

# Learning in the Context of Set Theoretic Estimation: an Efficient and Unifying Framework for Adaptive Machine Learning and Signal Processing

Sergios Theodoridis<sup>1</sup>

a joint work with  
K. Slavakis (Univ. of Peloponnese, Greece), and  
I. Yamada (Tokyo Institute of Technology, Japan)

<sup>1</sup>University of Athens, Athens, Greece.

Vienna,  
April 11th, 2012

“ΟΥΔΕΙΣ ΑΓΕΩΜΕΤΡΗΤΟΣ ΕΙΣΙΤΩ”

“ΟΥΔΕΙΣ ΑΓΕΩΜΕΤΡΗΤΟΣ ΕΙΣΙΤΩ”

(“Those who do not know geometry are not welcome here”)

*Plato's Academy of Philosophy*

## Part A

# Outline of Part A

- The set theoretic estimation approach and multiple intersecting closed convex sets.
- The fundamental tool of metric projections in Hilbert spaces.
- Online classification and regression.
- The concept of Reproducing Kernel Hilbert Spaces (RKHS) and nonlinear processing.
- Distributive learning in sensor networks.

## Problem Definition

Given

- A set of measurements  $(\mathbf{x}_n, y_n)_{n=1}^N$ , which are jointly distributed, and

## Problem Definition

Given

- A set of measurements  $(\mathbf{x}_n, y_n)_{n=1}^N$ , which are jointly distributed, and
- A parametric set of functions  $\mathcal{F} = \{f_{\alpha}(\mathbf{x}) : \alpha \in A \subset \mathbb{R}^k\}$ .

## Problem Definition

Given

- A set of measurements  $(\mathbf{x}_n, y_n)_{n=1}^N$ , which are jointly distributed, and
- A parametric set of functions  $\mathcal{F} = \{f_\alpha(\mathbf{x}) : \alpha \in A \subset \mathbb{R}^k\}$ .

Compute an  $f(\cdot)$ , within  $\mathcal{F}$ , that best approximates  $y$  given the value of  $\mathbf{x}$ :

$$y \approx f(\mathbf{x}).$$



## Problem Definition

Given

- A set of measurements  $(\mathbf{x}_n, y_n)_{n=1}^N$ , which are jointly distributed, and
- A parametric set of functions  $\mathcal{F} = \{f_\alpha(\mathbf{x}) : \alpha \in A \subset \mathbb{R}^k\}$ .

Compute an  $f(\cdot)$ , within  $\mathcal{F}$ , that best approximates  $y$  given the value of  $\mathbf{x}$ :

$$y \approx f(\mathbf{x}).$$

## Special Cases

Smoothing, prediction, curve-fitting, regression, classification, filtering, system identification, and beamforming.

## The More Classical Approach

Select a loss function  $\mathcal{L}(\cdot, \cdot)$  and estimate  $f(\cdot)$  so that

$$f(\cdot) \in \arg \min_{f_\alpha(\cdot): \alpha \in A} \sum_{n=1}^N \mathcal{L}(y_n, f_\alpha(\mathbf{x}_n)).$$

## The More Classical Approach

Select a loss function  $\mathcal{L}(\cdot, \cdot)$  and estimate  $f(\cdot)$  so that

$$f(\cdot) \in \arg \min_{f_\alpha(\cdot): \alpha \in A} \sum_{n=1}^N \mathcal{L}(y_n, f_\alpha(\mathbf{x}_n)).$$

## Drawbacks

- Most often, in practice, the choice of the cost is dictated **not by physical reasoning** but by **computational tractability**.

## The More Classical Approach

Select a loss function  $\mathcal{L}(\cdot, \cdot)$  and estimate  $f(\cdot)$  so that

$$f(\cdot) \in \arg \min_{f_\alpha(\cdot): \alpha \in A} \sum_{n=1}^N \mathcal{L}(y_n, f_\alpha(\mathbf{x}_n)).$$

## Drawbacks

- Most often, in practice, the choice of the cost is dictated **not by physical reasoning** but by **computational tractability**.
- The existence of **a-priori information** in the form of **constraints** makes the task even more difficult.

## The More Classical Approach

Select a loss function  $\mathcal{L}(\cdot, \cdot)$  and estimate  $f(\cdot)$  so that

$$f(\cdot) \in \arg \min_{f_\alpha(\cdot): \alpha \in A} \sum_{n=1}^N \mathcal{L}(y_n, f_\alpha(\mathbf{x}_n)).$$

## Drawbacks

- Most often, in practice, the choice of the cost is dictated **not by physical reasoning** but by **computational tractability**.
- The existence of **a-priori information** in the form of **constraints** makes the task even more difficult.
- The optimization task is solved iteratively, and iterations freeze after a **finite number of steps**. Thus, the obtained solution lies in a **neighborhood** of the optimal one.

## The More Classical Approach

Select a loss function  $\mathcal{L}(\cdot, \cdot)$  and estimate  $f(\cdot)$  so that

$$f(\cdot) \in \arg \min_{f_\alpha(\cdot): \alpha \in A} \sum_{n=1}^N \mathcal{L}(y_n, f_\alpha(\mathbf{x}_n)).$$

## Drawbacks

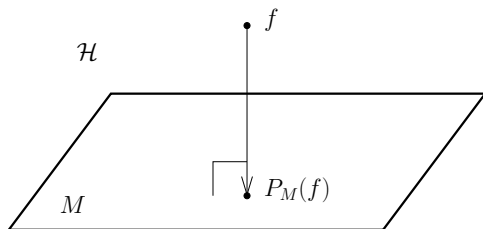
- Most often, in practice, the choice of the cost is dictated **not by physical reasoning** but by **computational tractability**.
- The existence of **a-priori information** in the form of **constraints** makes the task even more difficult.
- The optimization task is solved iteratively, and iterations freeze after a **finite number of steps**. Thus, the obtained solution lies in a **neighborhood** of the optimal one.
- The **stochastic nature** of the data and the existence of **noise** add another uncertainty to the optimality of the obtained solution.

- In this talk, we are concerned in finding a **set of solutions**, which are in **agreement** with all the available information.
- This will be achieved in the general context of
  - ▶ Set theoretic estimation.
  - ▶ Convexity.
  - ▶ Mappings or operators, e.g., projections, and their associated fixed point sets.

# Projection onto a Closed Subspace

## Theorem

Given a Euclidean  $\mathbb{R}^m$  or a Hilbert space  $\mathcal{H}$ , the projection of a point  $f$  onto a closed subspace  $M$  is the **unique** point  $P_M(f) \in M$  that lies **closest to  $f$**  (Pythagoras Theorem).





## Theorem

*Let  $C$  be a closed convex set in a Hilbert space  $\mathcal{H}$ . Then, for each  $f \in \mathcal{H}$ , there exists a **unique**  $f_* \in C$  such that*

$$\|f - f_*\| = \min_{g \in C} \|f - g\| =: d(f, C).$$

# Projection onto a Closed Convex Set

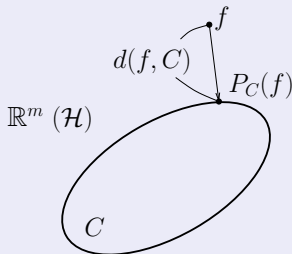
## Theorem

Let  $C$  be a closed convex set in a Hilbert space  $\mathcal{H}$ . Then, for each  $f \in \mathcal{H}$ , there exists a **unique**  $f_* \in C$  such that

$$\|f - f_*\| = \min_{g \in C} \|f - g\| =: d(f, C).$$

## Definition (Metric Projection Mapping)

The projection is the mapping  $P_C : \mathcal{H} \rightarrow C : f \mapsto P_C(f) := f_*$ .



# Projection onto a Closed Convex Set

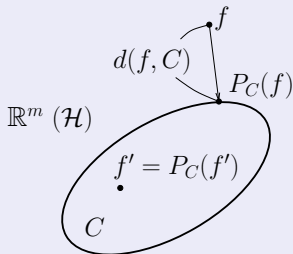
## Theorem

Let  $C$  be a closed convex set in a Hilbert space  $\mathcal{H}$ . Then, for each  $f \in \mathcal{H}$ , there exists a **unique**  $f_* \in C$  such that

$$\|f - f_*\| = \min_{g \in C} \|f - g\| =: d(f, C).$$

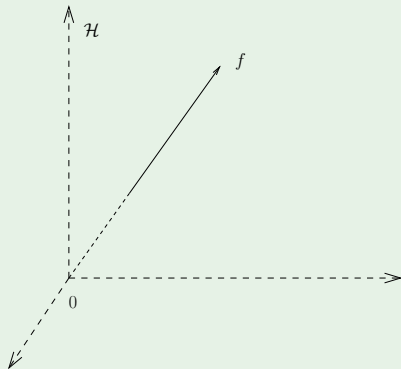
## Definition (Metric Projection Mapping)

The projection is the mapping  $P_C : \mathcal{H} \rightarrow C : f \mapsto P_C(f) := f_*$ .

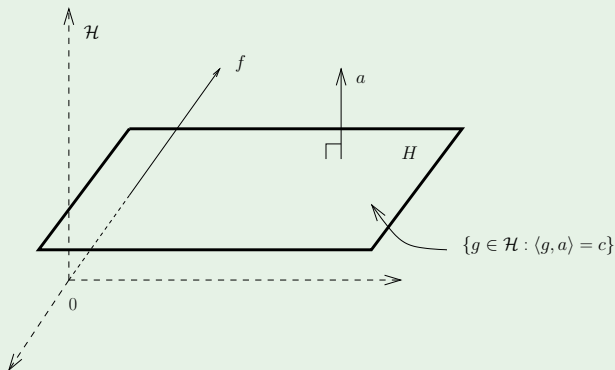


# Projection Mappings

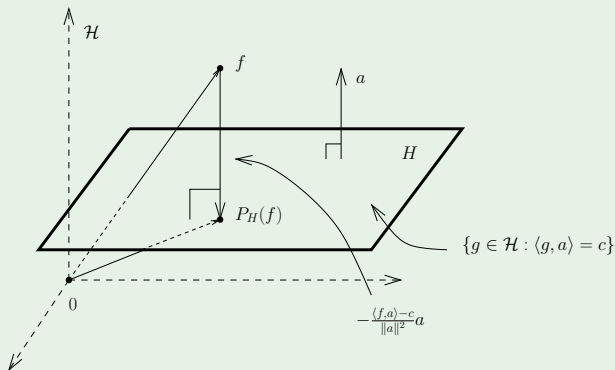
Example (Hyperplane  $H := \{g \in \mathcal{H} : \langle g, a \rangle = c\}$ )



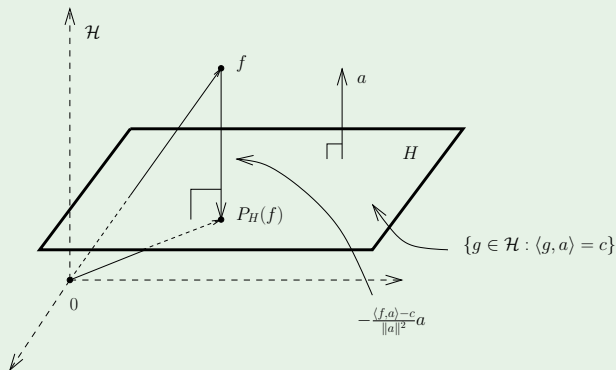
Example (Hyperplane  $H := \{g \in \mathcal{H} : \langle g, a \rangle = c\}$ )



Example (Hyperplane  $H := \{g \in \mathcal{H} : \langle g, a \rangle = c\}$ )

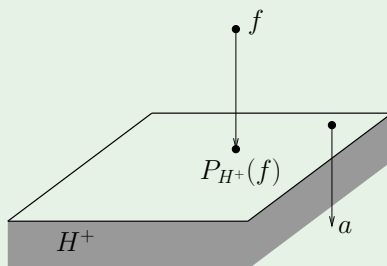


Example (Hyperplane  $H := \{g \in \mathcal{H} : \langle g, a \rangle = c\}$ )



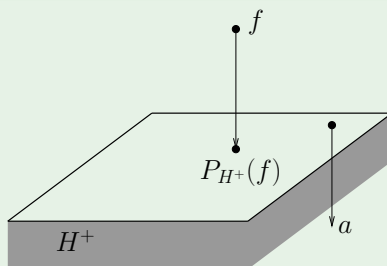
$$P_H(f) = f - \frac{\langle f, a \rangle - c}{\|a\|^2}a, \quad \forall f \in \mathcal{H}.$$

Example (Halfspace  $H^+ := \{g \in \mathcal{H} : \langle g, a \rangle \geq c\}$ )



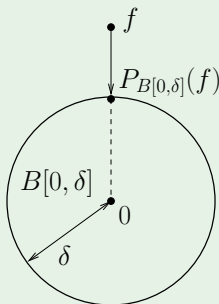


Example (Halfspace  $H^+ := \{g \in \mathcal{H} : \langle g, a \rangle \geq c\}$ )

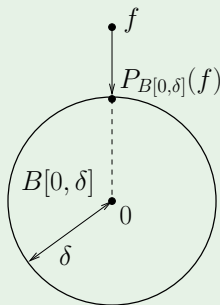


$$P_{H^+}(f) = f - \frac{\min\{0, \langle f, a \rangle - c\}}{\|a\|^2} a, \quad \forall f \in \mathcal{H}.$$

Example (Closed Ball  $B[0, \delta] := \{g \in \mathcal{H} : \|g\| \leq \delta\}$ )

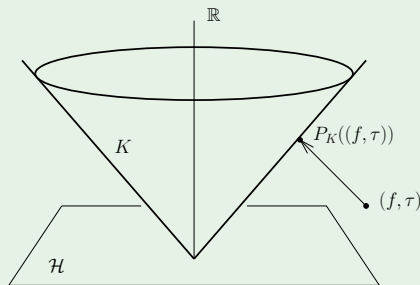


Example (Closed Ball  $B[0, \delta] := \{g \in \mathcal{H} : \|g\| \leq \delta\}$ )

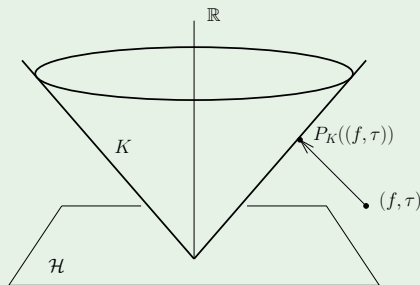


$$P_{B[0, \delta]}(f) := \frac{\delta}{\max\{\delta, \|f\|\}} f, \quad \forall f \in \mathcal{H}.$$

Example (Icecream Cone  $K := \{(f, \tau) \in \mathcal{H} \times \mathbb{R} : \|f\| \geq \tau\}$ )



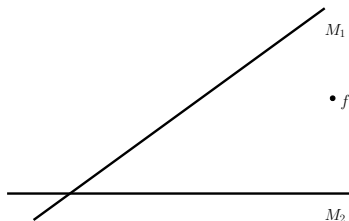
**Example (Icecream Cone**  $K := \{(f, \tau) \in \mathcal{H} \times \mathbb{R} : \|f\| \geq \tau\}$ )



$$P_K((f, \tau)) = \begin{cases} (f, \tau), & \text{if } \|f\| \leq \tau, \\ (0, 0), & \text{if } \|f\| \leq -\tau, \\ \frac{\|f\| + \tau}{2} \left( \frac{f}{\|f\|}, 1 \right), & \text{otherwise,} \end{cases} \quad \forall (f, \tau) \in \mathcal{H} \times \mathbb{R}.$$

# Alternating Projections

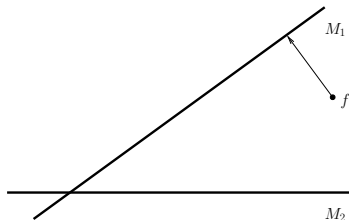
**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:



# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

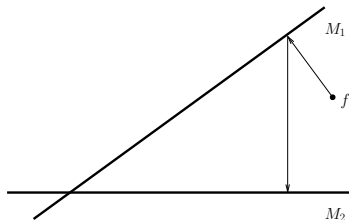
$$P_{M_1}(f).$$



# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

$$P_{M_2}P_{M_1}(f).$$

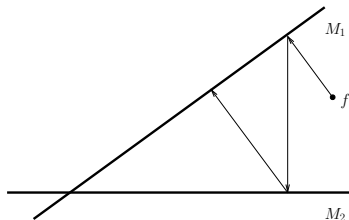




# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

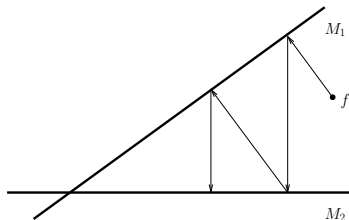
$$P_{M_1} P_{M_2} P_{M_1}(f).$$



# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

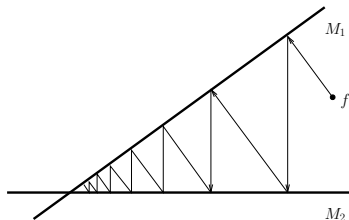
$$P_{M_2}P_{M_1}P_{M_2}P_{M_1}(f).$$



# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

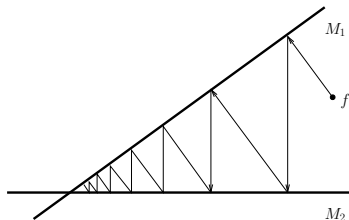
$$\cdots P_{M_2} P_{M_1} P_{M_2} P_{M_1}(f).$$



# Alternating Projections

**Composition of Projection Mappings:** Let  $M_1$  and  $M_2$  be closed subspaces in the Hilbert space  $\mathcal{H}$ . For any  $f \in \mathcal{H}$ , define the sequence of projections:

$$\cdots P_{M_2} P_{M_1} P_{M_2} P_{M_1}(f).$$



## Theorem ([von Neumann '33])

For any  $f \in \mathcal{H}$ ,  $\lim_{n \rightarrow \infty} (P_{M_2} P_{M_1})^n(f) = P_{M_1 \cap M_2}(f)$ .

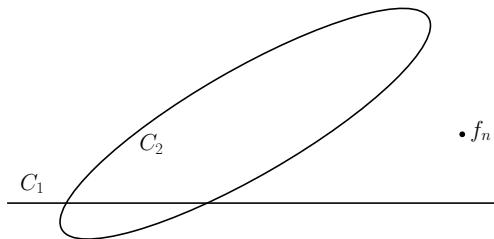
# Projections Onto Convex Sets (POCS)

## Theorem (POCS<sup>1</sup>)

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . For any  $f_0 \in \mathcal{H}$ , this defines the sequence of points

$$f_{n+1} := P_{C_p} \cdots P_{C_1}(f_n), \quad \forall n,$$

converges weakly to an  $f_* \in \bigcap_{i=1}^p C_i$ .



<sup>1</sup>[Bregman '65], [Gubin, Polyak, Raik '67].

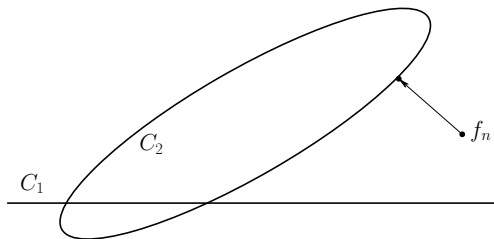
# Projections Onto Convex Sets (POCS)

## Theorem (POCS<sup>1</sup>)

Given a *finite* number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . For any  $f_0 \in \mathcal{H}$ , this defines the sequence of points

$$f_{n+1} := P_{C_p} \cdots P_{C_1}(f_n), \quad \forall n,$$

converges weakly to an  $f_* \in \bigcap_{i=1}^p C_i$ .



<sup>1</sup>[Bregman '65], [Gubin, Polyak, Raik '67].

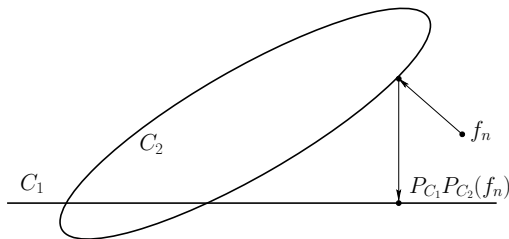
# Projections Onto Convex Sets (POCS)

## Theorem (POCS<sup>1</sup>)

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . For any  $f_0 \in \mathcal{H}$ , this defines the sequence of points

$$f_{n+1} := P_{C_p} \cdots P_{C_1}(f_n), \quad \forall n,$$

converges weakly to an  $f_* \in \bigcap_{i=1}^p C_i$ .



<sup>1</sup>[Bregman '65], [Gubin, Polyak, Raik '67].

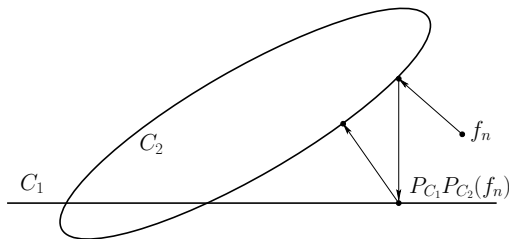
# Projections Onto Convex Sets (POCS)

## Theorem (POCS<sup>1</sup>)

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . For any  $f_0 \in \mathcal{H}$ , this defines the sequence of points

$$f_{n+1} := P_{C_p} \cdots P_{C_1}(f_n), \quad \forall n,$$

converges weakly to an  $f_* \in \bigcap_{i=1}^p C_i$ .



<sup>1</sup>[Bregman '65], [Gubin, Polyak, Raik '67].



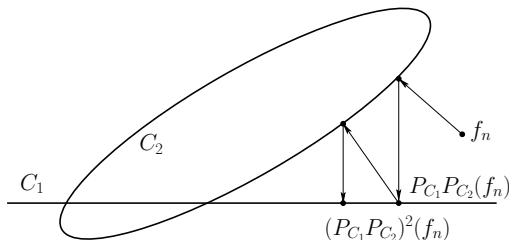
# Projections Onto Convex Sets (POCS)

## Theorem (POCS<sup>1</sup>)

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . For any  $f_0 \in \mathcal{H}$ , this defines the sequence of points

$$f_{n+1} := P_{C_p} \cdots P_{C_1}(f_n), \quad \forall n,$$

converges weakly to an  $f_* \in \bigcap_{i=1}^p C_i$ .



<sup>1</sup>[Bregman '65], [Gubin, Polyak, Raik '67].

# Extrapolated Parallel Projection Method (EPPM)

## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$

---

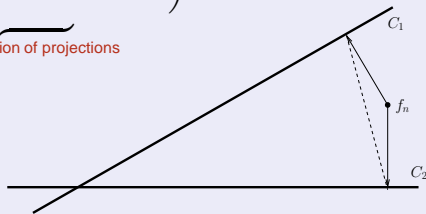
<sup>2</sup>[Pierra '84].

# Extrapolated Parallel Projection Method (EPPM)

## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$



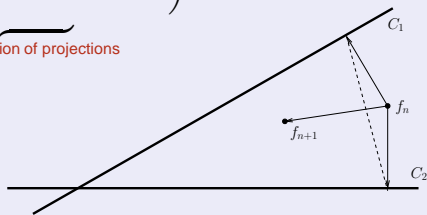
<sup>2</sup>[Pierra '84].

# Extrapolated Parallel Projection Method (EPPM)

## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$



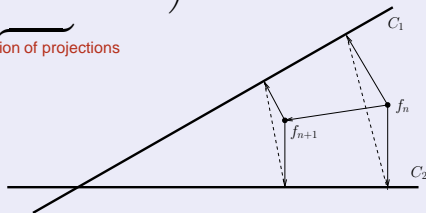
<sup>2</sup>[Pierra '84].

# Extrapolated Parallel Projection Method (EPPM)

## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$



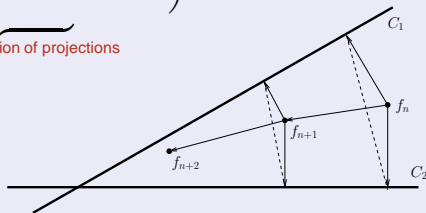
<sup>2</sup>[Pierra '84].

# Extrapolated Parallel Projection Method (EPPM)

## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$



<sup>2</sup>[Pierra '84].

# Extrapolated Parallel Projection Method (EPPM)

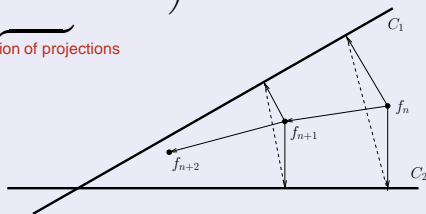
## EPPM<sup>2</sup>

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Let also a set of positive constants  $w_1, \dots, w_p$  such that  $\sum_{i=1}^p w_i = 1$ . Then for any  $f_0$ , the sequence

$$f_{n+1} = f_n + \mu_n \left( \underbrace{\sum_{i=1}^p w_i P_{C_i}(f_n)}_{\text{Convex combination of projections}} - f_n \right), \quad \forall n,$$

converges weakly to a point  $f_*$  in  $\bigcap_{i=1}^p C_i$ , where  $\mu_n \in (\epsilon, \mathcal{M}_n)$ , for  $\epsilon \in (0, 1)$ , and

$$\mathcal{M}_n := \frac{\sum_{i=1}^p w_i \|P_{C_i}(f_n) - f_n\|^2}{\|\sum_{i=1}^p w_i P_{C_i}(f_n) - f_n\|^2}.$$



<sup>2</sup>[Pierra '84].

