

# Methodologies and Algorithms for Structured Mixed-Integer Nonlinear Optimization

Andrés Gómez

Department of Industrial & Systems Engineering  
University of Southern California

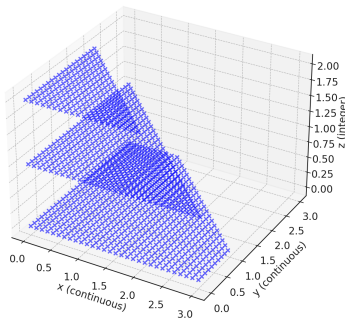
MIP Workshop 2025

# Agenda

- 1 Introduction
- 2 Branch and bound
- 3 Convexification

# Mixed-integer linear optimization (MILO)

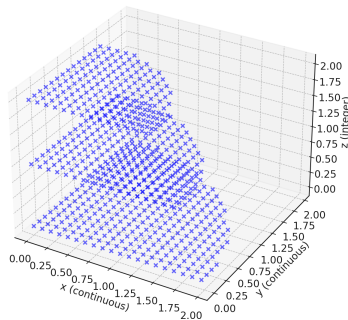
$$\begin{aligned} \min \quad & \mathbf{c}^\top \mathbf{x} + \mathbf{d}^\top \mathbf{z} \\ \text{s.t.} \quad & \mathbf{Ax} + \mathbf{Gz} \leq \mathbf{b} \\ & \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$



- Usually solved using branch-and-cut
- NP-hard in theory, *often* solvable in practice

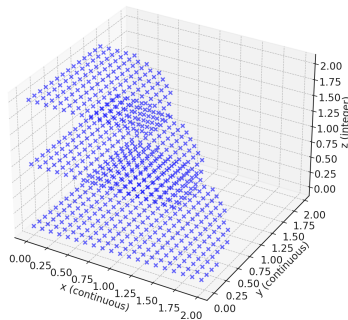
# Mixed-integer nonlinear optimization (MINLO)

$$\begin{aligned} \min f(\mathbf{x}, \mathbf{z}) \\ \text{s.t. } g_i(\mathbf{x}, \mathbf{z}) \leq 0 \quad i = 1, \dots, p \\ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$



# Mixed-integer nonlinear optimization (MINLO)

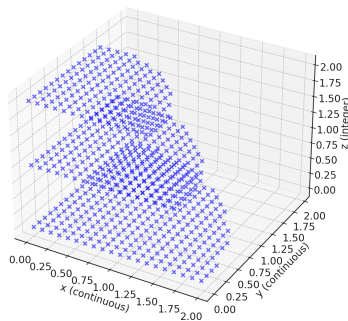
$$\begin{aligned} \min f(\mathbf{x}, \mathbf{z}) \\ \text{s.t. } g_i(\mathbf{x}, \mathbf{z}) \leq 0 \quad & i = 1, \dots, p \\ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$



- **Undecidable** Includes *Hilbert's 10th problem*
  - Given polynomial  $g$ , does there exist  $\mathbf{z} \in \mathbb{Z}^n$  satisfying  $g(\mathbf{z}) = 0$

# Mixed-integer nonlinear optimization (MINLO)

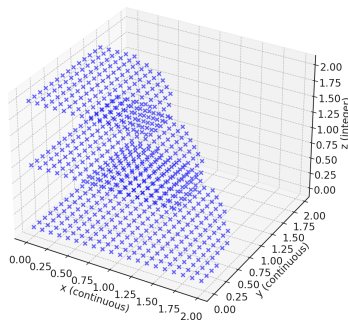
$$\begin{aligned} \min f(\mathbf{x}, \mathbf{z}) \\ \text{s.t. } g_i(\mathbf{x}, \mathbf{z}) \leq 0 \quad & i = 1, \dots, p \\ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$



- **Undecidable** Includes *Hilbert's 10th problem*
  - Given polynomial  $g$ , does there exist  $z \in \mathbb{Z}^n$  satisfying  $g(z) = 0$
- **Unstructured** Letting  $g(z) = z - z^2$ ,  $z \in \{0, 1\} \Leftrightarrow g(z) = 0$

# Mixed-integer nonlinear optimization (MINLO)

$$\begin{aligned} \min f(\mathbf{x}, \mathbf{z}) \\ \text{s.t. } g_i(\mathbf{x}, \mathbf{z}) \leq 0 \quad i = 1, \dots, p \\ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$



- **Undecidable** Includes *Hilbert's 10th problem*
  - Given polynomial  $g$ , does there exist  $z \in \mathbb{Z}^n$  satisfying  $g(z) = 0$
- **Unstructured** Letting  $g(z) = z - z^2$ ,  $z \in \{0, 1\} \Leftrightarrow g(z) = 0$

→ We assume continuous relaxation is “nice” (e.g.,  $f$  and  $g_i$  are convex)

# Mixed-integer nonlinear optimization (MINLO)

$$\begin{aligned} \min & f(\mathbf{x}, \mathbf{z}) \\ \text{s.t.} & g_i(\mathbf{x}, \mathbf{z}) \leq 0 && i = 1, \dots, p \\ & \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$

## Transformations

- Objective is linear:  $\min y$  s.t.  $f(\mathbf{x}, \mathbf{z}) \leq y$
- Single constraint:  $g(\mathbf{x}, \mathbf{z}) \leq 0$  with  $g(\mathbf{x}, \mathbf{z}) = \max_i g_i(\mathbf{x}, \mathbf{z})$
- Unconstrained:  $F(\mathbf{x}, \mathbf{z}) = \begin{cases} f(\mathbf{x}, \mathbf{z}) & \text{if } g(\mathbf{x}, \mathbf{z}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$



## Current “state-of-the-art” for MINLO

- Much less understood and mature than MILOs
- Concepts like number of variables/constraints are “uninformative”
- Most solvers and researchers are focused elsewhere
- Unlike MILOs, most of the heavy-lifting is left to the user

# Agenda

- 1 Introduction
- 2 Branch and bound**
  - Branch and bound for MILO
  - Branch and bound for MINLO
- 3 Convexification

# Agenda

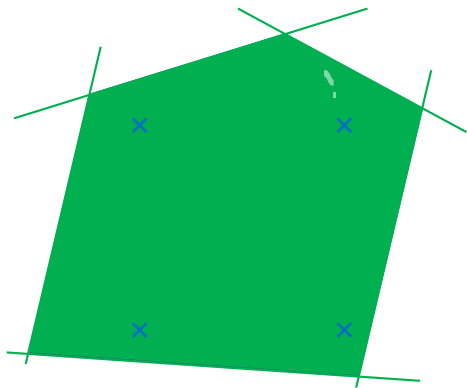
- 1 Introduction
- 2 Branch and bound
  - Branch and bound for MILO
  - Branch and bound for MINLO
- 3 Convexification

# Branch-and-cut for MILO



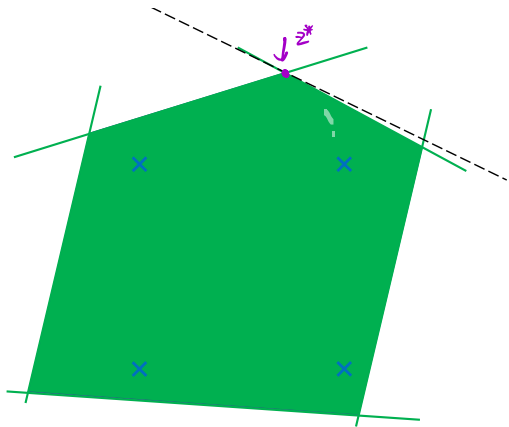
Discrete feasible region

# Branch-and-cut for MILO



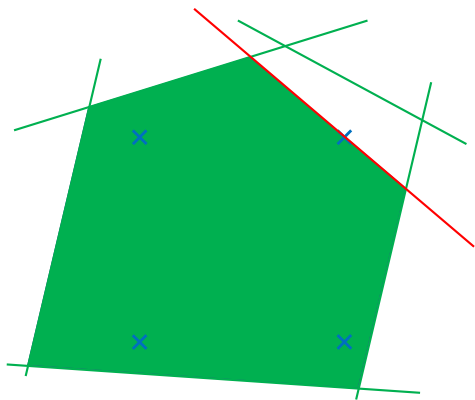
Linear programming relaxation

# Branch-and-cut for MILO



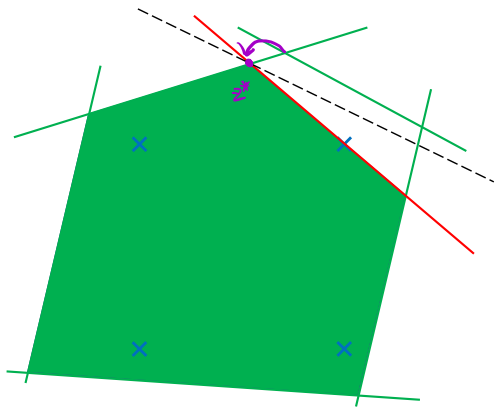
Solve (extreme point solution)

# Branch-and-cut for MILO



Improve relaxation (cutting plane)

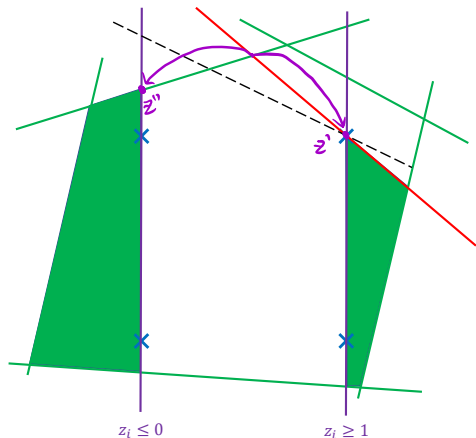
# Branch-and-cut for MILO



Solve (dual simplex)



# Branch-and-cut for MILO



Branch and resolve (dual simplex with two independent subproblems)

# Branch-and-cut for MILO

**Algorithm** To solve a mixed-integer **linear** program

- Start with a **linear** relaxation
- Dynamically refine using **cutting planes**
- Branch when needed
- Reoptimize using the **simplex method**
- When upper bound (best solution) = lower bound (relaxation), stop
  
- Other techniques
  - Heuristics, often based on rounding solutions from **linear** relaxations
  - Presolve, to improve the initial **linear** relaxation

# Branch-and-cut for MILO

**Algorithm** To solve a mixed-integer **linear** program

- Start with a **linear** relaxation
- Dynamically refine using **cutting planes**
- Branch when needed
- Reoptimize using the **simplex method**
- When upper bound (best solution) = lower bound (relaxation), stop
  
- Other techniques
  - Heuristics, often based on rounding solutions from **linear** relaxations
  - Presolve, to improve the initial **linear** relaxation

Algorithms revolve around deriving and exploiting linear relaxations

# Agenda

- 1 Introduction
- 2 **Branch and bound**
  - Branch and bound for MILO
  - **Branch and bound for MINLO**
- 3 Convexification

# Branch-and-cut for MINLO

The same algorithm works!

