

# Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher  
[volkan.cevher@epfl.ch](mailto:volkan.cevher@epfl.ch)

## *Lecture 11: Primal-dual optimization I*

Laboratory for Information and Inference Systems (LIONS)  
École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2020)



## License Information for Mathematics of Data Slides

- ▶ This work is released under a [Creative Commons License](#) with the following terms:
- ▶ **Attribution**
  - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- ▶ **Non-Commercial**
  - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- ▶ **Share Alike**
  - ▶ The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ▶ [Full Text of the License](#)

# Outline

- ▶ Today
  1. Min-max optimization (continued)
- ▶ Next week
  1. Algorithms for solving min-max optimization

## A minimax optimization template

### Minimax formulation

Consider the following problem that captures adversarial training, GANs, and robust reinforcement learning:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}), \quad (1)$$

where  $\Phi$  is differentiable and nonconvex in  $\mathbf{x}$  and nonconcave in  $\mathbf{y}$ .

o Key questions:

1. Where do the algorithms converge?
2. When do the algorithm converge?

## A minimax optimization template

### Minimax formulation

Consider the following problem that captures adversarial training, GANs, and robust reinforcement learning:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}), \quad (1)$$

where  $\Phi$  is differentiable and nonconvex in  $\mathbf{x}$  and nonconcave in  $\mathbf{y}$ .

o Key questions:

1. Where do the algorithms converge?
2. When do the algorithm converge?

### Recall: A buffet of negative results [5]

*“Even when the objective is a Lipschitz and smooth differentiable function, deciding whether a min-max point exists, in fact even deciding whether an approximate min-max point exists, is NP-hard. More importantly, an approximate local min-max point of large enough approximation is guaranteed to exist, but finding one such point is PPAD-complete. The same is true of computing an approximate fixed point of the (Projected) Gradient Descent/Ascent update dynamics.”*

## The difficulty of the nonconvex-nonconcave setting

### Minimax formulation

Consider the following problem that captures adversarial training, GANs, and robust reinforcement learning:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}), \quad (2)$$

where  $\Phi$  is differentiable and nonconvex in  $\mathbf{x}$  and nonconcave in  $\mathbf{y}$ .

### From minimax to minimization

Assume  $\Phi(\mathbf{x}, \mathbf{y}) = f(\mathbf{x})$  for all  $\mathbf{y}$ . The minimax optimization problem then seeks to find  $\mathbf{x}^*$  such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^p,$$

where  $\mathbf{x}^*$  is a global minimum of the nonconvex function  $f$ .

- ▶ Finding  $\mathbf{x}^*$  is NP-Hard even when  $f$  is smooth! (see the complexity supplementary material)
- ▶ Finding solutions to a nonconvex-nonconvex min-max problem is harder in general.

## Question 1 with a twist: Where do want the algorithms to converge?

### Definition (Saddle points & Local Nash equilibria)

The point  $(\mathbf{x}^*, \mathbf{y}^*)$  is called a saddle-point or a local Nash equilibrium (LNE) if it holds that

$$\Phi(\mathbf{x}^*, \mathbf{y}) \leq \Phi(\mathbf{x}^*, \mathbf{y}^*) \leq \Phi(\mathbf{x}, \mathbf{y}^*) \quad (\text{Saddle Point / LNE})$$

for all  $\mathbf{x}$  and  $\mathbf{y}$  within some neighborhood of  $\mathbf{x}^*$  and  $\mathbf{y}^*$ , i.e.,  $\|\mathbf{x} - \mathbf{x}^*\| \leq \delta$  and  $\|\mathbf{y} - \mathbf{y}^*\| \leq \delta$  for some  $\delta > 0$ .

### Necessary conditions

Through a Taylor expansion around  $\mathbf{x}^*$  and  $\mathbf{y}^*$  one can show that a LNE implies,

$$\nabla_{\mathbf{x}} \Phi(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} \Phi(\mathbf{x}, \mathbf{y}) = 0$$

$$\nabla_{\mathbf{x}\mathbf{x}} \Phi(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}\mathbf{y}} \Phi(\mathbf{x}, \mathbf{y}) \succeq 0$$

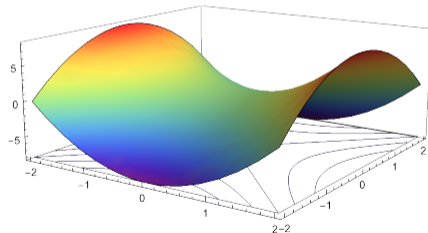


Figure:  $\Phi(x, y) = x^2 - y^2$

## Question 2 with a twist: When do generalized Robbins-Monro schemes converge?

- Given  $\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y})$ , define  $V(\mathbf{z}) = [-\nabla_{\mathbf{x}} \Phi(\mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} \Phi(\mathbf{x}, \mathbf{y})]$  with  $\mathbf{z} = [\mathbf{x}, \mathbf{y}]^\top$ .
- Given  $V(\mathbf{z})$ , define stochastic estimates of  $V(\mathbf{z}, \zeta) = V(\mathbf{z}) + U(\mathbf{z}, \zeta)$ , where
  - ▶  $U(\mathbf{z}, \zeta)$  is a bias term
  - ▶ We often have unbiasedness:  $EU(\mathbf{z}, \zeta) = 0$
  - ▶ The bias term can have bounded moments
  - ▶ We often have bounded variance:  $P(\|U(\mathbf{z}, \zeta)\| \geq t) \leq 2 \exp -\frac{t^2}{2\sigma^2}$  for  $\sigma > 0$ .
- An abstract template for generalized Robbins-Monro schemes, dubbed as  $\mathcal{A}$ :

$$\mathbf{z}^{k+1} = \mathbf{z}^k + \alpha_k V(\mathbf{z}^k, \zeta^k)$$

### The dessert section in the buffet of negative results: [12]

1. Bounded trajectories of  $\mathcal{A}$  always converge to an internally chain-transitive (ICT) set.
2. Trajectories of  $\mathcal{A}$  may converge with arbitrarily high probability to spurious attractors that contain no critical point of  $\Phi$ .



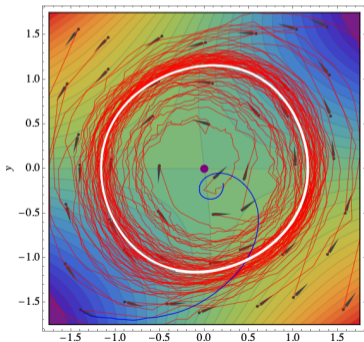
## Minimax is more difficult than just optimization [11]

○ Internally chain-transitive (ICT) sets characterize the convergence of dynamical systems [4].

- ▶ For optimization, {attracting ICT}  $\equiv$  {solutions}
- ▶ For minimax, {attracting ICT}  $\equiv$  {solutions}  $\cup$  {spurious sets}

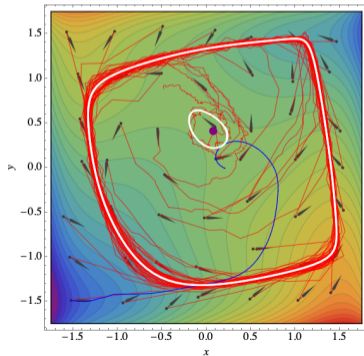
○ “Almost” bilinear  $\neq$  bilinear:

$$\Phi(x, y) = xy + \epsilon\phi(x), \phi(x) = \frac{1}{2}x^2 - \frac{1}{4}x^4$$



○ The “forsaken” solutions:

$$\Phi(y, x) = y(x-0.5) + \phi(y) - \phi(x), \phi(u) = \frac{1}{4}u^2 - \frac{1}{2}u^4 + \frac{1}{6}u^6$$



## A restricted minimax optimization template

### A restricted minimax formulation

Consider the following problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}), \quad (3)$$

where  $\Phi$  is convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ .

○ Key questions:

1. What problems does this template capture?
2. Where do the algorithms converge?
3. When do the algorithm converge?

## General nonsmooth problems

- We will show that the restricted template captures the familiar composite minimization:

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) + g(\mathbf{Ax}).$$

- ▶  $f, g$  are convex, nonsmooth functions; and  $\mathbf{A}$  is a linear operator.

### Examples

- ▶  $g(\mathbf{Ax}) = \|\mathbf{Ax} - \mathbf{b}\|_1$  or  $g(\mathbf{Ax}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$ .
- ▶  $g(\mathbf{Ax}) = \delta_{\{\mathbf{b}\}}(\mathbf{Ax})$ , where  $\delta_{\{\mathbf{b}\}}(\mathbf{Ax}) = \begin{cases} 0, & \text{if } \mathbf{Ax} = \mathbf{b}, \\ +\infty, & \text{if } \mathbf{Ax} \neq \mathbf{b}. \end{cases}$

- Observations:**
- The indicator example covers constrained problems, such as  $\min_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\}$ .
  - We need a tool, called Fenchel conjugation, to reveal the underlying minimax problem.

# Conjugation of functions

- o Idea: Represent a convex function in max-form:

## Definition

Let  $\mathcal{Q}$  be a Euclidean space and  $\mathcal{Q}^*$  be its dual space. Given a proper, closed and convex function  $f : \mathcal{Q} \rightarrow \mathbb{R} \cup \{+\infty\}$ , the function  $f^* : \mathcal{Q}^* \rightarrow \mathbb{R} \cup \{+\infty\}$  such that

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom}(f)} \{ \mathbf{y}^T \mathbf{x} - f(\mathbf{x}) \}$$

is called the Fenchel conjugate (or conjugate) of  $f$ .

- Observations:**
- o  $\mathbf{y}$  : slope of the hyperplane
  - o  $-f^*(\mathbf{y})$  : intercept of the hyperplane

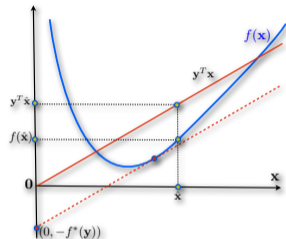


Figure: The conjugate function  $f^*(\mathbf{y})$  is the maximum gap between the linear function  $\mathbf{x}^T \mathbf{y}$  (red line) and  $f(\mathbf{x})$ .

## Conjugation of functions

### Definition

Given a **proper, closed and convex function**  $f : \mathcal{Q} \rightarrow \mathbb{R} \cup \{+\infty\}$ , the function  $f^* : \mathcal{Q}^* \rightarrow \mathbb{R} \cup \{+\infty\}$  such that

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom}(f)} \{ \mathbf{y}^T \mathbf{x} - f(\mathbf{x}) \}$$

is called the **Fenchel conjugate** (or conjugate) of  $f$ .

## Conjugation of functions

### Definition

Given a **proper, closed and convex function**  $f : \mathcal{Q} \rightarrow \mathbb{R} \cup \{+\infty\}$ , the function  $f^* : \mathcal{Q}^* \rightarrow \mathbb{R} \cup \{+\infty\}$  such that

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom}(f)} \{ \mathbf{y}^T \mathbf{x} - f(\mathbf{x}) \}$$

is called the **Fenchel conjugate** (or conjugate) of  $f$ .