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$

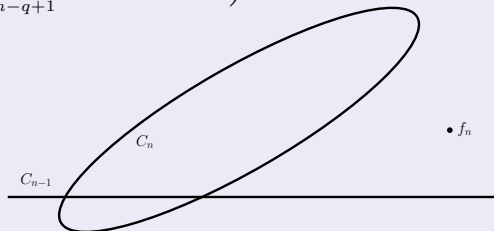
<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].



## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$

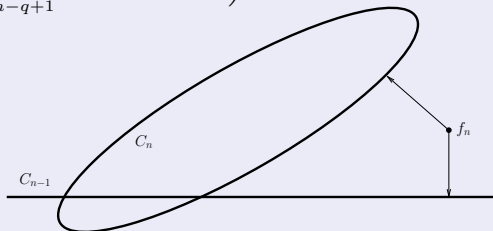


<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$

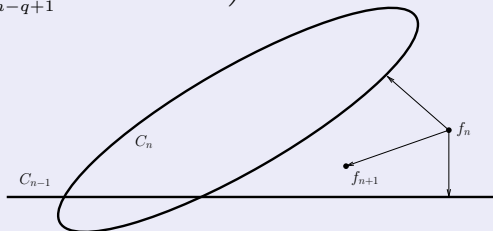


<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$



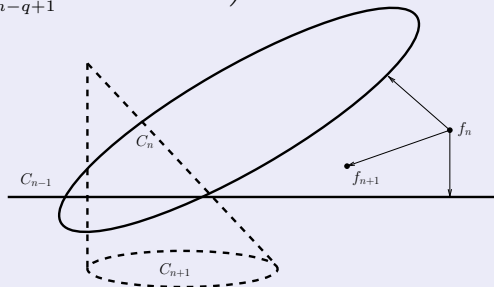
<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

# Infinite Number of Closed Convex Sets

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$



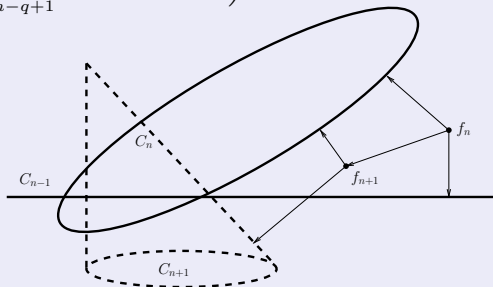
<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

# Infinite Number of Closed Convex Sets

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$



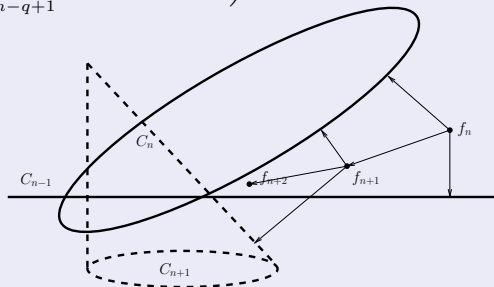
<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

# Infinite Number of Closed Convex Sets

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$



<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

# Infinite Number of Closed Convex Sets

## Adaptive Projected Subgradient Method (APSM)<sup>3</sup>

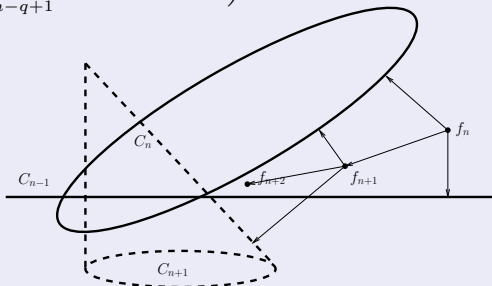
Given an **infinite number** of closed convex sets  $(C_n)_{n \geq 0}$ , let their associated projection mappings be  $(P_{C_n})_{n \geq 0}$ . For any starting point  $f_0$ , and an integer  $q > 0$ , let the sequence

$$f_{n+1} = f_n + \mu_n \left( \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right), \quad \forall n,$$

where  $\mu_n \in (0, 2\mathcal{M}_n)$ , and

$$\mathcal{M}_n := \frac{\sum_{j=n-q+1}^n w_j \|P_{C_j}(f_n) - f_n\|^2}{\left\| \sum_{j=n-q+1}^n w_j P_{C_j}(f_n) - f_n \right\|^2}.$$

Under certain constraints the above sequence converges strongly to a point  $f_* \in \text{clos}(\bigcup_{m \geq 0} \bigcap_{n \geq m} C_n)$ .



<sup>3</sup>[Yamada '03], [Yamada, Ogura '04], [Slavakis, Yamada, Ogura '06].

## The Task

Given a set of training samples  $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^m$  and a set of corresponding desired responses  $y_0, \dots, y_N$ , estimate a function  $f(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  that **fits the data**.



## The Task

Given a set of training samples  $\mathbf{x}_0, \dots, \mathbf{x}_N \in \mathbb{R}^m$  and a set of corresponding desired responses  $y_0, \dots, y_N$ , estimate a function  $f(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  that **fits the data**.

## The Expected / Empirical Risk Function approach

Estimate  $f$  so that the **expected risk** based on a loss function  $\mathcal{L}(\cdot, \cdot)$  is minimized:

$$\min_f \mathbb{E} \{ \mathcal{L}(f(\mathbf{x}), y) \},$$

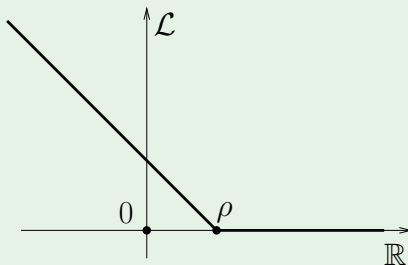
or, in practice, the **empirical risk** is minimized:

$$\min_f \sum_{n=0}^N \mathcal{L}(f(\mathbf{x}_n), y_n).$$

## Example (Classification)

For a given margin  $\rho \geq 0$ , and  $y_n \in \{+1, -1\}$ ,  $\forall n$ , define the **soft margin** loss function:

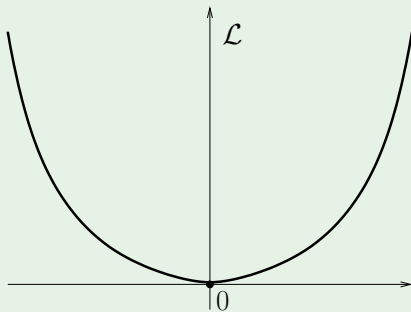
$$\mathcal{L}(f(\mathbf{x}_n), y_n) := \max\{0, \rho - y_n f(\mathbf{x}_n)\}, \quad \forall n.$$



## Example (Regression)

The square loss function:

$$\mathcal{L}(f(\mathbf{x}_n), y_n) := (y_n - f(\mathbf{x}_n))^2, \quad \forall n.$$



# The Set Theoretic Estimation Approach

## Main Idea

The goal here is to have a solution that is **in agreement with all the available information**, that resides in the data as well as in the available a-priori information.

# The Set Theoretic Estimation Approach

## Main Idea

The goal here is to have a solution that is **in agreement with all the available information**, that resides in the data as well as in the available a-priori information.

## The Means

- **Each piece of information**, associated with the training pair  $(x_n, y_n)$ , is represented in the solution space by a **set**.

# The Set Theoretic Estimation Approach

## Main Idea

The goal here is to have a solution that is **in agreement with all the available information**, that resides in the data as well as in the available a-priori information.

## The Means

- **Each piece of information**, associated with the training pair  $(x_n, y_n)$ , is represented in the solution space by a **set**.
- **Each piece of a-priori information**, i.e., each constraint, is also represented by a **set**.

# The Set Theoretic Estimation Approach

## Main Idea

The goal here is to have a solution that is **in agreement with all the available information**, that resides in the data as well as in the available a-priori information.

## The Means

- **Each piece of information**, associated with the training pair  $(x_n, y_n)$ , is represented in the solution space by a **set**.
- **Each piece of a-priori information**, i.e., each constraint, is also represented by a **set**.
- The **intersection** of all these sets constitutes the **family of solutions**.

# The Set Theoretic Estimation Approach

## Main Idea

The goal here is to have a solution that is **in agreement with all the available information**, that resides in the data as well as in the available a-priori information.

## The Means

- **Each piece of information**, associated with the training pair  $(x_n, y_n)$ , is represented in the solution space by a **set**.
- **Each piece of a-priori information**, i.e., each constraint, is also represented by a **set**.
- The **intersection** of all these sets constitutes the **family of solutions**.
- The family of solutions is known as the **feasibility set**.



That is, represent each **cost** and **constraint** by an equivalent **set**  $C_n$  and find the solution

$$f \in \bigcap_n C_n \subset \mathcal{H}.$$

# Classification: The Soft Margin Loss

## The Setting

Let the training data set  $(\mathbf{x}_n, y_n) \subset \mathbb{R}^m \times \{+1, -1\}$ ,  $n = 0, 1, \dots$   
Assume the two class task,

$$\begin{cases} y_n = +1, & \mathbf{x}_n \in W_1, \\ y_n = -1, & \mathbf{x}_n \in W_2. \end{cases}$$

Assume linear separable classes.

# Classification: The Soft Margin Loss

## The Setting

Let the training data set  $(\mathbf{x}_n, y_n) \subset \mathbb{R}^m \times \{+1, -1\}$ ,  $n = 0, 1, \dots$   
Assume the two class task,

$$\begin{cases} y_n = +1, & \mathbf{x}_n \in W_1, \\ y_n = -1, & \mathbf{x}_n \in W_2. \end{cases}$$

Assume linear separable classes.

## The Goal

# Classification: The Soft Margin Loss

## The Setting

Let the training data set  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \{+1, -1\}$ ,  $n = 0, 1, \dots$   
Assume the two class task,

$$\begin{cases} y_n = +1, & \mathbf{x}_n \in W_1, \\ y_n = -1, & \mathbf{x}_n \in W_2. \end{cases}$$

Assume linear separable classes.

## The Goal

Find  $f(\mathbf{x}) = \boldsymbol{\theta}^t \mathbf{x} + b$ , so that

$$\begin{cases} \boldsymbol{\theta}^t \mathbf{x}_n + b \geq \rho, & \text{if } y_n = +1, \\ \boldsymbol{\theta}^t \mathbf{x}_n + b \leq -\rho, & \text{if } y_n = -1. \end{cases}$$

# Classification: The Soft Margin Loss

## The Setting

Let the training data set  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \{+1, -1\}$ ,  $n = 0, 1, \dots$   
Assume the two class task,

$$\begin{cases} y_n = +1, & \mathbf{x}_n \in W_1, \\ y_n = -1, & \mathbf{x}_n \in W_2. \end{cases}$$

Assume linear separable classes.

## The Goal

Find  $f(\mathbf{x}) = \boldsymbol{\theta}^t \mathbf{x} + b$ , so that

$$\begin{cases} \boldsymbol{\theta}^t \mathbf{x}_n + b \geq \rho, & \text{if } y_n = +1, \\ \boldsymbol{\theta}^t \mathbf{x}_n + b \leq -\rho, & \text{if } y_n = -1. \end{cases} \quad \text{Hereafter, } (\boldsymbol{\theta} \leftarrow [\boldsymbol{\theta}^t], \quad \mathbf{x}_n \leftarrow [\mathbf{x}_n^t]).$$

## The Piece of Information

Find all those  $\theta$  so that  $y_n \theta^t x_n \geq \rho, \quad n = 0, 1, \dots$

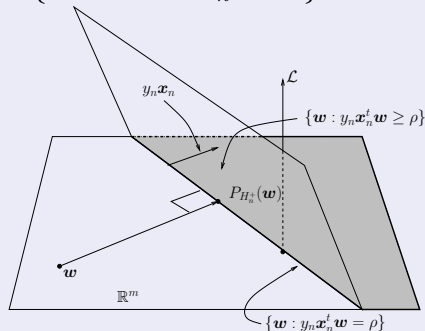
# Set Theoretic Estimation Approach to Classification

## The Piece of Information

Find all those  $\theta$  so that  $y_n \theta^t x_n \geq \rho$ ,  $n = 0, 1, \dots$

## The Equivalent Set

$$H_n^+ := \{\theta \in \mathbb{R}^m : y_n x_n^t \theta \geq \rho\}, n = 0, 1, \dots$$



## The feasibility set

For each pair  $(x_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

$$\text{find } \theta_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.



## The feasibility set

For each pair  $(x_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

$$\text{find } \theta_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\theta_0$ .

## The feasibility set

For each pair  $(x_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

$$\text{find } \theta_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\theta_0$ .
- Keep projecting as each  $H_n^+$  is formed.

## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n.$

## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.

$$\bullet \quad P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n.$$

$$\bullet \quad \boldsymbol{\theta}_{n-1}$$

## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

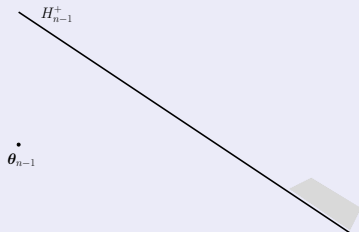
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n.$



## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

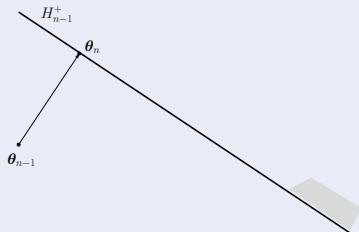
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n$ .



## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

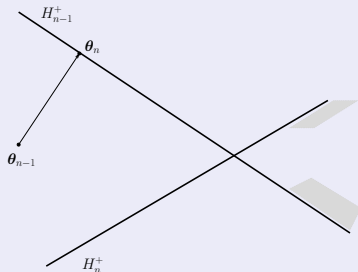
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n$ .



## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

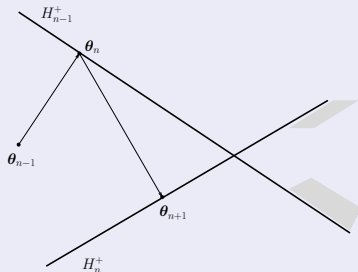
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $$P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n.$$





## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

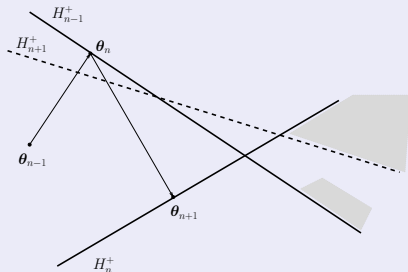
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n$ .



## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent halfspace  $H_n^+$ , and

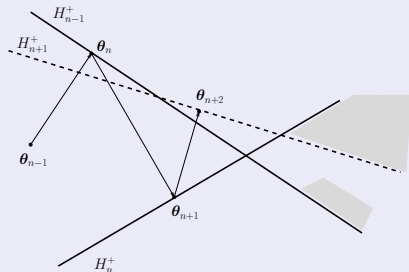
$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n H_n^+.$$

If linearly separable, the problem is feasible.

## The Algorithm

Each  $H_n^+$  is a convex set.

- Start from an arbitrary initial  $\boldsymbol{\theta}_0$ .
- Keep projecting as each  $H_n^+$  is formed.
- $P_{H_n^+}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \frac{\min\{0, \langle \boldsymbol{\theta}, y_n \mathbf{x}_n \rangle\}}{\|\mathbf{x}_n\|^2} y_n \mathbf{x}_n$ .



# Algorithmic Solution to Online Classification

$$\boldsymbol{\theta}_{n+1} := \boldsymbol{\theta}_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) \neq \boldsymbol{\theta}_n, \\ 1, & \text{otherwise.} \end{cases}$$

# Algorithmic Solution to Online Classification

$$\boldsymbol{\theta}_{n+1} := \boldsymbol{\theta}_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\boldsymbol{\theta}_n) \neq \boldsymbol{\theta}_n, \\ 1, & \text{otherwise.} \end{cases}$$

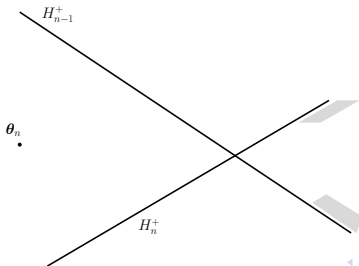
$\boldsymbol{\theta}_n$

# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

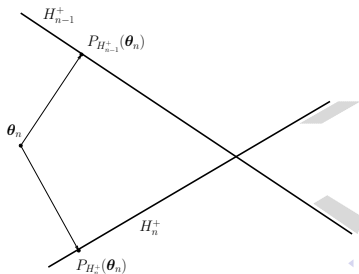


# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

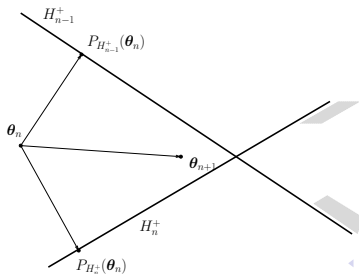


# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

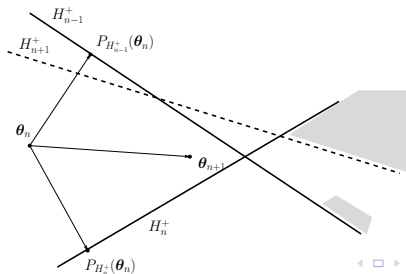


# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$



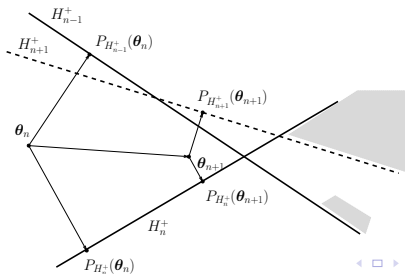


# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

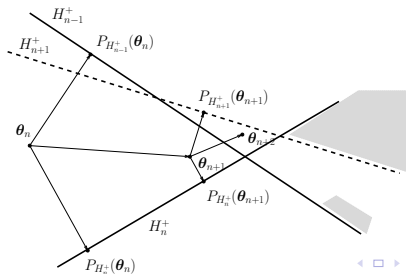


# Algorithmic Solution to Online Classification

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right),$$

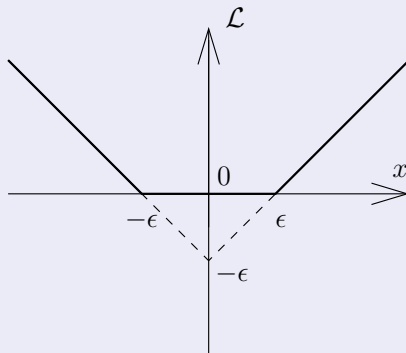
$$\mu_n \in (0, 2\mathcal{M}_n), \quad \text{and}$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{H_n^+}(\theta_n) - \theta_n\|^2}{\left\| \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) - \theta_n \right\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{H_n^+}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$



## The linear $\epsilon$ -insensitive loss function case

$$\mathcal{L}(x) := \max\{0, |x| - \epsilon\}, \quad x \in \mathbb{R}.$$



## The Piece of Information

Given  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \mathbb{R}$ , find  $\boldsymbol{\theta} \in \mathbb{R}^m$  such that

$$|\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon, \quad \forall n.$$

# Set Theoretic Estimation Approach to Regression

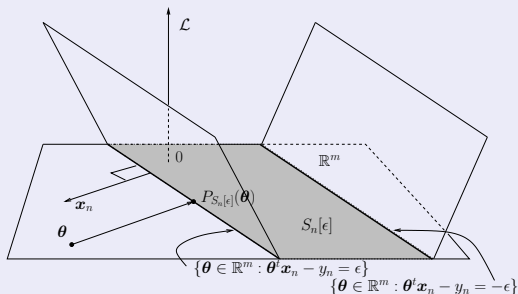
## The Piece of Information

Given  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \mathbb{R}$ , find  $\boldsymbol{\theta} \in \mathbb{R}^m$  such that

$$|\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon, \quad \forall n.$$

## The Equivalent Set (Hyperslab)

$$S_n[\epsilon] := \{\boldsymbol{\theta} \in \mathbb{R}^m : |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon\}, \quad \forall n.$$



## Projection onto a Hyperslab

$$P_{S_n[\epsilon]}(\boldsymbol{\theta}) = \boldsymbol{\theta} + \beta \mathbf{x}_n, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^m,$$

where

$$\beta := \begin{cases} \frac{y_n - \boldsymbol{\theta}^t \mathbf{x}_n - \epsilon}{\mathbf{x}_n^t \mathbf{x}_n}, & \text{if } \boldsymbol{\theta}^t \mathbf{x}_n - y_n < -\epsilon, \\ 0, & \text{if } |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon, \\ -\frac{\boldsymbol{\theta}^t \mathbf{x}_n - y_n - \epsilon}{\mathbf{x}_n^t \mathbf{x}_n}, & \text{if } \boldsymbol{\theta}^t \mathbf{x}_n - y_n > \epsilon. \end{cases}$$

## Projection onto a Hyperslab

$$P_{S_n[\epsilon]}(\boldsymbol{\theta}) = \boldsymbol{\theta} + \beta \mathbf{x}_n, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^m,$$

where

$$\beta := \begin{cases} \frac{y_n - \boldsymbol{\theta}^t \mathbf{x}_n - \epsilon}{\mathbf{x}_n^t \mathbf{x}_n}, & \text{if } \boldsymbol{\theta}^t \mathbf{x}_n - y_n < -\epsilon, \\ 0, & \text{if } |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon, \\ -\frac{\boldsymbol{\theta}^t \mathbf{x}_n - y_n - \epsilon}{\mathbf{x}_n^t \mathbf{x}_n}, & \text{if } \boldsymbol{\theta}^t \mathbf{x}_n - y_n > \epsilon. \end{cases}$$

## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent hyperslab  $S_n$ , and

$$\text{find } \boldsymbol{\theta}_* \in \bigcap_n S_n[\epsilon].$$

# Algorithm for the Online Regression

Assume weights  $\omega_j^{(n)} \geq 0$  such that  $\sum_{j=n-q+1}^n \omega_j^{(n)} = 1$ . For any  $\theta_0 \in \mathbb{R}^m$ ,

$$\theta_{n+1} := \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n \right), \quad \forall n \geq 0,$$

where the extrapolation coefficient  $\mu_n \in (0, 2\mathcal{M}_n)$  with

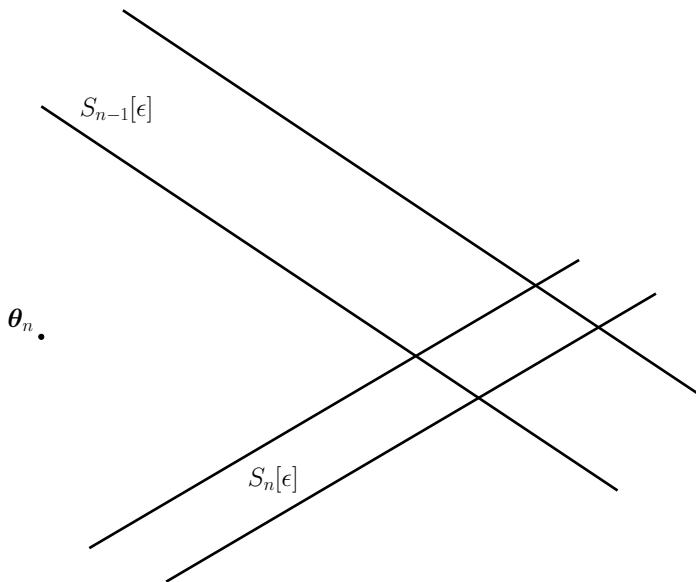
$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$