# Branch-and-cut for MINLO

The same algorithm works!

...but how to solve the continuous relaxations?

## Second order methods

- Hard to warm start (after branching or cuts)
  - Memory intensive
- Adds up when # nodes  $> 10^6$

## First order methods

- May struggle in heavily constrained problems
  - High-quality solutions difficult to obtain
- Numerical precision can be an issue

# Branch-and-cut for MINLO

Numerical precision is a very real issue in MINLO

```

C:\WINDOWS\system32\cmd.exe - run.bat
  0  0  0.0029  3  Cone: 235  847
  0  0  0.0029 161  0.0029  874
* 0+ 0  0.1800  0.0029  98.38%
  0  0  0.0029 200  0.1800  MIRcuts: 1  959  98.37%
  0  0  0.0029 200  0.1800  MIRcuts: 1  987  98.37%
* 0+ 0  0.0035  0.0029 16.90%
  0  2  0.0029 200  0.0035  0.0029 1016 16.65%
Elapsed time = 1.66 sec. (3548.12 ticks, tree = 0.01 MB, solutions = 2)
  3  5  0.0029 198  0.0035  0.0029 1104 16.36%
  6  8  0.0029 195  0.0035  0.0029 1189 16.36%
  9 11  0.0029 193  0.0035  0.0029 1276 16.36%
Integer feasible solution rejected --- infeasible on original model
 10 12  0.0029 192  0.0035  0.0029 1305 16.36%
 13 15  0.0029 190  0.0035  0.0029 1400 16.36%
 16 18  0.0029 187  0.0035  0.0029 1490 16.36%
 19 21  0.0030 185  0.0035  0.0029 1583 16.36%
Integer feasible solution rejected --- infeasible on original model
 20 22  0.0029 184  0.0035  0.0029 1616 16.36%
 23 25  0.0029 181  0.0035  0.0029 1713 16.36%
Integer feasible solution rejected --- infeasible on original model
 32 34  0.0029 173  0.0035  0.0029 1981 16.36%
Elapsed time = 3.25 sec. (7110.85 ticks, tree = 0.01 MB, solutions = 2)
Integer feasible solution rejected --- infeasible on original model
 42 44  0.0029 165  0.0035  0.0029 2281 16.36%
Integer feasible solution rejected --- infeasible on original model
Integer feasible solution rejected --- infeasible on original model
 50 52  0.0029 157  0.0035  0.0029 2517 16.36%
Integer feasible solution rejected --- infeasible on original model
 60 62  0.0029 147  0.0035  0.0029 2819 16.36%
Integer feasible solution rejected --- infeasible on original model

```

# Portfolio optimization

Given potential investments  $\{1, \dots, n\}$ , find a small portfolio maximizing return and minimizing risk



- **Decision variables**  $\mathbf{x} \in \mathbb{R}^n$ , where  $x_i = \%$  invested in security  $i$
- **Return**  $\boldsymbol{\mu} \in \mathbb{R}^n$ , where  $\mu_i =$  expected profit of investment  $i$   
→ Total return:  $\boldsymbol{\mu}^\top \mathbf{x}$
- **Risk**  $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ , where  $\Sigma_{ij} =$  covariance of returns from  $i$  and  $j$   
→ Variance of portfolio:  $\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}$
- **Size**  $\#$  of nonzero elements of  $\mathbf{x}$  is small



# Portfolio optimization

$$\max_{\mathbf{x}, \mathbf{z}} \boldsymbol{\mu}^\top \mathbf{x}$$

$$\text{s.t. } \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} \leq \alpha$$

$$\mathbf{1}^\top \mathbf{x} = 1$$

$$\mathbf{0} \leq \mathbf{x} \leq \mathbf{z}$$

$$\mathbf{1}^\top \mathbf{z} \leq k$$

$$\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n$$

$$\min_{\mathbf{x}, \mathbf{z}} \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}$$

$$\text{s.t. } \boldsymbol{\mu}^\top \mathbf{x} \geq \beta$$

$$\mathbf{1}^\top \mathbf{x} = 1$$

$$\mathbf{0} \leq \mathbf{x} \leq \mathbf{z}$$

$$\mathbf{1}^\top \mathbf{z} \leq k$$

$$\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n$$

Which formulation is preferable?

## Branch-and-cut for MIQO

$$\begin{aligned} \min & f(\mathbf{x}, \mathbf{z}) \\ \text{s.t.} & \mathbf{Ax} + \mathbf{Gz} \leq \mathbf{b} \\ & \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m \end{aligned}$$

where  $f$  is quadratic

Continuous relaxations can be solved via the simplex method<sup>12</sup>

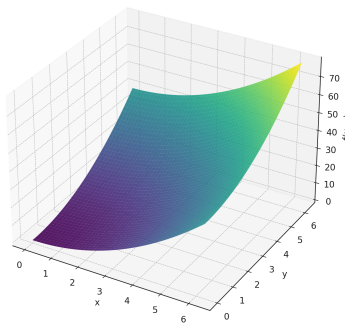
Keeping quadratic terms in the objective seems to help in MINLO

---

<sup>1</sup>Wolfe P (1959) The Simplex method for quadratic programming. *Econometrica*

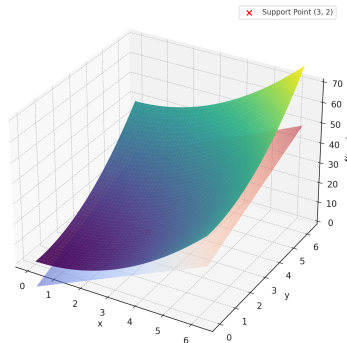
<sup>2</sup>Van de Panne C and Whinston A (1964) Simplicial methods for quadratic programming. *Naval Research Logistics*

# Linear outer approximations



Consider constraint  $f(\mathbf{x}) \leq t$

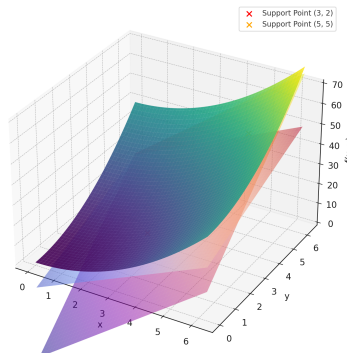
# Linear outer approximations



Consider constraint  $f(\mathbf{x}) \leq t$

Given  $\bar{\mathbf{x}}$ , can be relaxed as  $f(\bar{\mathbf{x}}) + \nabla f(\bar{\mathbf{x}})^\top (\mathbf{x} - \bar{\mathbf{x}}) \leq t$

# Linear outer approximations



Consider constraint  $f(\mathbf{x}) \leq t$

Given  $\bar{\mathbf{x}}$ , can be relaxed as  $f(\bar{\mathbf{x}}) + \nabla f(\bar{\mathbf{x}})^\top (\mathbf{x} - \bar{\mathbf{x}}) \leq t$

This process can be repeated for different support points

# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

- Assume  $UB=100$
- Construct an initial linear OA

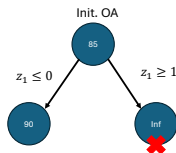
Init. OA



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

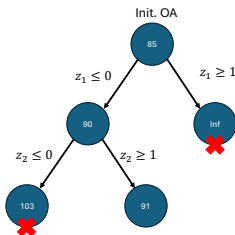
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual

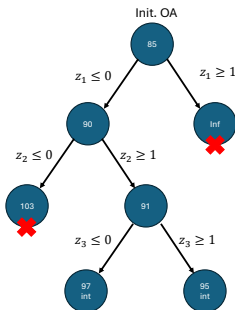




# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

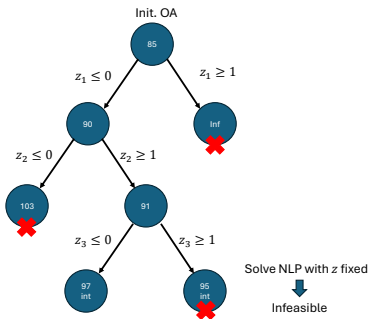
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

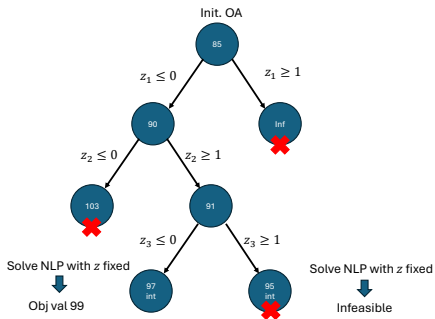
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual
- Integer nodes might be infeasible



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

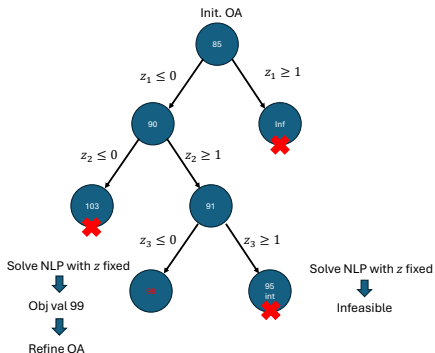
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual
- Integer nodes might be infeasible
- Incumbents obj values need to be handled carefully



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

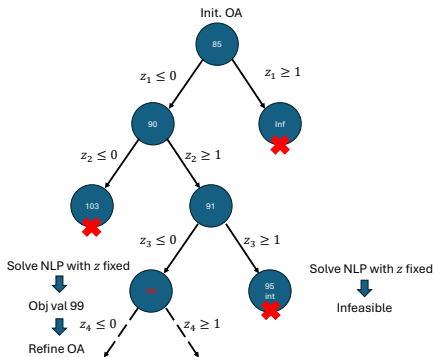
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual
- Integer nodes might be infeasible
- Incumbents obj values need to be handled carefully



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

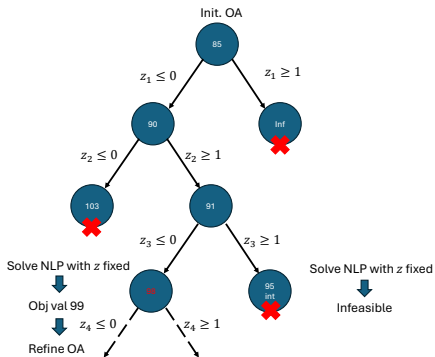
- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual
- Integer nodes might be infeasible
- Incumbents obj values need to be handled carefully
- No pruning at integer nodes



# Linear outer approximations in branch-and-bound

How to integrate in branch-and-bound?

