Aggregation and Empirical Risk Minimization

Guillaume Lecué (joint works with Shahar Mendelson)

CNRS, LATP Marseille

▲ロト ▲団ト ▲ヨト ▲ヨト 三ヨー のへで

Model



Introduction	Optimal aggregation via ERM	Lower bound for the ERM	Applications
Model			
		\mathcal{X} : measurable space μ : probability measure X_1, \ldots, X_n <i>n</i> i.i.d. r.v. $\sim \mu$	
	X ₂		
	(\mathcal{X},μ)		
	X,		
	•^	< ㅁ > < 团 > < 분 > < 분 >	ह १९९७



Aggregation and Empirical Risk Minimization

• Problem of prediction

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

• Problem of prediction

 $x \in \mathcal{X}$



CNRS, LATP Marseille

Sac

• Problem of prediction

 $x \in \mathcal{X} \longrightarrow y \in \mathbb{R}$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

• Problem of prediction

 $x \in \mathcal{X} \longrightarrow y \in \mathbb{R}$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small.

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Problem of prediction

$$\kappa \in \mathcal{X} \quad \longrightarrow \quad y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$ **a** Agnostic learning

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Problem of prediction

$$x \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

• Agnostic learning No assumption on Y

F class of functions : $\mathcal{X} \mapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

Agnostic learning No assumption on Y

F class of functions : $\mathcal{X} \mapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

$$\inf_{f\in F} \mathbb{E}(f(X)-Y)^2$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

 Agnostic learning No assumption on Y
E close of functions + X +

F class of functions : $\mathcal{X} \mapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

 $\inf_{f\in F}\mathbb{E}(f(X)-Y)^2$

• Problem of aggregation

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

• Agnostic learning No assumption on Y

F class of functions : $\mathcal{X} \mapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

 $\inf_{f\in F}\mathbb{E}(f(X)-Y)^2$

• Problem of aggregation |F| = M;

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

• Agnostic learning No assumption on Y

F class of functions : $\mathcal{X} \mapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

 $\inf_{f\in F}\mathbb{E}(f(X)-Y)^2$

• Problem of aggregation $|F| = M; F = \{f_1, \dots, f_M\};$

Problem of prediction

$$\kappa \in \mathcal{X} \longrightarrow y \in \mathbb{R}$$

construction of $\hat{f}(\cdot, (X_i, Y_i)_i)$ s.t. $\mathbb{E}(\hat{f}(X) - Y)^2$ small. $\mathbb{E}(\hat{f}(X) - Y)^2 = \mathbb{E}(\hat{f}(X) - \mathbb{E}(Y|X))^2 + \mathbb{E}(\mathbb{E}(Y|X) - Y)^2$ $\Rightarrow \mathbb{E}(Y|\cdot) \in \text{Class of functions (smooth, entropy,..)}$

• Agnostic learning No assumption on Y

F class of functions : $\mathcal{X} \longmapsto \mathbb{R} \Rightarrow$ best predictor of Y in F

 $\inf_{f\in F}\mathbb{E}(f(X)-Y)^2$

• Problem of aggregation |F| = M; $F = \{f_1, \dots, f_M\}$; construction of a procedure which has a risk as close as possible to

$$\min_{j=1,\dots,M} \mathbb{E}(f_j(X) - Y)^2$$

Problem of aggregation : Construction of a procedure $\tilde{f}(\cdot, D_n)$ such that

$$\mathbb{E}(\hat{f}(X) - Y)^2 \leq \min_{f \in F} \mathbb{E}(f(X) - Y)^2 + r_n(F)$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Problem of aggregation : Construction of a procedure $\tilde{f}(\cdot, D_n)$ such that

$$\mathbb{E}(\widehat{f}(X) - Y)^2 \leq \min_{f \in F} \mathbb{E}(f(X) - Y)^2 + r_n(F)$$

Empirical risk :

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

naa

Problem of aggregation : Construction of a procedure $\tilde{f}(\cdot, D_n)$ such that

$$\mathbb{E}(\widehat{f}(X) - Y)^2 \leq \min_{f \in F} \mathbb{E}(f(X) - Y)^2 + r_n(F)$$

Empirical risk :