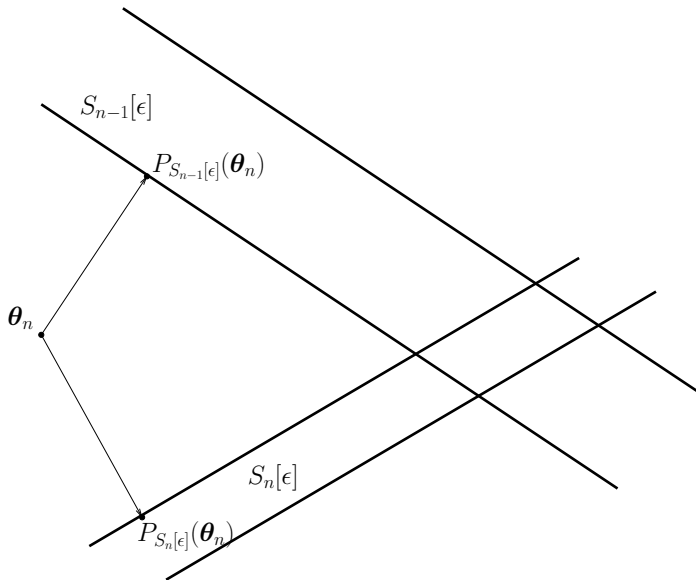
# Geometry of the Algorithm

$\theta_n \bullet$

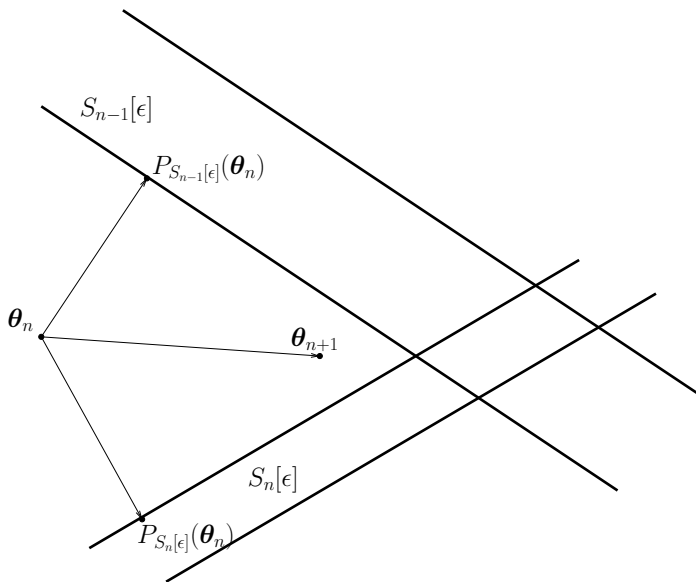
# Geometry of the Algorithm



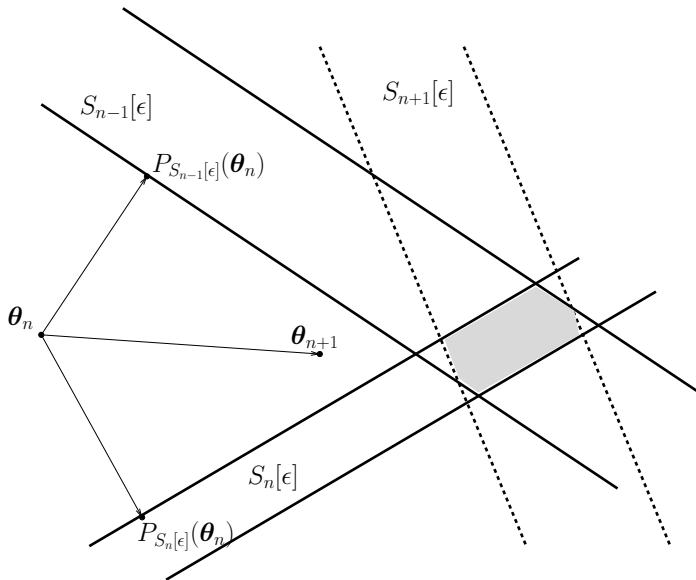
# Geometry of the Algorithm



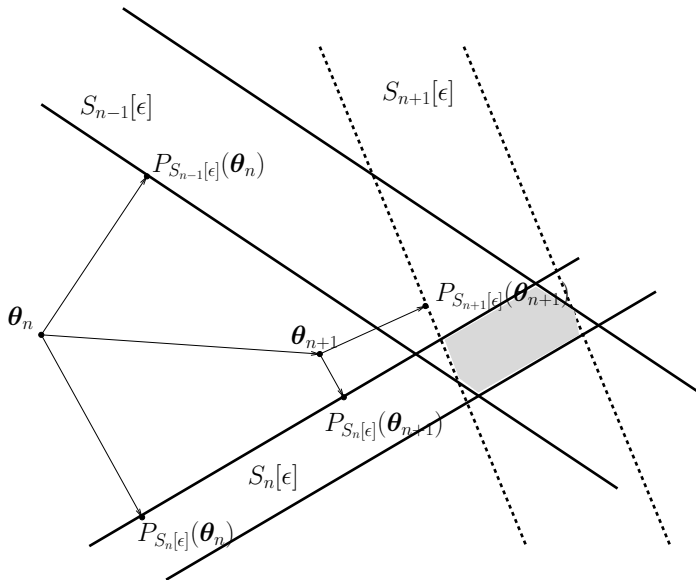
# Geometry of the Algorithm



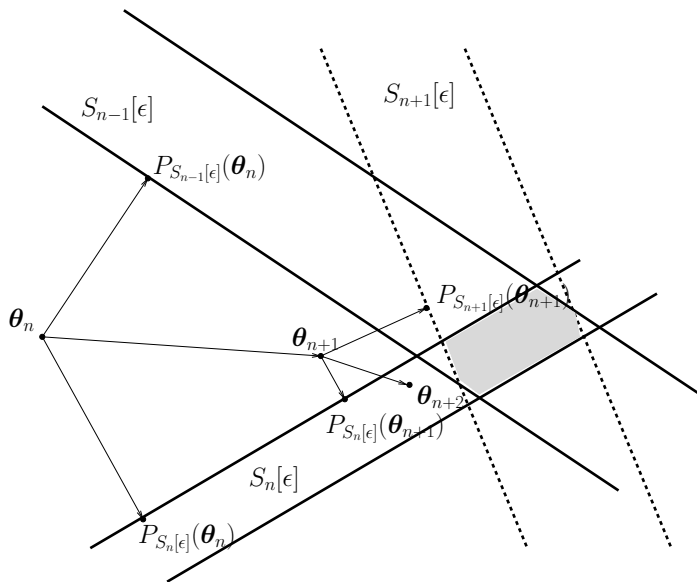
# Geometry of the Algorithm



# Geometry of the Algorithm



# Geometry of the Algorithm



## Definition

Consider a Hilbert space  $\mathcal{H}$  of functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ .



## Definition

Consider a Hilbert space  $\mathcal{H}$  of functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ .

Assume there exists a **kernel function**  $\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  such that

## Definition

Consider a Hilbert space  $\mathcal{H}$  of functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ .

Assume there exists a **kernel function**  $\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  such that

- $\kappa(\mathbf{x}, \cdot) \in \mathcal{H}, \forall \mathbf{x} \in \mathbb{R}^m,$

## Definition

Consider a Hilbert space  $\mathcal{H}$  of functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ .

Assume there exists a **kernel function**  $\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  such that

- $\kappa(\mathbf{x}, \cdot) \in \mathcal{H}, \forall \mathbf{x} \in \mathbb{R}^m$ ,
- $\langle f, \kappa(\mathbf{x}, \cdot) \rangle = f(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^m, \forall f \in \mathcal{H}$ , (**reproducing property**).

# Reproducing Kernel Hilbert Spaces (RKHS)

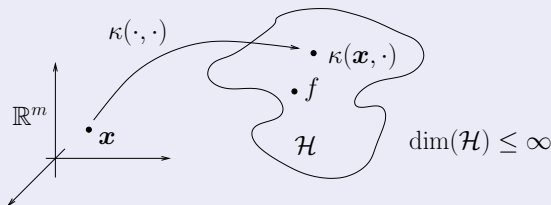
## Definition

Consider a Hilbert space  $\mathcal{H}$  of functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ .

Assume there exists a **kernel function**  $\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  such that

- $\kappa(\mathbf{x}, \cdot) \in \mathcal{H}, \forall \mathbf{x} \in \mathbb{R}^m$ ,
- $\langle f, \kappa(\mathbf{x}, \cdot) \rangle = f(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^m, \forall f \in \mathcal{H}$ , (**reproducing property**).

Then  $\mathcal{H}$  is called a Reproducing Kernel Hilbert Space (RKHS).



# Properties of the Kernel Function

- If such a kernel function exists, then it is a **symmetric and positive definite kernel**; for **any** real numbers  $a_0, a_1, \dots, a_N$ , **any**  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$ , and **any**  $N$ ,

$$\sum_{i=0}^N \sum_{j=0}^N a_i a_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

# Properties of the Kernel Function

- If such a kernel function exists, then it is a **symmetric and positive definite kernel**; for **any** real numbers  $a_0, a_1, \dots, a_N$ , **any**  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$ , and **any**  $N$ ,

$$\sum_{i=0}^N \sum_{j=0}^N a_i a_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

- The reverse is also true. Let

$$\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R},$$

be symmetric and positive definite. Then, there exists an RKHS of functions on  $\mathbb{R}^m$ , such that  $\kappa(\cdot, \cdot)$  is a reproducing kernel of  $\mathcal{H}$ .

# Properties of the Kernel Function

- If such a kernel function exists, then it is a **symmetric and positive definite kernel**; for **any** real numbers  $a_0, a_1, \dots, a_N$ , **any**  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$ , and **any**  $N$ ,

$$\sum_{i=0}^N \sum_{j=0}^N a_i a_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

- The reverse is also true. Let

$$\kappa(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R},$$

be symmetric and positive definite. Then, there exists an RKHS of functions on  $\mathbb{R}^m$ , such that  $\kappa(\cdot, \cdot)$  is a reproducing kernel of  $\mathcal{H}$ .

- Each RKHS is uniquely defined by a  $\kappa(\cdot, \cdot)$ , and each (symmetric) positive definite kernel,  $\kappa(\cdot, \cdot)$ , uniquely defines an RKHS<sup>4</sup>.

---

<sup>4</sup>[Aronszajn '50]

# Properties of the Kernel Function (cntd)

## The Kernel Trick

- The celebrated **kernel trick** is formed as follows.



# Properties of the Kernel Function (cntd)

## The Kernel Trick

- The celebrated **kernel trick** is formed as follows. Let

$$\mathbf{x} \mapsto \kappa(\mathbf{x}, \cdot) =: \phi(\mathbf{x}) \in \mathcal{H},$$

$$\mathbf{y} \mapsto \kappa(\mathbf{y}, \cdot) =: \phi(\mathbf{y}) \in \mathcal{H}.$$

- The celebrated **kernel trick** is formed as follows. Let

$$\mathbf{x} \mapsto \kappa(\mathbf{x}, \cdot) =: \phi(\mathbf{x}) \in \mathcal{H},$$

$$\mathbf{y} \mapsto \kappa(\mathbf{y}, \cdot) =: \phi(\mathbf{y}) \in \mathcal{H}.$$

Then,

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}).$$

- The celebrated **kernel trick** is formed as follows. Let

$$\mathbf{x} \mapsto \kappa(\mathbf{x}, \cdot) =: \phi(\mathbf{x}) \in \mathcal{H},$$

$$\mathbf{y} \mapsto \kappa(\mathbf{y}, \cdot) =: \phi(\mathbf{y}) \in \mathcal{H}.$$

Then,

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}).$$

- This is an important property since it leads to an easy, **black box** rule, which transforms a **nonlinear** task to a **linear** one; this is done by the following steps...

- Assume the **implicit** mapping

$$\mathbb{R}^m \ni \mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathcal{H}.$$

# Steps for Kernel Methods

- Assume the **implicit** mapping

$$\mathbb{R}^m \ni \mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathcal{H}.$$

- Solve the problem linearly in  $\mathcal{H}$ .

- Assume the **implicit** mapping

$$\mathbb{R}^m \ni \mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathcal{H}.$$

- Solve the problem linearly in  $\mathcal{H}$ .
- Use an algorithm that can be casted (modified) in terms of **inner products**.

- Assume the **implicit** mapping

$$\mathbb{R}^m \ni \mathbf{x} \mapsto \phi(\mathbf{x}) \in \mathcal{H}.$$

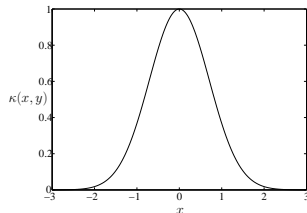
- Solve the problem linearly in  $\mathcal{H}$ .
- Use an algorithm that can be casted (modified) in terms of **inner products**.
- Replace inner product computations with kernel ones:

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}).$$

This is the step that brings the nonlinearity in the modeling.

- The Gaussian kernel:

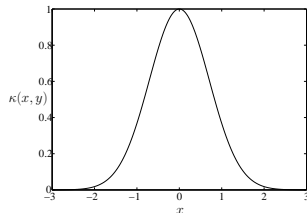
$$\kappa(\mathbf{x}, \mathbf{y}) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{\sigma^2}\right),$$





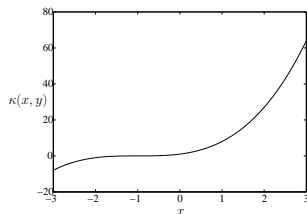
- The Gaussian kernel:

$$\kappa(\mathbf{x}, \mathbf{y}) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{\sigma^2}\right),$$



- The polynomial kernel:

$$\kappa(\mathbf{x}, \mathbf{y}) := (\mathbf{x}^t \mathbf{y} + 1)^d,$$



# The Representer Theorem

- Let a **strictly monotone increasing** function:  $\Omega : [0, \infty) \rightarrow \mathbb{R}$ ,

# The Representer Theorem

- Let a **strictly monotone increasing** function:  $\Omega : [0, \infty) \rightarrow \mathbb{R}$ ,
- and a (cost) function:  $\mathcal{L} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ .

# The Representer Theorem

- Let a **strictly monotone increasing** function:  $\Omega : [0, \infty) \rightarrow \mathbb{R}$ ,
- and a (cost) function:  $\mathcal{L} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ .
- Then, the solution of the task

$$\min_{f \in \mathcal{H}} \sum_{n=0}^N \mathcal{L}(y_n, f(\mathbf{x}_n)) + \Omega(\|f\|),$$

admits a representation of the form:

$$f_* = \sum_{n=0}^N a_n \kappa(\mathbf{x}_n, \cdot).$$

# The Representer Theorem

- Let a **strictly monotone increasing** function:  $\Omega : [0, \infty) \rightarrow \mathbb{R}$ ,
- and a (cost) function:  $\mathcal{L} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ .
- Then, the solution of the task

$$\min_{f \in \mathcal{H}} \sum_{n=0}^N \mathcal{L}(y_n, f(\mathbf{x}_n)) + \Omega(\|f\|),$$

admits a representation of the form:

$$f_* = \sum_{n=0}^N a_n \kappa(\mathbf{x}_n, \cdot).$$

## Example

$$\begin{aligned} \mathcal{L}(y_n, f(\mathbf{x}_n)) &:= (y_n - f(\mathbf{x}_n))^2, \\ \Omega(\|f\|) &:= \|f\|^2 = \langle f, f \rangle. \end{aligned}$$

## The Goal

Let the training data set  $(\mathbf{x}_n, y_n) \subset \mathbb{R}^m \times \mathbb{R}$ ,  $n = 0, 1, \dots$

- $\mathbf{x}_n \mapsto \kappa(\mathbf{x}_n, \cdot)$ , which is a function of one variable.

## The Goal

Let the training data set  $(\mathbf{x}_n, y_n) \subset \mathbb{R}^m \times \mathbb{R}$ ,  $n = 0, 1, \dots$

- $\mathbf{x}_n \mapsto \kappa(\mathbf{x}_n, \cdot)$ , which is a function of one variable.
- Find  $f \in \mathcal{H}$  such that

$$|f(\mathbf{x}_n) - y_n| \leq \epsilon, \quad \forall n.$$

## The Piece of Information

Given  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \mathbb{R}$ ,  $n = 0, 1, 2, \dots$ , find  $f \in \mathcal{H}$  such that

$$|\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n| \leq \epsilon, \quad \forall n.$$



# Set Theoretic Estimation Approach to Regression

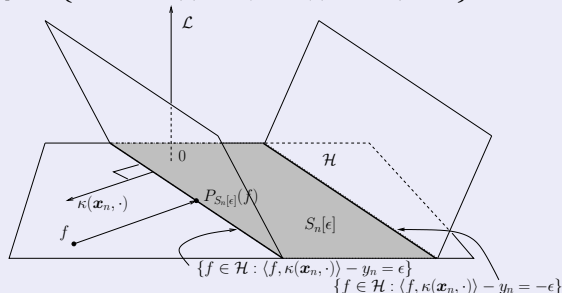
## The Piece of Information

Given  $(\mathbf{x}_n, y_n) \in \mathbb{R}^m \times \mathbb{R}$ ,  $n = 0, 1, 2, \dots$ , find  $f \in \mathcal{H}$  such that

$$|\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n| \leq \epsilon, \quad \forall n.$$

## The Equivalent Set (Hyperslab)

$$S_n[\epsilon] := \{f \in \mathcal{H} : |\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n| \leq \epsilon\}, \quad \forall n.$$



## Projection onto a Hyperslab

$$P_{S_n[\epsilon]}(f) = f + \beta \kappa(\mathbf{x}_n, \cdot), \forall f \in \mathcal{H},$$

where

$$\beta := \begin{cases} \frac{y_n - \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - \epsilon}{\kappa(\mathbf{x}_n, \mathbf{x}_n)}, & \text{if } \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n < -\epsilon, \\ 0, & \text{if } |\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n| \leq \epsilon, \\ -\frac{\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n - \epsilon}{\kappa(\mathbf{x}_n, \mathbf{x}_n)}, & \text{if } \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n > \epsilon. \end{cases}$$

## Projection onto a Hyperslab

$$P_{S_n[\epsilon]}(f) = f + \beta \kappa(\mathbf{x}_n, \cdot), \forall f \in \mathcal{H},$$

where

$$\beta := \begin{cases} \frac{y_n - \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - \epsilon}{\kappa(\mathbf{x}_n, \mathbf{x}_n)}, & \text{if } \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n < -\epsilon, \\ 0, & \text{if } |\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n| \leq \epsilon, \\ -\frac{\langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n - \epsilon}{\kappa(\mathbf{x}_n, \mathbf{x}_n)}, & \text{if } \langle f, \kappa(\mathbf{x}_n, \cdot) \rangle - y_n > \epsilon. \end{cases}$$

## The feasibility set

For each pair  $(\mathbf{x}_n, y_n)$ , form the equivalent hyperslab  $S_n$ , and

$$\text{find } f_* \in \bigcap_{n \geq n_0} S_n[\epsilon].$$

# Algorithm for Online Regression in RKHS

For  $f_0 \in \mathcal{H}$ , execute the following algorithm<sup>5</sup>

$$f_{n+1} := f_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(f_n) - f_n \right), \quad \forall n \geq 0,$$

where the extrapolation coefficient  $\mu_n \in (0, 2\mathcal{M}_n)$  with

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(f_n) - f_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(f_n) - f_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(f_n) \neq f_n, \\ 1, & \text{otherwise.} \end{cases}$$

---

<sup>5</sup>[Slavakis, Theodoridis, Yamada '09].

As time goes by:

$$f_n := \sum_{i=0}^{n-1} \gamma_i^{(n)} \kappa(\mathbf{x}_i, \cdot).$$

# Sparsification

As time goes by:

$$f_n := \sum_{i=0}^{n-1} \gamma_i^{(n)} \kappa(\mathbf{x}_i, \cdot).$$

Memory and computational load grows unbounded as  $n \rightarrow \infty$ !

# Sparsification

As time goes by:

$$f_n := \sum_{i=0}^{n-1} \gamma_i^{(n)} \kappa(\mathbf{x}_i, \cdot).$$

Memory and computational load grows unbounded as  $n \rightarrow \infty$ !

To cope with the problem, we additionally constrain the norm of  $f_n$  by a predefined  $\delta > 0$ <sup>6</sup>:

$$\forall n \geq 0, \quad f_n \in B[0, \delta] := \{f \in \mathcal{H} : \|f\| \leq \delta\} : \text{Closed Ball.}$$

---

<sup>6</sup>[Slavakis, Theodoridis, Yamada '08], [Slavakis, Theodoridis '08].

As time goes by:

$$f_n := \sum_{i=0}^{n-1} \gamma_i^{(n)} \kappa(\mathbf{x}_i, \cdot).$$

Memory and computational load grows unbounded as  $n \rightarrow \infty$ !

To cope with the problem, we additionally constrain the norm of  $f_n$  by a predefined  $\delta > 0^6$ :

$$\forall n \geq 0, \quad f_n \in B[0, \delta] := \{f \in \mathcal{H} : \|f\| \leq \delta\} : \text{Closed Ball.}$$

## Goal

Thus, we are looking for a classifier  $f \in \mathcal{H}$  such that

$$f \in B[0, \delta] \cap \left( \bigcap_{n \geq n_0} S_n[\epsilon] \right).$$

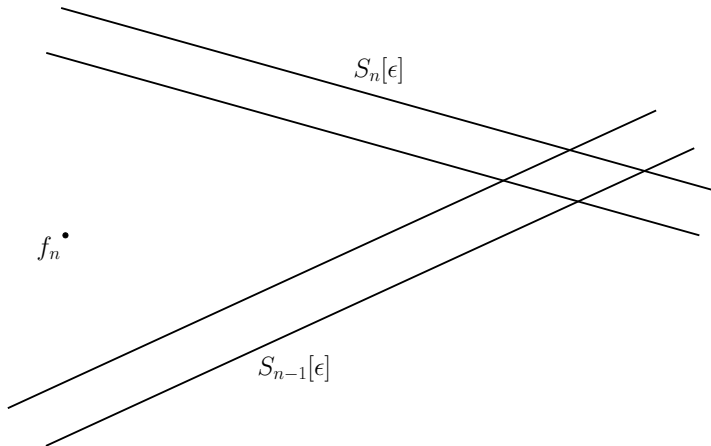
<sup>6</sup>[Slavakis, Theodoridis, Yamada '08], [Slavakis, Theodoridis '08].



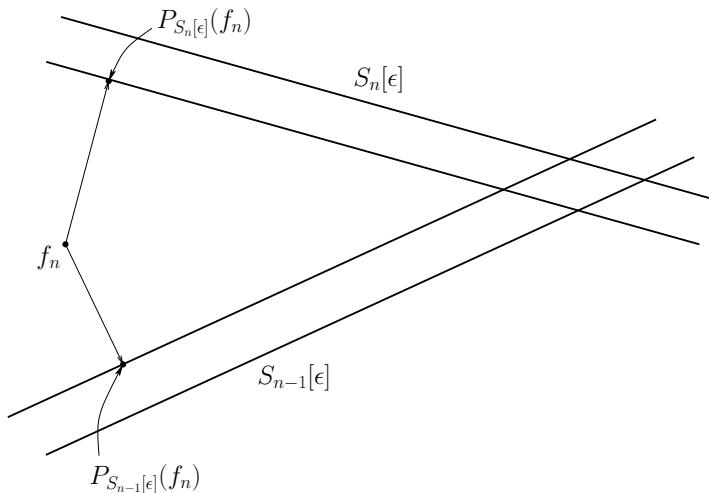
# Geometric Illustration of the Algorithm

$f_n^\bullet$

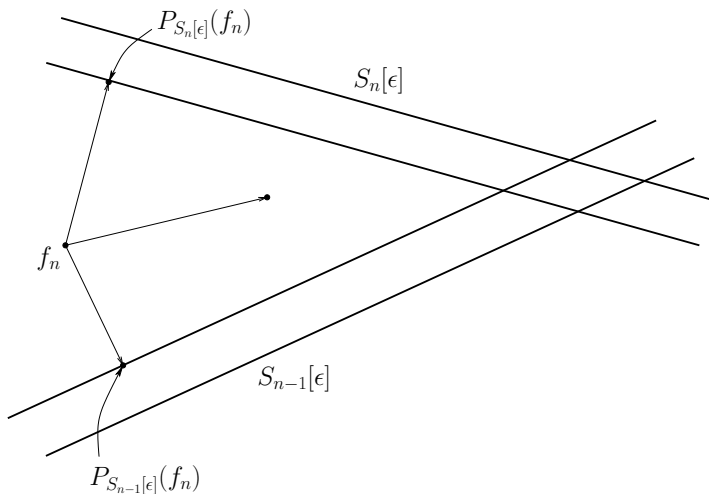
# Geometric Illustration of the Algorithm



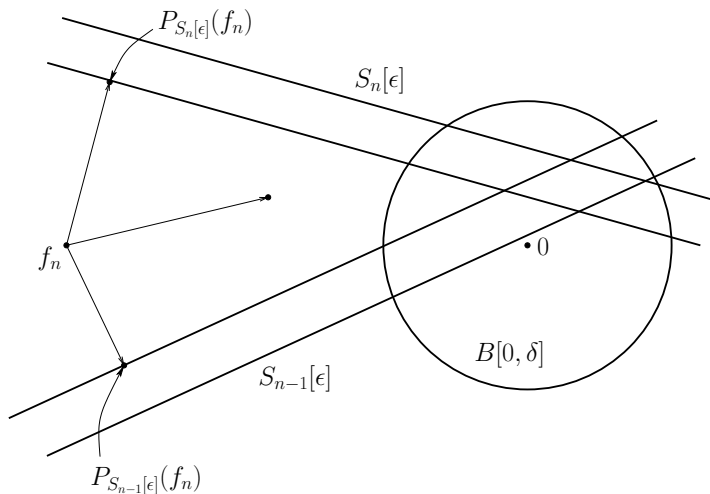
# Geometric Illustration of the Algorithm



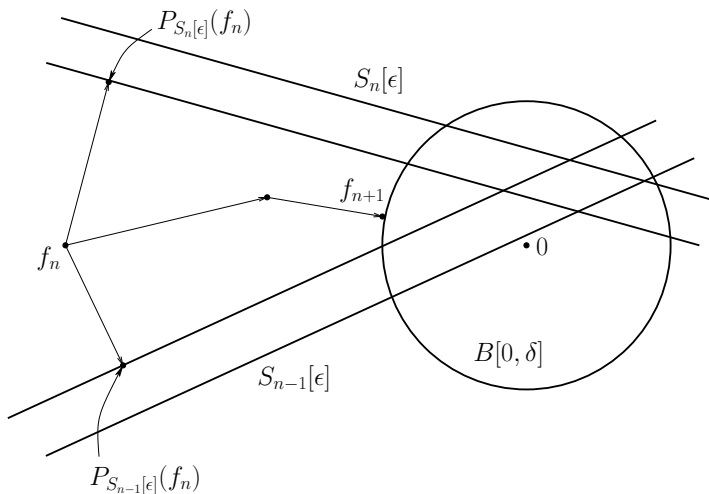
# Geometric Illustration of the Algorithm



# Geometric Illustration of the Algorithm



# Geometric Illustration of the Algorithm



## Problem Definition

- In a distributed network, the nodes are tasked to collect information and estimate a parameter of interest. The general concept can be summarized as follows:

## Problem Definition

- In a distributed network, the nodes are tasked to collect information and estimate a parameter of interest. The general concept can be summarized as follows:
  - ▶ The nodes sense an amount of data from the environment.