- Assume  $UB=100$
- Construct an initial linear OA
- Branching, pruning by bound/infeasibility as usual
- Integer nodes might be infeasible
- Incumbents obj values need to be handled carefully
- No pruning at integer nodes



How to best construct linear outer approximations?

# Constructing effective linear outer approximations

Assume support points  $\{\bar{\mathbf{x}}^j\}_{j=1}^r$  and approximate<sup>3</sup>

$$f(\mathbf{x}) = \sum_{i=1}^n h_i(x_i)?$$

**Direct** Add  $r$  linear inequalities

$$f(\mathbf{x}) \geq f(\bar{\mathbf{x}}^j) + \nabla f(\bar{\mathbf{x}}^j)^\top (\mathbf{x} - \bar{\mathbf{x}}^j), \quad \forall j = 1, \dots, r$$

---

<sup>3</sup>Tawarmalani M and Sahinidis N (2005) A polyhedral branch-and-cut approach to global optimization. *Mathematical Programming*

# Constructing effective linear outer approximations

Assume support points  $\{\bar{\mathbf{x}}^j\}_{j=1}^r$  and approximate<sup>3</sup>

$$f(\mathbf{x}) = \sum_{i=1}^n h_i(x_i)?$$

**Direct** Add  $r$  linear inequalities

$$f(\mathbf{x}) \geq f(\bar{\mathbf{x}}^j) + \nabla f(\bar{\mathbf{x}}^j)^\top (\mathbf{x} - \bar{\mathbf{x}}^j), \quad \forall j = 1, \dots, r$$

**Extended** Add  $n$  variables and  $nr$  linear inequalities

$$f(\mathbf{x}) \geq \sum_{i=1}^n y_i$$

$$y_i \geq h_i(\bar{x}_i^j) + h_i'(\bar{x}_i^j)(x_i^j - \bar{x}_i^j), \quad \forall i = 1, \dots, n, j = 1, \dots, r$$

---

<sup>3</sup>Tawarmalani M and Sahinidis N (2005) A polyhedral branch-and-cut approach to global optimization. *Mathematical Programming*



# Constructing effective linear outer approximations

**Example** Outer approximate function

$$f(\mathbf{x}) = |x_1| + |x_2| + |x_3| + |x_4|$$

# Constructing effective linear outer approximations

**Example** Outer approximate function

$$f(\mathbf{x}) = |x_1| + |x_2| + |x_3| + |x_4|$$

**Direct** Add  $2^n = 16$  linear inequalities

$$f(\mathbf{x}) \geq x_1 + x_2 + x_3 + x_4, \quad f(\mathbf{x}) \geq x_1 + x_2 + x_3 - x_4, \quad f(\mathbf{x}) \geq x_1 + x_2 - x_3 + x_4$$

$$f(\mathbf{x}) \geq x_1 + x_2 - x_3 - x_4, \quad f(\mathbf{x}) \geq x_1 - x_2 + x_3 + x_4, \quad f(\mathbf{x}) \geq x_1 - x_2 + x_3 - x_4$$

$\vdots$

# Constructing effective linear outer approximations

**Example** Outer approximate function

$$f(\mathbf{x}) = |x_1| + |x_2| + |x_3| + |x_4|$$

**Direct** Add  $2^n = 16$  linear inequalities

$$f(\mathbf{x}) \geq x_1 + x_2 + x_3 + x_4, \quad f(\mathbf{x}) \geq x_1 + x_2 + x_3 - x_4, \quad f(\mathbf{x}) \geq x_1 + x_2 - x_3 + x_4$$

$$f(\mathbf{x}) \geq x_1 + x_2 - x_3 - x_4, \quad f(\mathbf{x}) \geq x_1 - x_2 + x_3 + x_4, \quad f(\mathbf{x}) \geq x_1 - x_2 + x_3 - x_4$$

⋮

**Extended** Add  $n = 4$  variables and  $2n = 8$  linear inequalities

$$f(\mathbf{x}) \geq \sum_{i=1}^n y_i$$

$$y_i \geq x_i, \quad y_i \geq -x_i \quad i = 1, \dots, 4$$

## Constructing effective linear outer approximations

### Proposition (Tawarmalani and Sahinidis 2005)

*For separable functions, the extended formulation with support points  $\{\bar{x}^j\}_{j=1}^r$  is equivalent to the direct linear outer approximation supported at every  $x$  such that for every  $i \in [n]$  there exists  $j \in [r]$  such that  $x_i = \bar{x}_i^j$ .*

- Polynomial extended formulations  $\Leftrightarrow$  Exponential direct OA
- Linear ineqs in extended space  $\Leftrightarrow$  Nonlinear ineqs in original space

# Outer approximations of nonlinear functions?

## Quadratic functions

$$f(\mathbf{x}) = 5x_1^2 + 4x_2^2 + 9x_3^2 + 4x_1x_2 + 6x_1x_3 + 12x_2x_3$$

# Outer approximations of nonlinear functions?

## Quadratic functions

$$f(\mathbf{x}) = 5x_1^2 + 4x_2^2 + 9x_3^2 + 4x_1x_2 + 6x_1x_3 + 12x_2x_3$$

$$f(\mathbf{x}) = (x_1 + 2x_2 + 3x_3)^2 + 4x_1^2$$

# Outer approximations of nonlinear functions?

## Quadratic functions

$$f(\mathbf{x}) = 5x_1^2 + 4x_2^2 + 9x_3^2 + 4x_1x_2 + 6x_1x_3 + 12x_2x_3$$

$$f(\mathbf{x}) = (x_1 + 2x_2 + 3x_3)^2 + 4x_1^2$$

$$f(\mathbf{x}) = x_4^2 + 4x_1^2 \text{ with } x_4 = x_1 + 2x_2 + 3x_3$$

# Outer approximations of nonlinear functions?

## Quadratic functions

$$f(\mathbf{x}) = 5x_1^2 + 4x_2^2 + 9x_3^2 + 4x_1x_2 + 6x_1x_3 + 12x_2x_3$$

$$f(\mathbf{x}) = (x_1 + 2x_2 + 3x_3)^2 + 4x_1^2$$

$$f(\mathbf{x}) = x_4^2 + 4x_1^2 \text{ with } x_4 = x_1 + 2x_2 + 3x_3$$

Any convex quadratic function of rank  $k$  can be written as a separable function with  $k$  additional variables

→ Cholesky decomposition, eigendecomposition...



# Outer approximations of nonlinear functions?

Conic quadratic functions <sup>4</sup> (Handling the Lorentz cone)

$$x_0 \geq \sqrt{\sum_{i=1}^n x_i^2}$$

---

<sup>4</sup>Vielma JP et al (2017) Extended formulations in mixed-integer conic quadratic programming. *Mathematical Programming Computation*

# Outer approximations of nonlinear functions?

Conic quadratic functions <sup>4</sup> (Handling the Lorentz cone)

$$x_0 \geq \sqrt{\sum_{i=1}^n x_i^2}$$
$$\Leftrightarrow x_0^2 \geq \sum_{i=1}^n x_i^2 \text{ and } x_0 \geq 0$$

---

<sup>4</sup>Vielma JP et al (2017) Extended formulations in mixed-integer conic quadratic programming. *Mathematical Programming Computation*

# Outer approximations of nonlinear functions?

Conic quadratic functions <sup>4</sup> (Handling the Lorentz cone)

$$x_0 \geq \sqrt{\sum_{i=1}^n x_i^2}$$

$$\Leftrightarrow x_0^2 \geq \sum_{i=1}^n x_i^2 \text{ and } x_0 \geq 0$$

$$\Leftrightarrow x_0 \geq \sum_{i=1}^n x_i^2 / x_0 \text{ and } x_0 \geq 0$$

$$\Leftrightarrow x_0 \geq \sum_{i=1}^n y_i \text{ and } x_0 \geq 0, y_i \geq x_i^2 / x_0, \forall i \in [n]$$

---

<sup>4</sup>Vielma JP et al (2017) Extended formulations in mixed-integer conic quadratic programming. *Mathematical Programming Computation*

# Outer approximations of nonlinear functions?

Conic quadratic functions <sup>4</sup> (Handling the Lorentz cone)

$$\begin{aligned}x_0 &\geq \sqrt{\sum_{i=1}^n x_i^2} \\ \Leftrightarrow x_0^2 &\geq \sum_{i=1}^n x_i^2 \text{ and } x_0 \geq 0 \\ \Leftrightarrow x_0 &\geq \sum_{i=1}^n x_i^2 / x_0 \text{ and } x_0 \geq 0 \\ \Leftrightarrow x_0 &\geq \sum_{i=1}^n y_i \text{ and } x_0 \geq 0, y_i \geq x_i^2 / x_0, \forall i \in [n]\end{aligned}$$

20x speedup when first implemented

---

<sup>4</sup>Vielma JP et al (2017) Extended formulations in mixed-integer conic quadratic programming. *Mathematical Programming Computation*

# Summary

- Lack of dual simplex hampers algorithms
- Several approach exist in the literature <sup>5</sup>
- Several popular approaches rely on linear outer approximations
- Effective implementations: integrated with branch-and-bound, calls to interior point method, addition of variables, reformulations...
- In practice, varying degrees of success

---

<sup>5</sup>Kronqvist J et al (2019) A review and comparison of solvers for convex MINLP.  
*Optimization and Engineering*

# Agenda

- 1 Introduction
- 2 Branch and bound
- 3 Convexification**
  - Convexification for MILO
  - Convexification for MINLO in sparse regression

# Agenda

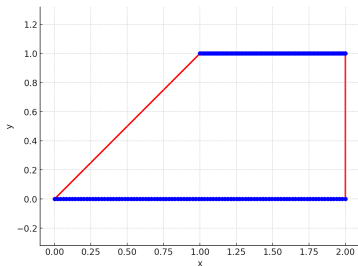
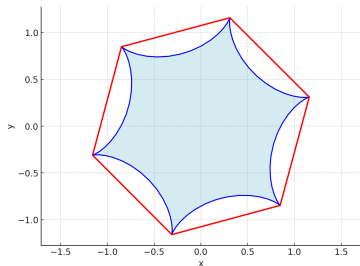
- 1 Introduction
- 2 Branch and bound
- 3 Convexification
  - Convexification for MILO
  - Convexification for MINLO in sparse regression

# Convex hull

## Definition

Given a set  $X \subseteq \mathbb{R}^n$ , the convex hull of  $X$ , denoted as  $\text{conv}(X)$ , is

- The smallest convex set containing  $X$
- The set of all convex combinations of points in  $X$ .