### Properties

- $f^*$  is a **convex** and lower semicontinuous function by construction as the supremum of affine functions of  $\mathbf{y}$ .
- The **conjugate** of the **conjugate** of a convex function  $f$  is the same function  $f$ ; i.e.,  $f^{**} = f$  for  $f \in \mathcal{F}(\mathcal{Q})$ .
- The **conjugate** of the **conjugate** of a non-convex function  $f$  is its lower convex envelope when  $\mathcal{Q}$  is compact:
  - ▶  $f^{**}(\mathbf{x}) = \sup \{ g(\mathbf{x}) : g \text{ is convex and } g \leq f, \forall \mathbf{x} \in \mathcal{Q} \}$ .
- For closed convex  $f$ ,  $\mu$ -strong convexity w.r.t.  $\|\cdot\|$  is equivalent to  $\frac{1}{\mu}$  smoothness of  $f^*$  w.r.t.  $\|\cdot\|_*$ .
  - ▶ Recall dual norm:  $\|\mathbf{y}\|_* = \sup_{\mathbf{x}} \{ \langle \mathbf{x}, \mathbf{y} \rangle : \|\mathbf{x}\| \leq 1 \}$ .
  - ▶ See for example Theorem 3 in [16].

## Examples

### $\ell_2$ -norm-squared

$$f(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x}\|^2 \Rightarrow f^*(\mathbf{y}) = \max_{\mathbf{x}} \langle \mathbf{y}, \mathbf{x} \rangle - \frac{\lambda}{2} \|\mathbf{x}\|^2.$$

- Take the derivative and equate to 0:  $0 = \mathbf{y} - \lambda \mathbf{x} \iff \mathbf{x} = \frac{1}{\lambda} \mathbf{y} \iff f^*(\mathbf{y}) = \frac{1}{\lambda} \|\mathbf{y}\|^2 - \frac{1}{2\lambda} \|\mathbf{y}\|^2 = \frac{1}{2\lambda} \|\mathbf{y}\|^2.$

### $\ell_1$ -norm

$$f(\mathbf{x}) = \lambda \|\mathbf{x}\|_1 \Rightarrow f^*(\mathbf{y}) = \max_{\mathbf{x}} \langle \mathbf{y}, \mathbf{x} \rangle - \lambda \|\mathbf{x}\|_1.$$

- By definition of the  $\ell_1$ -norm:  $f^*(\mathbf{y}) = \max_{\mathbf{x}} \sum_{i=1}^n y_i x_i - \lambda |x_i| = \max_{\mathbf{x}} \sum_{i=1}^n y_i \text{sign}(x_i) |x_i| - \lambda |x_i|.$
- By inspection:
  - ▶ If all  $|y_i| \leq \lambda$ , then  $\forall i, (y_i \text{sign}(x_i) - \lambda) |x_i| \leq 0$ . Taking  $\mathbf{x} = 0$  gives the maximum value:  $f^*(\mathbf{y}) = 0$ .
  - ▶ If for at least one  $i, |y_i| > \lambda$ ,  $(y_i \text{sign}(x_i) - \lambda) |x_i| \rightarrow +\infty$  as  $|x_i| \rightarrow +\infty$ .
- $f^*(\mathbf{y}) = \delta_{\mathbf{y}: \|\cdot\|_\infty \leq \lambda}(\mathbf{y}) = \begin{cases} 0, & \text{if } \|\mathbf{y}\|_\infty \leq \lambda \\ +\infty, & \text{if } \|\mathbf{y}\|_\infty > \lambda \end{cases}$

**Remark:** ○ See advanced material at the end for non-convex examples, such as  $f(\mathbf{x}) = \|\mathbf{x}\|_0$ .

## General nonsmooth problems

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) + g(\mathbf{Ax})$$

- By Fenchel-conjugation, we have  $g(\mathbf{Ax}) = \max_{\mathbf{y}} \langle \mathbf{Ax}, \mathbf{y} \rangle - g^*(\mathbf{y})$ , where  $g^*$  is the conjugate of  $g$ .
- Min-max formulation:

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) + g(\mathbf{Ax}) = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y}} \{\Phi(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{Ax}, \mathbf{y} \rangle - g^*(\mathbf{y})\}$$

### An example with linear constraints

- If  $g(\mathbf{Ax}) = \delta_{\{\mathbf{b}\}}(\mathbf{Ax}) = \begin{cases} 0, & \text{if } \mathbf{Ax} = \mathbf{b}, \\ +\infty, & \text{if } \mathbf{Ax} \neq \mathbf{b}, \end{cases}$

$$\Rightarrow g^*(\mathbf{y}) = \max_{\mathbf{x}} \langle \mathbf{y}, \mathbf{x} \rangle - \delta_{\{\mathbf{b}\}}(\mathbf{x}) = \max_{\mathbf{x}: \mathbf{x}=\mathbf{b}} \langle \mathbf{y}, \mathbf{x} \rangle = \langle \mathbf{y}, \mathbf{b} \rangle.$$

- We reach the minimax formulation (or the so-called “Lagrangian”) via conjugation:

$$\min_{\mathbf{x}} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\} = \min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{Ax}) = \min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}) + \langle \mathbf{Ax} - \mathbf{b}, \mathbf{y} \rangle.$$



## A special case in minimax optimization

### Bilinear min-max template

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}) + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle - h(\mathbf{y}),$$

where  $\mathcal{X} \subseteq \mathbb{R}^p$  and  $\mathcal{Y} \subseteq \mathbb{R}^n$ .

- ▶  $f: \mathcal{X} \rightarrow \mathbb{R}$  is convex.
- ▶  $h: \mathcal{Y} \rightarrow \mathbb{R}$  is convex.

## Example: Sparse recovery

An example from sparseland  $\mathbf{b} = \mathbf{A}\mathbf{x}^\dagger + \mathbf{w}$ : constrained formulation

The basis pursuit denoising (BPDN) formulation is given by

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_1 : \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 \leq \|\mathbf{w}\|_2, \|\mathbf{x}\|_\infty \leq 1 \right\}. \quad (\text{BPDN})$$

A **primal problem** prototype

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{A}\mathbf{x} - \mathbf{b} \in \mathcal{K}, \mathbf{x} \in \mathcal{X} \right\},$$

The above template captures BPDN formulation with

- ▶  $f(\mathbf{x}) = \|\mathbf{x}\|_1$ .
- ▶  $\mathcal{K} = \{\|\mathbf{u}\| \in \mathbb{R}^n : \|\mathbf{u}\| \leq \|\mathbf{w}\|_2\}$ .
- ▶  $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_\infty \leq 1\}$ .

## An alternative formulation

### A primal problem prototype

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{Ax} - \mathbf{b} \in \mathcal{K}, \mathbf{x} \in \mathcal{X} \right\}, \quad (4)$$

- ▶  $f$  is a proper, closed and **convex** function
- ▶  $\mathcal{X}$  and  $\mathcal{K}$  are nonempty, closed **convex** sets
- ▶  $\mathbf{A} \in \mathbb{R}^{n \times p}$  and  $\mathbf{b} \in \mathbb{R}^n$  are known
- ▶ An optimal solution  $\mathbf{x}^*$  to (4) satisfies  $f(\mathbf{x}^*) = f^*$ ,  $\mathbf{Ax}^* - \mathbf{b} \in \mathcal{K}$  and  $\mathbf{x}^* \in \mathcal{X}$

### A simplified template without loss of generality

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b} \right\}, \quad (5)$$

- ▶  $f$  is a proper, closed and **convex** function
- ▶  $\mathbf{A} \in \mathbb{R}^{n \times p}$  and  $\mathbf{b} \in \mathbb{R}^n$  are known
- ▶ An optimal solution  $\mathbf{x}^*$  to (5) satisfies  $f(\mathbf{x}^*) = f^*$ ,  $\mathbf{Ax}^* = \mathbf{b}$

## Reformulation between templates

### A primal problem template

$$\min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{Ax} - \mathbf{b} \in \mathcal{K}, \mathbf{x} \in \mathcal{X} \right\}.$$

First step: Let  $\mathbf{r}_1 = \mathbf{Ax} - \mathbf{b} \in \mathbb{R}^n$  and  $\mathbf{r}_2 = \mathbf{x} \in \mathbb{R}^p$ .

$$\min_{\mathbf{x}, \mathbf{r}_1, \mathbf{r}_2} \left\{ f(\mathbf{x}) : \mathbf{r}_1 \in \mathcal{K}, \mathbf{r}_2 \in \mathcal{X}, \mathbf{Ax} - \mathbf{b} = \mathbf{r}_1, \mathbf{x} = \mathbf{r}_2 \right\}.$$

- Define  $\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{r}_1 \\ \mathbf{r}_2 \end{bmatrix} \in \mathbb{R}^{2p+n}$ ,  $\bar{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & -\mathbf{I}_{n \times n} & \mathbf{0}_{n \times p} \\ \mathbf{I}_{p \times p} & \mathbf{0}_{p \times n} & -\mathbf{I}_{p \times p} \end{bmatrix}$ ,  $\bar{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix}$ ,  $\bar{f}(\mathbf{z}) = f(\mathbf{x}) + \delta_{\mathcal{K}}(\mathbf{r}_1) + \delta_{\mathcal{X}}(\mathbf{r}_2)$ ,  
where  $\delta_{\mathcal{X}}(\mathbf{x}) = 0$ , if  $\mathbf{x} \in \mathcal{X}$ , and  $\delta_{\mathcal{X}}(\mathbf{x}) = +\infty$ , o/w.

### The simplified template

$$\min_{\mathbf{z} \in \mathbb{R}^{2p+n}} \left\{ \bar{f}(\mathbf{z}) : \bar{\mathbf{A}}\mathbf{z} = \bar{\mathbf{b}} \right\}.$$

## From constrained formulation back to minimax

### A general template

$$\min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\}.$$

Other examples:

- ▶ **Standard convex optimization** formulations: *linear programming, convex quadratic programming, second order cone programming, semidefinite programming and geometric programming.*
- ▶ **Reformulations** of existing unconstrained problems via **convex splitting**: *composite convex minimization, consensus optimization, ...*

### Formulating as min-max

$$\max_{\mathbf{y} \in \mathbb{R}^n} \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle = \begin{cases} 0, & \text{if } \mathbf{Ax} = \mathbf{b}, \\ +\infty, & \text{if } \mathbf{Ax} \neq \mathbf{b}. \end{cases}$$

$$\min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\} = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \{ \Phi(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle \}$$

## Dual problem

$$\min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\} = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \{\Phi(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle\}$$

◦ We define the dual problem

$$\max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}) := \max_{\mathbf{y} \in \mathbb{R}^n} \underbrace{\left\{ \min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle \right\}}_{d(\mathbf{y})}.$$

### Concavity of dual problem

Even if  $f(\mathbf{x})$  is not convex,  $d(\mathbf{y})$  is concave:

- ▶ For each  $\mathbf{x}$ ,  $d(\mathbf{y})$  is linear; i.e., it is both convex and concave.
- ▶ Pointwise minimum of concave functions is still concave.