## Problem Definition

- In a distributed network, the nodes are tasked to collect information and estimate a parameter of interest. The general concept can be summarized as follows:
  - ▶ The nodes sense an amount of data from the environment.
  - ▶ Computations are performed **locally** in each node.

## Problem Definition

- In a distributed network, the nodes are tasked to collect information and estimate a parameter of interest. The general concept can be summarized as follows:
  - ▶ The nodes sense an amount of data from the environment.
  - ▶ Computations are performed **locally** in each node.
  - ▶ Each node transmits the **locally** obtained estimate to a **neighborhood** of nodes.

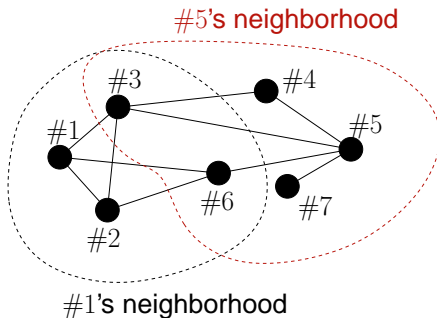
## Problem Definition

- In a distributed network, the nodes are tasked to collect information and estimate a parameter of interest. The general concept can be summarized as follows:
  - ▶ The nodes sense an amount of data from the environment.
  - ▶ Computations are performed **locally** in each node.
  - ▶ Each node transmits the **locally** obtained estimate to a **neighborhood** of nodes.

The goal is to drive the **locally** computed estimates to converge to the **same** value. This is known as **consensus**.

# The Diffusion Topology

- The most commonly used topology is the **diffusion** network:



## Problem Formulation

- Let a node set denoted as  $\mathcal{N} := \{1, 2, \dots, N\}$  and **each node**,  $k$ , **at time**,  $n$ , has access to the measurements

$$y_k(n) \in \mathbb{R}, \quad \mathbf{x}_{k,n} \in \mathbb{R}^m,$$

## Problem Formulation

- Let a node set denoted as  $\mathcal{N} := \{1, 2, \dots, N\}$  and **each node**,  $k$ , **at time**,  $n$ , has access to the measurements

$$y_k(n) \in \mathbb{R}, \quad \mathbf{x}_{k,n} \in \mathbb{R}^m,$$

we assume that there exists a **linear system**,  $\boldsymbol{\theta}_*$ , such that

$$y_k(n) = \mathbf{x}_{k,n}^t \boldsymbol{\theta}_* + v_k(n),$$

where  $v_k(n)$  is the noise.

## Problem Formulation

- Let a node set denoted as  $\mathcal{N} := \{1, 2, \dots, N\}$  and **each node**,  $k$ , **at time**,  $n$ , has access to the measurements

$$y_k(n) \in \mathbb{R}, \quad \mathbf{x}_{k,n} \in \mathbb{R}^m,$$

we assume that there exists a **linear system**,  $\boldsymbol{\theta}_*$ , such that

$$y_k(n) = \mathbf{x}_{k,n}^t \boldsymbol{\theta}_* + v_k(n),$$

where  $v_k(n)$  is the noise.

The task is to estimate the **common**  $\boldsymbol{\theta}_*$ .

- **Combine** estimates received from the neighborhood  $\mathcal{N}_k$ :

$$\phi_k(n) := \sum_{l \in \mathcal{N}_k} c_{k,l}(n+1) \theta_l(n).$$



- **Combine** estimates received from the neighborhood  $\mathcal{N}_k$ :

$$\phi_k(n) := \sum_{l \in \mathcal{N}_k} c_{k,l}(n+1) \theta_l(n).$$

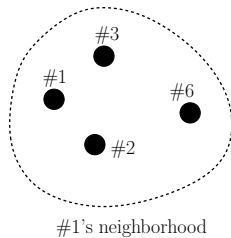
- Perform the **adaptation** step<sup>7</sup>:

$$\theta_k(n+1) := \phi_k(n) + \mu_k(n+1) \left( \sum_{j=n-q+1}^n \omega_{k,j} P_{S_{k,j}}(\phi_k(n)) - \phi_k(n) \right).$$

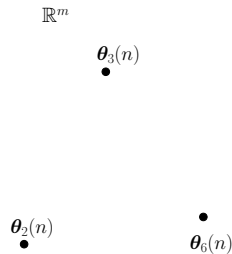
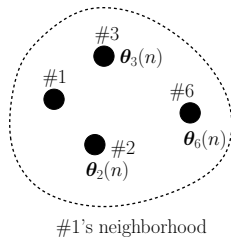
---

<sup>7</sup>[Chouvardas, Slavakis, Theodoridis, '11].

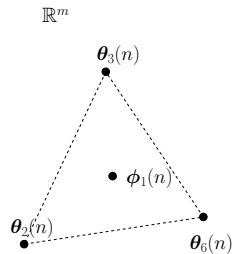
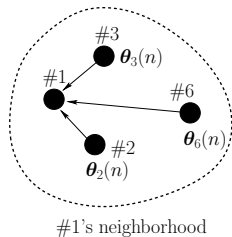
# The Geometry of the Algorithm



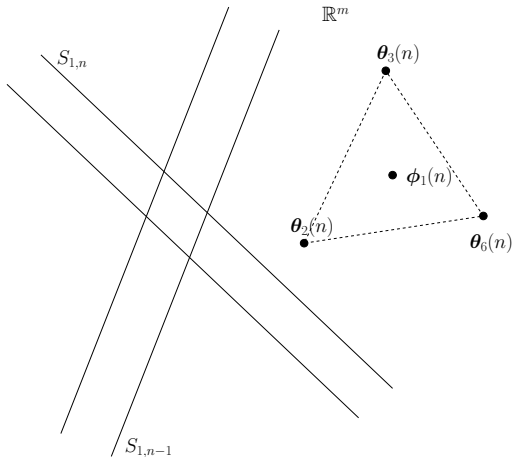
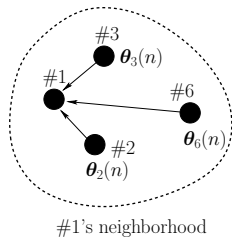
# The Geometry of the Algorithm



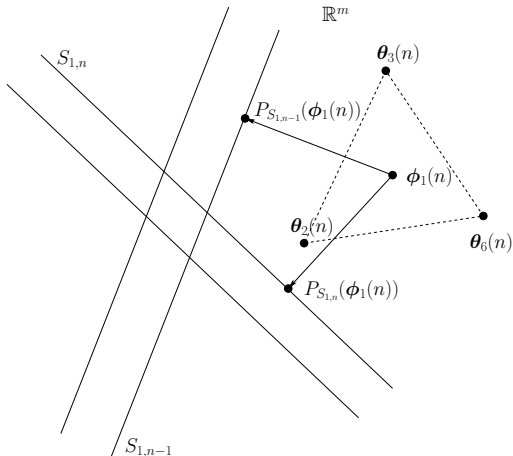
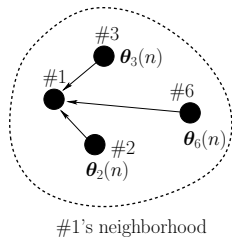
# The Geometry of the Algorithm



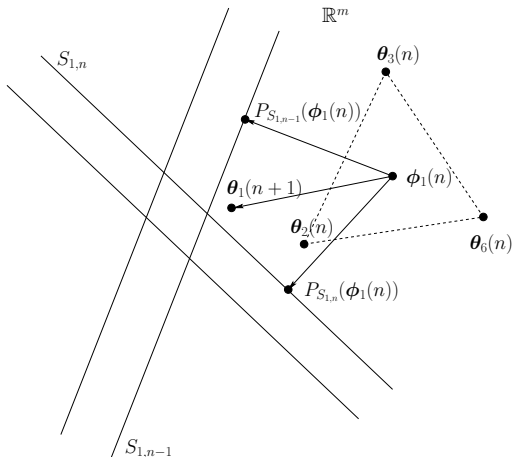
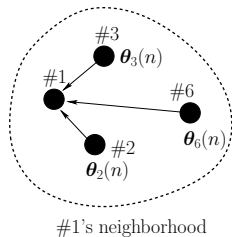
# The Geometry of the Algorithm



# The Geometry of the Algorithm



# The Geometry of the Algorithm



## Part B



- Incorporate a-priori information into our algorithmic framework.

- Incorporate a-priori information into our algorithmic framework.
- An operator theoretic approach will be followed.

- Incorporate a-priori information into our algorithmic framework.
- An operator theoretic approach will be followed.
- Such an approach will be illustrated through two paradigms:

- Incorporate a-priori information into our algorithmic framework.
- An operator theoretic approach will be followed.
- Such an approach will be illustrated through two paradigms:
  - ▶ Beamforming task.

- Incorporate a-priori information into our algorithmic framework.
- An operator theoretic approach will be followed.
- Such an approach will be illustrated through two paradigms:
  - ▶ Beamforming task.
  - ▶ Sparsity-aware learning problem.

- Incorporate a-priori information into our algorithmic framework.
- An operator theoretic approach will be followed.
- Such an approach will be illustrated through two paradigms:
  - ▶ Beamforming task.
  - ▶ Sparsity-aware learning problem.
- Our objective is to show that a large variety of constrained online learning tasks can be unified under a common umbrella; the **Adaptive Projected Subgradient Method (APSM)**.

# The Underlying Concepts

## A Mapping and its Fixed Point Set

- A **mapping** defined in a Hilbert space  $\mathcal{H}$ :

$$T : \mathcal{H} \rightarrow \mathcal{H}.$$

# The Underlying Concepts

## A Mapping and its Fixed Point Set

- A **mapping** defined in a Hilbert space  $\mathcal{H}$ :

$$T : \mathcal{H} \rightarrow \mathcal{H}.$$

- Given a mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , its **fixed point set** is defined as

$$\text{Fix}(T) := \{f \in \mathcal{H} : T(f) = f\},$$

i.e., all those points which are **unaffected** by  $T$ .



# The Underlying Concepts

## A Mapping and its Fixed Point Set

- A **mapping** defined in a Hilbert space  $\mathcal{H}$ :

$$T : \mathcal{H} \rightarrow \mathcal{H}.$$

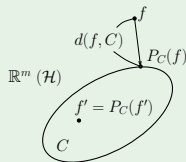
- Given a mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , its **fixed point set** is defined as

$$\text{Fix}(T) := \{f \in \mathcal{H} : T(f) = f\},$$

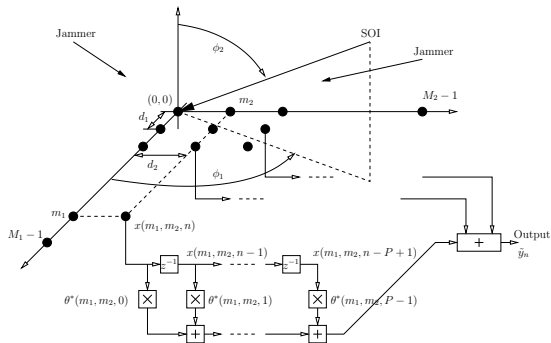
i.e., all those points which are **unaffected** by  $T$ .

### Example

If  $C$  is a closed convex set in  $\mathcal{H}$ , then  $\text{Fix}(P_C) = C$ .

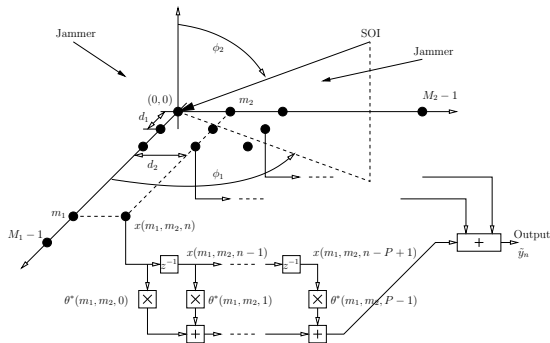


- Antenna arrays are vastly utilized for space-time filtering:



# Beamforming

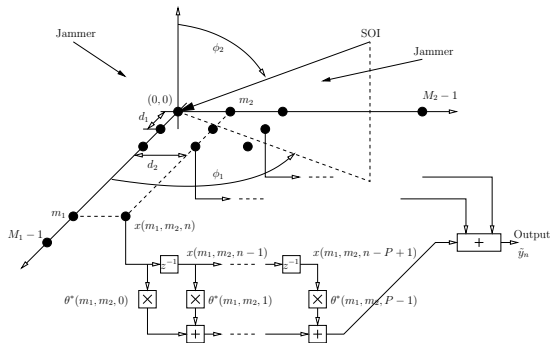
- Antenna arrays are vastly utilized for space-time filtering:



- The superscript  $*$  stands for complex conjugation.

# Beamforming

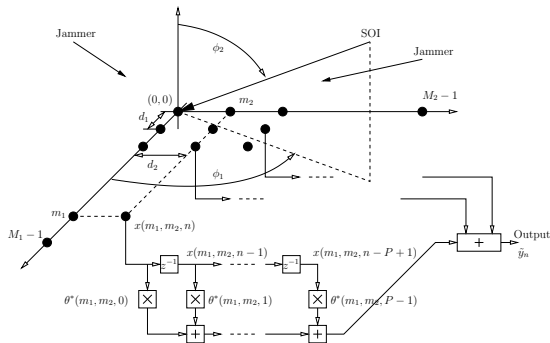
- Antenna arrays are vastly utilized for space-time filtering:



- ▶ The superscript  $*$  stands for complex conjugation.
- ▶ SOI: Signal Of Interest.

# Beamforming

- Antenna arrays are vastly utilized for space-time filtering:



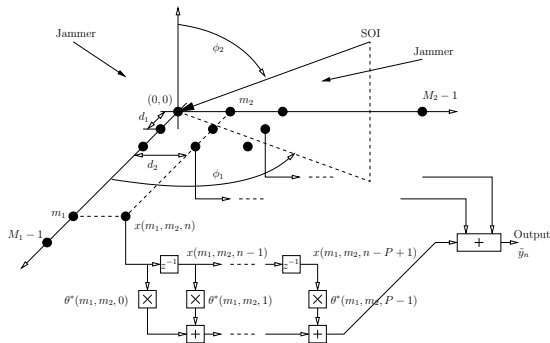
- ▶ The superscript  $*$  stands for complex conjugation.
- ▶ SOI: Signal Of Interest.

- After some re-arrangements, the output of the array is given by

$$\tilde{y}_n := \boldsymbol{\theta}^t \mathbf{x}_n, \quad n = 0, 1, 2, \dots$$

# Beamforming

- Antenna arrays are vastly utilized for space-time filtering:



- ▶ The superscript  $*$  stands for complex conjugation.
  - ▶ SOI: Signal Of Interest.
- After some re-arrangements, the output of the array is given by

$$\tilde{y}_n := \boldsymbol{\theta}^t \mathbf{x}_n, \quad n = 0, 1, 2, \dots$$

The **beamformer** is the vector  $\boldsymbol{\theta}$ .

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

## A-priori information



## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:
  - ▶ Knowledge of the approximate location of the SOIs and jammers.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:
  - ▶ Knowledge of the approximate location of the SOIs and jammers.
  - ▶ Array calibration errors.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:
  - ▶ Knowledge of the approximate location of the SOIs and jammers.
  - ▶ Array calibration errors.
  - ▶ Inoperative array elements.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:
  - ▶ Knowledge of the approximate location of the SOIs and jammers.
  - ▶ Array calibration errors.
  - ▶ Inoperative array elements.
  - ▶ Bounds on the weights of the array elements.

## The Goal of Beamforming

By utilizing all the available **a-priori knowledge**, reconstruct the SOIs, while, in the meantime, suppress the jamming signals.

### A-priori information

- Known locations of the SOIs and/or the jammers.
- Robustness against erroneous information and array imperfections:
  - ▶ Knowledge of the approximate location of the SOIs and jammers.
  - ▶ Array calibration errors.
  - ▶ Inoperative array elements.
  - ▶ Bounds on the weights of the array elements.

Given the previous a-priori info, and the set of data  $(y_n, \mathbf{x}_n)$ ,  $n = 0, 1, 2, \dots$ , compute  $\theta$  such that

$$\theta^t \mathbf{x}_n \approx y_n, \quad \forall n.$$

## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if



## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if

- only the source of interest transmits a signal of value 1,

## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if

- only the source of interest transmits a signal of value 1,
- and there is no noise in the system.

## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if

- only the source of interest transmits a signal of value 1,
- and there is no noise in the system.

**Remark:** The steering vector comprises information like the location of the associated source, and the geometry of the array.

## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if

- only the source of interest transmits a signal of value 1,
- and there is no noise in the system.

**Remark:** The steering vector comprises information like the location of the associated source, and the geometry of the array.

## Distortionless constraint

If  $\mathbf{s}_{\text{SOI}}$  is the steering vector associated to a SOI, then we would like to have:

$$\mathbf{s}_{\text{SOI}}^t \boldsymbol{\theta} = 1.$$

# Distortionless and Null Constraints

## Definition (Steering vector)

Each transmitting source is associated to a **steering vector**,  $\mathbf{s}$ , defined as the vector which collects all the signal values in the array if

- only the source of interest transmits a signal of value 1,
- and there is no noise in the system.

**Remark:** The steering vector comprises information like the location of the associated source, and the geometry of the array.

## Distortionless constraint

If  $\mathbf{s}_{\text{SOI}}$  is the steering vector associated to a SOI, then we would like to have:

$$\mathbf{s}_{\text{SOI}}^t \boldsymbol{\theta} = 1.$$

## Nulls

If  $\mathbf{s}_{\text{jam}}$  is the steering vector associated to a jammer, then we would like to have:

$$\mathbf{s}_{\text{jam}}^t \boldsymbol{\theta} = 0.$$

# Affinely Constrained Beamforming

A large variety of a-priori knowledge in beamforming problems can be cast by means of affine constraints; given a matrix  $C$  and a vector  $g$ :

$$C^t \theta = g.$$

# Affinely Constrained Beamforming

A large variety of a-priori knowledge in beamforming problems can be cast by means of affine constraints; given a matrix  $C$  and a vector  $g$ :

$$C^t \theta = g.$$

## Example

Let  $C := [s_{\text{SOL}}, s_{\text{jam}}]$ , and  $g := [1, 0]^t$ .

# Affinely Constrained Beamforming

A large variety of a-priori knowledge in beamforming problems can be cast by means of affine constraints; given a matrix  $C$  and a vector  $g$ :

$$C^t \theta = g.$$

## Example

Let  $C := [s_{\text{SOL}}, s_{\text{jam}}]$ , and  $g := [1, 0]^t$ .

Define the following **affine** set  $V := \arg \min_{\theta \in \mathbb{R}^m} \|C^t \theta - g\|$ ,



# Affinely Constrained Beamforming

A large variety of a-priori knowledge in beamforming problems can be cast by means of affine constraints; given a matrix  $C$  and a vector  $g$ :

$$C^t \theta = g.$$

## Example

Let  $C := [s_{\text{SOL}}, s_{\text{jam}}]$ , and  $g := [1, 0]^t$ .

Define the following **affine** set  $V := \arg \min_{\theta \in \mathbb{R}^m} \|C^t \theta - g\|$ , which contains, in general, an infinite number of points, and covers also the case of **inconsistent** a-priori constraints, i.e., the case:

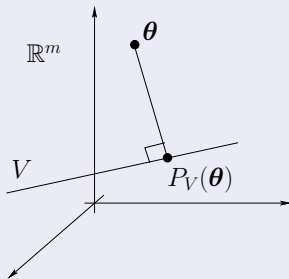
$$\forall \theta, \quad C^t \theta \neq g.$$

## Projection onto the affine set $V$

Given  $V := \arg \min_{\theta \in \mathbb{R}^m} \|C^t \theta - g\|$ , the metric projection mapping onto  $V$  is given by

$$P_V(\theta) = \theta - C^{t\dagger}(C^t \theta - g), \quad \forall \theta \in \mathbb{R}^m,$$

where  $(\cdot)^\dagger$  denotes the Moore-Penrose pseudoinverse of a matrix.



# Affinely Constrained Algorithm

- At time  $n$ , given the training data  $(y_n, \mathbf{x}_n)$ , define the hyperslab:

$$S_n[\epsilon] := \{\boldsymbol{\theta} \in \mathbb{R}^m : |\mathbf{x}_n^t \boldsymbol{\theta} - y_n| \leq \epsilon\}.$$

# Affinely Constrained Algorithm

- At time  $n$ , given the training data  $(y_n, \mathbf{x}_n)$ , define the hyperslab:

$$S_n[\epsilon] := \{\boldsymbol{\theta} \in \mathbb{R}^m : |\mathbf{x}_n^t \boldsymbol{\theta} - y_n| \leq \epsilon\}.$$

- For any initial point  $\boldsymbol{\theta}_0$ , and  $\forall n$ ,

$$\boldsymbol{\theta}_{n+1} := P_V \left( \boldsymbol{\theta}_n + \mu_n \left( \sum_{i=n-q+1}^n \omega_i^{(n)} P_{S_i[\epsilon]}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right) \right),$$

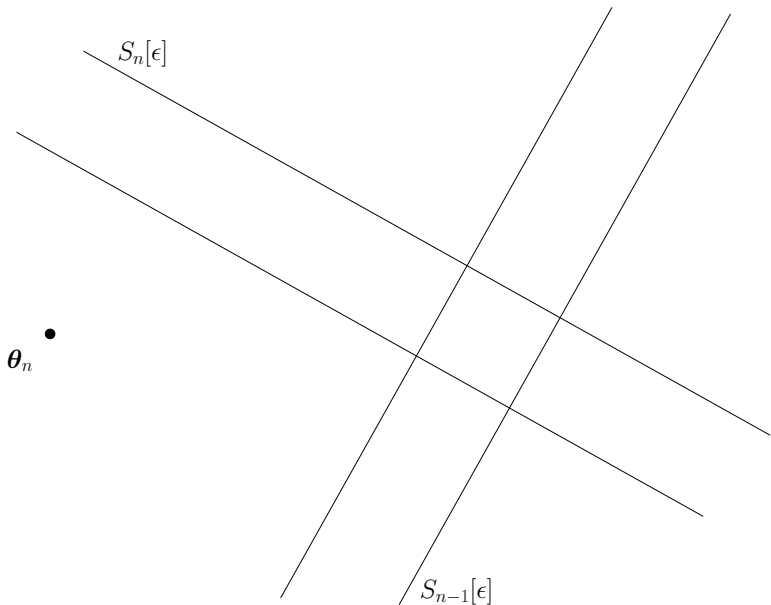
$$\mu_n \in (0, 2\mathcal{M}_n),$$

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\boldsymbol{\theta}_n) \neq \boldsymbol{\theta}_n, \\ 1, & \text{otherwise.} \end{cases}$$