## Convex optimization

Consider the optimization  $\min_{\mathbf{x} \in X} \mathbf{a}^\top \mathbf{x}$

### Proposition

*If set  $X$  is convex, then any local minimum is a global minimum.*

Intuition: Optimization over set  $X$  is “easy” under convexity

## Convex optimization

Consider the optimization 
$$\min_{\mathbf{x} \in X} \mathbf{a}^\top \mathbf{x}$$

### Proposition

*If set  $X$  is convex, then any local minimum is a global minimum.*

Intuition: Optimization over set  $X$  is “easy” under convexity

### Proposition

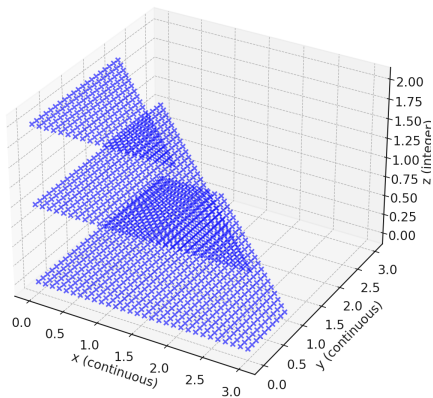
*The optimization problem is equivalent to*

$$\min_{\mathbf{x} \in \text{conv}(X)} \mathbf{a}^\top \mathbf{x},$$

*i.e., there exist a solution that is optimal for both.*

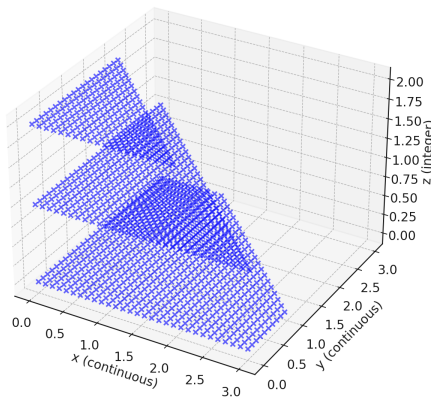
Intuition: Any optimization problem can be reduced to a convex problem

# Convexification in mixed-integer linear optimization



What is the convex hull?

# Convexification in mixed-integer linear optimization



$$x + y + z \leq 4$$

$$0 \leq x \leq 3$$

$$0 \leq y \leq 3$$

$$0 \leq z \leq 2, z \in \mathbb{Z}$$

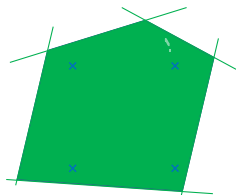
What is the convex hull?

# Convexification in mixed-integer linear optimization

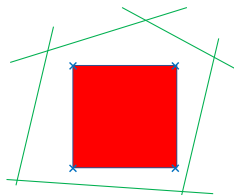
## Proposition (Meyer 1974)

The convex hull of set  $\{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m : \mathbf{Ax} + \mathbf{Gz} \leq \mathbf{b}\}$  is a polyhedron.

Linear relaxation



Convex hull

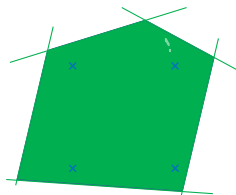


# Convexification in mixed-integer linear optimization

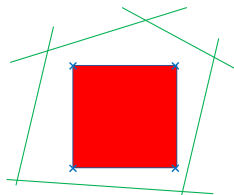
## Proposition (Meyer 1974)

The convex hull of set  $\{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{Z}^m : \mathbf{Ax} + \mathbf{Gz} \leq \mathbf{b}\}$  is a polyhedron.

Linear relaxation

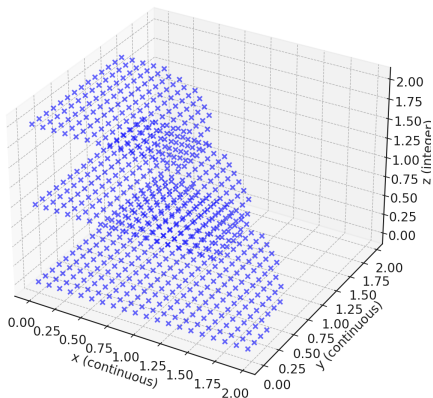


Convex hull



Instead of computing exact convex hulls, convexifications are dynamically added to branch-and-bound algorithms via cutting planes

# Convexification in mixed-integer nonlinear optimization



$$x^2 + y^2 + z \leq 4$$

$$0 \leq x \leq 3$$

$$0 \leq y \leq 3$$

$$0 \leq z \leq 2, z \in \mathbb{Z}$$

What is the convex hull? How to implement in practice?

# Agenda

- 1 Introduction
- 2 Branch and bound
- 3 Convexification
  - Convexification for MILO
  - Convexification for MINLO in sparse regression



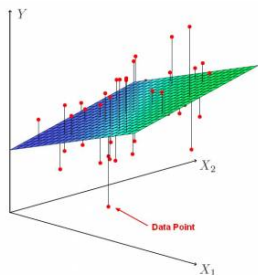
# Least squares regression

Consider dataset<sup>6</sup>  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  where  $\mathbf{x}_i \in \mathbb{R}^n$

## Least squares with ridge regularization

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m (y_i - \mathbf{a}_i^T \mathbf{x})^2 + \lambda \sum_{j=1}^n x_j^2$$

for some  $\lambda \geq 0$



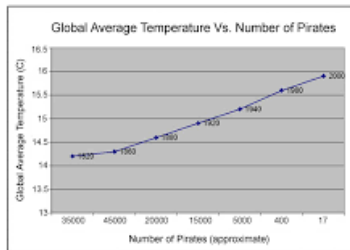
---

<sup>6</sup>Hoerl AE and Kennard RW (1970) Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*

# Shortcomings of ordinary least squares

## Prone to overfitting

### STOP GLOBAL WARMING: BECOME A PIRATE

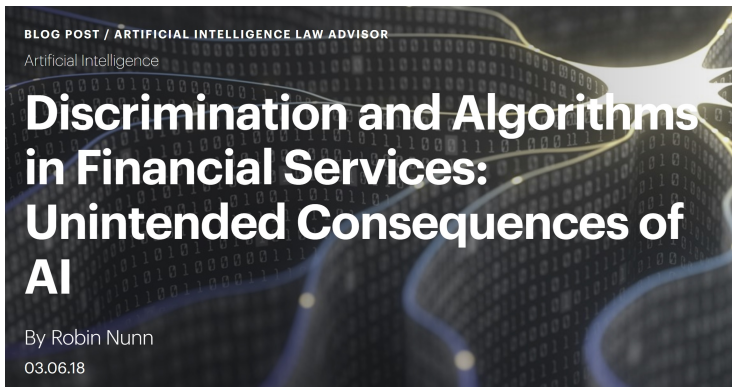


WWW.VENGANZA.ORG



Can fail to make meaningful predictions out-of-sample

## Shortcomings of ordinary least squares



In some cases, interpretability is far more important than accuracy

# Linear regression

Least squares in action with the “Communities and crime” dataset

- Data with socio-economic data, law enforcement data and crime data (US census, LEMAS survey and FBI)
- $n = 100$  features,  $m = 1993$  cities

# Linear regression

Least squares in action with the “Communities and crime” dataset

- Data with socio-economic data, law enforcement data and crime data (US census, LEMAS survey and FBI)
- $n = 100$  features,  $m = 1993$  cities

Solution metrics Optimal solution found in milliseconds,  $R^2 = 0.84$

# Linear regression

Least squares in action with the “Communities and crime” dataset

- Data with socio-economic data, law enforcement data and crime data (US census, LEMAS survey and FBI)
- $n = 100$  features,  $m = 1993$  cities

Solution metrics Optimal solution found in milliseconds,  $R^2 = 0.84$

## Solution

BlspFecCap	0.0260832	FctBchCrimImp	0.0911841
BlspFecCurt	0.209001	FctBchCrimTweq	-0.0212925
LandArea	0.0212184	FctDomeCitySt	0.0346204
LeasFecOfficCntr	0.0260337	FctDomeHouseSt	-0.0138617
MalFecDIrect	0.119487	FctDomeRateSt	0.00120485
MalFecDmWear	0.071392	FctDomeRegIncl	-0.0136794
MedBumSt	0.024658	FctDomePwr	-0.0147045
MedDmCntrFctC	-0.025073	FctDmEmply	-0.0148241
MedDmCntrFctIncBmtg	-0.0740048	FctDmPwrFract	-0.0204822
MedDmCntrFctIncBmtg	0.002048	FctDmRateMm	-0.0420264
MedDmCntrFctIncBmtg	0.025481	FctDmRateMm	0.0537563
MedDmCntrFctIncBmtg	0.015325	FctDmRateMm	-0.0080557
MedDmCntrFctIncBmtg	0.012478	FctDmRateMm	-0.0754108
MedDmCntrFctIncBmtg	-0.05884	FctDmRateMm	-0.0911742
MedDmCntrFctIncBmtg	-0.012478	FctDmRateMm	-0.098852
MedDmCntrFctIncBmtg	0.14224	FctDmRateMm	0.039868
MedDmCntrFctIncBmtg	0.18077	FctDmRateMm	0.087588
MedDmCntrFctIncBmtg	0.02822	FctDmRateMm	-0.0114222
MedDmCntrFctIncBmtg	0.040795	FctDmRateMm	0.0131009
MedDmCntrFctIncBmtg	0.025084	FctDmRateMm	0.018131
MedDmCntrFctIncBmtg	-0.085235	FctDmRateMm	0.0202614
MedDmCntrFctIncBmtg	0.002048	FctDmRateMm	-0.12062
MedDmCntrFctIncBmtg	0.023084	FctDmRateMm	0.0208774
MedDmCntrFctIncBmtg	-0.029284	FctDmRateMm	0.022203
MedDmCntrFctIncBmtg	-0.0080359	FctDmRateMm	-0.0442128
MedDmCntrFctIncBmtg	0.0736016	FctDmRateMm	-0.0444774
MedDmCntrFctIncBmtg	-0.0736016	FctDmRateMm	-0.0332621
MedDmCntrFctIncBmtg	0.001135	FctDmRateMm	0.0782709
MedDmCntrFctIncBmtg	0.072027	FctDmRateMm	-0.0148107
MedDmCntrFctIncBmtg	0.0242781	FctDmRateMm	0.0176021
MedDmCntrFctIncBmtg	-0.092774	FctDmRateMm	-0.0244771
MedDmCntrFctIncBmtg	0.0080359	FctDmRateMm	0.0114271
MedDmCntrFctIncBmtg	0.186242	FctDmRateMm	-0.030747
MedDmCntrFctIncBmtg	-0.015325	FctDmRateMm	-0.032702
MedDmCntrFctIncBmtg	-0.017715	FctDmRateMm	0.039124

# Use parsimony

Occam's razor / Principle of parsimony (William of Ockham  $\approx$  1300)



Why did the tree fall?