**Remark:**

- If we can exchange min and max, we obtain a **concave** maximization problem.

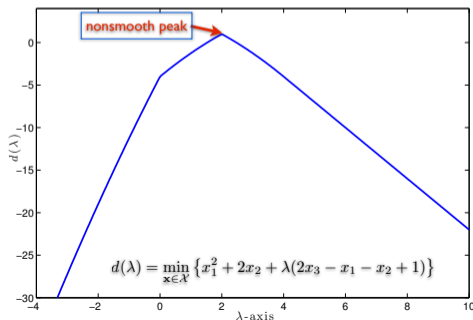
## Example: Nonsmoothness of the dual function

- Consider a constrained convex problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^3} \quad & \{ f(\mathbf{x}) := x_1^2 + 2x_2 \}, \\ \text{s.t.} \quad & 2x_3 - x_1 - x_2 = 1, \\ & \mathbf{x} \in \mathcal{X} := [-2, 2] \times [-2, 2] \times [0, 2]. \end{aligned}$$

- The **dual function** is **concave** and **nonsmooth** as written and then illustrated below.

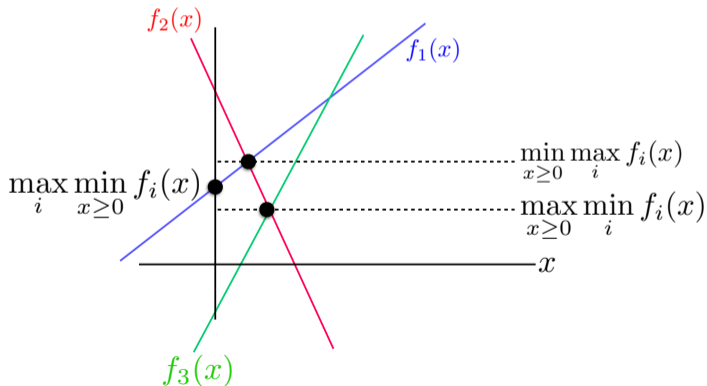
$$d(\lambda) := \min_{\mathbf{x} \in \mathcal{X}} \{ x_1^2 + 2x_2 + \lambda(2x_3 - x_1 - x_2 - 1) \}$$



## Exchanging min and max: A dangerous proposal

- Weak duality:

$$\underbrace{\max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y})}_{\text{Dual problem}} =: \boxed{\max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^p} \Phi(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \Phi(\mathbf{x}, \mathbf{y})} = \underbrace{\min_{\mathbf{x} \in \mathbb{R}^p} \{f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b}\}}_{\text{Primal problem}} = \begin{cases} f^*, & \text{if } \mathbf{Ax} = \mathbf{b} \\ +\infty, & \text{if } \mathbf{Ax} \neq \mathbf{b} \end{cases}$$





## A proof of weak duality

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b} \right\} = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \left\{ \Phi(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle \right\}$$

- Since  $\mathbf{Ax}^* = \mathbf{b}$ , it holds for any  $\mathbf{y}$

$$\begin{aligned} \Phi(\mathbf{x}^*, \mathbf{y}) &= f^* = f(\mathbf{x}^*) + \langle \mathbf{y}, \mathbf{Ax}^* - \mathbf{b} \rangle \\ &\geq \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle \right\} \\ &= \min_{\mathbf{x} \in \mathbb{R}^p} \Phi(\mathbf{x}, \mathbf{y}). \end{aligned}$$

- Take maximum of both sides in  $\mathbf{y}$  and note that  $f^*$  is independent of  $\mathbf{y}$ :

$$f^* = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \Phi(\mathbf{x}, \mathbf{y}) \geq \max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^p} \Phi(\mathbf{x}, \mathbf{y}) =: \max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}) = d^*.$$

## Strong duality and saddle points

### Strong duality

$$f^* = f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \Phi(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^p} \Phi(\mathbf{x}, \mathbf{y}) =: \max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}) = d^*.$$

Under strong duality and assuming existence of  $\mathbf{x}^*$ ,  $\Phi(\mathbf{x}, \mathbf{y})$  has a saddle point. We have primal and dual optimal values coincide, i.e.,  $f^* = d^*$ .

## Strong duality and saddle points

### Strong duality

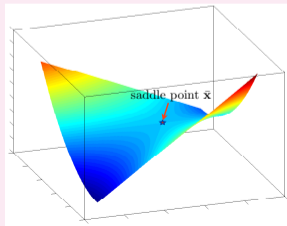
$$f^* = f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \Phi(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x} \in \mathbb{R}^p} \Phi(\mathbf{x}, \mathbf{y}) =: \max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}) = d^*.$$

Under strong duality and assuming existence of  $\mathbf{x}^*$ ,  $\Phi(\mathbf{x}, \mathbf{y})$  has a saddle point. We have primal and dual optimal values coincide, i.e.,  $f^* = d^*$ .

### Recall saddle point / LNE

A point  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathbb{R}^p \times \mathbb{R}^n$  is called a **saddle point** of  $\Phi$  if

$$\Phi(\mathbf{x}^*, \mathbf{y}) \leq \Phi(\mathbf{x}^*, \mathbf{y}^*) \leq \Phi(\mathbf{x}, \mathbf{y}^*), \quad \forall \mathbf{x} \in \mathbb{R}^p, \mathbf{y} \in \mathbb{R}^n.$$



## Toy example: Strong duality

### Primal problem

- Consider the following primal minimization problem:  $\min_{\mathbf{x}} P(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) := \frac{1}{2}\|\mathbf{x}\|^2 + \|\mathbf{x}\|_1$
- Using conjugation and strong duality

$$\begin{aligned} P(\mathbf{x}^*) &= \min_{\mathbf{x}} P(\mathbf{x}) = \min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}) + \langle \mathbf{x}, \mathbf{y} \rangle - g^*(\mathbf{y}), && \text{by conjugation} \\ &= \max_{\mathbf{y}} -g^*(\mathbf{y}) + \min_{\mathbf{x}} f(\mathbf{x}) + \langle \mathbf{x}, \mathbf{y} \rangle, && \text{by changing min-max} \\ &= \max_{\mathbf{y}} -g^*(\mathbf{y}) - \max_{\mathbf{x}} \langle \mathbf{x}, -\mathbf{y} \rangle - f(\mathbf{x}), && \text{by } \min f = -\max -f \\ &= \max_{\mathbf{y}} -g^*(\mathbf{y}) - f^*(-\mathbf{y}), && \text{by conjugation.} \end{aligned}$$

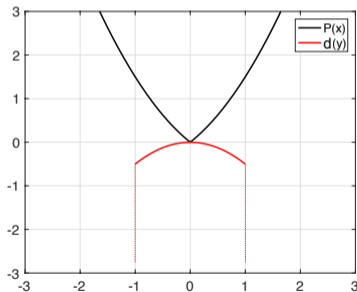
### Dual problem

- Dual problem:  $d^* = \max_{\mathbf{y}} d(\mathbf{y}) = -g^*(\mathbf{y}) - f^*(-\mathbf{y})$
- Recall  $f^*(-\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|^2$  and  $g^*(\mathbf{y}) = \delta_{\mathbf{y}: \|\mathbf{y}\|_{\infty} \leq 1}(\mathbf{y})$ .

## Toy example: Strong duality

$$\text{Primal problem: } \min_{\mathbf{x}} P(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2 + \|\mathbf{x}\|_1$$

$$\text{Dual problem: } \max_{\mathbf{y}} -\frac{1}{2} \|\mathbf{y}\|^2 - \delta_{\mathbf{y}: \|\mathbf{y}\|_\infty \leq 1}(\mathbf{y})$$



$$d(\mathbf{y}) = \begin{cases} -\frac{1}{2} \|\mathbf{y}\|^2, & \text{if } \|\mathbf{y}\|_\infty \leq 1 \\ -\infty, & \text{if } \|\mathbf{y}\|_\infty > 1 \end{cases}$$

## Back to convex-concave: Necessary and sufficient condition for strong duality

- Existence of a saddle point is not automatic even in convex-concave setting!
- Recall the minimax template:

$$\min_{\mathbf{x} \in \mathbb{R}^p} \max_{\mathbf{y} \in \mathbb{R}^n} \{\Phi(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{Ax} - \mathbf{b} \rangle\}$$

### Theorem (Necessary and sufficient optimality condition)

Under the *Slater's condition*:  $\text{relint}(\text{dom } f) \cap \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\} \neq \emptyset$ , strong duality holds, where the primal and dual problems are given by

$$f^* := \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^p} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{Ax} = \mathbf{b}, \end{cases} \quad \text{and} \quad d^* := \max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}).$$

- Remarks:**
- By definition of  $f^*$  and  $d^*$ , we always have  $d^* \leq f^*$  (**weak duality**).
  - If a primal solution exists and the Slater's condition holds, we have  $d^* = f^*$  (**strong duality**).

## Slater's qualification condition

- Denote  $\text{relint}(\text{dom } f)$  the **relative interior** of the domain.
- The **Slater condition** requires

$$\text{relint}(\text{dom } f) \cap \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\} \neq \emptyset. \quad (6)$$

### Special cases

- ▶ If  $\text{dom } f = \mathbb{R}^p$ , then (6)  $\Leftrightarrow \boxed{\exists \bar{\mathbf{x}} : \mathbf{A}\bar{\mathbf{x}} = \mathbf{b}}$ .
- ▶ If  $\text{dom } f = \mathbb{R}^p$  and instead of  $\mathbf{Ax} = \mathbf{b}$ , we have the feasible set  $\{\mathbf{x} : h(\mathbf{x}) \leq 0\}$ , where  $h$  is  $\mathbb{R}^p \rightarrow \mathbb{R}^q$  is convex, then

$$(6) \Leftrightarrow \boxed{\exists \bar{\mathbf{x}} : h(\bar{\mathbf{x}}) < 0.}$$

## Example: Slater's condition

### Example

Let us consider solving  $\min_{\mathbf{x} \in \mathcal{D}_\alpha} f(\mathbf{x})$  and so the feasible set is  $\mathcal{D}_\alpha := \mathcal{X} \cap \mathcal{A}_\alpha$ , where

$$\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 \leq 1\}, \mathcal{A}_\alpha := \{\mathbf{x} \in \mathbb{R}^2 : x_1 + x_2 = \alpha\},$$

where  $\alpha \in \mathbb{R}$ .