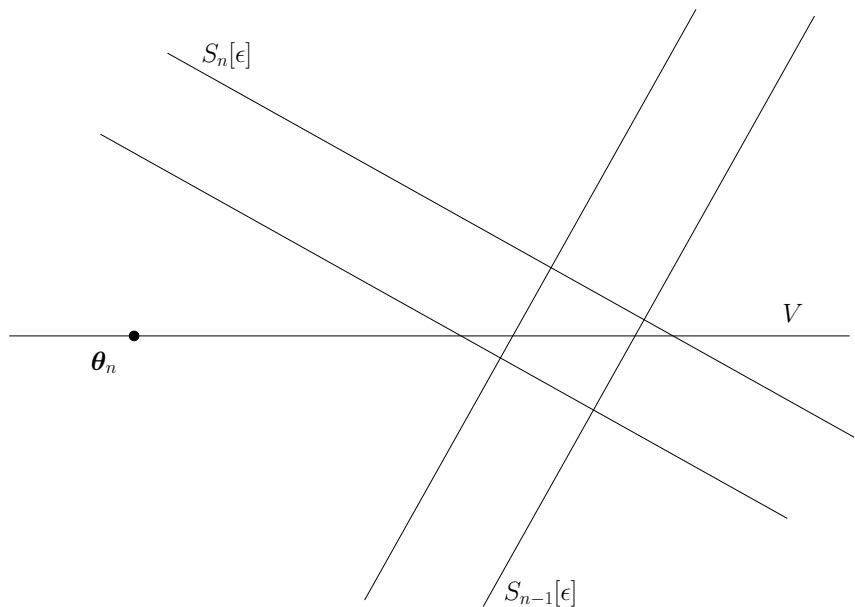
# Geometry of the Algorithm

$\theta_n$  •

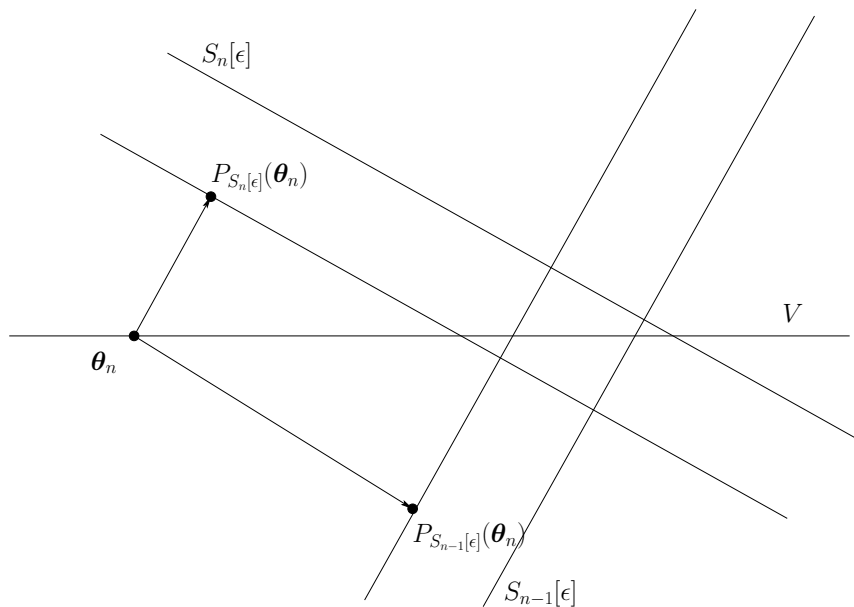
# Geometry of the Algorithm



# Geometry of the Algorithm

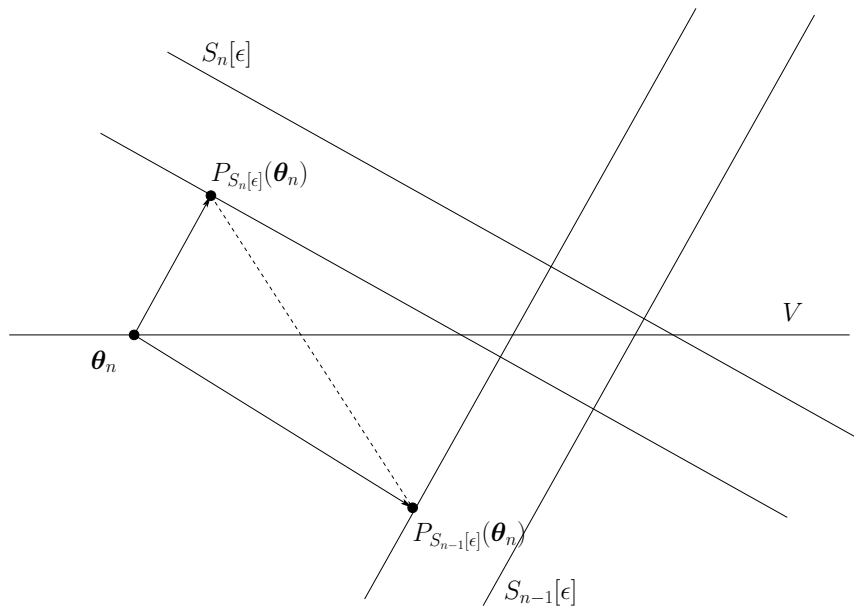


# Geometry of the Algorithm

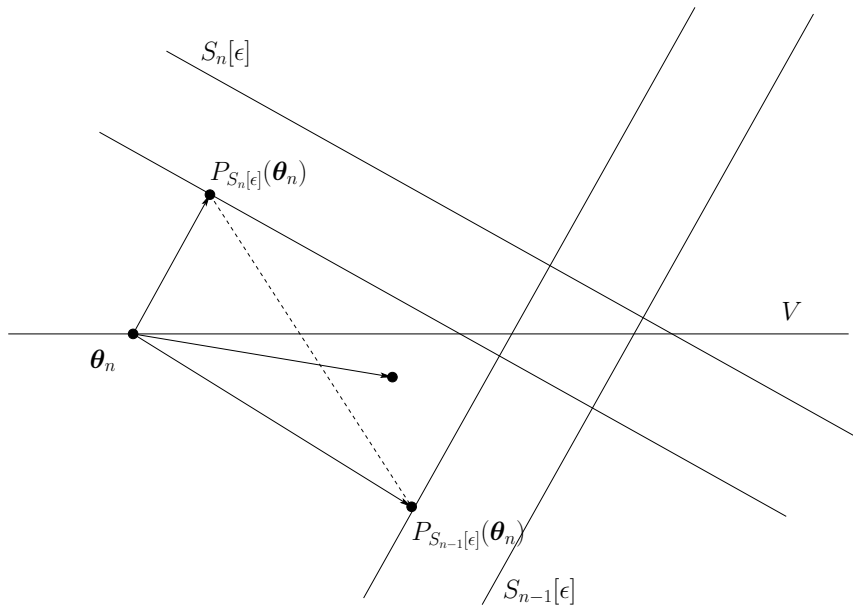




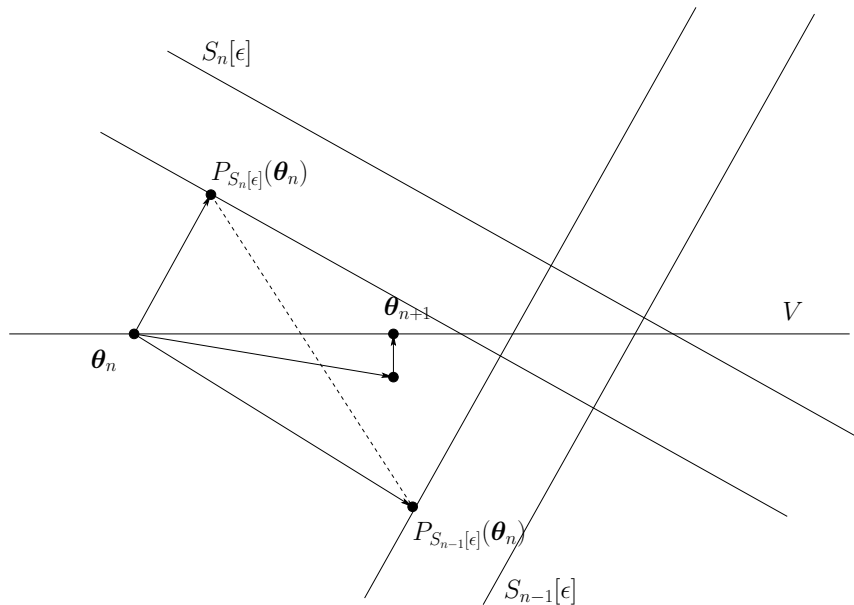
# Geometry of the Algorithm



# Geometry of the Algorithm



# Geometry of the Algorithm



# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

- Robustness is a key design issue in beamforming.

# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

- Robustness is a key design issue in beamforming.
- There are cases, for example, where the location of the SOI is known **approximately**.

# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

- Robustness is a key design issue in beamforming.
- There are cases, for example, where the location of the SOI is known **approximately**.
- A mathematical formulation for such a scenario is as follows;

# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

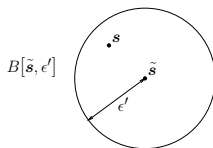
- Robustness is a key design issue in beamforming.
- There are cases, for example, where the location of the SOI is known **approximately**.
- A mathematical formulation for such a scenario is as follows;
  - ▶ given the approximate steering vector  $\tilde{s}$ ,



# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

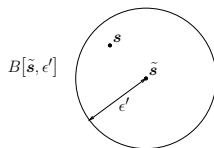
- Robustness is a key design issue in beamforming.
- There are cases, for example, where the location of the SOI is known **approximately**.
- A mathematical formulation for such a scenario is as follows;
  - ▶ given the approximate steering vector  $\tilde{s}$ ,
  - ▶ and a **ball of uncertainty**  $B[\tilde{s}, \epsilon']$ , of radius  $\epsilon'$  around  $\tilde{s}$ :



# Robustness in Beamforming

## Towards More Elaborated Constrained Learning

- Robustness is a key design issue in beamforming.
- There are cases, for example, where the location of the SOI is known **approximately**.
- A mathematical formulation for such a scenario is as follows;
  - ▶ given the approximate steering vector  $\tilde{s}$ ,
  - ▶ and a **ball of uncertainty**  $B[\tilde{s}, \epsilon']$ , of radius  $\epsilon'$  around  $\tilde{s}$ :



- ▶ calculate those  $\theta$  such that, for some user-defined  $\epsilon'' \geq 0$ ,

$$\theta^t s \in [1 - \epsilon'', 1 + \epsilon''], \quad \forall s \in B[\tilde{s}, \epsilon'].$$

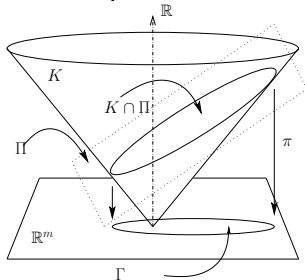
# The Icecream Cone

- The previous task breaks down to a number of more fundamental problems of the following type; find a vector that belongs to

$$\Gamma := \left\{ \boldsymbol{\theta} \in \mathbb{R}^m : \boldsymbol{\theta}^t \mathbf{s} \geq \gamma, \forall \mathbf{s} \in B[\tilde{\mathbf{s}}, \epsilon'] \right\} = \left\{ \begin{array}{l} \text{all vectors that satisfy an} \\ \text{infinite number of inequalities} \end{array} \right\}.$$

- If  $\Gamma \neq \emptyset$ , then the previous problem is equivalent to<sup>8</sup>

finding a point in  $K \cap \Pi$ ,  
 $K$ : an icecream cone,  
 $\Pi$ : a hyperplane.



<sup>8</sup>[Slavakis, Yamada' 07], [Slavakis, Theodoridis, Yamada' 09].

Given  $(\mathbf{x}_n, y_n)$ , find a  $\boldsymbol{\theta} \in \mathbb{R}^m$  such that

$$|\boldsymbol{\theta}^t \mathbf{x}_n - y_n| \leq \epsilon,$$

Given  $(\mathbf{x}_n, y_n)$ , find a  $\boldsymbol{\theta} \in \mathbb{R}^m$  such that

$$\begin{aligned} |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| &\leq \epsilon, \\ \boldsymbol{\theta}^t \mathbf{s} &\geq \gamma, \quad \forall \mathbf{s} \in B[\tilde{\mathbf{s}}, \epsilon'], \quad (\text{Robustness}). \end{aligned}$$

# Algorithm for Robust Regression

Assume weights  $\omega_j^{(n)} \geq 0$  such that  $\sum_{j=n-q+1}^n \omega_j^{(n)} = 1$ . For any  $\theta_0 \in \mathbb{R}^m$ ,

$$\theta_{n+1} := P_{\Pi} P_K \left( \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n \right) \right), \quad \forall n \geq 0,$$

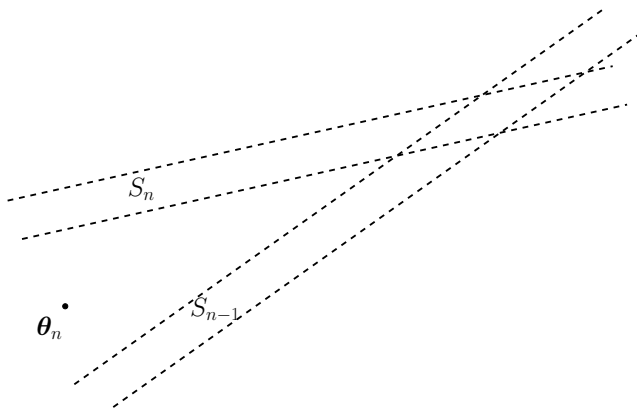
where the extrapolation coefficient  $\mu_n \in (0, 2\mathcal{M}_n)$  with

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

# Geometry of the Algorithm

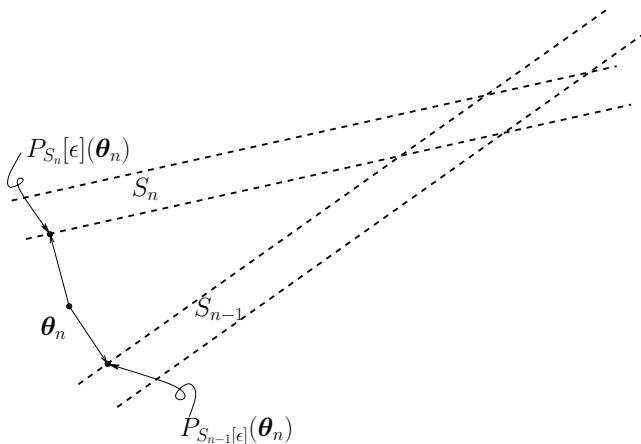
$\theta_n^\bullet$

# Geometry of the Algorithm

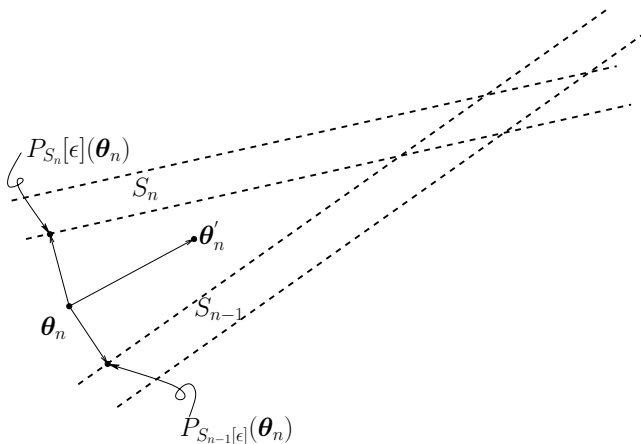




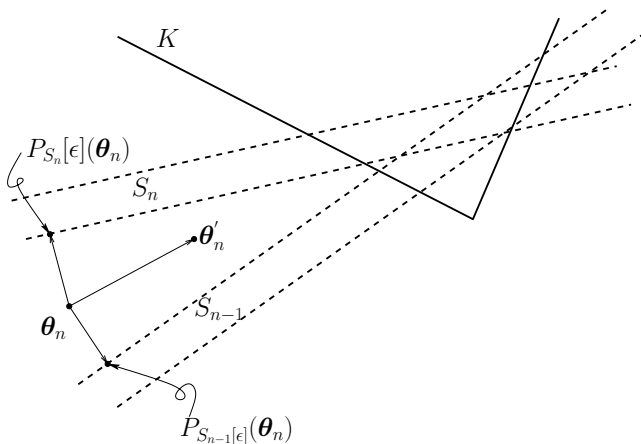
# Geometry of the Algorithm



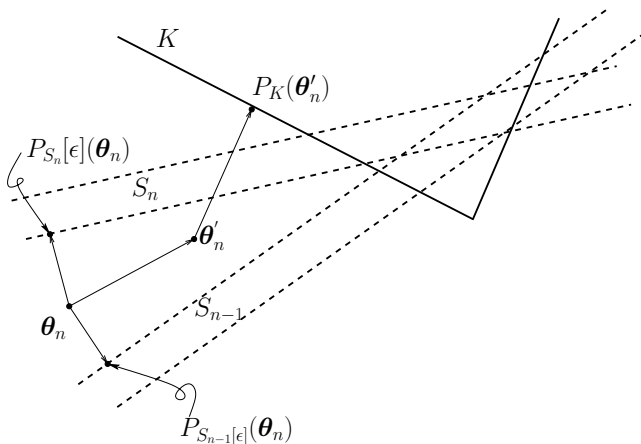
# Geometry of the Algorithm



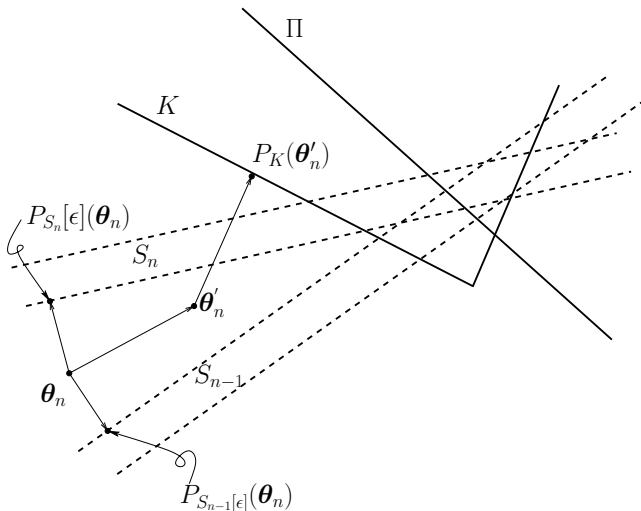
# Geometry of the Algorithm



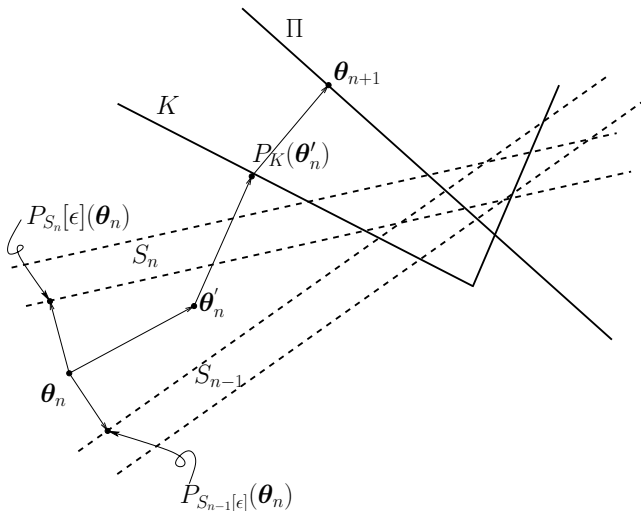
# Geometry of the Algorithm



# Geometry of the Algorithm



# Geometry of the Algorithm



# Handling A-Priori Information

How did we handle a-priori information?

## How did we handle a-priori information?

- Each piece of a-priori info was represented by a closed convex set, e.g.,  $K$ ,  $\Pi$ .



## How did we handle a-priori information?

- Each piece of a-priori info was represented by a closed convex set, e.g.,  $K$ ,  $\Pi$ .
- In order to be in agreement with **all** of the pieces of the a-priori info, we looked for a point into the intersection of the closed convex sets, e.g.,  $K \cap \Pi$ .

## How did we handle a-priori information?

- Each piece of a-priori info was represented by a closed convex set, e.g.,  $K$ ,  $\Pi$ .
- In order to be in agreement with **all** of the pieces of the a-priori info, we looked for a point into the intersection of the closed convex sets, e.g.,  $K \cap \Pi$ .
- In algorithmic terms, we utilized the **composition mapping**  $P_{\Pi}P_K : \mathbb{R}^m \rightarrow \mathbb{R}^m$ .

# Handling A-Priori Information

## How did we handle a-priori information?

- Each piece of a-priori info was represented by a closed convex set, e.g.,  $K$ ,  $\Pi$ .
- In order to be in agreement with **all** of the pieces of the a-priori info, we looked for a point into the intersection of the closed convex sets, e.g.,  $K \cap \Pi$ .
- In algorithmic terms, we utilized the **composition mapping**  $P_{\Pi}P_K : \mathbb{R}^m \rightarrow \mathbb{R}^m$ .

This strategy reminds us of POCS:

## POCS

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Then,

$$\forall \boldsymbol{\theta} \in \mathbb{R}^m, \quad (P_{C_p} \cdots P_{C_1})^n(\boldsymbol{\theta}) \xrightarrow[n \rightarrow \infty]{w} \exists \boldsymbol{\theta}_* \in \bigcap_{i=1}^p C_i.$$

# Handling A-Priori Information

## How did we handle a-priori information?

- Each piece of a-priori info was represented by a closed convex set, e.g.,  $K$ ,  $\Pi$ .
- In order to be in agreement with **all** of the pieces of the a-priori info, we looked for a point into the intersection of the closed convex sets, e.g.,  $K \cap \Pi$ .
- In algorithmic terms, we utilized the **composition mapping**  $P_{\Pi}P_K : \mathbb{R}^m \rightarrow \mathbb{R}^m$ .

This strategy reminds us of POCS:

## POCS

Given a **finite** number of closed convex sets  $C_1, \dots, C_p$ , with  $\bigcap_{i=1}^p C_i \neq \emptyset$ , let their associated projection mappings be  $P_{C_1}, \dots, P_{C_p}$ . Then,

$$\forall \theta \in \mathbb{R}^m, \quad (P_{C_p} \cdots P_{C_1})^n(\theta) \xrightarrow[n \rightarrow \infty]{w} \exists \theta_* \in \bigcap_{i=1}^p C_i.$$

## Key assumption

The a-priori info is **consistent**, i.e.,  $\bigcap_{i=1}^p C_i \neq \emptyset$ .

# Inconsistent A-Priori Information

- Is the case of **inconsistent** a-priori info possible in practice?

# Inconsistent A-Priori Information

- Is the case of **inconsistent** a-priori info possible in practice?
- **Yes**, in highly constrained learning tasks; the more constraints we add to the problem, the smaller the intersection of the associated convex sets becomes.

# Inconsistent A-Priori Information

- Is the case of **inconsistent** a-priori info possible in practice?
- **Yes**, in highly constrained learning tasks; the more constraints we add to the problem, the smaller the intersection of the associated convex sets becomes.

## Example

A beamforming problem where there is erroneous info on SOI and jammer steering vectors, array calibration errors, info on inoperative array elements, and stringent bounds on the weights of the array.

# Inconsistent A-Priori Information

- Is the case of **inconsistent** a-priori info possible in practice?
- **Yes**, in highly constrained learning tasks; the more constraints we add to the problem, the smaller the intersection of the associated convex sets becomes.

## Example

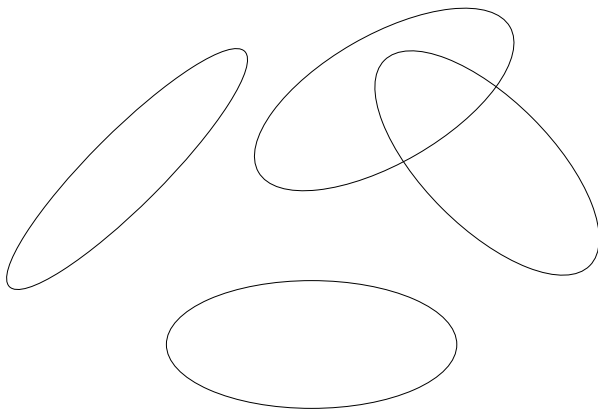
A beamforming problem where there is erroneous info on SOI and jammer steering vectors, array calibration errors, info on inoperative array elements, and stringent bounds on the weights of the array.

How do we deal with the case of inconsistent a-priori info, i.e.,

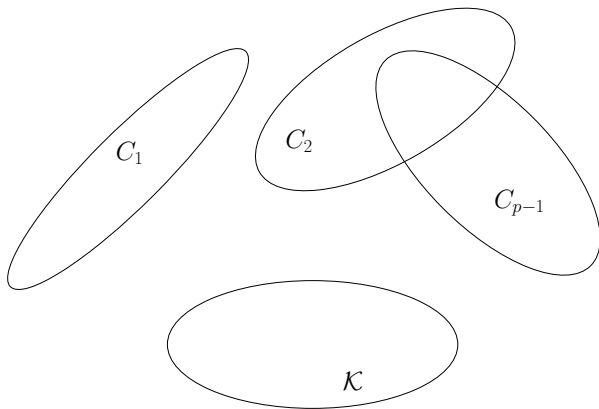
$$\bigcap_{i=1}^p C_i = \emptyset?$$



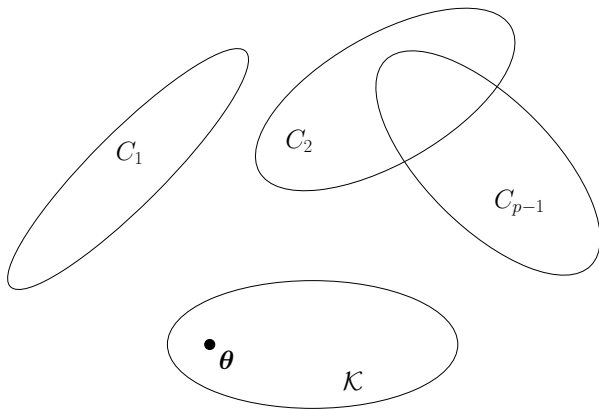
# An Intuitive Suggestion



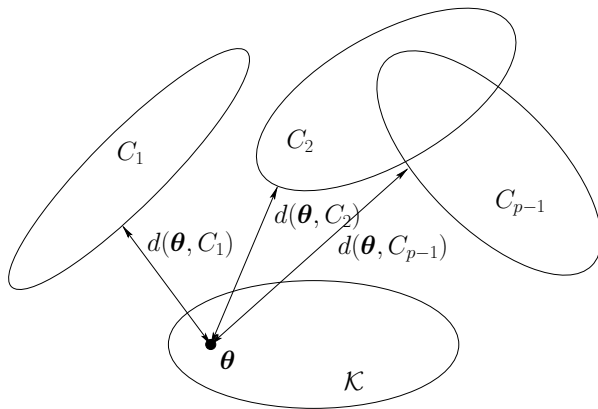
# An Intuitive Suggestion



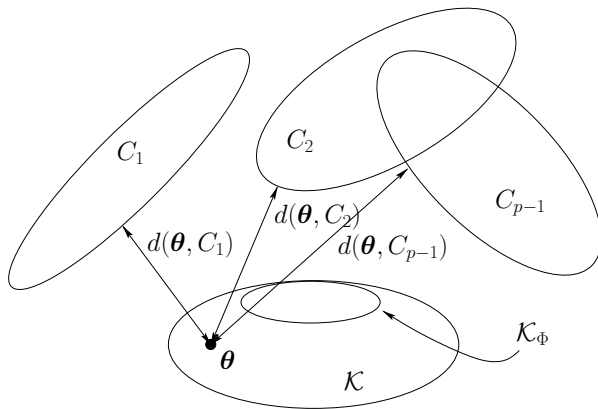
# An Intuitive Suggestion



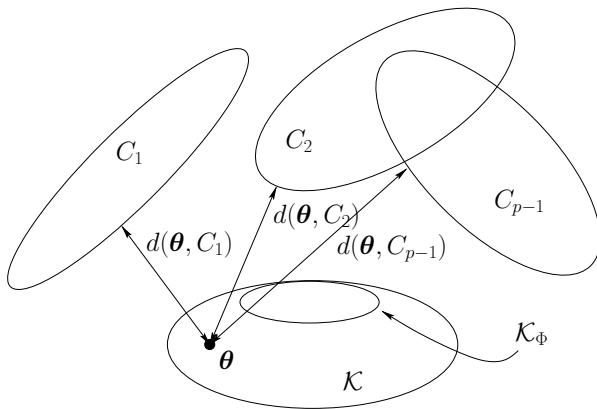
# An Intuitive Suggestion



# An Intuitive Suggestion



# An Intuitive Suggestion



## Definition ( $\mathcal{K}_\Phi$ )

All those points of  $\mathcal{K}$  which minimize a function  $\Phi$  of the distances  $\{d(\cdot, C_i)\}_{i=1}^{p-1}$ .