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

$$\mathbb{E}[R_n(f)] = \mathbb{E}(f(X) - Y)^2 := R(f)$$

Empirical risk minimization algorithm :

$$\hat{f}^{ERM} \in \operatorname{Arg}\min_{f \in F} R_n(f)$$

Aggregation and Empirical Risk Minimization

Empirical risk minimization algorithm over the convex hull of F:

 $\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f \in \operatorname{Conv}(F)} R_n(f),$

Empirical risk minimization algorithm over the convex hull of F:

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in \operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$

Empirical risk minimization algorithm over the convex hull of F:

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in \operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ Aggregation with exponential weights :

$$\tilde{f} := \sum_{j=1}^{M} w_{j,n} f_j$$

Empirical risk minimization algorithm over the convex hull of F:

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in \operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ Aggregation with exponential weights :

$$\tilde{f} := \sum_{j=1}^{M} w_{j,n} f_j \text{ where } w_{j,n} = \frac{\exp(-nR_n(f_j)/T)}{\sum_{k=1}^{M} \exp(-nR_n(f_k)/T)}$$

T : temperature parameter.

Empirical risk minimization algorithm over the convex hull of F:

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ Aggregation with exponential weights :

$$\tilde{f} := \sum_{j=1}^{M} w_{j,n} f_j \text{ where } w_{j,n} = \frac{\exp(-nR_n(f_j)/T)}{\sum_{k=1}^{M} \exp(-nR_n(f_k)/T)}$$

T : temperature parameter. Cumulative Aggregation with exponential weights :

$$\tilde{f} := \sum_{j=1}^{M} \bar{w}_j f_j$$

Empirical risk minimization algorithm over the convex hull of F:

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ Aggregation with exponential weights :

$$\tilde{f} := \sum_{j=1}^{M} w_{j,n} f_j \text{ where } w_{j,n} = \frac{\exp(-nR_n(f_j)/T)}{\sum_{k=1}^{M} \exp(-nR_n(f_k)/T)}$$

T : temperature parameter. Cumulative Aggregation with exponential weights :

$$ilde{f} := \sum_{j=1}^M ar{w}_j f_j ext{ where } ar{w}_j = rac{1}{n} \sum_{p=1}^n w_{j,p}.$$

Optimal rate of aggregation

Definition

 $\psi(n, M)$: optimal rate of aggregation with confidence $0 < \delta < 1/2$ \tilde{f} : optimal aggregation procedure with confidence δ

Optimal rate of aggregation

Definition

 $\psi(n, M)$: optimal rate of aggregation with confidence $0 < \delta < 1/2$ \tilde{f} : optimal aggregation procedure with confidence δ

• $\forall n, M, \forall F$ of cardinality M and any target Y (all bounded by b), with ν^n -probability at least $1 - \delta$,

$$R(\tilde{f}) \leq \min_{f \in F} R(f) + c_1(\delta) \psi(n, M),$$

Optimal rate of aggregation

Definition

 $\psi(n, M)$: optimal rate of aggregation with confidence $0 < \delta < 1/2$ \tilde{f} : optimal aggregation procedure with confidence δ

• $\forall n, M, \forall F$ of cardinality M and any target Y (all bounded by b), with ν^n -probability at least $1 - \delta$,

$$R(\tilde{f}) \leq \min_{f \in F} R(f) + c_1(\delta)\psi(n, M),$$

∃c₂ > 0 s.t. ∀M, n, ∀f̄, ∃F of cardinality M and a target Y (all bounded by b) such that, with νⁿ-probability at least 1/2,

$$R(\bar{f}) \geq \min_{f \in F} R(f) + c_2 \psi(n, M).$$

Lower bound

 $\exists c_2 > 0 \text{ s.t. } \forall M, n, \forall \overline{f}, \exists F \text{ of cardinality } M \text{ and a target } Y \text{ (all bounded by } b) \text{ such that, with } \nu^n\text{-probability at least } 1/2,$