## Example: Slater's condition

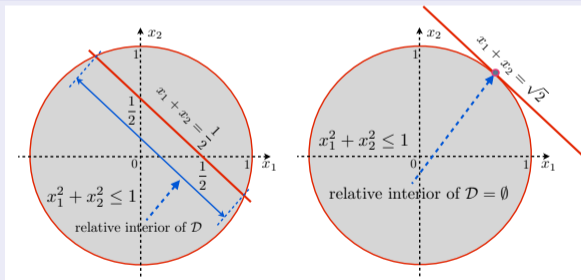
### Example

Let us consider solving  $\min_{\mathbf{x} \in \mathcal{D}_\alpha} f(\mathbf{x})$  and so the feasible set is  $\mathcal{D}_\alpha := \mathcal{X} \cap \mathcal{A}_\alpha$ , where

$$\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 \leq 1\}, \quad \mathcal{A}_\alpha := \{\mathbf{x} \in \mathbb{R}^2 : x_1 + x_2 = \alpha\},$$

where  $\alpha \in \mathbb{R}$ .

### Two cases where Slater's condition holds and does not hold



$\mathcal{D}_{1/2}$  satisfies Slater's condition –  $\mathcal{D}_{\sqrt{2}}$ -does not satisfy Slater's condition

## Performance of optimization algorithms

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{A}\mathbf{x} = \mathbf{b}, \right\}, \quad (\text{Affine-Constrained})$$

### Exact vs. approximate solutions

- ▶ Computing an **exact solution**  $\mathbf{x}^*$  to (Affine-Constrained) is **impracticable**
- ▶ Algorithms seek  $\mathbf{x}_\epsilon^*$  that **approximates**  $\mathbf{x}^*$  up to  $\epsilon$  in some sense

### A performance metric: Time-to-reach $\epsilon$

time-to-reach  $\epsilon$  = number of iterations to reach  $\epsilon$   $\times$  per iteration time

A key issue: Number of iterations to reach  $\epsilon$

**The notion of  $\epsilon$ -accuracy is elusive in constrained optimization!**

## Numerical $\epsilon$ -accuracy

- **Unconstrained case:** All iterates are feasible (no advantage from infeasibility)!

$$f(\mathbf{x}_\epsilon^*) - f^* \leq \epsilon$$

$$f^* = \min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x})$$

- **Constrained case:** We need to also measure the infeasibility of the iterates!

$$f^* - f(\mathbf{x}_\epsilon^*) \leq \epsilon \quad !!!$$

$$f^* = \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{Ax} = \mathbf{b} \right\} \quad (7)$$

### Our definition of $\epsilon$ -accurate solutions [20]

Given a numerical tolerance  $\epsilon \geq 0$ , a point  $\mathbf{x}_\epsilon^* \in \mathbb{R}^p$  is called an  $\epsilon$ -solution of (7) if

$$\begin{cases} f(\mathbf{x}_\epsilon^*) - f^* & \leq \epsilon \text{ (objective residual),} \\ \|\mathbf{Ax}_\epsilon^* - \mathbf{b}\| & \leq \epsilon \text{ (feasibility gap),} \end{cases}$$

- ▶ When  $\mathbf{x}^*$  is unique, we can also obtain  $\|\mathbf{x}_\epsilon^* - \mathbf{x}^*\| \leq \epsilon$  (iterate residual).

## Numerical $\epsilon$ -accuracy

### Constrained problems

Given a numerical tolerance  $\epsilon \geq 0$ , a point  $\mathbf{x}_\epsilon^* \in \mathbb{R}^p$  is called an  $\epsilon$ -solution of (7) if

$$\begin{cases} f(\mathbf{x}_\epsilon^*) - f^* & \leq \epsilon \text{ (objective residual),} \\ \|\mathbf{A}\mathbf{x}_\epsilon^* - \mathbf{b}\| & \leq \epsilon \text{ (feasibility gap),} \end{cases}$$

- ▶ When  $\mathbf{x}^*$  is unique, we can also obtain  $\|\mathbf{x}_\epsilon^* - \mathbf{x}^*\| \leq \epsilon$  (iterate residual).

### General minimax problems

Since duality gap is 0 at the solution, we measure the primal-dual gap

$$\text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\bar{\mathbf{x}}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}) \leq \epsilon. \quad (8)$$

#### Remarks:

- $\epsilon$  can be different for the objective, feasibility gap, or the iterate residual.
- It is easy to show  $\text{Gap}(\mathbf{x}, \mathbf{y}) \geq 0$  and  $\text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = 0$  iff  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is a saddle point.

## Primal-dual gap function for nonsmooth minimization

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + g(\mathbf{Ax}) = \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{f(\mathbf{x}) + \langle \mathbf{Ax}, \mathbf{y} \rangle - g^*(\mathbf{y})}_{\Phi(\mathbf{x}, \mathbf{y})} = \max_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathbf{Ax}, \mathbf{y} \rangle - g^*(\mathbf{y})$$

- Primal problem:  $\min_{\mathbf{x} \in \mathcal{X}} P(\mathbf{x})$  where

$$P(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}).$$

- Dual problem:  $\max_{\mathbf{y} \in \mathcal{Y}} d(\mathbf{y})$  where

$$d(\mathbf{y}) = \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \mathbf{y}).$$

- The primal-dual gap, i.e.,  $\text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ , is literally (primal value at  $\bar{\mathbf{x}}$ ) – (dual value at  $\bar{\mathbf{y}}$ ):

$$\text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = P(\bar{\mathbf{x}}) - d(\bar{\mathbf{y}}) = \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\bar{\mathbf{x}}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}).$$

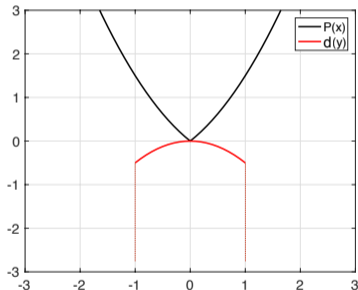
## Toy example for nonnegativity of gap

○  $P(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2 + \|\mathbf{x}\|_1$

○  $d(\mathbf{y}) = -\frac{1}{2}\|\mathbf{y}\|^2 - \delta_{\mathbf{y}: \|\mathbf{y}\|_\infty \leq 1}(\mathbf{y})$

Recall the indicator function

$$\delta_{\mathbf{y}: \|\mathbf{y}\|_\infty \leq 1}(\mathbf{y}) = \begin{cases} 0, & \text{if } \|\mathbf{y}\|_\infty \leq 1 \\ +\infty, & \text{if } \|\mathbf{y}\|_\infty > 1 \end{cases}$$



$$d(\mathbf{y}) = \begin{cases} -\frac{1}{2}\|\mathbf{y}\|^2, & \text{if } \|\mathbf{y}\|_\infty \leq 1 \\ -\infty, & \text{if } \|\mathbf{y}\|_\infty > 1 \end{cases}$$

## Primal-dual gap function in the general case

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \mathbf{y})$$

- Saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$  is such that  $\forall \mathbf{x} \in \mathbb{R}^p, \forall \mathbf{y} \in \mathbb{R}^n$ :

$$\Phi(\mathbf{x}^*, \mathbf{y}) \stackrel{(*)}{\leq} \Phi(\mathbf{x}^*, \mathbf{y}^*) \stackrel{(**)}{\leq} \Phi(\mathbf{x}, \mathbf{y}^*).$$

- Nonnegativity of Gap:

$$\begin{aligned} \text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) &= \max_{\mathbf{y} \in \mathcal{X}} \Phi(\bar{\mathbf{x}}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}) \\ &\geq \Phi(\bar{\mathbf{x}}, \mathbf{y}^*) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}), && \text{by the definition of maximization} \\ &\geq \Phi(\mathbf{x}^*, \mathbf{y}^*) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}), && \text{by the inequality (**)} \\ &\geq \Phi(\mathbf{x}^*, \bar{\mathbf{y}}) - \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x}, \bar{\mathbf{y}}), && \text{by the inequality (*)} \\ &\geq 0, && \text{by the definition of minimization.} \end{aligned}$$

- If  $(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = (\mathbf{x}^*, \mathbf{y}^*)$ , then all the inequalities will be equalities and  $\text{Gap}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = 0$ .

## Optimality conditions for minimax

### Saddle point

We say  $(\mathbf{x}^*, \mathbf{y}^*)$  is a primal-dual solution corresponding to primal and dual problems

$$f^* := \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^p} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{A}\mathbf{x} = \mathbf{b}, \end{cases} \quad \text{and} \quad d^* := \max_{\mathbf{y} \in \mathbb{R}^n} d(\mathbf{y}) = \max_{\mathbf{y} \in \mathbb{R}^n} \min_{\mathbf{x}} \Phi(\mathbf{x}, \mathbf{y}).$$

if it is a saddle point of  $\Phi(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \langle \mathbf{y}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle$ :

$$\Phi(\mathbf{x}^*, \mathbf{y}) \leq \Phi(\mathbf{x}^*, \mathbf{y}^*) \leq \Phi(\mathbf{x}, \mathbf{y}^*), \quad \forall \mathbf{x} \in \mathbb{R}^p, \mathbf{y} \in \mathbb{R}^n.$$

### Karush-Khun-Tucker (KKT) conditions

Under our assumptions, an equivalent characterization of  $(\mathbf{x}^*, \mathbf{y}^*)$  is via the KKT conditions of the problem

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) : \mathbf{A}\mathbf{x} = \mathbf{b},$$

which reads

$$\begin{cases} 0 \in \partial_{\mathbf{x}} \Phi(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{A}^T \mathbf{y}^* + \partial f(\mathbf{x}^*), \\ 0 = \nabla_{\mathbf{y}} \Phi(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{A}\mathbf{x}^* - \mathbf{b}. \end{cases}$$



## A naive proposal: Gradient descent-ascent (GDA)

### Towards algorithms for minimax optimization

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}).$$

We assume that

- ▶  $\Phi(\cdot, \mathbf{y})$  is convex,
- ▶  $\Phi(\mathbf{x}, \cdot)$  is concave,
- ▶  $\Phi$  is smooth in the following sense:

$$\left\| \begin{bmatrix} \nabla_{\mathbf{x}} \Phi(\mathbf{x}_1, \mathbf{y}_1) \\ -\nabla_{\mathbf{y}} \Phi(\mathbf{x}_1, \mathbf{y}_1) \end{bmatrix} - \begin{bmatrix} \nabla_{\mathbf{x}} \Phi(\mathbf{x}_2, \mathbf{y}_2) \\ -\nabla_{\mathbf{y}} \Phi(\mathbf{x}_2, \mathbf{y}_2) \end{bmatrix} \right\| \leq L \left\| \begin{bmatrix} \mathbf{x}_1 - \mathbf{x}_2 \\ \mathbf{y}_1 - \mathbf{y}_2 \end{bmatrix} \right\|. \quad (9)$$

- Let us try to use gradient descent for  $\mathbf{x}$ , gradient ascent for  $\mathbf{y}$  to obtain a solution