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .



# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .
- Assign to each  $C_i$  a convex weight  $\beta_i$ , i.e.,  $\beta_i \in (0, 1]$  and  $\sum_{i=1}^{p-1} \beta_i = 1$ .

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .
- Assign to each  $C_i$  a convex weight  $\beta_i$ , i.e.,  $\beta_i \in (0, 1]$  and  $\sum_{i=1}^{p-1} \beta_i = 1$ .
- Define the function:

$$\Phi(\boldsymbol{\theta}) := \frac{1}{2} \sum_{i=1}^{p-1} \beta_i d^2(\boldsymbol{\theta}, C_i), \quad \forall \boldsymbol{\theta} \in \mathcal{K}.$$

Our objective is to look for the **minimizers**  $\mathcal{K}_\Phi$  of this function.

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .
- Assign to each  $C_i$  a convex weight  $\beta_i$ , i.e.,  $\beta_i \in (0, 1]$  and  $\sum_{i=1}^{p-1} \beta_i = 1$ .
- Define the function:

$$\Phi(\boldsymbol{\theta}) := \frac{1}{2} \sum_{i=1}^{p-1} \beta_i d^2(\boldsymbol{\theta}, C_i), \quad \forall \boldsymbol{\theta} \in \mathcal{K}.$$

Our objective is to look for the **minimizers**  $\mathcal{K}_\Phi$  of this function.

Notice that  $\Phi' = I - \sum_{i=1}^{p-1} \beta_i P_{C_i}$ .

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .
- Assign to each  $C_i$  a convex weight  $\beta_i$ , i.e.,  $\beta_i \in (0, 1]$  and  $\sum_{i=1}^{p-1} \beta_i = 1$ .
- Define the function:

$$\Phi(\boldsymbol{\theta}) := \frac{1}{2} \sum_{i=1}^{p-1} \beta_i d^2(\boldsymbol{\theta}, C_i), \quad \forall \boldsymbol{\theta} \in \mathcal{K}.$$

Our objective is to look for the **minimizers**  $\mathcal{K}_\Phi$  of this function.

$$\text{Notice that } \Phi' = I - \sum_{i=1}^{p-1} \beta_i P_{C_i}.$$

- Define the **mapping**  $T : \mathbb{R}^m \rightarrow \mathbb{R}^m$  as

$$T := P_{\mathcal{K}} \left( I - \lambda \left( I - \sum_{i=1}^{p-1} \beta_i P_{C_i} \right) \right), \quad \lambda \in (0, 2).$$

# A Way to Handle Inconsistent A-Priori Information

Given a number of a-priori constraints, represented as  $p$  closed convex sets,

- Identify one of them as the **absolute constraint**  $\mathcal{K}$ , and rename the other ones as  $C_1, C_2, \dots, C_{p-1}$ .
- Assign to each  $C_i$  a convex weight  $\beta_i$ , i.e.,  $\beta_i \in (0, 1]$  and  $\sum_{i=1}^{p-1} \beta_i = 1$ .
- Define the function:

$$\Phi(\theta) := \frac{1}{2} \sum_{i=1}^{p-1} \beta_i d^2(\theta, C_i), \quad \forall \theta \in \mathcal{K}.$$

Our objective is to look for the **minimizers**  $\mathcal{K}_\Phi$  of this function.

$$\text{Notice that } \Phi' = I - \sum_{i=1}^{p-1} \beta_i P_{C_i}.$$

- Define the **mapping**  $T : \mathbb{R}^m \rightarrow \mathbb{R}^m$  as

$$T := P_{\mathcal{K}} \left( I - \lambda \left( I - \sum_{i=1}^{p-1} \beta_i P_{C_i} \right) \right), \quad \lambda \in (0, 2).$$

- Then,  $\text{Fix}(T) = \mathcal{K}_\Phi$ .

For any  $\theta_0 \in \mathbb{R}^m$ ,

$$\theta_{n+1} := \textcolor{red}{T} \left( \theta_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n \right) \right), \quad \forall n \geq 0,$$

where the extrapolation coefficient  $\mu_n \in (0, 2\mathcal{M}_n)$  with

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) - \theta_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\theta_n) \neq \theta_n, \\ 1, & \text{otherwise.} \end{cases}$$

## Problem definition

## Problem definition

- In a number of applications, many of the parameters to be estimated are a-priori known to be **zero**. That is, the parameter vector,  $\theta$ , is sparse.

$$\theta^t = [*, *, 0, 0, 0, *, 0, \dots].$$



## Problem definition

- In a number of applications, many of the parameters to be estimated are a-priori known to be **zero**. That is, the parameter vector,  $\theta$ , is sparse.

$$\theta^t = [*, *, 0, 0, 0, *, 0, \dots].$$

If the locations of the zeros were known, the problem would be trivial.

However, the locations **of the zeros are not known a-priori**. This makes the task challenging.

## Problem definition

- In a number of applications, many of the parameters to be estimated are a-priori known to be **zero**. That is, the parameter vector,  $\theta$ , is sparse.

$$\theta^t = [*, *, 0, 0, 0, *, 0, \dots].$$

If the locations of the zeros were known, the problem would be trivial.

However, the locations **of the zeros are not known a-priori**. This makes the task challenging.

- Typical applications include echo cancellation in Internet telephony, MIMO channel estimation, Compressed Sensing (CS), etc.

## Problem definition

- In a number of applications, many of the parameters to be estimated are a-priori known to be **zero**. That is, the parameter vector,  $\theta$ , is sparse.

$$\theta^t = [*, *, 0, 0, 0, *, 0, \dots].$$

If the locations of the zeros were known, the problem would be trivial.

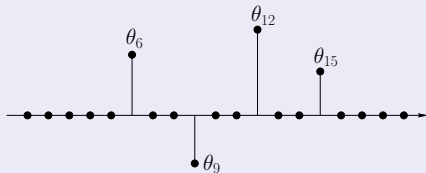
However, the locations **of the zeros are not known a-priori**. This makes the task challenging.

- Typical applications include echo cancellation in Internet telephony, MIMO channel estimation, Compressed Sensing (CS), etc.
- Sparsity promotion is achieved via  **$\ell_1$ -norm regularization** of a loss function:

$$\min_{\theta \in \mathbb{R}^m} \sum_{n=0}^N \mathcal{L}(y_n, \mathbf{x}_n^t \theta) + \lambda \|\theta\|_1, \quad \lambda > 0.$$

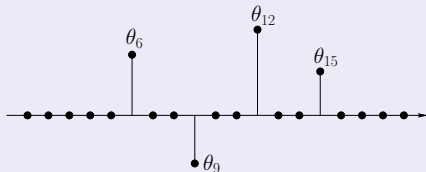
## The $\ell_0$ norm

$$\|\boldsymbol{\theta}\|_0 := \text{card}\{i : \theta_i \neq 0\}.$$



## The $\ell_0$ norm

$$\|\boldsymbol{\theta}\|_0 := \text{card}\{i : \theta_i \neq 0\}.$$



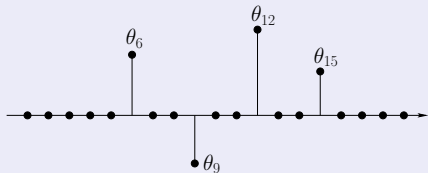
- Consider the linear model:

$$y_n := \mathbf{x}_n^t \boldsymbol{\theta} + v_n, \quad \forall n,$$

where  $(v_n)_{n \geq 0}$  denotes the noise process.

## The $\ell_0$ norm

$$\|\boldsymbol{\theta}\|_0 := \text{card}\{i : \theta_i \neq 0\}.$$



- Consider the linear model:

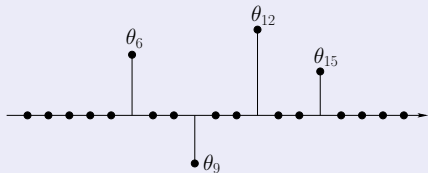
$$y_n := \mathbf{x}_n^t \boldsymbol{\theta} + v_n, \quad \forall n,$$

where  $(v_n)_{n \geq 0}$  denotes the noise process.

- Define  $\mathbf{X}_N := [\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N]$ ,  $\mathbf{y}_N := [y_0, y_1, \dots, y_N]^t$ , and  $\epsilon \geq 0$ .

## The $\ell_0$ norm

$$\|\boldsymbol{\theta}\|_0 := \text{card}\{i : \theta_i \neq 0\}.$$



- Consider the linear model:

$$y_n := \mathbf{x}_n^t \boldsymbol{\theta} + v_n, \quad \forall n,$$

where  $(v_n)_{n \geq 0}$  denotes the noise process.

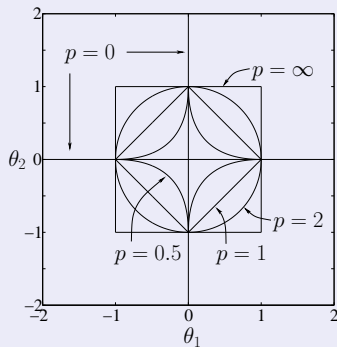
- Define  $\mathbf{X}_N := [\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N]$ ,  $\mathbf{y}_N := [y_0, y_1, \dots, y_N]^t$ , and  $\epsilon \geq 0$ .
- A typical Compressed Sensing task is formulated as follows:

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^m} \quad & \|\boldsymbol{\theta}\|_0 \\ \text{s.t.} \quad & \|\mathbf{X}_N^t \boldsymbol{\theta} - \mathbf{y}_N\| \leq \epsilon. \end{aligned}$$

# Alternatives to the $\ell_0$ Norm

## The $\ell_p$ norm ( $0 < p \leq 1$ )

$$\|\boldsymbol{\theta}\|_p := \left( \sum_{i=1}^m |\theta_i|^p \right)^{\frac{1}{p}}.$$





## The $\ell_1$ -ball case

## The $\ell_1$ -ball case

- Given  $(\mathbf{x}_n, y_n)$ ,  $n = 0, 1, 2, \dots$ , find  $\boldsymbol{\theta}$  such that

$$\begin{aligned} |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| &\leq \epsilon, \quad n = 0, 1, 2, \dots \\ \boldsymbol{\theta} &\in B_{\ell_1}[\delta] := \{\boldsymbol{\theta}' \in \mathbb{R}^m : \|\boldsymbol{\theta}'\|_1 \leq \delta\}. \end{aligned}$$

## The $\ell_1$ -ball case

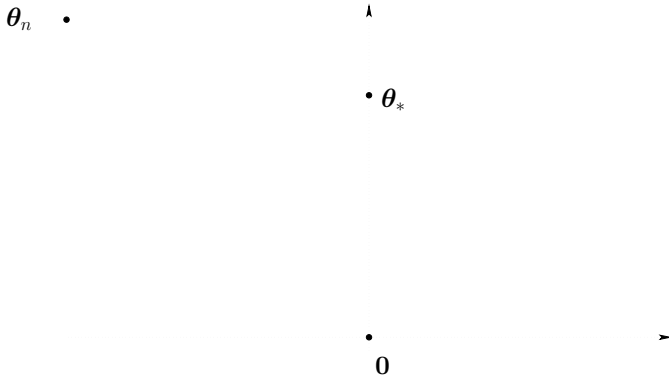
- Given  $(\mathbf{x}_n, y_n)$ ,  $n = 0, 1, 2, \dots$ , find  $\boldsymbol{\theta}$  such that

$$\begin{aligned} |\boldsymbol{\theta}^t \mathbf{x}_n - y_n| &\leq \epsilon, \quad n = 0, 1, 2, \dots \\ \boldsymbol{\theta} &\in B_{\ell_1}[\delta] := \{\boldsymbol{\theta}' \in \mathbb{R}^m : \|\boldsymbol{\theta}'\|_1 \leq \delta\}. \end{aligned}$$

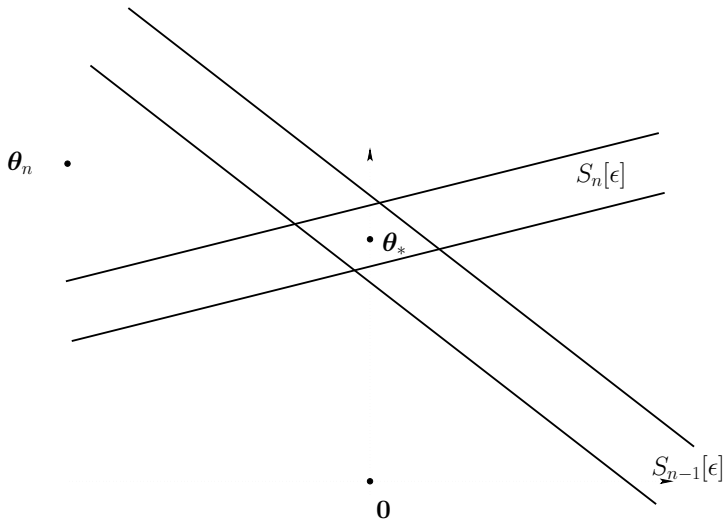
- The recursion:

$$\boldsymbol{\theta}_{n+1} := P_{B_{\ell_1}[\delta]} \left( \boldsymbol{\theta}_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right) \right).$$

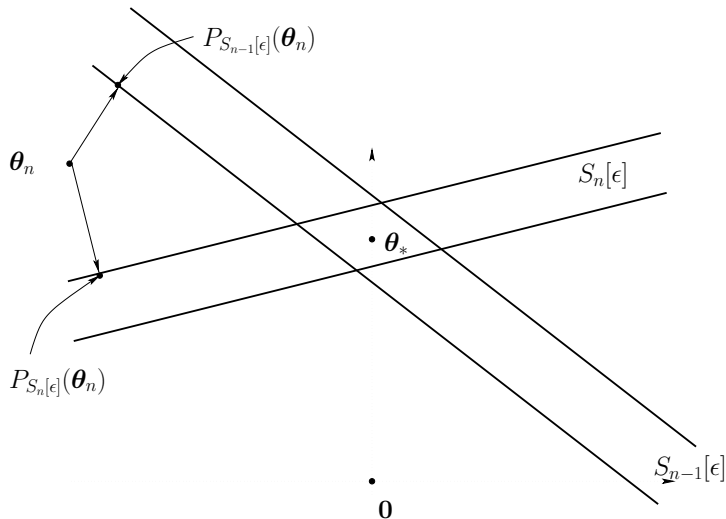
# Geometric Illustration of the Algorithm



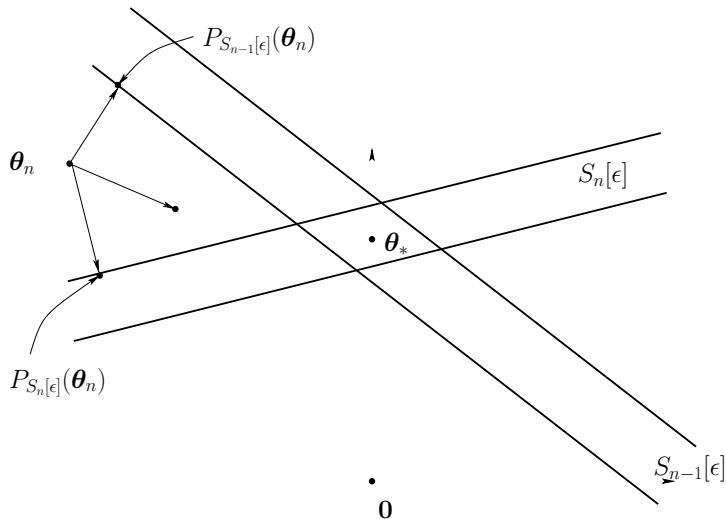
# Geometric Illustration of the Algorithm



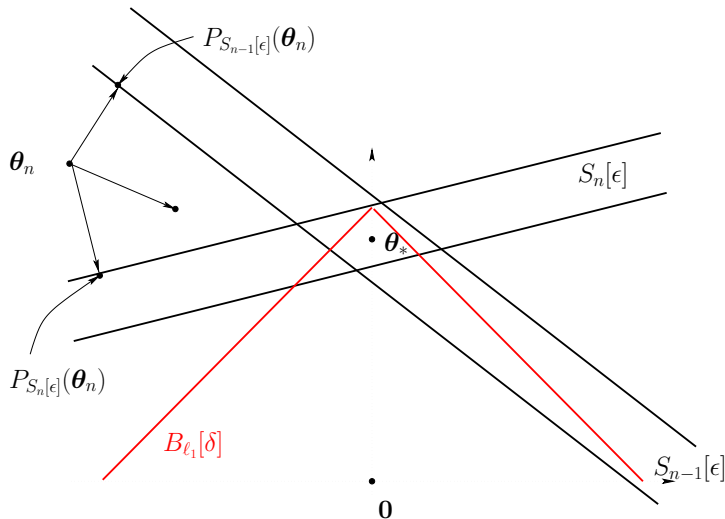
# Geometric Illustration of the Algorithm



# Geometric Illustration of the Algorithm

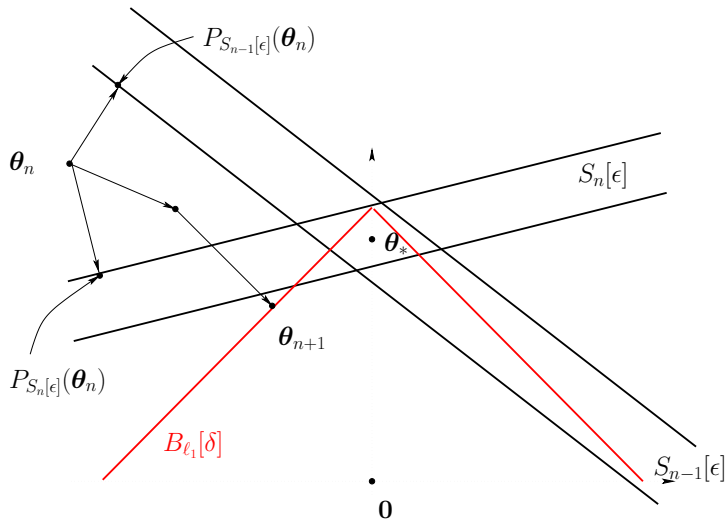


# Geometric Illustration of the Algorithm





# Geometric Illustration of the Algorithm



**Remark:** The convergence can be **significantly** speeded up, if in place of the  $\ell_1$ -ball a **weighted**  $\ell_1$ -ball is used to constrain the solutions.

**Remark:** The convergence can be **significantly** speeded up, if in place of the  $\ell_1$ -ball a **weighted**  $\ell_1$ -ball is used to constrain the solutions.

- Definition:

$$\|\boldsymbol{\theta}\|_{1,w} := \sum_{i=1}^m w_i |\theta_i|,$$
$$B_{\ell_1}[\boldsymbol{w}_n, \delta] := \{\boldsymbol{\theta} \in \mathbb{R}^m : \|\boldsymbol{\theta}\|_{1,w} \leq \delta\}.$$

**Remark:** The convergence can be **significantly** speeded up, if in place of the  $\ell_1$ -ball a **weighted**  $\ell_1$ -ball is used to constrain the solutions.

- Definition:

$$\|\boldsymbol{\theta}\|_{1,w} := \sum_{i=1}^m w_i |\theta_i|,$$
$$B_{\ell_1}[\mathbf{w}_n, \delta] := \{\boldsymbol{\theta} \in \mathbb{R}^m : \|\boldsymbol{\theta}\|_{1,w} \leq \delta\}.$$

- Time-adaptive weighted norm:

$$w_{n,i} := \frac{1}{|\theta_{n,i}| + \epsilon'_n}.$$

**Remark:** The convergence can be **significantly** speeded up, if in place of the  $\ell_1$ -ball a **weighted**  $\ell_1$ -ball is used to constrain the solutions.

- Definition:

$$\|\boldsymbol{\theta}\|_{1,w} := \sum_{i=1}^m w_i |\theta_i|,$$
$$B_{\ell_1}[\mathbf{w}_n, \delta] := \{\boldsymbol{\theta} \in \mathbb{R}^m : \|\boldsymbol{\theta}\|_{1,w} \leq \delta\}.$$

- Time-adaptive weighted norm:

$$w_{n,i} := \frac{1}{|\theta_{n,i}| + \epsilon'_n}.$$

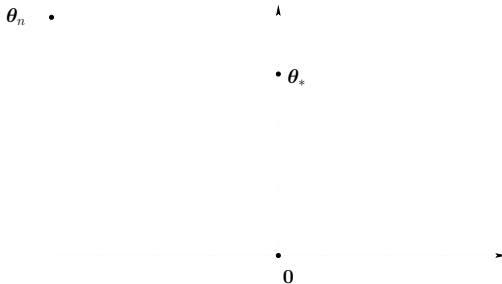
- The recursion<sup>9</sup>:

$$\boldsymbol{\theta}_{n+1} := P_{B_{\ell_1}[\mathbf{w}_n, \delta]} \left( \boldsymbol{\theta}_n + \mu_n \left( \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(\boldsymbol{\theta}_n) - \boldsymbol{\theta}_n \right) \right).$$

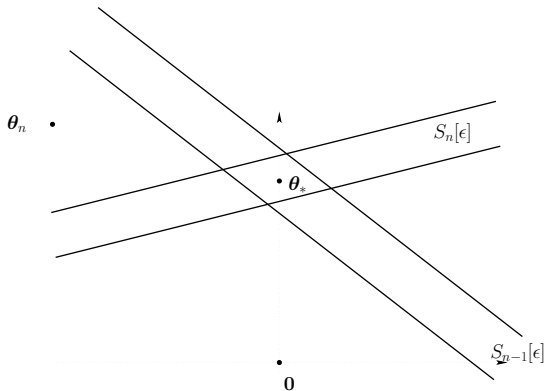
---

<sup>9</sup>[Kopsinis, Slavakis, Theodoridis, '11].