- The wind
- Two meteorites crashed into earth. One hit the tree, the other hit the first meteorite, obliterating both and destroying the evidence

# Use parsimony

Occam's razor / Principle of parsimony (William of Ockham  $\approx$  1300)



Why did the tree fall?

- The wind
- Two meteorites crashed into earth. One hit the tree, the other hit the first meteorite, obliterating both and destroying the evidence

Given two competing explanations, the simplest one is often right.



# Use parsimony

## Best subset selection

- Let  $k$  be the target complexity of the model. Among all  $\binom{n}{k}$  subsets of  $k$  features, find the one that best fits the model

---

<sup>7</sup>Furnival G and Wilson R (1974) Regressions by leaps and bounds. *Technometrics*

# Use parsimony

## Best subset selection

- Let  $k$  be the target complexity of the model. Among all  $\binom{n}{k}$  subsets of  $k$  features, find the one that best fits the model
- Solve

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n x_j^2 \\ \text{s.t.} \quad & \sum_{j=1}^n \mathbb{1}_{\{x_j \neq 0\}} \leq k \end{aligned}$$

- Implemented<sup>7</sup> in R packages for  $n < 30$

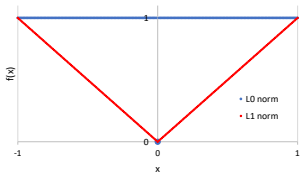
---

<sup>7</sup>Furnival G and Wilson R (1974) Regressions by leaps and bounds. *Technometrics*

# Relaxations

Lasso/ elastic net (Tibshirani 1996, Zou and Hastie 2005)

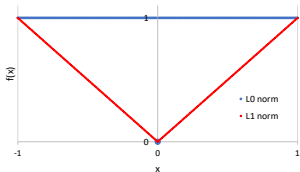
- The best convex underestimator of the “ $\ell_0$ -norm” function  $f(x) = \mathbb{1}_{\{x \neq 0\}}$  on  $-1 \leq x \leq 1$  is the  $\ell_1$ -norm  $|x|$



# Relaxations

Lasso/ elastic net (Tibshirani 1996, Zou and Hastie 2005)

- The best convex underestimator of the “ $\ell_0$ -norm” function  $f(x) = \mathbb{1}_{\{x \neq 0\}}$  on  $-1 \leq x \leq 1$  is the  $\ell_1$ -norm  $|x|$



- Solve

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m (y_i - \mathbf{a}_i^T \mathbf{x})^2 + \lambda \sum_{j=1}^n x_j^2$$

$$\text{s.t. } \sum_{j=1}^n |x_j| \leq \kappa$$

where  $\kappa$  is a parameter (to be tuned) controlling sparsity vs accuracy

# Relaxations



## Robert Tibshirani

[FOLLOW](#)

Professor of Biomedical Data Sciences, and of Statistics, [Stanford University](#)

Verified email at stanford.edu - [Homepage](#)

[Statistics](#) [Applied Statistics](#) [Statistical learning](#) [machine learning](#) [data science](#)

TITLE

CITED BY YEAR

[Unsupervised learning](#)

T Hastie, R Tibshirani, J Friedman  
The elements of statistical learning, 485-585

42343 2009

[An introduction to the bootstrap](#)

B Efron, RJ Tibshirani  
CRC press

39319 1994

[Regression shrinkage and selection via the lasso](#)

R Tibshirani  
Journal of the Royal Statistical Society. Series B (Methodological), 267-288

26819 1996

[Generalized additive models](#)

TJ Hastie  
Statistical models in S, 249-307

46474 2017

[Generalized Additive Models](#)

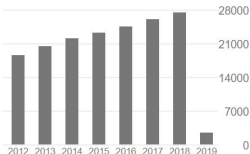
TJ Hastie, RJ Tibshirani  
CRC Press

16171 <sup>\*</sup> 1990

Cited by

[VIEW ALL](#)

	All	Since 2014
Citations	289715	126679
h-index	143	96
i10-index	382	301



Co-authors

[VIEW ALL](#)

- Trevor Hastie  
Professor of Statistics, Stanford ...
>
- B Efron  
Professor of statistics, Stanford ...
>
- Jerome Friedman
>

# Relaxations




## Albert Einstein

[FOLLOW](#)

Institute of Advanced Studies, Princeton

No verified email

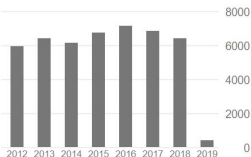
Physics

Cited by

[VIEW ALL](#)

	All	Since 2014
Citations	124162	33880
h-index	113	64
i10-index	370	206

TITLE	CITED BY	YEAR
Can quantum-mechanical description of physical reality be considered complete? A Einstein, B Podolsky, N Rosen Physical review 47 (10), 777	17493	1935
Über einen die Erzeugung und Verwandlung des Lichtes betreffenden heuristischen Gesichtspunkt A Einstein Ann. Phys. 17, 132-148	11362*	1905
On the movement of small particles suspended in stationary liquids required by the molecular-kinetic theory of heat A Einstein Annalen der Physik 17, 549-560	9690*	1905
Zur Elektrodynamik bewegter Körper A Einstein	5504*	
Sitzungsber. K A Einstein Preuss. Akad. Wiss., Phys. Math. Kl 3, 18	4950*	1925
Graviton Mass and Inertia Mass A Einstein Ann Physik 35, 898	4814*	1911
	-	



# Relaxations

Lasso in action with the “Communities and crime” dataset ( $n = 100$ )

Solution metrics Optimal solution found in milliseconds

## Solutions

Large  $\kappa$  ( $R^2 = 0.81$ )

HousVacant	0.237656
LemasPctOfficDrugUn	0.000206104
MalePctDivorce	0.0260223
NumStreet	0.14165
PctHousNoPhone	0.0266273
PctIlleg	0.311418
PctPersDenseHous	0.197454
PctVacantBoarded	0.0405917
PopDens	0.0193849
pctWPubAsst	0.0445904
racepctblack	0.186306

Small  $\kappa$  ( $R^2 = 0.25$ )

LandArea	0.121248
NumIlleg	0.685344
NumImmig	0.183812
PctPersDenseHous	0.000258902

# Mixed-integer optimization

## Mixed-integer optimization <sup>89</sup>

- Best subset selection can be formulated as a MIO
- Letting binary variable  $z_j = 1$  iff feature  $j$  is included, solve

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n x_j^2 \\ \text{s.t.} \quad & \sum_{j=1}^n z_j \leq k \\ & -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n \end{aligned}$$

<sup>8</sup>Bertsimas D et al (2016) Best subset selection via a modern optimization lens. *The Annals of Statistics*

<sup>9</sup>Cozad A et al (2014) Learning surrogate models for simulation-based optimization. *AICHE*.



# Mixed-integer optimization

MIO in action with the “Communities and crime” dataset ( $n = 100$ )

Solution ( $R^2 = 0.81$ )

HousVacant	0.250896
MalePctDivorce	0.135992
PctIlleg	0.524062
PctPersDenseHous	0.175159

# Mixed-integer optimization

MIO in action with the “Communities and crime” dataset ( $n = 100$ )

Solution ( $R^2 = 0.81$ )

HousVacant	0.250896
MalePctDivorce	0.135992
PctIlleg	0.524062
PctPersDenseHous	0.175159

Solution metrics 20 hours to optimality, millions of nodes (Gurobi, 2022)

# Mixed-integer optimization



Is best subset selection really worth it?<sup>10</sup>

- Best subset is slow
- Lasso is better in some situations, and can be

improved otherwise



Of course!<sup>11</sup>?

- Solution times are appropriate in many cases
- Lasso is better in very low SNR regimes, and best

subset can be adapted



Both methods have merits!<sup>12</sup>?

<sup>10</sup> Hastie T, Tibshirani R, Tibshirani R (2020) Best subset, forward stepwise or Lasso? Analysis and recommendations based on extensive comparisons. *Statistical Science*

<sup>11</sup> Mazumder R, Radchenko P, Dedieu A (2023) Subset selection with shrinkage: Sparse linear modeling when the SNR is low. *Operations Research*

<sup>12</sup> Chen Y, Taeb A, Bühlmann P (2020) A look at robustness and stability of  $\ell_1$ - versus  $\ell_0$ -regularization: Discussion of papers by Bertsimas et al. and Hastie et al. *Statistical Science*

# Improving the formulation

How good is the convex relaxation?  $\mathbf{z} \in \{0, 1\}^n \rightarrow \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n x_j^2 \\ \text{s.t.} \quad & \sum_{j=1}^n z_j \leq k \\ & -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n \end{aligned}$$

## Improving the formulation

How good is the convex relaxation?  $\mathbf{z} \in \{0, 1\}^n \rightarrow \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n x_j^2 \\ \text{s.t.} \quad & \sum_{j=1}^n z_j \leq k \\ & -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n \end{aligned}$$

In an optimal solution,  $z_j^* = |x_j|/M$

$\Rightarrow$  The continuous relaxation is in fact lasso!