#### GDA

1. Choose  $\mathbf{x}^0, \mathbf{y}^0$  and  $\tau$ .
2. For  $k = 0, 1, \dots$ , perform:  
 $\mathbf{x}^{k+1} := \mathbf{x}^k - \tau \nabla_{\mathbf{x}} \Phi(\mathbf{x}^k, \mathbf{y}^k).$   
 $\mathbf{y}^{k+1} := \mathbf{y}^k + \tau \nabla_{\mathbf{y}} \Phi(\mathbf{x}^k, \mathbf{y}^k).$

## GDA on a simple problem

### Min-max problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y}).$$

#### SimGDA

1. Choose  $\mathbf{x}^0, \mathbf{y}^0$  and  $\tau$ .
2. For  $k = 0, 1, \dots$ , perform:  
 $\mathbf{x}^{k+1} := \mathbf{x}^k - \tau \nabla_{\mathbf{x}} \Phi(\mathbf{x}^k, \mathbf{y}^k).$   
 $\mathbf{y}^{k+1} := \mathbf{y}^k + \tau \nabla_{\mathbf{y}} \Phi(\mathbf{x}^k, \mathbf{y}^k).$

#### AltGDA

1. Choose  $\mathbf{x}^0, \mathbf{y}^0$  and  $\tau$ .
2. For  $k = 0, 1, \dots$ , perform:  
 $\mathbf{x}^{k+1} := \mathbf{x}^k - \tau \nabla_{\mathbf{x}} \Phi(\mathbf{x}^k, \mathbf{y}^k).$   
 $\mathbf{y}^{k+1} := \mathbf{y}^k + \tau \nabla_{\mathbf{y}} \Phi(\mathbf{x}^{k+1}, \mathbf{y}^k).$

### Example [9]

Let  $\Phi(x, y) = xy$ ,  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ , then,

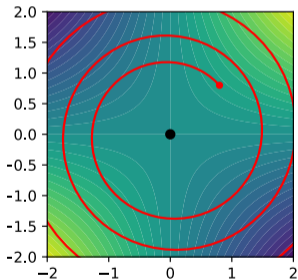
- ▶ for the iterates of SimGDA:  $x_{k+1}^2 + y_{k+1}^2 = (1 + \eta^2)(x_k^2 + y_k^2),$
- ▶ for the iterates of AltGDA:  $x_{k+1}^2 + y_{k+1}^2 = C(x_0^2 + y_0^2).$

o SimGDA diverges and AltGDA does not converge!

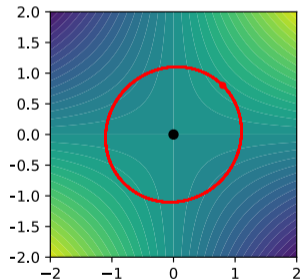
## Practical performance

$$\min_{x \in \mathbb{R}} \max_{y \in \mathbb{R}} xy$$

○ Simultaneous GDA



○ Alternating GDA



## Between convex-concave and nonconvex-nonconcave

### Nonconvex-concave problems

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y})$$

- $\Phi(\mathbf{x}, \mathbf{y})$  is nonconvex in  $\mathbf{x}$ , concave in  $\mathbf{y}$ , smooth in  $\mathbf{x}$  and  $\mathbf{y}$ .

### Recall

Define  $f(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\mathbf{x}, \mathbf{y})$ .

- Gradient descent applied to nonconvex  $f$  requires  $\mathcal{O}(\epsilon^{-2})$  iterations to give an  $\epsilon$ -stationary point.
- (Sub)gradient of  $f$  can be computed using Danskin's theorem:

$$\nabla_{\mathbf{x}} \Phi(\cdot, \mathbf{y}^*(\cdot)) \in \partial f(\cdot), \text{ where } \mathbf{y}^*(\cdot) \in \arg \max_{\mathbf{y} \in \mathcal{Y}} \Phi(\cdot, \mathbf{y}),$$

which is **tractable since  $\Phi$  is concave in  $\mathbf{y}$  [17]**.

- Remark:**
- “Conceptually” much easier than nonconvex-nonconcave case.

# Epilogue

	Gradient complexity	Optimality measure	Reference
convex-concave	$\mathcal{O}(\epsilon^{-1})^1$	$\epsilon$ optimality w.r.t. duality gap	Nemirovski, 2004; Chambolle & Pock, 2011;
nonconvex-concave	$\tilde{\mathcal{O}}(\epsilon^{-2.5})^3$	$\epsilon$ -stationarity w.r.t. gradient mapping norm	Tran-Dinh & Cevher, 2014. <sup>2</sup>
nonconvex-nonconcave	HARD	HARD	Lin, Jin, & Jordan, 2020. <sup>4</sup> Daskalakis, Stratis, & Zampetakis, 2020; Hsieh, Mertikopoulos, & Cevher, 2020. <sup>5</sup>

<sup>1</sup>Rates are not directly comparable as duality gap and gradient mapping norm are not necessarily of the same order!

<sup>2</sup>Arkadi Nemirovski, "Prox-method with rate of convergence  $\mathcal{O}(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems." *SIAM Journal on Optimization* 15.1 (2004): 229-251.

Antonin Chambolle, and Thomas Pock, "A first-order primal-dual algorithm for convex problems with applications to imaging." *Journal of mathematical imaging and vision* 40.1 (2011): 120-145.

Quoc Tran-Dinh, and Volkan Cevher, "Constrained convex minimization via model-based excessive gap." *Advances in Neural Information Processing Systems*. 2014.

<sup>3</sup>The rate is  $\tilde{\mathcal{O}}(\epsilon^{-2})$  for strongly concave problems.

<sup>4</sup>Tianyi Lin, Chi Jin, and Michael Jordan, "Near-optimal algorithms for minimax optimization." *arXiv preprint arXiv:2002.02417* (2020).

<sup>5</sup>Constantinos Daskalakis, Stratis Skoulakis, and Manolis Zampetakis, "The complexity of constrained min-max optimization." *arXiv preprint arXiv:2009.09623* (2020).

Ya-Ping Hsieh, Panayotis Mertikopoulos, and Volkan Cevher, "The limits of min-max optimization algorithms: convergence to spurious non-critical sets." *arXiv preprint arXiv:2006.09065* (2020).

## A new hope

$$\min_{x \in \mathbb{R}} \max_{y \in \mathbb{R}} xy$$

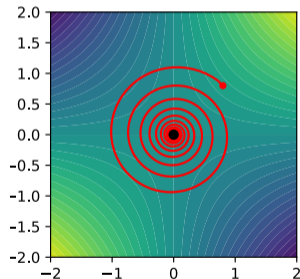
o Next lecture: Some algorithms that actually **converge**!

o Convergence of the sequence:

There exists  $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$ , such that  $\mathbf{z}_k \rightarrow \mathbf{z}^*$ .

o Convergence rate:

$$\text{Gap} \left( \frac{1}{K} \sum_{k=1}^K \mathbf{x}^k, \frac{1}{K} \sum_{k=1}^K \mathbf{y}^k \right) \leq \mathcal{O} \left( \frac{1}{K} \right).$$



## Wrap up!

- Recitation on Friday is on Homework #2...

## A *convex* proto-problem for *structured* sparsity

A combinatorial approach for estimating  $\mathbf{x}^\natural$  from  $\mathbf{b} = \mathbf{A}\mathbf{x}^\natural + \mathbf{w}$

We may consider the sparsest estimator or its surrogate with a valid sparsity pattern:

$$\hat{\mathbf{x}} \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_s : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa, \|\mathbf{x}\|_\infty \leq 1 \right\} \quad (\mathcal{P}_s)$$

with some  $\kappa \geq 0$ . If  $\kappa = \|\mathbf{w}\|_2$ , then the structured sparse  $\mathbf{x}^\natural$  is a feasible solution.

## Sparsity and structure together [7]

Given some weights  $\mathbf{d} \in \mathbb{R}^d$ ,  $\mathbf{e} \in \mathbb{R}^p$  and an integer input  $c \in \mathbb{Z}^l$ , we define

$$\|\mathbf{x}\|_s := \min_{\omega} \{ \mathbf{d}^T \omega + \mathbf{e}^T \mathbf{s} : M \begin{bmatrix} \omega \\ \mathbf{s} \end{bmatrix} \leq \mathbf{c}, \mathbb{1}_{\text{supp}(\mathbf{x})} = \mathbf{s}, \omega \in \{0, 1\}^d \}$$

for all feasible  $\mathbf{x}$ ,  $\infty$  otherwise. The parameter  $\omega$  is useful for *latent* modeling.



## A *convex* proto-problem for *structured* sparsity

A combinatorial approach for estimating  $\mathbf{x}^\natural$  from  $\mathbf{b} = \mathbf{A}\mathbf{x}^\natural + \mathbf{w}$

We may consider the sparsest estimator or its surrogate with a valid sparsity pattern:

$$\hat{\mathbf{x}} \in \arg \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ \|\mathbf{x}\|_s : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa, \|\mathbf{x}\|_\infty \leq 1 \right\} \quad (\mathcal{P}_s)$$

with some  $\kappa \geq 0$ . If  $\kappa = \|\mathbf{w}\|_2$ , then the structured sparse  $\mathbf{x}^\natural$  is a feasible solution.

## Sparsity and structure together [7]

Given some weights  $\mathbf{d} \in \mathbb{R}^d, \mathbf{e} \in \mathbb{R}^p$  and an integer input  $\mathbf{c} \in \mathbb{Z}^l$ , we define

$$\|\mathbf{x}\|_s := \min_{\boldsymbol{\omega}} \{ \mathbf{d}^T \boldsymbol{\omega} + \mathbf{e}^T \mathbf{s} : M \begin{bmatrix} \boldsymbol{\omega} \\ \mathbf{s} \end{bmatrix} \leq \mathbf{c}, \mathbf{1}_{\text{supp}(\mathbf{x})} = \mathbf{s}, \boldsymbol{\omega} \in \{0, 1\}^d \}$$

for all feasible  $\mathbf{x}$ ,  $\infty$  otherwise. The parameter  $\boldsymbol{\omega}$  is useful for *latent* modeling.

A convex candidate solution for  $\mathbf{b} = \mathbf{A}\mathbf{x}^\natural + \mathbf{w}$

We use the *convex* estimator based on the *tightest* convex relaxation of  $\|\mathbf{x}\|_s$ :

$$\hat{\mathbf{x}} \in \arg \min_{\mathbf{x} \in \text{dom}(\|\cdot\|_s)} \left\{ \|\mathbf{x}\|_s^{**} : \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 \leq \kappa \right\} \text{ with some } \kappa \geq 0, \text{dom}(\|\cdot\|_s) := \{\mathbf{x} : \|\mathbf{x}\|_s < \infty\}.$$

## Tractability & tightness of biconjugation

### Proposition (Hardness of conjugation)

Let  $F(s) : 2^{\mathfrak{B}} \rightarrow \mathbb{R} \cup \{+\infty\}$  be a set function defined on the support  $s = \text{supp}(\mathbf{x})$ . Conjugate of  $F$  over the unit infinity ball  $\|\mathbf{x}\|_{\infty} \leq 1$  is given by

$$g^*(\mathbf{y}) = \sup_{s \in \{0,1\}^{\mathfrak{P}}} |\mathbf{y}|^T s - F(s).$$

### Observations:

- ▶  $F(s)$  is general set function

**Computation:** NP-Hard

- ▶  $F(s) = \|\mathbf{x}\|_s$

**Computation:** Integer Linear Program (ILP) in general. However, if

- ▶  $M$  is Totally Unimodular **TU**
- ▶  $(M, \mathbf{c})$  is Total Dual Integral **TDI**

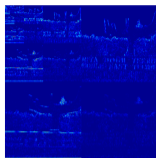
then tight convex relaxations with a linear program (LP, which is “usually” tractable)

Otherwise, relax to LP anyway!

