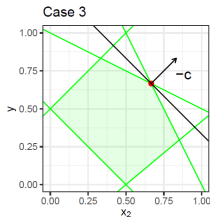
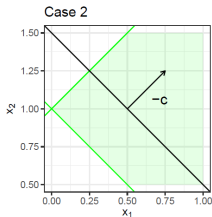
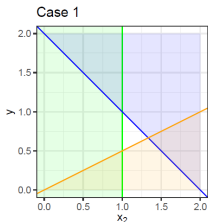
- So we move the objective function **in the opposite direction** of the normal vector until the line just touches the polygon.



SOLUTIONS TO LP

There are three ways to solve Linear programming:

- 1 LP is **infeasible**, the feasible set is empty ($\mathcal{S} = \emptyset$)
- 2 LP is unconstrained
- 3 LP has at least one optimal solution



SOLUTIONS TO LP

- If LP is solvable and constrained (neither case 1 nor case 2), there is always an optimal point that can **not** be convexly combined from other points in the polytope.
- The optimal solution is then a corner, edge or side of the polytope.