## Improving the formulation

How good is the convex relaxation?  $\mathbf{z} \in \{0, 1\}^n \rightarrow \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{0} \leq \mathbf{z} \leq \mathbf{1}} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^T \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n x_j^2 \\ \text{s.t.} \quad & \sum_{j=1}^n z_j \leq k \\ & -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n \end{aligned}$$

In an optimal solution,  $z_j^* = |x_j|/M$

$\Rightarrow$  The continuous relaxation is in fact lasso!

Since lasso is not a good approximation, this formulation is slow...

# Improving the formulation

Improve the convex relaxation    Need to exploit nonlinearities

$$\begin{aligned}
 \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n, \mathbf{t} \in \mathbb{R}_+^n} & \quad \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n t_j \\
 \text{s.t.} & \quad x_j^2 \leq t_j \quad \forall j = 1, \dots, n \\
 & \quad \sum_{j=1}^n z_j \leq k \\
 & \quad -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n
 \end{aligned}$$

# Improving the formulation

Improve the convex relaxation    Need to exploit nonlinearities

$$\begin{aligned}
 \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n, \mathbf{t} \in \mathbb{R}_+^n} & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n t_j \\
 \text{s.t.} & \quad x_j^2 \leq t_j \quad \forall j = 1, \dots, n \\
 & \quad \sum_{j=1}^n z_j \leq k \\
 & \quad -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n
 \end{aligned}$$

What is the convex hull of

$$S = \{x \in \mathbb{R}, z \in \{0,1\}, t \in \mathbb{R} : x^2 \leq t, x(1-z) = 0\}$$



## Improving the formulation

$$S = \underbrace{\{(x, z, t) \in \mathbb{R}^3 : 0 \leq t, x = z = 0\}}_{S_1} \cup \underbrace{\{(x, z, t) \in \mathbb{R}^3 : x^2 \leq t, z = 1\}}_{S_2}?$$

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists (x_i, z_i, t_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \lambda_1 x_1 + \lambda_2 x_2, \quad z = \lambda_1 z_1 + \lambda_2 z_2, \quad t = \lambda_1 t_1 + \lambda_2 t_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$(x_1, z_1, t_1) \in S_1 \Leftrightarrow x_1 = z_1 = 0, t_1 \geq 0$$

$$(x_2, z_2, t_2) \in S_2 \Leftrightarrow x_2^2 \leq t_2, z_2 = 1$$

## Improving the formulation

$$S = \underbrace{\{(x, z, t) \in \mathbb{R}^3 : 0 \leq t, x = z = 0\}}_{S_1} \cup \underbrace{\{(x, z, t) \in \mathbb{R}^3 : x^2 \leq t, z = 1\}}_{S_2}?$$

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists (x_i, z_i, t_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \lambda_1 x_1 + \lambda_2 x_2, \quad z = \lambda_1 z_1 + \lambda_2 z_2, \quad t = \lambda_1 t_1 + \lambda_2 t_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$(x_1, z_1, t_1) \in S_1 \Leftrightarrow x_1 = z_1 = 0, t_1 \geq 0$$

$$(x_2, z_2, t_2) \in S_2 \Leftrightarrow x_2^2 \leq t_2, z_2 = 1$$

Change of variables:  $\tilde{x}_i = x_i \lambda_i, \tilde{z}_i = z_i \lambda_i, \tilde{t}_i = t_i \lambda_i$

## Improving the formulation

$$S = \underbrace{\{(x, z, t) \in \mathbb{R}^3 : 0 \leq t, x = z = 0\}}_{S_1} \cup \underbrace{\{(x, z, t) \in \mathbb{R}^3 : x^2 \leq t, z = 1\}}_{S_2}?$$

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists (x_i, z_i, t_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_1 + \tilde{x}_2, \quad z = \tilde{z}_1 + \tilde{z}_2, \quad t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$(x_1, z_1, t_1) \in S_1 \Leftrightarrow \tilde{x}_1 = \tilde{z}_1 = 0, \quad \tilde{t}_1 \geq 0$$

$$(x_2, z_2, t_2) \in S_2 \Leftrightarrow (\tilde{x}_2/\lambda_2)^2 \leq \tilde{t}_2/\lambda_2, \quad \tilde{z}_2/\lambda_2 = 1$$

Change of variables:  $\tilde{x}_i = x_i \lambda_i$ ,  $\tilde{z}_i = z_i \lambda_i$ ,  $\tilde{t}_i = t_i \lambda_i$

## Improving the formulation

$$S = \underbrace{\{(x, z, t) \in \mathbb{R}^3 : 0 \leq t, x = z = 0\}}_{S_1} \cup \underbrace{\{(x, z, t) \in \mathbb{R}^3 : x^2 \leq t, z = 1\}}_{S_2}?$$

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists (x_i, z_i, t_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_1 + \tilde{x}_2, \quad z = \tilde{z}_1 + \tilde{z}_2, \quad t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$(\tilde{x}_1, \tilde{z}_1, \tilde{t}_1) \in \lambda_1 S_1 \Leftrightarrow \tilde{x}_1 = \tilde{z}_1 = 0, \tilde{t}_1 \geq 0$$

$$(\tilde{x}_2, \tilde{z}_2, \tilde{t}_2) \in \lambda_2 S_2 \Leftrightarrow \tilde{x}_2^2 / \lambda_2 \leq \tilde{t}_2, \tilde{z}_2 = \lambda_2$$

Change of variables:  $\tilde{x}_i = x_i \lambda_i, \tilde{z}_i = z_i \lambda_i, \tilde{t}_i = t_i \lambda_i$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_1 + \tilde{x}_2, \quad z = \tilde{z}_1 + \tilde{z}_2, \quad t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$\tilde{x}_1 = \tilde{z}_1 = 0, \quad \tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2 / \lambda_2 \leq \tilde{t}_2, \quad \tilde{z}_2 = \lambda_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_1 + \tilde{x}_2, \quad z = \tilde{z}_1 + \tilde{z}_2, \quad t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \quad \lambda_1 \geq 0, \quad \lambda_2 \geq 0$$

$$\tilde{x}_1 = \tilde{z}_1 = 0, \quad \tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2 / \lambda_2 \leq \tilde{t}_2, \quad \tilde{z}_2 = \lambda_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_2, z = \tilde{z}_2, t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \lambda_1 \geq 0, \lambda_2 \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2/\lambda_2 \leq \tilde{t}_2, \tilde{z}_2 = \lambda_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_2, z = \tilde{z}_2, t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \lambda_1 \geq 0, \lambda_2 \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2/\lambda_2 \leq \tilde{t}_2, \tilde{z}_2 = \lambda_2$$



## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_2, z = \lambda_2, t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \lambda_1 \geq 0, \lambda_2 \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2/\lambda_2 \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$x = \tilde{x}_2, z = \lambda_2, t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + \lambda_2 = 1, \lambda_1 \geq 0, \lambda_2 \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$\tilde{x}_2^2/\lambda_2 \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + z_2 = 1, \lambda_1 \geq 0, z \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$x^2/z \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$t = \tilde{t}_1 + \tilde{t}_2$$

$$\lambda_1 + z_2 = 1, \lambda_1 \geq 0, z \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$x^2/z \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$t = \tilde{t}_1 + \tilde{t}_2$$

$$z \leq 1, z \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$x^2/z \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$t = \tilde{t}_1 + \tilde{t}_2$$

$$z \leq 1, z \geq 0$$

$$\tilde{t}_1 \geq 0$$

$$x^2/z \leq \tilde{t}_2$$

## Improving the formulation

$(x, z, t) \in \text{conv}(S)$  if and only<sup>13</sup>  $\exists(\tilde{x}_i, \tilde{z}_i, \tilde{t}_i) \in S_i$  and  $\lambda_i \in \mathbb{R}^i$  such that

$$t \geq x^2/z, 0 \leq z \leq 1$$

### Proposition (Frangioni and Gentile 2006)

*The convex hull of set*

$$S = \{x \in \mathbb{R}, z \in \{0, 1\}, t \in \mathbb{R} : x^2 \leq t, x(1 - z) = 0\}$$

*is*

$$\text{conv}(S) = \{(x, z, t) \in \mathbb{R}^3 : x^2 \leq tz, 0 \leq z \leq 1\}$$

---

<sup>13</sup>Frangioni A and Gentile C (2006) Perspective cuts for a class of convex 0-1 mixed-integer programs. *Mathematical Programming*

# Improving the formulation

Improve the convex relaxation <sup>1415</sup>

$$\begin{aligned}
 \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n, \mathbf{t} \in \mathbb{R}_+^n} \quad & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n t_j \\
 \text{s.t.} \quad & x_j^2 \leq t_j \quad \forall j = 1, \dots, n \\
 & \sum_{j=1}^n z_j \leq k \\
 & -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n
 \end{aligned}$$

<sup>14</sup>Dong H et al (2018) Regularization vs relaxation: A convexification perspective of statistical variable selection. *Optimization Online*

<sup>15</sup>Xie W and Deng X (2020) Scalable algorithms for the sparse ridge regression. *SIAM Journal on Optimization*



# Improving the formulation

Improve the convex relaxation <sup>1415</sup>

$$\begin{aligned}
 \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n, \mathbf{t} \in \mathbb{R}_+^n} & \sum_{i=1}^m \left( y_i - \mathbf{a}_i^\top \mathbf{x} \right)^2 + \lambda \sum_{j=1}^n t_j \\
 \text{s.t.} & \quad x_j^2 \leq t_j z_j \quad \forall j = 1, \dots, n \\
 & \quad \sum_{j=1}^n z_j \leq k \\
 & \quad -Mz_j \leq x_j \leq Mz_j \quad \forall j = 1, \dots, n
 \end{aligned}$$

The perspective reformulation! ([Disjunctive programming](#))

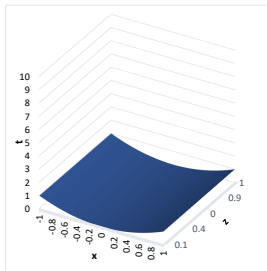
<sup>14</sup>Dong H et al (2018) Regularization vs relaxation: A convexification perspective of statistical variable selection. *Optimization Online*

<sup>15</sup>Xie W and Deng X (2020) Scalable algorithms for the sparse ridge regression. *SIAM Journal on Optimization*