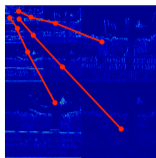
- ▶  $F(s)$  is submodular

**Computation:** Polynomial-time

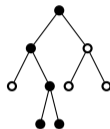
## Tree sparsity [15, 6, 3, 21]



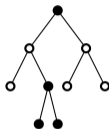
Wavelet coefficients



Wavelet tree



Valid selection of nodes



Invalid selection of nodes

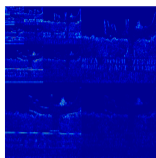
**Structure:** *We seek the sparsest signal with a rooted connected subtree support.*

**Linear description:** A **valid** support satisfy  $s_{\text{parent}} \geq s_{\text{child}}$  over tree  $\mathcal{T}$

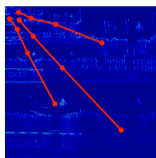
$$\mathbf{T}\mathbf{1}_{\text{supp}(\mathbf{x})} := \mathbf{T}\mathbf{s} \geq 0$$

where  $\mathbf{T}$  is the directed edge-node incidence matrix, which is **TU**.

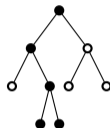
## Tree sparsity [15, 6, 3, 21]



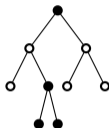
Wavelet coefficients



Wavelet tree



Valid selection of nodes



Invalid selection of nodes

**Structure:** *We seek the sparsest signal with a rooted connected subtree support.*

**Linear description:** A **valid** support satisfy  $s_{\text{parent}} \geq s_{\text{child}}$  over tree  $\mathcal{T}$

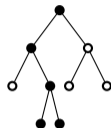
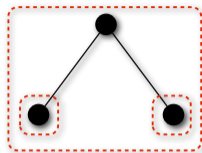
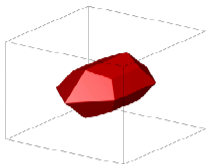
$$\mathbf{T}\mathbf{1}_{\text{supp}(\mathbf{x})} := \mathbf{T}\mathbf{s} \geq 0$$

where  $\mathbf{T}$  is the directed edge-node incidence matrix, which is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_s^{**} = \min_{\mathbf{s} \in [0,1]^p} \{\mathbf{1}^T \mathbf{s} : \mathbf{T}\mathbf{s} \geq 0, |\mathbf{x}| \leq \mathbf{s}\}$

for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

## Tree sparsity [15, 6, 3, 21]



$$\mathcal{G}_H = \{\{1, 2, 3\}, \{2\}, \{3\}\} \quad \text{valid selection of nodes}$$

**Structure:** *We seek the sparsest signal with a rooted connected subtree support.*

**Linear description:** A **valid** support satisfy  $s_{\text{parent}} \geq s_{\text{child}}$  over tree  $\mathcal{T}$

$$\mathbf{T} \mathbf{1}_{\text{supp}(\mathbf{x})} := \mathbf{T} \mathbf{s} \geq 0$$

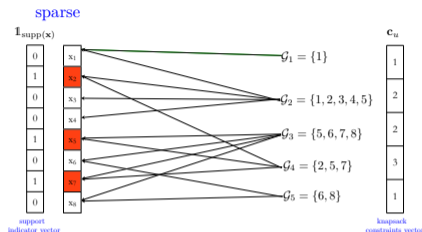
where  $\mathbf{T}$  is the directed edge-node incidence matrix, which is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_s^{**} = \min_{\mathbf{s} \in [0, 1]^p} \{\mathbf{1}^T \mathbf{s} : \mathbf{T} \mathbf{s} \geq 0, |\mathbf{x}| \leq \mathbf{s}\} \stackrel{*}{=} \sum_{\mathcal{G} \in \mathcal{G}_H} \|x_{\mathcal{G}}\|_{\infty}$

for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

The set  $\mathcal{G} \in \mathcal{G}_H$  are defined as each node and all its descendants.

## Group knapsack sparsity [23, 10, 8]



**Structure:** *We seek the sparsest signal with group allocation constraints.*

**Linear description:** A **valid** support obeys budget constraints over  $\mathcal{G}$

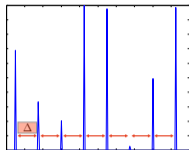
$$\mathfrak{B}^T \mathbf{s} \leq c_u$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix or  $\mathcal{G}$  has a *loopless* group intersection graph, it is **TU**.

**Remark:** We can also budget a lowerbound  $c_\ell \leq \mathfrak{B}^T \mathbf{s} \leq c_u$ .

## Group knapsack sparsity [23, 10, 8]



$$\mathfrak{B}^T = \begin{bmatrix} 1 & 1 & \cdots & 1 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 1 & \cdots & 1 & 1 & 0 & \cdots & 0 \\ & & & & & & & \ddots & \\ & & & & & & & & \ddots & \\ 0 & \cdots & 0 & 0 & 1 & 1 & \cdots & 1 & 1 \end{bmatrix}_{(p-\Delta+1) \times p}$$

**Structure:** We seek the sparsest signal with group allocation constraints.

**Linear description:** A **valid** support obeys budget constraints over  $\mathfrak{G}$

$$\mathfrak{B}^T \mathbf{s} \leq \mathbf{c}_u$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathfrak{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix or  $\mathfrak{G}$  has a **loopless** group intersection graph, it is **TU**.

**Remark:** We can also budget a lowerbound  $\mathbf{c}_\ell \leq \mathfrak{B}^T \mathbf{s} \leq \mathbf{c}_u$ .

$$\text{Biconjugate: } \|\mathbf{x}\|_s^{**} = \begin{cases} \|\mathbf{x}\|_1 & \text{if } \mathbf{x} \in [-1, 1]^p, \mathfrak{B}^T |\mathbf{x}| \leq \mathbf{c}_u, \\ \infty & \text{otherwise} \end{cases}$$

For the neuronal spike example, we have  $\mathbf{c}_u = \mathbf{1}$ .

## Group knapsack sparsity [23, 10, 8]

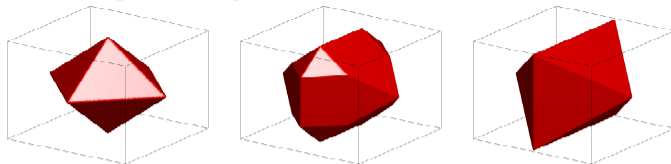


Figure: \*

(left)  $\|\mathbf{x}\|_s^{**} \leq 1$  (middle)  $\|\mathbf{x}\|_s^{**} \leq 1.5$  (right)  $\|\mathbf{x}\|_s^{**} \leq 2$  for  $\mathcal{G} = \{\{1, 2\}, \{2, 3\}\}$

**Structure:** We seek the sparsest signal with group allocation constraints.

**Linear description:** A valid support obeys budget constraints over  $\mathcal{G}$

$$\mathfrak{B}^T \mathbf{s} \leq \mathbf{c}_u$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix or  $\mathcal{G}$  has a *loopless* group intersection graph, it is **TU**.

**Remark:** We can also budget a lowerbound  $\mathbf{c}_l \leq \mathfrak{B}^T \mathbf{s} \leq \mathbf{c}_u$ .

**Biconjugate:**  $\|\mathbf{x}\|_s^{**} = \begin{cases} \|\mathbf{x}\|_1 & \text{if } \mathbf{x} \in [-1, 1]^p, \mathfrak{B}^T |\mathbf{x}| \leq \mathbf{c}_u, \\ \infty & \text{otherwise} \end{cases}$



# Group knapsack sparsity example: A stylized spike train

- ▶ Basis pursuit (BP):  $\|\mathbf{x}\|_1$
- ▶ TU-relax (TU):

$$\|\mathbf{x}\|_s^{**} = \begin{cases} \|\mathbf{x}\|_1 & \text{if } \mathbf{x} \in [-1, 1]^p, \mathfrak{B}^T |\mathbf{x}| \leq \mathbf{c}_u, \\ \infty & \text{otherwise} \end{cases}$$

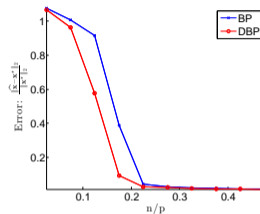
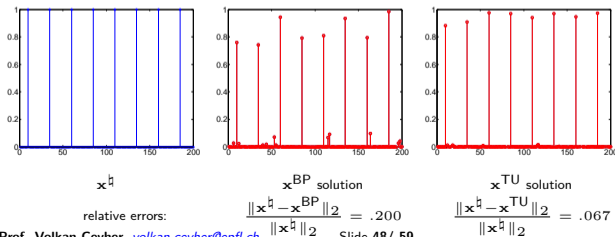
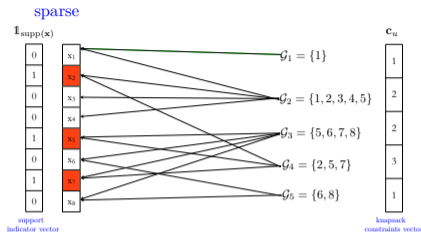


Figure: Recovery for  $n = 0.18p$ .