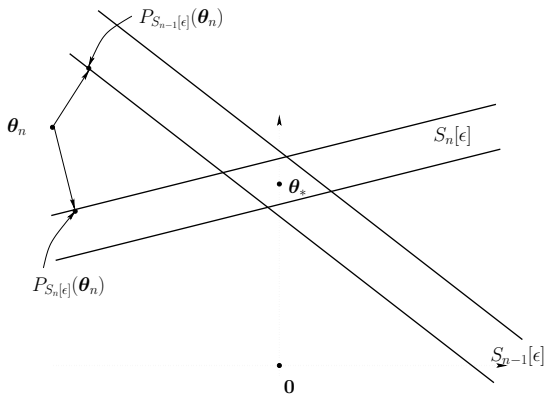
# Geometric Illustration of the Algorithm



# Geometric Illustration of the Algorithm

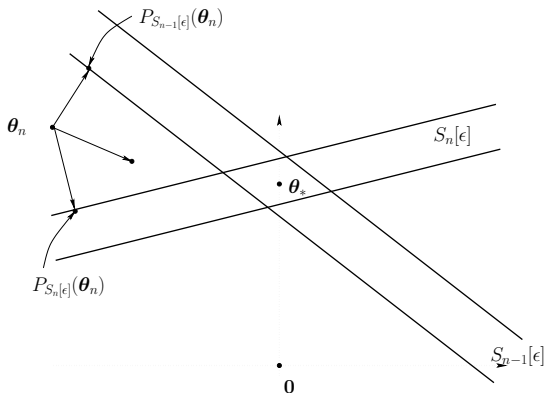


# Geometric Illustration of the Algorithm

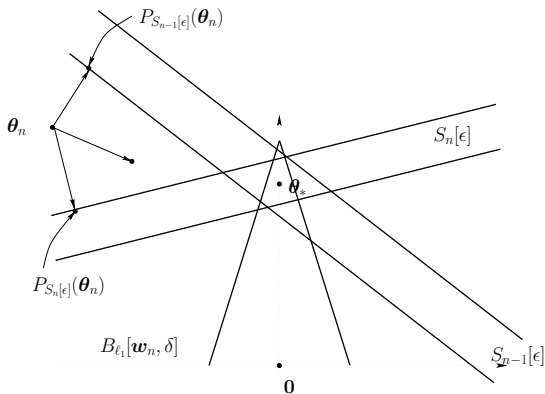




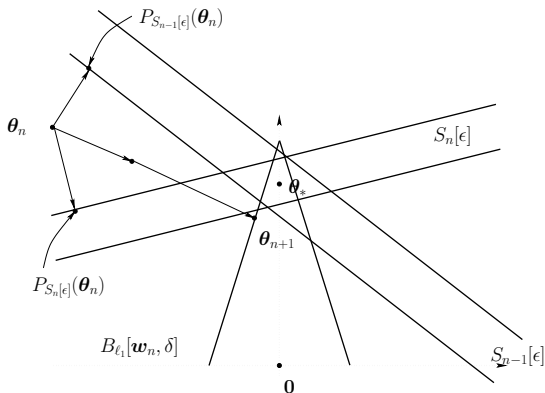
# Geometric Illustration of the Algorithm



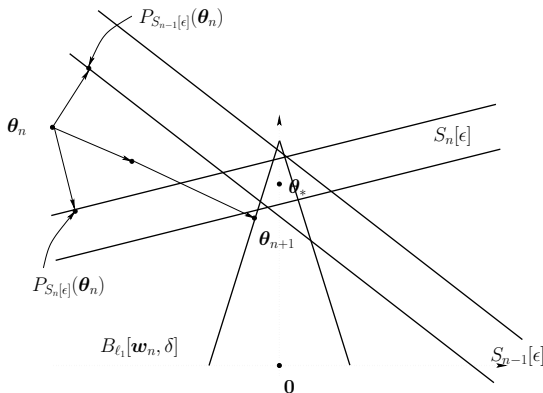
# Geometric Illustration of the Algorithm



# Geometric Illustration of the Algorithm



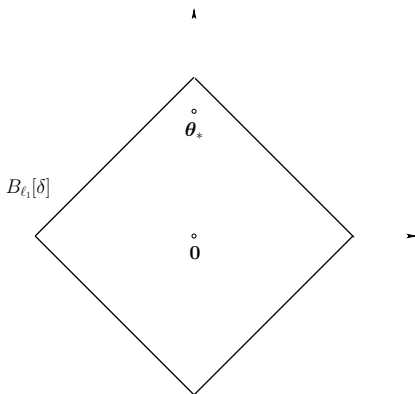
# Geometric Illustration of the Algorithm



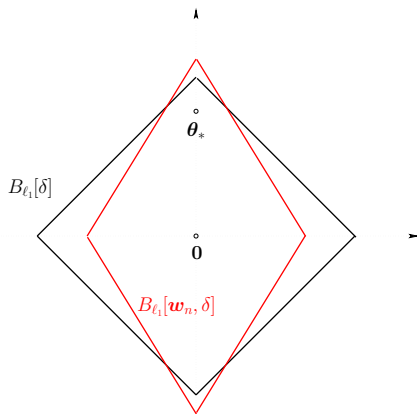
Projecting onto  $B_{\ell_1}[\mathbf{w}_n, \delta]$  is equivalent to a specific **soft thresholding** operation.

Note that our constraint, i.e., the weighted  $\ell_1$ -ball is a **time-varying constraint**.

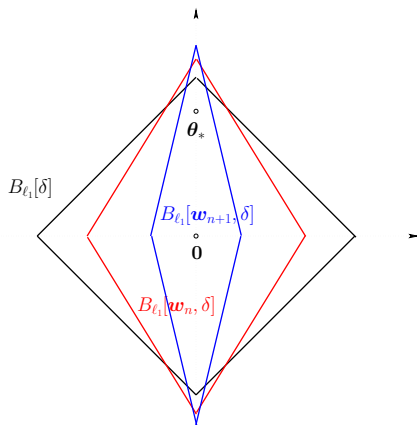
Note that our constraint, i.e., the weighted  $\ell_1$ -ball is a **time-varying constraint**.



Note that our constraint, i.e., the weighted  $\ell_1$ -ball is a **time-varying constraint**.

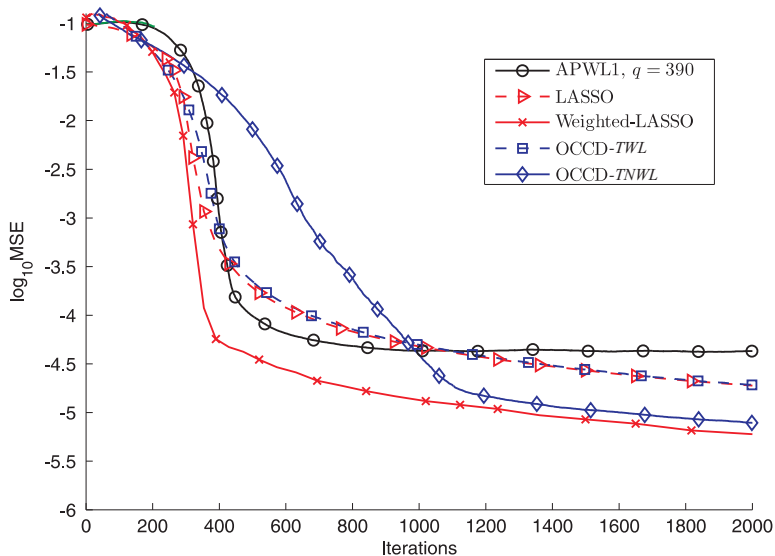


Note that our constraint, i.e., the weighted  $\ell_1$ -ball is a **time-varying constraint**.



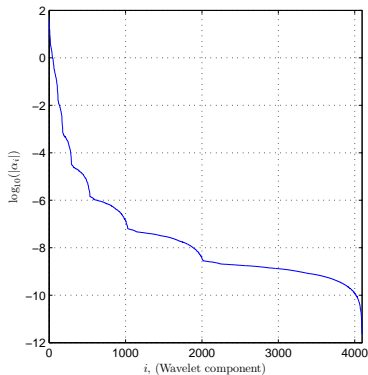
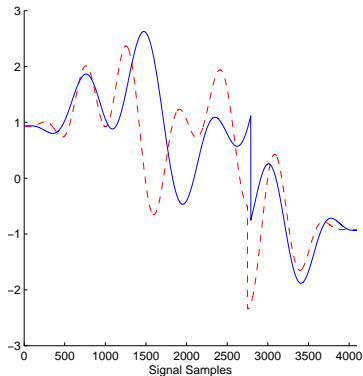


# Time Invariant Signal



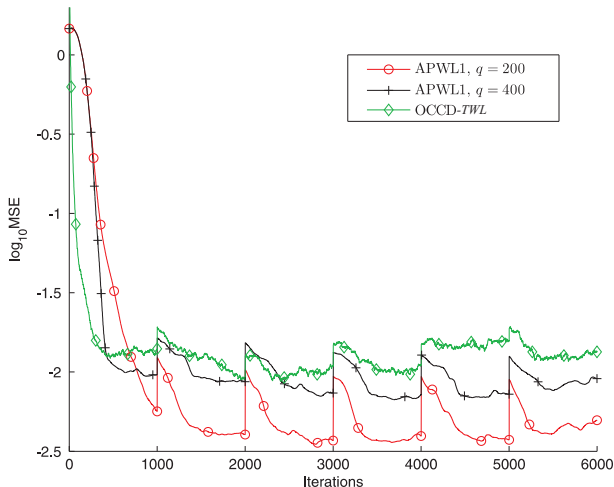
$m := 1024$ ,  $\|\theta_*\|_0 := 100$  wavelet coefficients. The radius of the  $\ell_1$ -ball is set to  $\delta := 101$ .

# Time Varying Signal



$m := 4096$ . The radius of the  $\ell_1$ -ball is set to  $\delta := 40$ .  
The sum of two [chirp signals](#).

# Time Varying Signal



Movies of the [OCCD](#), and the [APWL1sub](#).

# Generalized Thresholding

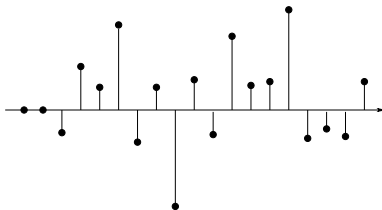
What we have seen so far corresponds to **soft thresholding** operations.

# Generalized Thresholding

What we have seen so far corresponds to **soft thresholding** operations.

## Hard thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .

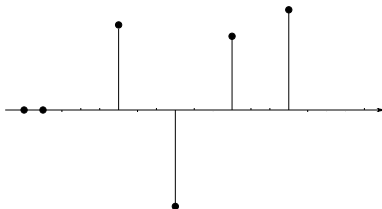


# Generalized Thresholding

What we have seen so far corresponds to **soft thresholding** operations.

## Hard thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Keep those as they are, while nullify the rest of them.

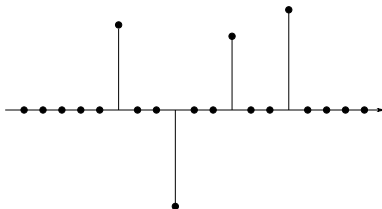


# Generalized Thresholding

What we have seen so far corresponds to **soft thresholding** operations.

## Hard thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Keep those as they are, while nullify the rest of them.

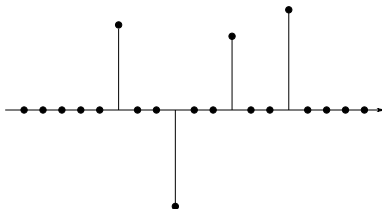


# Generalized Thresholding

What we have seen so far corresponds to **soft thresholding** operations.

## Hard thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Keep those as they are, while nullify the rest of them.



## Generalized thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .

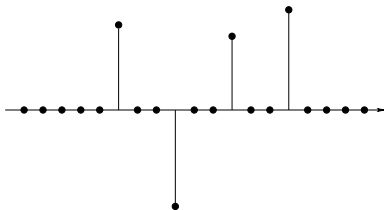


# Generalized Thresholding

What we have seen so far corresponds to **soft thresholding** operations.

## Hard thresholding

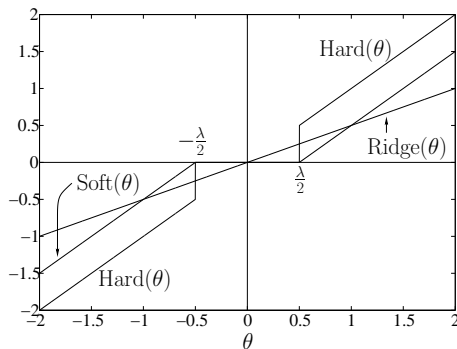
- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Keep those as they are, while nullify the rest of them.



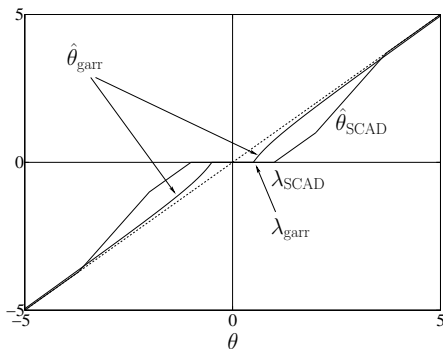
## Generalized thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- **Shrink**, under some rule, the rest of the components.

# Examples of Generalized Thresholding Mappings



(a) Hard, soft thresholding, and the ridge regression estimate.



(b) The SCAD and garrote thresholding.

## Penalized Least-Squares Thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .

## Penalized Least-Squares Thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Let  $\theta_i$  be one of the rest of the components.

## Penalized Least-Squares Thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Let  $\theta_i$  be one of the rest of the components.
- In order to **shrink**  $\theta_i$ , solve the optimization task:

$$\min_{\hat{\theta}_i \in \mathbb{R}} \frac{1}{2}(\hat{\theta}_i - \theta_i)^2 + \lambda p(|\hat{\theta}_i|), \quad \lambda > 0,$$

where  $p(\cdot)$  stands for a user-defined penalty function, which might be **non-convex**.

## Penalized Least-Squares Thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Let  $\theta_i$  be one of the rest of the components.
- In order to **shrink**  $\theta_i$ , solve the optimization task:

$$\min_{\hat{\theta}_i \in \mathbb{R}} \frac{1}{2} (\hat{\theta}_i - \theta_i)^2 + \lambda p(|\hat{\theta}_i|), \quad \lambda > 0,$$

where  $p(\cdot)$  stands for a user-defined penalty function, which might be **non-convex**.

- Under some mild conditions, the previous optimization task possesses a **unique solution**  $\hat{\theta}_{i*}$ .

## Penalized Least-Squares Thresholding

- Identify the  $K$  largest, in magnitude, components of a vector  $\theta$ .
- Let  $\theta_i$  be one of the rest of the components.
- In order to **shrink**  $\theta_i$ , solve the optimization task:

$$\min_{\hat{\theta}_i \in \mathbb{R}} \frac{1}{2} (\hat{\theta}_i - \theta_i)^2 + \lambda p(|\hat{\theta}_i|), \quad \lambda > 0,$$

where  $p(\cdot)$  stands for a user-defined penalty function, which might be **non-convex**.

- Under some mild conditions, the previous optimization task possesses a **unique solution**  $\hat{\theta}_{i*}$ .

## Definition (Generalized Thresholding Mapping)

The **Generalized Thresholding mapping** is defined as follows:

$$T_{\text{GT}} : \theta_i \mapsto \hat{\theta}_{i*}.$$

- Given  $K$ , define the set of all tuples of length  $K$ :

$$\mathcal{I} := \{(i_1, i_2, \dots, i_K) : 1 \leq i_1 < i_2 < \dots < i_K \leq m\}.$$



# Fixed Point Set of $T_{\text{GT}}$

- Given  $K$ , define the set of all tuples of length  $K$ :

$$\mathcal{T} := \{(i_1, i_2, \dots, i_K) : 1 \leq i_1 < i_2 < \dots < i_K \leq m\}.$$

- Given a tuple  $J \in \mathcal{T}$ , define the subspace:

$$M_J := \{\boldsymbol{\theta} \in \mathbb{R}^m : \theta_i = 0, \forall i \notin J\}.$$

# Fixed Point Set of $T_{\text{GT}}$

- Given  $K$ , define the set of all tuples of length  $K$ :

$$\mathcal{T} := \{(i_1, i_2, \dots, i_K) : 1 \leq i_1 < i_2 < \dots < i_K \leq m\}.$$

- Given a tuple  $J \in \mathcal{T}$ , define the subspace:

$$M_J := \{\boldsymbol{\theta} \in \mathbb{R}^m : \theta_i = 0, \forall i \notin J\}.$$

- Then, the fixed point set of  $T_{\text{GT}}$  is a **union of subspaces**:

$$\text{Fix}(T_{\text{GT}}) = \bigcup_{J \in \mathcal{T}} M_J, \quad (\text{non-convex set}).$$

# Fixed Point Set of $T_{GT}$

- Given  $K$ , define the set of all tuples of length  $K$ :

$$\mathcal{T} := \{(i_1, i_2, \dots, i_K) : 1 \leq i_1 < i_2 < \dots < i_K \leq m\}.$$

- Given a tuple  $J \in \mathcal{T}$ , define the subspace:

$$M_J := \{\boldsymbol{\theta} \in \mathbb{R}^m : \theta_i = 0, \forall i \notin J\}.$$

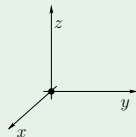
- Then, the fixed point set of  $T_{GT}$  is a **union of subspaces**:

$$\text{Fix}(T_{GT}) = \bigcup_{J \in \mathcal{T}} M_J, \quad (\text{non-convex set}).$$

## Example

For the 3-dimensional case  $\mathbb{R}^3$ , and if  $K := 2$ ,

$$\text{Fix}(T_{GT}) = xy\text{-plane} \cup yz\text{-plane} \cup xz\text{-plane}.$$



## Definition (Nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$  is called **nonexpansive** if

$$\|T(f_1) - T(f_2)\| \leq \|f_1 - f_2\|, \quad \forall f_1, f_2 \in \mathcal{H}.$$

The fixed point set of a nonexpansive mapping is **closed and convex**.

# First Steps Towards a Unifying Framework

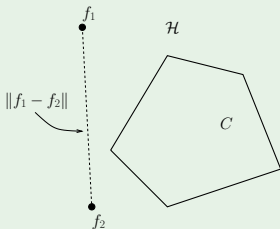
## Definition (Nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$  is called **nonexpansive** if

$$\|T(f_1) - T(f_2)\| \leq \|f_1 - f_2\|, \quad \forall f_1, f_2 \in \mathcal{H}.$$

The fixed point set of a nonexpansive mapping is **closed and convex**.

## Example (Projection Mapping)



# First Steps Towards a Unifying Framework

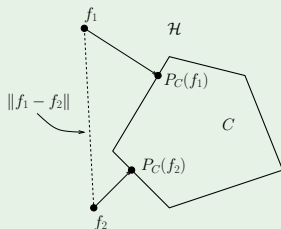
## Definition (Nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$  is called **nonexpansive** if

$$\|T(f_1) - T(f_2)\| \leq \|f_1 - f_2\|, \quad \forall f_1, f_2 \in \mathcal{H}.$$

The fixed point set of a nonexpansive mapping is **closed and convex**.

## Example (Projection Mapping)



# First Steps Towards a Unifying Framework

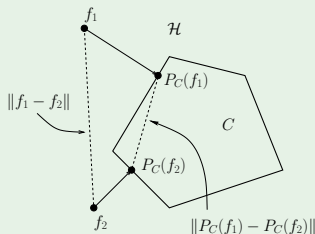
## Definition (Nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$  is called **nonexpansive** if

$$\|T(f_1) - T(f_2)\| \leq \|f_1 - f_2\|, \quad \forall f_1, f_2 \in \mathcal{H}.$$

The fixed point set of a nonexpansive mapping is **closed and convex**.

## Example (Projection Mapping)



$$\text{Fix}(P_C) = C.$$

# Quasi-nonexpansive Mapping

## Definition (Quasi-nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , with  $\text{Fix}(T) \neq \emptyset$ , is called **quasi-nonexpansive**, if

$$\|T(f) - h\| \leq \|f - h\|, \quad \forall f \in \mathcal{H}, \forall h \in \text{Fix}(T).$$

The fixed point set of  $T$  is **convex**.

$\mathcal{H}$

$f$   
•

$\text{Fix}(T)$



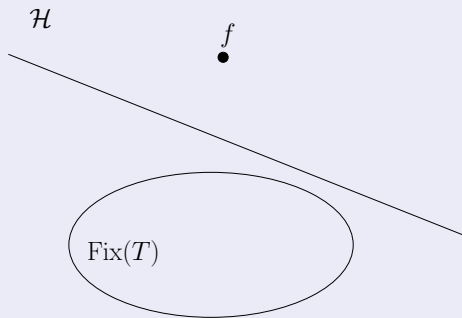
# Quasi-nonexpansive Mapping

## Definition (Quasi-nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , with  $\text{Fix}(T) \neq \emptyset$ , is called **quasi-nonexpansive**, if

$$\|T(f) - h\| \leq \|f - h\|, \quad \forall f \in \mathcal{H}, \forall h \in \text{Fix}(T).$$

The fixed point set of  $T$  is **convex**.



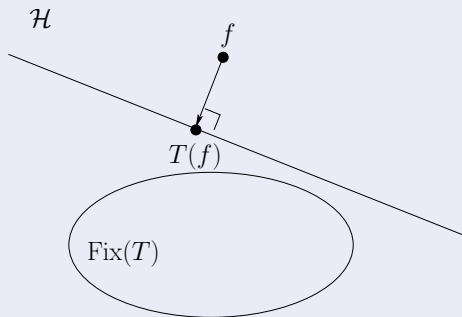
# Quasi-nonexpansive Mapping

## Definition (Quasi-nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , with  $\text{Fix}(T) \neq \emptyset$ , is called **quasi-nonexpansive**, if

$$\|T(f) - h\| \leq \|f - h\|, \quad \forall f \in \mathcal{H}, \forall h \in \text{Fix}(T).$$

The fixed point set of  $T$  is **convex**.



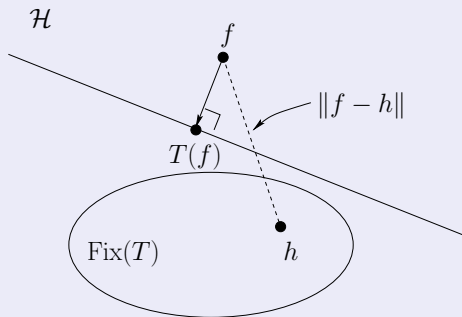
# Quasi-nonexpansive Mapping

## Definition (Quasi-nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , with  $\text{Fix}(T) \neq \emptyset$ , is called **quasi-nonexpansive**, if

$$\|T(f) - h\| \leq \|f - h\|, \quad \forall f \in \mathcal{H}, \forall h \in \text{Fix}(T).$$

The fixed point set of  $T$  is **convex**.



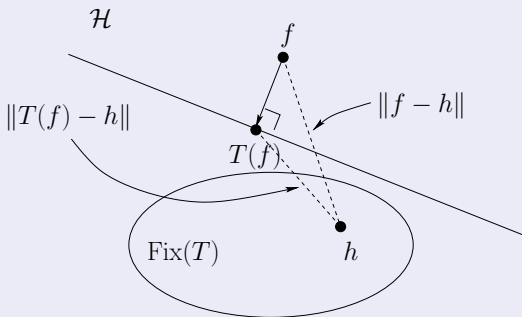
# Quasi-nonexpansive Mapping

## Definition (Quasi-nonexpansive Mapping)

A mapping  $T : \mathcal{H} \rightarrow \mathcal{H}$ , with  $\text{Fix}(T) \neq \emptyset$ , is called **quasi-nonexpansive**, if

$$\|T(f) - h\| \leq \|f - h\|, \quad \forall f \in \mathcal{H}, \forall h \in \text{Fix}(T).$$

The fixed point set of  $T$  is **convex**.

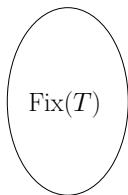


Every nonexpansive mapping is quasi-nonexpansive.

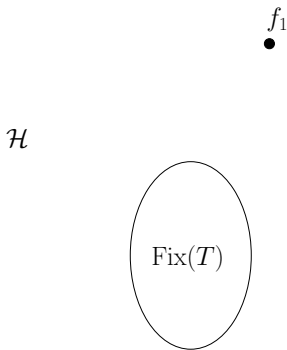
Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.

Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.

$\mathcal{H}$

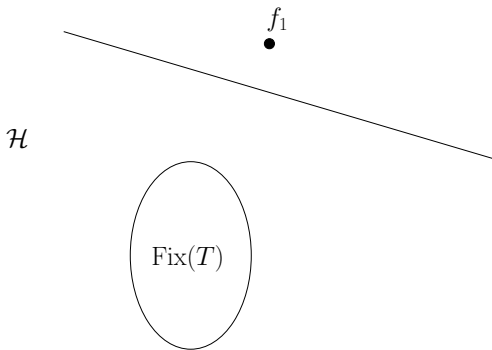


Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.