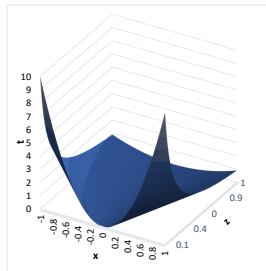
## Improving the formulation

Constraint  $tz \geq x^2$  with  $t, z \geq 0$  is convex and SOCP representable

It represents a substantial improvement in the relaxation quality



Graph of  $t = x^2$  (big-M)

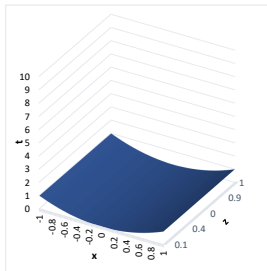


Graph of  $t = x^2/z$

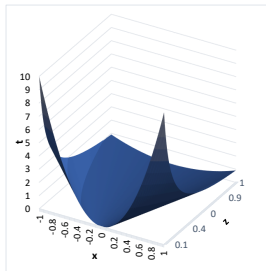
## Improving the formulation

Constraint  $tz \geq x^2$  with  $t, z \geq 0$  is convex and SOCP representable

It represents a substantial improvement in the relaxation quality



Graph of  $t = x^2$  (big-M)



Graph of  $t = x^2/z$

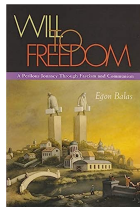
Solution times in “Communities and crime”: 2s (**15,000 $\times$  speedup**)

# Disjunctive programming

Disjunctive programming was invented by Egon Balas in the 80s

<https://www.wsj.com/articles/>

egon-balas-jailed-and-tortured-in-romania-found-salvation-in-math-11553869800



Generalizes to nonlinear optimization<sup>16</sup>

---

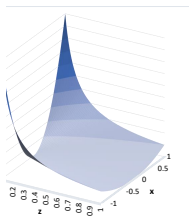
<sup>16</sup>Ceria S and Soares J (1999) Convex programming for disjunctive convex optimization. *Mathematical Programming*

# Disjunctive Programming

Given a convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , consider

$$f^\pi(\mathbf{x}, \lambda) = \begin{cases} \lambda f(\mathbf{x}/\lambda) & \text{if } \lambda > 0 \\ \lim_{\lambda \rightarrow 0^+} \lambda f(\mathbf{x}/\lambda) & \text{if } \lambda = 0 \\ +\infty & \text{otherwise.} \end{cases}$$

- $f^\pi$  is convex and homogeneous
- If  $f(\mathbf{x}) = a_0 + \mathbf{a}^\top \mathbf{x}$ , then  $f^\pi(\mathbf{x}, \lambda) = a_0 \lambda + \mathbf{a}^\top \mathbf{x}$
- If  $f(x) = x^2$ , then  $f^\pi(x, z) = x^2/z$  with  $0/0 = 0$  and  $x^2/0 = +\infty$  if  $x \neq 0$



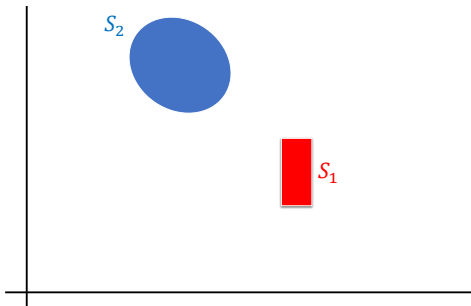
# Disjunctive programming

For  $i \in \{1, \dots, k\}$ , let  $S_i = \{x \in \mathbb{R}^n : g_{ij}(x) \leq 0, j = 1 \dots, m\}$  convex.

Then  $x \in \text{cl conv} \left( \bigcup_{i=1}^k S_i \right)$  iff  $\exists x^i \in \mathbb{R}^n$  and  $\lambda^i \in \mathbb{R}_+$  such that

$$x^i \in S_i^\pi(\lambda) = \{x \in \mathbb{R}^n : g_{ij}^\pi(x, \lambda^i) \leq 0, j = 1 \dots, m\}$$

$$x = \sum_{i=1}^k x^i, \text{ and } 1 = \sum_{i=1}^k \lambda^i.$$



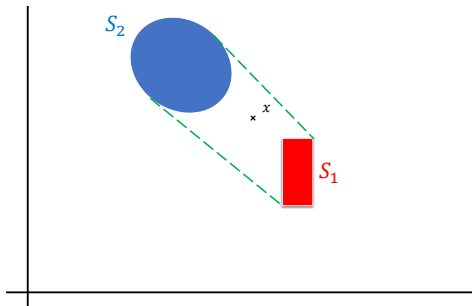
# Disjunctive programming

For  $i \in \{1, \dots, k\}$ , let  $S_i = \{x \in \mathbb{R}^n : g_{ij}(x) \leq 0, j = 1 \dots, m\}$  convex.

Then  $x \in \text{cl conv} \left( \bigcup_{i=1}^k S_i \right)$  iff  $\exists x^i \in \mathbb{R}^n$  and  $\lambda^i \in \mathbb{R}_+$  such that

$$x^i \in S_i^\pi(\lambda) = \{x \in \mathbb{R}^n : g_{ij}^\pi(x, \lambda^i) \leq 0, j = 1 \dots, m\}$$

$$x = \sum_{i=1}^k x^i, \text{ and } 1 = \sum_{i=1}^k \lambda^i.$$



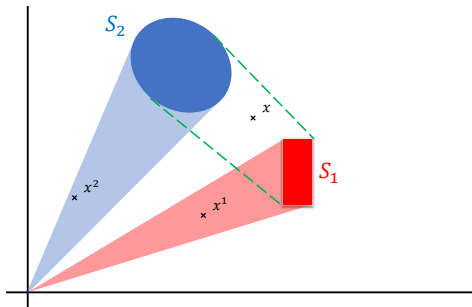
# Disjunctive programming

For  $i \in \{1, \dots, k\}$ , let  $S_i = \{x \in \mathbb{R}^n : g_{ij}(x) \leq 0, j = 1 \dots, m\}$  convex.

Then  $x \in \text{cl conv} \left( \bigcup_{i=1}^k S_i \right)$  iff  $\exists x^i \in \mathbb{R}^n$  and  $\lambda^i \in \mathbb{R}_+$  such that

$$x^i \in S_i^\pi(\lambda) = \{x \in \mathbb{R}^n : g_{ij}^\pi(x, \lambda^i) \leq 0, j = 1 \dots, m\}$$

$$x = \sum_{i=1}^k x^i, \text{ and } 1 = \sum_{i=1}^k \lambda^i.$$





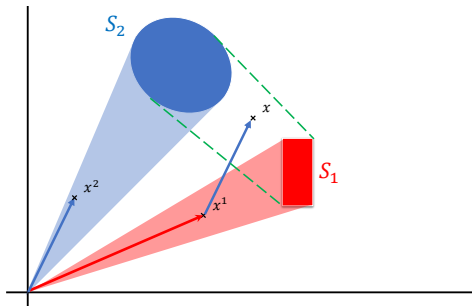
# Disjunctive programming

For  $i \in \{1, \dots, k\}$ , let  $S_i = \{x \in \mathbb{R}^n : g_{ij}(x) \leq 0, j = 1 \dots, m\}$  convex.

Then  $x \in \text{cl conv} \left( \bigcup_{i=1}^k S_i \right)$  iff  $\exists x^i \in \mathbb{R}^n$  and  $\lambda^i \in \mathbb{R}_+$  such that

$$x^i \in S_i^\pi(\lambda) = \{x \in \mathbb{R}^n : g_{ij}^\pi(x, \lambda^i) \leq 0, j = 1 \dots, m\}$$

$$x = \sum_{i=1}^k x^i, \text{ and } 1 = \sum_{i=1}^k \lambda^i.$$



## Disjunctive programming

**Implications of disjunctive programming** Given any disjunctive set, we can create an equivalent convex (conic-representable) representation by creating  $k$  copies  $x^i$  of variables  $x$ , and  $mk$  constraints  $g_{ij}^{\pi}(x^i, \lambda^i) \leq 0$ . In other words, formulation increases by a factor of  $k$ .

Is it useful?

## Disjunctive programming

**Implications of disjunctive programming** Given any disjunctive set, we can create an equivalent convex (conic-representable) representation by creating  $k$  copies  $x^i$  of variables  $x$ , and  $mk$  constraints  $g_{ij}^{\pi}(x^i, \lambda^i) \leq 0$ . In other words, formulation increases by a factor of  $k$ .

**Is it useful?** If used carefully

- Number of additional variables can grow exponentially
- Fourier–Motzkin elimination can be difficult in closed form
- Cuts from disjunctive programming may be hard to implement

# Implementation of disjunctive programming

$$S = \bigcup_{i=1}^{\ell} \{\mathbf{x} \in \mathbb{R}^n : f_i(\mathbf{x}) \leq 0\}$$

**Implementation 1** Add variables  $\{(\mathbf{x}^i, \lambda_i) \in \mathbb{R}^{n+1}\}_{i=1}^{\ell}$  and formulate as

$$\begin{aligned} \mathbf{x} &= \sum_{i=1}^{\ell} \mathbf{x}^i, \quad 1 = \sum_{i=1}^{\ell} \lambda_i, \quad \boldsymbol{\lambda} \geq \mathbf{0} \\ f_i^{\pi}(\mathbf{x}^i, \lambda_i) &\leq 0 \quad \forall i \in \{1, \dots, \ell\} \end{aligned}$$

- Can be effective when  $\ell$  is small (e.g.,  $\ell = 2$ )
- But number of variables and constraints may be prohibitive...

# Implementation of disjunctive programming

$$S = \bigcup_{i=1}^{\ell} \{\mathbf{x} \in \mathbb{R}^n : f_i(\mathbf{x}) \leq 0\}$$

Implementation 2 Add linear cuts using representation

$$0 \geq \max_{\alpha \in \mathbb{R}^n, \beta \in \mathbb{R}, \gamma \in \mathbb{R}_+^{\ell}} \alpha^{\top} \mathbf{x} + \beta$$

$$+ \min_{\{(\mathbf{x}^i, \lambda_i)\}} \left\{ - \sum_{i=1}^{\ell} \alpha^{\top} \mathbf{x}^i - \sum_{i=1}^{\ell} \beta \lambda_i + \sum_{i=1}^{\ell} \gamma_i f_i^{\pi}(\mathbf{x}^i, \lambda_i) \right\}$$

- Given fixed  $\mathbf{x}$ , solve separation (max) and add linear cut
- Requires computing Fenchel conjugates (min)
- But linear cuts can be ineffective...