## Group knapsack sparsity: A simple variation



**Structure:** We seek the signal with the minimal overall group allocation.

$$\text{Objective: } \mathbf{1}^T \mathbf{s} \rightarrow \|\mathbf{x}\|_{\omega} = \begin{cases} \min_{\omega \in \mathbb{Z}_{++}} \omega & \text{if } \mathbf{x} \in [-1, 1]^p, \mathfrak{B}^T \mathbf{s} \leq \omega \mathbf{1}, \\ \infty & \text{otherwise} \end{cases}$$

**Linear description:** A valid support obeys budget constraints over  $\mathfrak{G}$

$$\mathfrak{B}^T \mathbf{s} \leq \omega \mathbf{1}$$

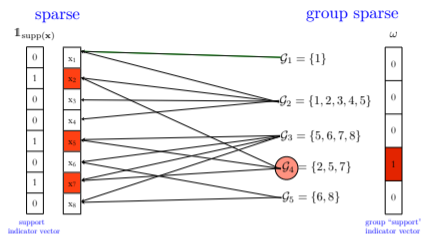
where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathfrak{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix or  $\mathfrak{G}$  has a *loopless* group intersection graph, it is **TU**.

$$\text{Biconjugate: } \|\mathbf{x}\|_s^{**} = \begin{cases} \max_{\mathcal{G} \in \mathfrak{G}} \|\mathbf{x}^{\mathcal{G}}\|_1 & \text{if } \mathbf{x} \in [-1, 1]^p, \\ \infty & \text{otherwise} \end{cases}$$

**Remark:** The regularizer is known as *exclusive Lasso* [23, 19]

## Group cover sparsity: Minimal group cover [2, 18, 13]



**Structure:** We seek the signal covered by a minimal number of groups.

$$\text{Objective: } \mathbf{1}^T s \rightarrow d^T w$$

**Linear description:** At least one group containing a sparse coefficient is selected

$$\mathfrak{B}w \geq s$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathfrak{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix, or  $\mathfrak{G}$  has a *loopless* group intersection graph it is TU.

## Group cover sparsity: **Minimal group cover** [2, 18, 13]

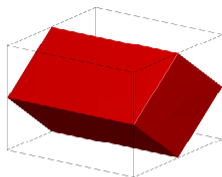


Figure:  $\mathcal{G} = \{\{1, 2\}, \{2, 3\}\}$ , unit group weights  $\mathbf{d} = \mathbf{1}$ .

**Structure:** *We seek the signal covered by a minimal number of groups.*

$$\text{Objective: } \mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega}$$

**Linear description:** **At least one** group containing a sparse coefficient is selected

$$\mathfrak{B}\boldsymbol{\omega} \geq \mathbf{s}$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix, or  $\mathcal{G}$  has a *loopless* group intersection graph it is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}^*} = \min_{\boldsymbol{\omega} \in [0,1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathfrak{B}\boldsymbol{\omega} \geq |\mathbf{x}|\}$  for  $\mathbf{x} \in [-1, 1]^P$ ,  $\infty$  otherwise

## Group cover sparsity: **Minimal group cover** [2, 18, 13]

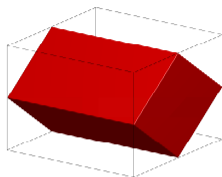


Figure:  $\mathcal{G} = \{\{1, 2\}, \{2, 3\}\}$ , unit group weights  $\mathbf{d} = \mathbf{1}$ .

**Structure:** *We seek the signal covered by a minimal number of groups.*

$$\text{Objective: } \mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega}$$

**Linear description:** **At least one** group containing a sparse coefficient is selected

$$\mathfrak{B}\boldsymbol{\omega} \geq \mathbf{s}$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix, or  $\mathcal{G}$  has a *loopless* group intersection graph it is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}}^{**} = \min_{\boldsymbol{\omega} \in [0,1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathfrak{B}\boldsymbol{\omega} \geq |\mathbf{x}|\}$  for  $\mathbf{x} \in [-1, 1]^P$ ,  $\infty$  otherwise  
 $\stackrel{*}{=} \min_{\mathbf{v}_i \in \mathbb{R}^P} \{\sum_{i=1}^M d_i \|\mathbf{v}_i\|_{\infty} : \mathbf{x} = \sum_{i=1}^M \mathbf{v}_i, \forall \text{supp}(\mathbf{v}_i) \subseteq \mathcal{G}_i\}$ ,

## Group cover sparsity: **Minimal group cover** [2, 18, 13]

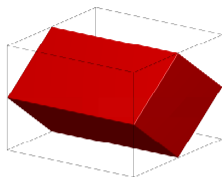


Figure:  $\mathcal{G} = \{\{1, 2\}, \{2, 3\}\}$ , unit group weights  $\mathbf{d} = \mathbf{1}$ .

**Structure:** *We seek the signal covered by a minimal number of groups.*

$$\text{Objective: } \mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega}$$

**Linear description:** **At least one** group containing a sparse coefficient is selected

$$\mathfrak{B}\boldsymbol{\omega} \geq \mathbf{s}$$

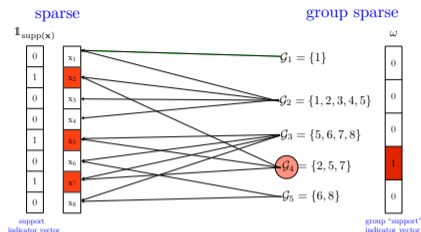
where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\mathfrak{B}$  is an interval matrix, or  $\mathcal{G}$  has a *loopless* group intersection graph it is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}}^{**} = \min_{\boldsymbol{\omega} \in [0,1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathfrak{B}\boldsymbol{\omega} \geq |\mathbf{x}|\}$  for  $\mathbf{x} \in [-1, 1]^P$ ,  $\infty$  otherwise  
 $\stackrel{*}{=} \min_{\mathbf{v}_i \in \mathbb{R}^P} \{\sum_{i=1}^M d_i \|\mathbf{v}_i\|_{\infty} : \mathbf{x} = \sum_{i=1}^M \mathbf{v}_i, \forall \text{supp}(\mathbf{v}_i) \subseteq \mathcal{G}_i\}$ ,

**Remark:** Weights  $\mathbf{d}$  can depend on the **sparsity** within each groups (**not TU**) [7].

## Budgeted group cover sparsity



**Structure:** We seek the sparsest signal covered by  $G$  groups.

$$\text{Objective: } d^T \omega \rightarrow \mathbb{1}^T s$$

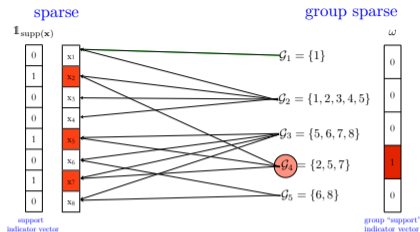
**Linear description:** At least one of the  $G$  selected groups cover each sparse coefficient.

$$\mathfrak{B}\omega \geq s, \mathbb{1}^T \omega \leq G$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathcal{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\begin{bmatrix} \mathfrak{B} \\ \mathbb{1} \end{bmatrix}$  is an interval matrix, it is TU.

## Budgeted group cover sparsity



**Structure:** We seek the sparsest signal covered by  $G$  groups.

$$\text{Objective: } \mathbf{d}^T \omega \rightarrow \mathbb{1}^T \mathbf{s}$$

**Linear description:** At least one of the  $G$  selected groups cover each sparse coefficient.

$$\mathfrak{B}\omega \geq \mathbf{s}, \mathbb{1}^T \omega \leq G$$

where  $\mathfrak{B}$  is the biadjacency matrix of  $\mathfrak{G}$ , i.e.,  $\mathfrak{B}_{ij} = 1$  iff  $i$ -th coefficient is in  $\mathcal{G}_j$ .

When  $\begin{bmatrix} \mathfrak{B} \\ \mathbb{1} \end{bmatrix}$  is an interval matrix, it is TU.

**Biconjugate:**  $\|\mathbf{x}\|_{\omega}^{**} = \min_{\omega \in [0,1]^M} \{\|\mathbf{x}\|_1 : \mathfrak{B}\omega \geq |\mathbf{x}|, \mathbb{1}^T \omega \leq G\}$   
for  $\mathbf{x} \in [-1, 1]^p, \infty$  otherwise.



## Budgeted group cover example: Interval overlapping groups

- ▶ Basis pursuit (BP):  $\|\mathbf{x}\|_1$
- ▶ Sparse group Lasso (SGL<sub>q</sub>):

$$(1 - \alpha) \sum_{\mathcal{G} \in \mathbb{G}} \sqrt{|\mathcal{G}|} \|\mathbf{x}^{\mathcal{G}}\|_q + \alpha \|\mathbf{x}^{\mathcal{G}}\|_1$$

- ▶ TU-relax (TU):

$$\|\mathbf{x}\|_{\omega}^{**} = \min_{\omega \in [0,1]^M} \{ \|\mathbf{x}\|_1 : \mathfrak{B}\omega \geq |\mathbf{x}|, \mathbf{1}^T \omega \leq G \}$$

for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

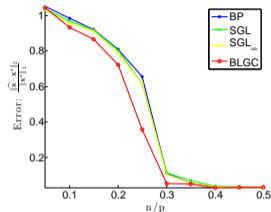
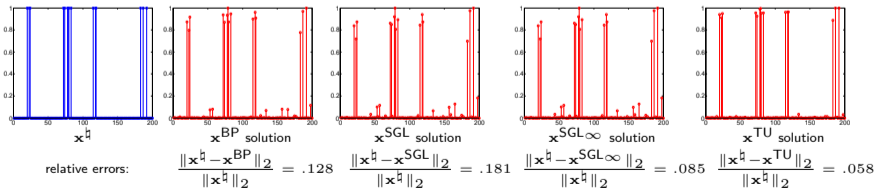
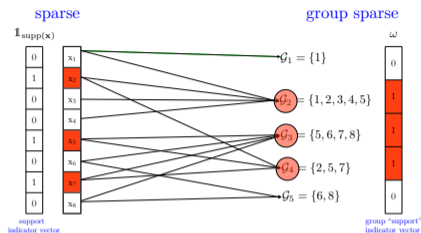


Figure: Recovery for  $n = 0.25p$ ,  $s = 15$ ,  $p = 200$ ,  $G = 5$  out of  $M = 29$  groups.



## Group intersection sparsity [14, 22, 1]



**Structure:** We seek the signal intersecting with minimal number of groups.

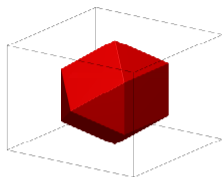
$$\text{Objective: } \mathbf{1}^T s \rightarrow d^T w$$

**Linear description:** All groups containing a sparse coefficient are selected

$$\mathbf{H}_k s \leq w, \forall k \in \mathfrak{F}$$

where  $\mathbf{H}_k(i, j) = \begin{cases} 1 & \text{if } j = k, j \in \mathcal{G}_i \\ 0 & \text{otherwise} \end{cases}$ , which is TU.

## Group intersection sparsity [14, 22, 1]



$\mathfrak{G} = \{\{1, 2\}, \{2, 3\}\}$ , unit group weights  $\mathbf{d} = \mathbf{1}$   
(left) intersection (right) cover.