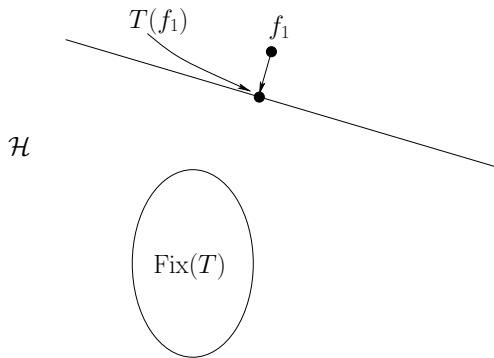




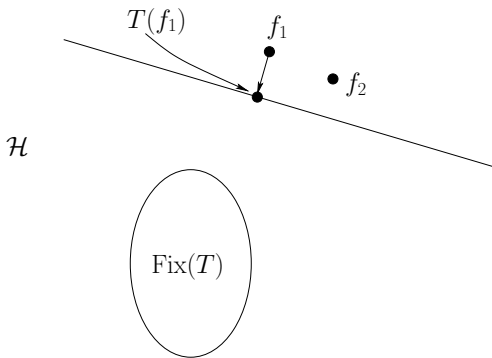
Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



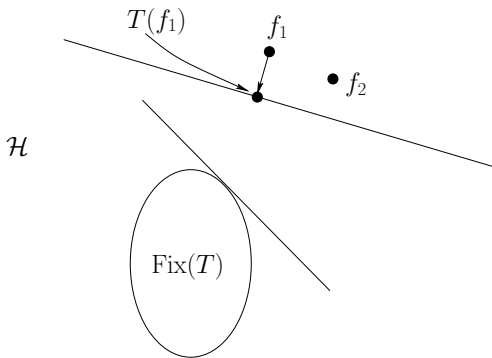
Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



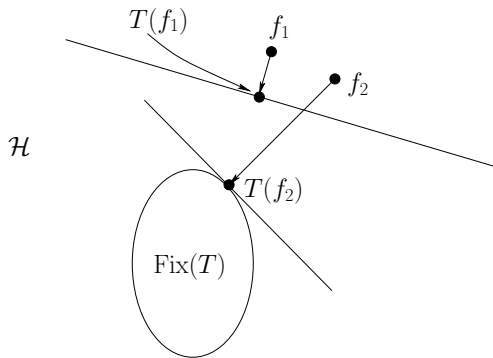
Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



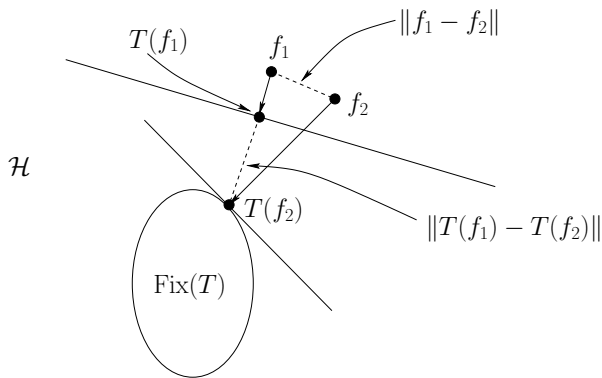
Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



Projecting onto **arbitrary separating hyperplanes** generates a quasi-nonexpansive mapping which is **not** nonexpansive.



# The Subgradient

## Definition (Subgradient)

Given a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , the subgradient,  $\Theta'(f)$ , is an element of  $\mathcal{H}$  such that

$$\langle g - f, \Theta'(f) \rangle + \Theta(f) \leq \Theta(g), \quad \forall g \in \mathcal{H}.$$

In other words, the **hyperplane**  $\{(g, \langle g - f, \Theta'(f) \rangle + \Theta(f)) : g \in \mathcal{H}\}$ , **supports** the graph of  $\Theta$  at the point  $(f, \Theta(f))$ .

## Definition (Level set)

$$\text{lev}_{\leq 0}(\Theta) := \{f \in \mathcal{H} : \Theta(f) \leq 0\}.$$



# The Subgradient

## Definition (Subgradient)

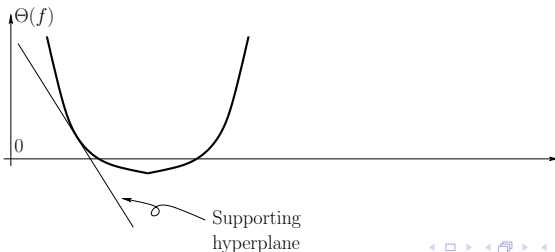
Given a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , the subgradient,  $\Theta'(f)$ , is an element of  $\mathcal{H}$  such that

$$\langle g - f, \Theta'(f) \rangle + \Theta(f) \leq \Theta(g), \quad \forall g \in \mathcal{H}.$$

In other words, the **hyperplane**  $\{(g, \langle g - f, \Theta'(f) \rangle + \Theta(f)) : g \in \mathcal{H}\}$ , **supports** the graph of  $\Theta$  at the point  $(f, \Theta(f))$ .

## Definition (Level set)

$$\text{lev}_{\leq 0}(\Theta) := \{f \in \mathcal{H} : \Theta(f) \leq 0\}.$$





# The Subgradient

## Definition (Subgradient)

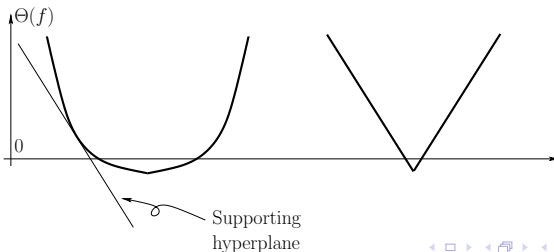
Given a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , the subgradient,  $\Theta'(f)$ , is an element of  $\mathcal{H}$  such that

$$\langle g - f, \Theta'(f) \rangle + \Theta(f) \leq \Theta(g), \quad \forall g \in \mathcal{H}.$$

In other words, the **hyperplane**  $\{(g, \langle g - f, \Theta'(f) \rangle + \Theta(f)) : g \in \mathcal{H}\}$ , **supports** the graph of  $\Theta$  at the point  $(f, \Theta(f))$ .

## Definition (Level set)

$$\text{lev}_{\leq 0}(\Theta) := \{f \in \mathcal{H} : \Theta(f) \leq 0\}.$$



# The Subgradient

## Definition (Subgradient)

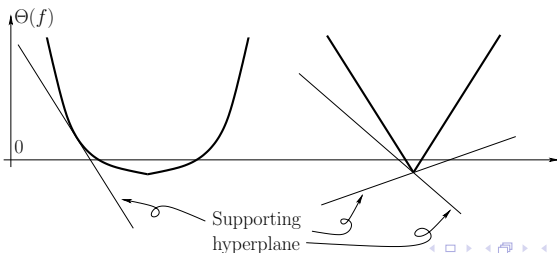
Given a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , the subgradient,  $\Theta'(f)$ , is an element of  $\mathcal{H}$  such that

$$\langle g - f, \Theta'(f) \rangle + \Theta(f) \leq \Theta(g), \quad \forall g \in \mathcal{H}.$$

In other words, the **hyperplane**  $\{(g, \langle g - f, \Theta'(f) \rangle + \Theta(f)) : g \in \mathcal{H}\}$ , **supports** the graph of  $\Theta$  at the point  $(f, \Theta(f))$ .

## Definition (Level set)

$$\text{lev}_{\leq 0}(\Theta) := \{f \in \mathcal{H} : \Theta(f) \leq 0\}.$$



# The Subgradient

## Definition (Subgradient)

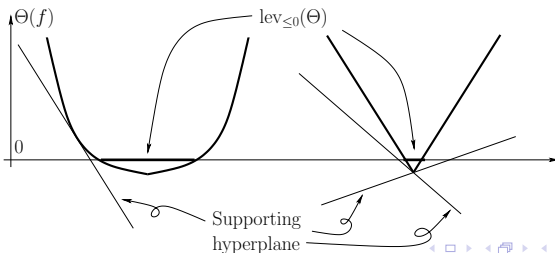
Given a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , the subgradient,  $\Theta'(f)$ , is an element of  $\mathcal{H}$  such that

$$\langle g - f, \Theta'(f) \rangle + \Theta(f) \leq \Theta(g), \quad \forall g \in \mathcal{H}.$$

In other words, the **hyperplane**  $\{(g, \langle g - f, \Theta'(f) \rangle + \Theta(f)) : g \in \mathcal{H}\}$ , **supports** the graph of  $\Theta$  at the point  $(f, \Theta(f))$ .

## Definition (Level set)

$$\text{lev}_{\leq 0}(\Theta) := \{f \in \mathcal{H} : \Theta(f) \leq 0\}.$$



# The Subgradient Projection Mapping

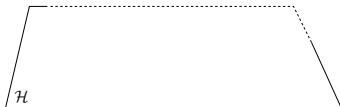
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# The Subgradient Projection Mapping

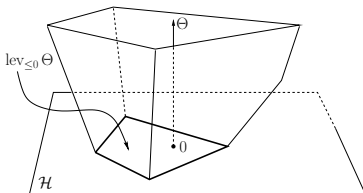
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# The Subgradient Projection Mapping

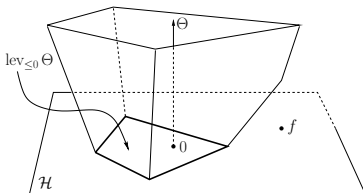
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# The Subgradient Projection Mapping

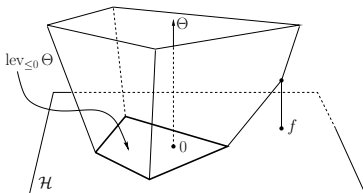
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# The Subgradient Projection Mapping

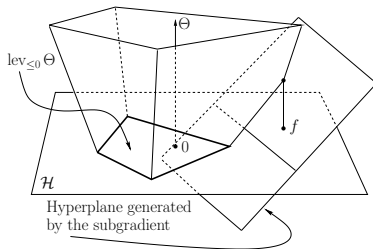
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.





# The Subgradient Projection Mapping

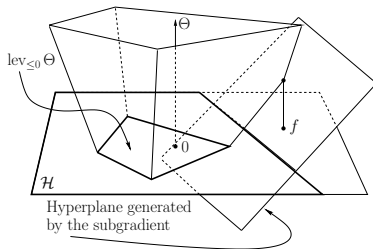
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# The Subgradient Projection Mapping

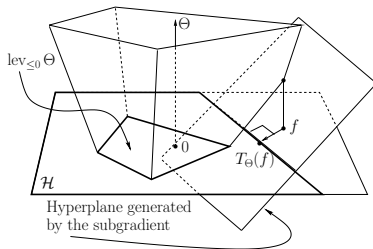
A Quasi-nonexpansive mapping

## Definition (Subgradient projection mapping)

Let a convex function  $\Theta : \mathcal{H} \rightarrow \mathbb{R}$ , with  $\text{lev}_{\leq 0}(\Theta) \neq \emptyset$ . Then, the subgradient projection mapping  $T_{\Theta} : \mathcal{H} \rightarrow \mathcal{H}$  is defined as follows:

$$T_{\Theta}(f) := \begin{cases} f - \frac{\Theta(f)}{\|\Theta'(f)\|^2} \Theta'(f), & \text{if } f \notin \text{lev}_{\leq 0}(\Theta), \\ f, & \text{if } f \in \text{lev}_{\leq 0}(\Theta). \end{cases}$$

The mapping  $T_{\Theta}$  is a **quasi-nonexpansive** one.



# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1), \lambda \in (0, 2),$

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1), \lambda \in (0, 2),$
- and any initial point  $f_0 \in \mathcal{H},$

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

APSM

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

APSM

$(f_n)$

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1), \lambda \in (0, 2),$
- and any initial point  $f_0 \in \mathcal{H},$

### APSM

$$\overbrace{\left( (1 - \lambda)I + \lambda T_{\Theta_n} \right)}^{\text{loss function info}} (f_n)$$

= relaxed subgradient  
projection  $T_{\Theta_n}$

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

### APSM

$$\underbrace{(\alpha R_n + (1 - \alpha)I)}_{= T_n: \text{averaged quasi-nonexpansive mapping}} \underbrace{((1 - \lambda)I + \lambda T_{\Theta_n})}_{= \text{relaxed subgradient projection } T_{\Theta_n}} (f_n)$$

The diagram shows the APSM iteration formula. It consists of two main parts in brackets, each with a descriptive label below it. The first part is  $(\alpha R_n + (1 - \alpha)I)$ , labeled "a-priori info" above and " $= T_n: \text{averaged quasi-nonexpansive mapping}$ " below. The second part is  $((1 - \lambda)I + \lambda T_{\Theta_n})$ , labeled "loss function info" above and " $= \text{relaxed subgradient projection } T_{\Theta_n}$ " below. The entire expression is applied to  $(f_n)$ .



# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

### APSM

$$f_{n+1} := \underbrace{(\alpha R_n + (1 - \alpha)I)}_{\substack{= T_n: \text{averaged} \\ \text{quasi-nonexpansive mapping}}} \underbrace{((1 - \lambda)I + \lambda T_{\Theta_n})}_{\substack{= \text{relaxed subgradient} \\ \text{projection } T_{\Theta_n}}}(f_n), \quad \forall n,$$

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

### APSM

$$f_{n+1} := \underbrace{\left( \alpha R_n + (1 - \alpha)I \right)}_{\substack{= T_n: \text{averaged} \\ \text{quasi-nonexpansive mapping}}} \underbrace{\left( (1 - \lambda)I + \lambda T_{\Theta_n} \right)}_{\substack{= \text{relaxed subgradient} \\ \text{projection } T_{\Theta_n}}} (f_n), \quad \forall n,$$

where

- $(R_n)_{n=0,1,\dots}$  is a sequence of quasi-nonexpansive mappings. This sequence of mappings comprises the a-priori information.

# Generalizing the Previous Methodologies

## Adaptive Projected Subgradient Method (APSM)

- For some user-defined
  - ▶  $\alpha \in (0, 1)$ ,  $\lambda \in (0, 2)$ ,
- and any initial point  $f_0 \in \mathcal{H}$ ,

### APSM

$$f_{n+1} := \underbrace{(\alpha R_n + (1 - \alpha)I)}_{\substack{= T_n: \text{averaged} \\ \text{quasi-nonexpansive mapping}}} \underbrace{((1 - \lambda)I + \lambda T_{\Theta_n})}_{\substack{= \text{relaxed subgradient} \\ \text{projection } T_{\Theta_n}}}(f_n), \quad \forall n,$$

where

- $(R_n)_{n=0,1,\dots}$  is a sequence of quasi-nonexpansive mappings. This sequence of mappings comprises the a-priori information.
- $(\Theta_n)_{n=0,1,\dots}$  is a sequence of loss/penalty function which quantifies the deviation of the sequential training data from the underlying model.

# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \sum_{i=n-q+1}^n \frac{\omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \frac{\sum_{i=n-q+1}^n \omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

Then, the APSM becomes:  $\forall n$ ,

# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \frac{\sum_{i=n-q+1}^n \omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

Then, the APSM becomes:  $\forall n$ ,

$$\left( f_n + \mu_n \left( \sum_{i=n-q+1}^n \omega_i^{(n)} P_{S_i[\epsilon]}(f_n) - f_n \right) \right),$$

# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \frac{\sum_{i=n-q+1}^n \omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

Then, the APSM becomes:  $\forall n$ ,

$$\underbrace{(\alpha R_n + (1 - \alpha)I)}_{\substack{T_n: \text{ averaged} \\ \text{quasi-nonexpansive mapping}}} \left( f_n + \mu_n \left( \sum_{i=n-q+1}^n \omega_i^{(n)} P_{S_i[\epsilon]}(f_n) - f_n \right) \right),$$

# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \frac{\sum_{i=n-q+1}^n \omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

Then, the APSM becomes:  $\forall n$ ,

$$f_{n+1} = \underbrace{(\alpha R_n + (1 - \alpha)I)}_{\substack{T_n: \text{averaged} \\ \text{quasi-nonexpansive mapping}}} \left( f_n + \mu_n \left( \sum_{i=n-q+1}^n \omega_i^{(n)} P_{S_i[\epsilon]}(f_n) - f_n \right) \right),$$



# Candidates for the Loss Functions $(\Theta_n)_{n=0,1,\dots}$

Given the current estimate  $f_n$ , define  $\forall f \in \mathcal{H}$ ,

$$\Theta_n(f) := \begin{cases} \frac{\sum_{i=n-q+1}^n \omega_i^{(n)} d(f_n, S_i[\epsilon])}{\sum_{j=n-q+1}^n \omega_j^{(n)} d(f_n, S_j[\epsilon])} d(f, S_i[\epsilon]), & \text{if } f \notin \bigcap_{i=n-q+1}^n S_i[\epsilon], \\ 0, & \text{otherwise.} \end{cases}$$

Then, the APSM becomes:  $\forall n$ ,

$$f_{n+1} = \underbrace{(\alpha R_n + (1 - \alpha)I)}_{\substack{T_n: \text{ averaged} \\ \text{quasi-nonexpansive mapping}}} \left( f_n + \mu_n \left( \sum_{i=n-q+1}^n \omega_i^{(n)} P_{S_i[\epsilon]}(f_n) - f_n \right) \right),$$

where the extrapolation coefficient  $\mu_n \in (0, 2\mathcal{M}_n)$  with

$$\mathcal{M}_n := \begin{cases} \frac{\sum_{j=n-q+1}^n \omega_j^{(n)} \|P_{S_j[\epsilon]}(f_n) - f_n\|^2}{\|\sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(f_n) - f_n\|^2}, & \text{if } \sum_{j=n-q+1}^n \omega_j^{(n)} P_{S_j[\epsilon]}(f_n) \neq f_n, \\ 1, & \text{otherwise.} \end{cases}$$

## Example (Examples of averaged quasi-nonexpansive mappings)

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).
  - ▶ The projection  $P_{B_{\ell_1}[\delta]}$  onto the  $\ell_1$  ball, (sparsity-aware learning).

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).
  - ▶ The projection  $P_{B_{\ell_1}[\delta]}$  onto the  $\ell_1$  ball, (sparsity-aware learning).
  - ▶ The projections  $(P_{B_{\ell_1}[\mathbf{w}_n, \delta]})_{n=0,1,\dots}$  onto a sequence of weighted  $\ell_1$  balls, (sparsity-aware learning).

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).
  - ▶ The projection  $P_{B_{\ell_1}[\delta]}$  onto the  $\ell_1$  ball, (sparsity-aware learning).
  - ▶ The projections  $(P_{B_{\ell_1}[\mathbf{w}_n, \delta]})_{n=0,1,\dots}$  onto a sequence of weighted  $\ell_1$  balls, (sparsity-aware learning).
- The composition of projections  $P_{C_1} \cdots P_{C_p}$ , where  $C_1, \dots, C_p$  are closed convex sets with  $\bigcap_{i=1}^p C_i \neq \emptyset$ .

## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).
  - ▶ The projection  $P_{B_{\ell_1}[\delta]}$  onto the  $\ell_1$  ball, (sparsity-aware learning).
  - ▶ The projections  $(P_{B_{\ell_1}[\mathbf{w}_n, \delta]})_{n=0,1,\dots}$  onto a sequence of weighted  $\ell_1$  balls, (sparsity-aware learning).
- The composition of projections  $P_{C_1} \cdots P_{C_p}$ , where  $C_1, \dots, C_p$  are closed convex sets with  $\bigcap_{i=1}^p C_i \neq \emptyset$ .
  - ▶ The composition of projections  $P_{\Pi}P_K$ , where  $\Pi$  is a hyperplane and  $K$  is an icecream cone, (beamforming).



## Example (Examples of averaged quasi-nonexpansive mappings)

- The projection  $P_C$  onto a closed convex set  $C$  of  $\mathcal{H}$ .
  - ▶ The projection  $P_V$  onto an affine set  $V$  of  $\mathbb{R}^m$ , (beamforming).
  - ▶ The projection  $P_{B_{\ell_1}[\delta]}$  onto the  $\ell_1$  ball, (sparsity-aware learning).
  - ▶ The projections  $(P_{B_{\ell_1}[\mathbf{w}_n, \delta]})_{n=0,1,\dots}$  onto a sequence of weighted  $\ell_1$  balls, (sparsity-aware learning).
- The composition of projections  $P_{C_1} \cdots P_{C_p}$ , where  $C_1, \dots, C_p$  are closed convex sets with  $\bigcap_{i=1}^p C_i \neq \emptyset$ .
  - ▶ The composition of projections  $P_{\Pi}P_K$ , where  $\Pi$  is a hyperplane and  $K$  is an icecream cone, (beamforming).
- The composition  $P_{\mathcal{K}}\left(I - \lambda\left(I - \sum_{i=1}^{p-1} \beta_i P_{C_i}\right)\right)$ ,  $\lambda \in (0, 2)$ , where  $\mathcal{K} \cap \left(\bigcap_{i=1}^{p-1} C_i\right) = \emptyset$ , (beamforming).

- Surprisingly, the APSM retains its performance and theoretical properties in the case where the Generalized Thresholding mapping  $T_{GT}$  is used in the place of  $T_n$ !

- Surprisingly, the APSM retains its performance and theoretical properties in the case where the Generalized Thresholding mapping  $T_{GT}$  is used in the place of  $T_n$ !
- Recall that  $\text{Fix}(T_{GT})$  is a union of subspaces, which is a non-convex set.

- Surprisingly, the APSM retains its performance and theoretical properties in the case where the Generalized Thresholding mapping  $T_{GT}$  is used in the place of  $T_n$ !
- Recall that  $\text{Fix}(T_{GT})$  is a union of subspaces, which is a non-convex set.
- Such an application motivates the extension of the concept of a quasi-nonexpansive mapping to that of a partially quasi-nonexpansive one<sup>10</sup>.

---

<sup>10</sup>[Kopsinis et al '11a].

# Theoretical Properties

Define at  $n \geq 0$ ,  $\Omega_n := \text{Fix}(T_n) \cap \text{lev}_{\leq 0} \Theta_n$ . Let  $\Omega := \bigcap_{n \geq n_0} \Omega_n \neq \emptyset$ , for some nonnegative integer  $n_0$ . Assume also that  $\frac{\mu_n}{M_n} \in [\epsilon_1, 2 - \epsilon_2]$ ,  $\forall n \geq n_0$ , for some sufficiently small  $\epsilon_1, \epsilon_2 > 0$ . Under the addition of some mild assumptions, the following statements hold true<sup>11</sup>.

- **Monotone approximation.**  $d(f_{n+1}, \Omega) \leq d(f_n, \Omega)$ ,  $\forall n \geq n_0$ .
- **Asymptotic minimization.**  $\lim_{n \rightarrow \infty} \Theta_n(f_n) = 0$ .
- **Cluster points.** If we assume that the set of all sequential strong cluster points  $\mathfrak{S}((f_n)_{n=0,1,\dots})$  is nonempty, then

$$\mathfrak{S}((f_n)_{n=0,1,\dots}) \subset \limsup_{n \rightarrow \infty} \text{Fix}(T_n) \cap \limsup_{n \rightarrow \infty} \text{lev}_{\leq 0}(\Theta_n),$$

where  $\limsup_{n \rightarrow \infty} A_n := \bigcap_{r>0} \bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} (A_k + B[0, r])$ , and  $B[0, r]$  is a closed ball of center 0 and radius  $r$ .

- **Strong convergence.** Assume that there exists a hyperplane  $\Pi \subset \mathcal{H}$  such that  $\text{ri}_{\Pi}(\Omega) \neq \emptyset$ . Then, there exists an  $f_* \in \mathcal{H}$  such that  $\lim_{n \rightarrow \infty} f_n = f_*$ .

---

<sup>11</sup>[Slavakis, Yamada, '11].

- K. Slavakis, P. Bouboulis, and S. Theodoridis. Adaptive multiregression in reproducing kernel Hilbert spaces: the multiaccess MIMO channel case. *IEEE Trans. Neural Networks and Learning Systems*, 23(2): 260–276, Feb. 2012,
- Y. Kopsinis, K. Slavakis, S. Theodoridis, and S. McLaughlin. Generalized thresholding sparsity-aware online learning in a union of subspaces, 2011, <http://arxiv.org/abs/1112.0665>.
- K. Slavakis and I. Yamada. The adaptive projected subgradient method constrained by families of quasi-nonexpansive mappings and its application to online learning, 2011, <http://arxiv.org/abs/1008.5231>.
- S. Theodoridis, K. Slavakis, and I. Yamada. Adaptive learning in a world of projections: a unifying framework for linear and nonlinear classification and regression tasks. *IEEE Signal Processing Magazine*, 28(1):97–123, 2011.
- Y. Kopsinis, K. Slavakis, and S. Theodoridis. Online sparse system identification and signal reconstruction using projections onto weighted  $\ell_1$  balls. *IEEE Trans. Signal Processing*, 59(3):936–952, 2011.
- S. Chouvardas, K. Slavakis, and S. Theodoridis. Adaptive robust distributed learning in diffusion sensor networks. *IEEE Trans. Signal Processing*, 59(10): 4692–4707, 2011.

- K. Slavakis, S. Theodoridis, and I. Yamada. Adaptive constrained learning in reproducing kernel Hilbert spaces: the robust beamforming case. *IEEE Trans. Signal Processing*, 57(12):4744–4764, Dec. 2009.
- K. Slavakis, S. Theodoridis, and I. Yamada. Online kernel-based classification using adaptive projection algorithms. *IEEE Trans. Signal Processing*, 56(7), Part 1: 2781–2796, July 2008.
- K. Slavakis and I. Yamada. Robust wideband beamforming by the hybrid steepest descent method. *IEEE Trans. Signal Processing*, 55(9): 4511–4522, Sept. 2007.
- K. Slavakis, I. Yamada, and N. Ogura. The adaptive projected subgradient method over the fixed point set of strongly attracting nonexpansive mappings. *Numerical Functional Analysis and Optimization*, 27(7&8):905–930, 2006.
- I. Yamada and N. Ogura. Adaptive projected subgradient method for asymptotic minimization of sequence of nonnegative convex functions. *Numerical Functional Analysis and Optimization*, 25(7&8):593–617, 2004.
- I. Yamada, K. Slavakis, and K. Yamada. An efficient robust adaptive filtering algorithm based on parallel subgradient projection techniques. *IEEE Trans. Signal Processing*, 50(5): 1091–1101, May 2002.

## Matlab code

`http://users.uop.gr/~slavakis/publications.htm`