# Implementation of disjunctive programming

$$S = \bigcup_{i=1}^{\ell} \{\mathbf{x} \in \mathbb{R}^n : f_i(\mathbf{x}) \leq 0\}$$

**Implementation 3** Add nonlinear cuts (e.g., Fourier-Motzkin with duality)

- May achieve a good compromise between convex hull and linear cuts
- But adding nonlinear cuts in branch-and-bound is not easy
  - Not supported in OA branch-and-bound solvers
  - Could require column generation to implement effectively
  - **May be of different classes than original function**

## Low-rank functions

What is the convex hull of

$$S = \left\{ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : \left( \mathbf{a}^\top \mathbf{x} \right)^2 \leq t, \mathbf{x} \circ (\mathbf{1} - \mathbf{z}) \right\}$$

where “ $\circ$ ” is the entrywise product, i.e.,  $\mathbf{x} \circ (\mathbf{1} - \mathbf{z}) \Leftrightarrow x_i(1 - z_i) = 0$

## Low-rank functions

What is the convex hull of

$$S = \left\{ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : (\mathbf{a}^\top \mathbf{x})^2 \leq t, \mathbf{x} \circ (\mathbf{1} - \mathbf{z}) \right\}$$

where “ $\circ$ ” is the entrywise product, i.e.,  $\mathbf{x} \circ (\mathbf{1} - \mathbf{z}) \Leftrightarrow x_i(1 - z_i) = 0$

Disjunctive programming

$$S = \bigcup_{\bar{\mathbf{z}} \in \{0, 1\}^n} \left\{ (\mathbf{x}, \mathbf{z}, t) \in \mathbb{R}^{2n+1} : (\mathbf{a}^\top \mathbf{x})^2 \leq t, \mathbf{x} \circ (\mathbf{1} - \bar{\mathbf{z}}) = \mathbf{0} \right\}$$

$\Rightarrow$  Exponential number of variables/constraints



## Low-rank functions

Consider optimization over set  $S$ ,

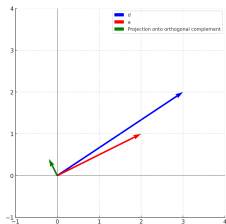
$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n} \mathbf{c}^\top \mathbf{z} + \mathbf{d}^\top \mathbf{x} + (\mathbf{a}^\top \mathbf{x})^2 \quad \text{s.t. } \mathbf{x} \circ (\mathbf{1} - \mathbf{z})$$

# Low-rank functions

Consider optimization over set  $S$ ,

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n} \mathbf{c}^\top \mathbf{z} + \mathbf{d}^\top \mathbf{x} + \left( \mathbf{a}^\top \mathbf{x} \right)^2 \quad \text{s.t. } \mathbf{x} \circ (\mathbf{1} - \mathbf{z})$$

$\mathbf{d} \neq \mu \mathbf{a}$  for some  $\mu \in \mathbb{R} \implies \exists \mathbf{h} \in \mathbb{R}^n$  such that  $\mathbf{h}^\top \mathbf{a} = 0$  and  $\mathbf{h}^\top \mathbf{d} < 0$   
 $\implies$  Unbounded, letting  $\mathbf{z} = \mathbf{1}$  and  $\mathbf{x} = \gamma \mathbf{h}$  with  $\gamma \rightarrow \infty$



## Low-rank functions

Consider optimization over set  $S$ ,

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n} \mathbf{c}^\top \mathbf{z} + \mathbf{d}^\top \mathbf{x} + (\mathbf{a}^\top \mathbf{x})^2 \quad \text{s.t. } \mathbf{x} \circ (\mathbf{1} - \mathbf{z})$$

If  $\mathbf{d} = \mu \mathbf{a}$ , then optimization

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0,1\}^n, y \in \mathbb{R}} \mathbf{c}^\top \mathbf{z} + \mu y + y^2 \quad \text{s.t. } y = \mathbf{a}^\top \mathbf{x}, \mathbf{x} \circ (\mathbf{1} - \mathbf{z})$$

has an optimal solution with at most one non-zero  $x_i$

$\implies$  All extreme points of  $\text{conv}(S)$  have at most one non-zero  $x_i$

## Low-rank functions

Given any convex function<sup>17</sup>  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  and  $\mathbf{A} \in \mathbb{R}^{k \times n}$ , define

$$S = \{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : t \geq f(\mathbf{A}\mathbf{x}), \mathbf{x} \circ (\mathbf{1} - \mathbf{z}) = \mathbf{0}\}$$

Proposition (Han and Gómez 2024)

$$cl \ conv(S) = cl \ conv \left( \bigcup_{\mathcal{I} \subseteq [n]: |\mathcal{I}| \leq k} V(\mathcal{I}) \cup R \right)$$

where

$$V(\mathcal{I}) = \{\{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : t \geq f(\mathbf{A}\mathbf{x}), x_i = 0 \ \forall i \notin \mathcal{I}, z_i = 1 \ \forall i \in \mathcal{I}\}\}$$

$$R = \{\{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : t \geq 0, \mathbf{A}\mathbf{x} = \mathbf{0}, \mathbf{z} = \mathbf{1}\}\}$$

<sup>17</sup>Han S and Gómez A (2024) Compact extended formulations for low rank functions with indicators. *Mathematics of Operations Research*

# Low-rank functions

Computing

$$\text{cl conv} \left( \bigcup_{\mathcal{I} \subseteq [n]: |\mathcal{I}| \leq k} V(\mathcal{I}) \cup R \right) \text{ vs. } \text{cl conv} \left( \bigcup_{\mathcal{I} \subseteq [n]} V(\mathcal{I}) \right)$$

involves  $\mathcal{O}(n^k)$  vs  $2^n$  disjunctions

## Low-rank functions

Computing

$$\text{cl conv} \left( \bigcup_{\mathcal{I} \subseteq [n]: |\mathcal{I}| \leq k} V(\mathcal{I}) \cup R \right) \text{ vs. } \text{cl conv} \left( \bigcup_{\mathcal{I} \subseteq [n]} V(\mathcal{I}) \right)$$

involves  $\mathcal{O}(n^k)$  vs  $2^n$  disjunctions

For special case of  $k = 1$ ,

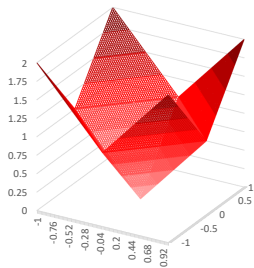
$$S = \left\{ \mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^n, t \in \mathbb{R} : (\mathbf{a}^\top \mathbf{x})^2 \leq t, \mathbf{x} \circ (\mathbf{1} - \mathbf{z}) = \mathbf{0} \right\},$$

we find after Fourier-Motzkin elimination that

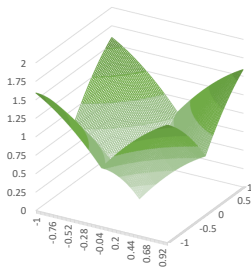
$$\text{cl conv}(S) = \left\{ (\mathbf{x}, \mathbf{z}, t) \in \mathbb{R}^{2n+1} : (\mathbf{a}^\top \mathbf{x})^2 / \min\{1, \mathbf{1}^\top \mathbf{z}\} \leq t, \mathbf{0} \leq \mathbf{z} \leq \mathbf{1} \right\}$$

# Rank-one convexification in sparse regression

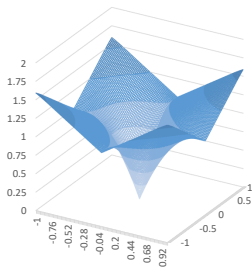
Can be interpreted as strong regularization<sup>18</sup>



Lasso as regularization



Perspective as regularization



Rank-one as regularization

- In tall instances ( $n \ll m$ ), solution from relaxation is integral in practice
- But relaxation is more sophisticated (SOCP  $\rightarrow$  SDP)

<sup>18</sup>Atamtürk A and Gómez A (2025) Rank-one convexifications for sparse regression. *Journal of Machine Learning Research*.

# Full implementation

Tailored branch-and-bound algorithm based on perspective relaxation<sup>19</sup>

- Project out unnecessary variables
- Coordinate descent to solve relaxations
- Active sets
- Dual bounds
- Strong branching

---

<sup>19</sup>Hazimeh H et al (2022) Sparse regression at scale: Branch-and-bound rooted in first-order optimization. *Mathematical Programming*



# Full implementation

Tailored branch-and-bound algorithm based on perspective relaxation<sup>19</sup>

- Project out unnecessary variables
- Coordinate descent to solve relaxations
- Active sets
- Dual bounds
- Strong branching

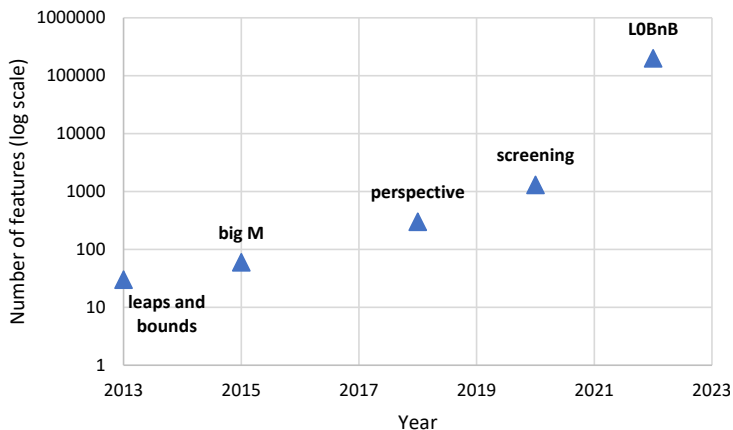
p	LOBnB	GRB	MSK	B
$10^3$	<b>0.7</b>	70	92	(4%)
$10^4$	<b>3</b>	(15%)	1697	–
$10^5$	<b>34</b>	–	–	–
$10^6$	<b>1112</b>	–	–	–

Time comparison (in seconds) with Gurobi, Mosek and Baron

---

<sup>19</sup>Hazimeh H et al (2022) Sparse regression at scale: Branch-and-bound rooted in first-order optimization. *Mathematical Programming*

# The journey so far...



Dimension of problems that can be comfortably solved

# Conclusion

- Convexification is harder than in MILO
- Some methods extend (less intuitive)
  - Disjunctive programming
  - RLTs
  - Lifting
- Implementation is non-trivial
- ... but it can work