**Structure:** *We seek the signal intersecting with minimal number of groups.*

**Objective:**  $\mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega}$

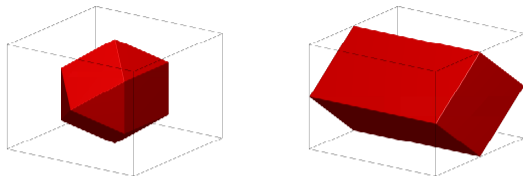
**Linear description:** All groups containing a sparse coefficient are selected

$$\mathbf{H}_k \mathbf{s} \leq \boldsymbol{\omega}, \forall k \in \mathfrak{F}$$

where  $\mathbf{H}_k(i, j) = \begin{cases} 1 & \text{if } j = k, j \in \mathcal{G}_i \\ 0 & \text{otherwise} \end{cases}$ , which is TU.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}^*} = \min_{\boldsymbol{\omega} \in [0, 1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathbf{H}_k |\mathbf{x}| \leq \boldsymbol{\omega}, \forall k \in \mathfrak{F}\}$   
for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

## Group intersection sparsity [14, 22, 1]



$\mathfrak{G} = \{\{1, 2\}, \{2, 3\}\}$ , unit group weights  $\mathbf{d} = \mathbf{1}$   
(left) intersection (right) cover.

**Structure:** We seek the signal intersecting with minimal number of groups.

**Objective:**  $\mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega}$  (submodular)

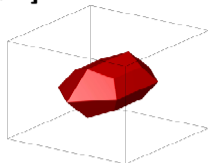
**Linear description:** All groups containing a sparse coefficient are selected

$$\mathbf{H}_k \mathbf{s} \leq \boldsymbol{\omega}, \forall k \in \mathfrak{F}$$

where  $\mathbf{H}_k(i, j) = \begin{cases} 1 & \text{if } j = k, j \in \mathcal{G}_i \\ 0 & \text{otherwise} \end{cases}$ , which is TU.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}^*} = \min_{\boldsymbol{\omega} \in [0, 1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathbf{H}_k |\mathbf{x}| \leq \boldsymbol{\omega}, \forall k \in \mathfrak{F}\}^* = \sum_{\mathcal{G} \in \mathfrak{G}} \|x_{\mathcal{G}}\|_{\infty}$   
for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

## Group intersection sparsity [14, 22, 1]



$$\mathfrak{G} = \{\{1, 2, 3\}, \{2\}, \{3\}\}, \text{ unit group weights } \mathbf{d} = \mathbf{1}.$$

**Structure:** *We seek the signal intersecting with minimal number of groups.*

$$\text{Objective: } \mathbf{1}^T \mathbf{s} \rightarrow \mathbf{d}^T \boldsymbol{\omega} \quad (\text{submodular})$$

**Linear description:** All groups containing a sparse coefficient are selected

$$\mathbf{H}_k \mathbf{s} \leq \omega, \forall k \in \mathfrak{F}$$

where  $\mathbf{H}_k(i, j) = \begin{cases} 1 & \text{if } j = k, j \in \mathcal{G}_i \\ 0 & \text{otherwise} \end{cases}$ , which is **TU**.

**Biconjugate:**  $\|\mathbf{x}\|_{\boldsymbol{\omega}^*} = \min_{\boldsymbol{\omega} \in [0,1]^M} \{\mathbf{d}^T \boldsymbol{\omega} : \mathbf{H}_k |\mathbf{x}| \leq \omega, \forall k \in \mathfrak{F}\}^* = \sum_{\mathcal{G} \in \mathfrak{G}} \|x_{\mathcal{G}}\|_{\infty}$   
for  $\mathbf{x} \in [-1, 1]^p$ ,  $\infty$  otherwise.

**Remark:** For hierarchical  $\mathfrak{G}_H$ , group intersection and tree sparsity models coincide.

## Beyond linear costs: Graph dispersiveness

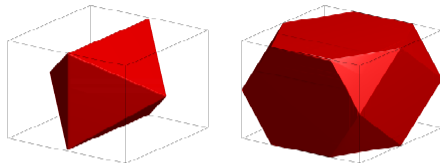


Figure: (left)  $\|\mathbf{x}\|_s^{**} = 0$  (right)  $\|\mathbf{x}\|_s^{**} \leq 1$  for  $\mathcal{E} = \{\{1, 2\}, \{2, 3\}\}$  (chain graph)

**Structure:** *We seek a signal dispersive over a given graph  $\mathcal{G}(\mathfrak{V}, \mathcal{E})$*

**Objective:**  $\mathbf{1}^T \mathbf{s} \rightarrow \sum_{(i,j) \in \mathcal{E}} s_i s_j$  (non-linear, supermodular function)

**Linearization:**

$$\|\mathbf{x}\|_s = \min_{\mathbf{z} \in \{0,1\}^{|\mathcal{E}|}} \left\{ \sum_{(i,j) \in \mathcal{E}} z_{ij} : z_{ij} \geq s_i + s_j - 1 \right\}$$

When edge-node incidence matrix of  $\mathcal{G}(\mathfrak{V}, \mathcal{E})$  is TU (e.g., bipartite graphs), it is **TU**.

## Beyond linear costs: Graph dispersiveness

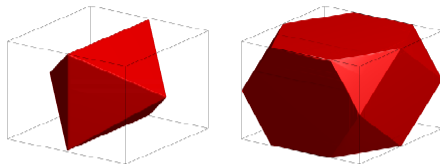


Figure: (left)  $\|\mathbf{x}\|_s^{**} = 0$  (right)  $\|\mathbf{x}\|_s^{**} \leq 1$  for  $\mathcal{E} = \{\{1, 2\}, \{2, 3\}\}$  (chain graph)

**Structure:** We seek a signal dispersive over a given graph  $\mathcal{G}(\mathfrak{V}, \mathcal{E})$

**Objective:**  $\mathbf{1}^T \mathbf{s} \rightarrow \sum_{(i,j) \in \mathcal{E}} s_i s_j$  (non-linear, supermodular function)

**Linearization:**

$$\|\mathbf{x}\|_s = \min_{\mathbf{z} \in \{0,1\}^{|\mathcal{E}|}} \left\{ \sum_{(i,j) \in \mathcal{E}} z_{ij} : z_{ij} \geq s_i + s_j - 1 \right\}$$

When edge-node incidence matrix of  $\mathcal{G}(\mathfrak{V}, \mathcal{E})$  is TU (e.g., bipartite graphs), it is TU.

**Biconjugate:**  $\|\mathbf{x}\|_s^{**} = \sum_{(i,j) \in \mathcal{E}} (|x_i| + |x_j| - 1)_+$  for  $\mathbf{x} \in [-1, 1]^p, \infty$  otherwise.

## References I

- [1] Francis Bach.  
Structured sparsity-inducing norms through submodular functions.  
*Adv. Neur. Inf. Proc. Sys. (NIPS)*, pages 118–126, 2010.
- [2] L. Baldassarre, N. Bhan, V. Cevher, and A. Kyrillidis.  
Group-sparse model selection: Hardness and relaxations.  
*arXiv preprint arXiv:1303.3207*, 2013.
- [3] R.G. Baraniuk, V. Cevher, M.F. Duarte, and C. Hegde.  
Model-based compressive sensing.  
*IEEE Trans. Inf. Theory*, 56(4):1982–2001, April 2010.
- [4] Michel Benaïm and Morris W. Hirsch.  
Asymptotic pseudotrajectories and chain recurrent flows, with applications.  
*Journal of Dynamics and Differential Equations*, 8(1):141–176, 1996.
- [5] Constantinos Daskalakis, Stratis Skoulakis, and Manolis Zampetakis.  
The complexity of constrained min-max optimization.  
*arXiv preprint arXiv:2009.09623*, 2020.



## References II

- [6] Marco F. Duarte, Dharmpal Davenport, Mark A. adn Takhar, Jason N. Laska, Ting Sun, Kevin F. Kelly, and Richard G. Baraniuk.  
Single-pixel imaging via compressive sampling.  
*IEEE Sig. Proc. Mag.*, 25(2):83–91, March 2008.
- [7] Marwa El Halabi and Volkan Cevher.  
A totally unimodular view of structured sparsity.  
*preprint*, 2014.  
arXiv:1411.1990v1 [cs.LG].
- [8] W Gerstner and W. Kistler.  
*Spiking neuron models: Single neurons, populations, plasticity*.  
Cambridge university press, 2002.
- [9] Gauthier Gidel, Hugo Berard, Gaëtan Vignoud, Pascal Vincent, and Simon Lacoste-Julien.  
A variational inequality perspective on generative adversarial networks.  
*In International Conference on Learning Representations*, 2018.
- [10] C. Hegde, M. Duarte, and V. Cevher.  
Compressive sensing recovery of spike trains using a structured sparsity model.  
*In Sig. Proc. with Adapative Sparse Struct. Rep. (SPARS)*, 2009.

## References III

- [11] Ya-Ping Hsieh.  
Convergence without convexity: Sampling, optimization, and games.  
Technical report, EPFL, 2020.
- [12] Ya-Ping Hsieh, Panayotis Mertikopoulos, and Volkan Cevher.  
The limits of min-max optimization algorithms: convergence to spurious non-critical sets.  
*arXiv preprint arXiv:2006.09065*, 2020.
- [13] J. Huang, T. Zhang, and D. Metaxas.  
Learning with structured sparsity.  
*J. Mach. Learn. Res.*, 12:3371–3412, 2011.
- [14] R. Jenatton, A. Gramfort, V. Michel, G. Obozinski, F. Bach, and B. Thirion.  
Multi-scale mining of fmri data with hierarchical structured sparsity.  
In *Pattern Recognition in NeuroImaging (PRNI)*, 2011.
- [15] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach.  
Proximal methods for hierarchical sparse coding.  
*J. Mach. Learn. Res.*, 12:2297–2334, 2011.

## References IV

- [16] Sham M Kakade, Shai Shalev-Shwartz, and Ambuj Tewari.  
Regularization techniques for learning with matrices.  
*The Journal of Machine Learning Research*, 13(1):1865–1890, 2012.
- [17] Tianyi Lin, Chi Jin, Michael Jordan, et al.  
Near-optimal algorithms for minimax optimization.  
*arXiv preprint arXiv:2002.02417*, 2020.
- [18] G. Obozinski, L. Jacob, and J.P. Vert.  
Group lasso with overlaps: The latent group lasso approach.  
*arXiv preprint arXiv:1110.0413*, 2011.
- [19] G. Obozinski, B. Taskar, and M.I. Jordan.  
Joint covariate selection and joint subspace selection for multiple classification problems.  
*Statistics and Computing*, 20(2):231–252, 2010.
- [20] Q. Tran-Dinh and V. Cevher.  
Constrained convex minimization via model-based excessive gap.  
In *Proc. the Neural Information Processing Systems Foundation conference (NIPS2014)*, pages 1–9, Montreal, Canada, December 2014.

## References V

- [21] Peng Zhao, Guilherme Rocha, and Bin Yu.  
Grouped and hierarchical model selection through composite absolute penalties.  
*Department of Statistics, UC Berkeley, Tech. Rep, 703, 2006.*
- [22] Peng Zhao and Bin Yu.  
On model selection consistency of Lasso.  
*J. Mach. Learn. Res., 7:2541–2563, 2006.*
- [23] H. Zhou, M.E. Sehl, J.S. Sinsheimer, and K. Lange.  
Association screening of common and rare genetic variants by penalized regression.  
*Bioinformatics, 26(19):2375, 2010.*