 $R(\overline{f}) \geq \min_{f \in F} R(f) + c_2 \frac{\log M}{n}.$

Lower bound

 $\exists c_2 > 0 \text{ s.t. } \forall M, n, \forall \overline{f}, \exists F \text{ of cardinality } M \text{ and a target } Y \text{ (all bounded by } b) \text{ such that, with } \nu^n\text{-probability at least } 1/2,$

$$R(\overline{f}) \geq \min_{f \in F} R(f) + c_2 \frac{\log M}{n}.$$

Given n observations and a dictionary F of cardinality M, the minimum price that one has to pay to mimic the best function in F is at least

 $\frac{\log M}{n}$

Lower bound

 $\exists c_2 > 0 \text{ s.t. } \forall M, n, \forall \overline{f}, \exists F \text{ of cardinality } M \text{ and a target } Y \text{ (all bounded by } b) \text{ such that, with } \nu^n\text{-probability at least } 1/2,$

$$R(\overline{f}) \geq \min_{f \in F} R(f) + c_2 \frac{\log M}{n}.$$

Given n observations and a dictionary F of cardinality M, the minimum price that one has to pay to mimic the best function in F is at least

$\frac{\log M}{n}$

Is it possible to achieve this rate?

Suboptimality of ERM over F and Conv(F)

• $\forall n, M, \exists F \text{ of cardinality } M \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERM}) \ge \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

Suboptimality of ERM over F and Conv(F)

• $\forall n, M, \exists F \text{ of cardinality } M \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERM}) \ge \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

2 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}.$$

Suboptimality of ERM over F and Conv(F)

• $\forall n, M, \exists F \text{ of cardinality } M \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERM}) \ge \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

2 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

1

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}$$

 $(1/\sqrt{n}) >> (\log M)/n \Longrightarrow \hat{f}^{ERM}$ and $\hat{f}^{ERMconv}$ cannot achieve the rate $(\log M)/n$.

Suboptimality of ERM over Conv(F)

First remark :

A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F).

naa

Suboptimality of ERM over Conv(F)

First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$
First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

 \implies the ERM over F is optimal among all the selectors.

First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

 \implies the ERM over F is optimal among all the selectors.

There exists some convex combinations of the elements in F which are optimal aggregation procedures.

First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

 \implies the ERM over F is optimal among all the selectors.

There exists some convex combinations of the elements in F which are optimal aggregation procedures.

idea :

First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

 \implies the ERM over F is optimal among all the selectors.

There exists some convex combinations of the elements in F which are optimal aggregation procedures.

idea :

 \implies the ERM over Conv(F) is optimal among all the convex aggregates;

First remark : A Selector is an aggregation procedure taking its values in the dictionary F (for instance the ERM over F). $\forall \bar{f}_n$ (selector), $\forall n, M, \exists F$ of cardinality M such that, w.p.g. 1/2

$$R(\overline{f}_n) \geq \min_{f \in F} R(f) + c \sqrt{\frac{\log M}{n}}$$

 \implies the ERM over F is optimal among all the selectors.

There exists some convex combinations of the elements in F which are optimal aggregation procedures.

idea :

 \implies the ERM over Conv(F) is optimal among all the convex aggregates;

 \implies the ERM over Conv(F) is an optimal aggregation procedure.

< □ > < 同

Suboptimality of ERM over Conv(F)

Second remark :

 $\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f \in \operatorname{Conv}(F)} R_n(f),$

글 🕨 🖌 글

Suboptimality of ERM over $\underline{Conv}(F)$

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in \operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ \implies the target of the ERM over $\operatorname{Conv}(F)$ is $\min_{f \in \operatorname{Conv}(F)} R(f)$.

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ \implies the target of the ERM over $\operatorname{Conv}(F)$ is $\min_{f \in \operatorname{Conv}(F)} R(f)$.

$$R(\hat{f}^{ERMconv}) \leq \min_{f \in \text{Conv}(F)} R(f) + c\phi(M, n)$$

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ \implies the target of the ERM over $\operatorname{Conv}(F)$ is $\min_{f \in \operatorname{Conv}(F)} R(f)$.

$$R(\hat{f}^{ERMconv}) \leq \min_{f \in \text{Conv}(F)} R(f) + c\phi(M, n)$$

$$\leq \min_{f \in F} R(f) + \left[c\phi(M, n) + \min_{f \in F} R(f) - \min_{f \in \text{Conv}(F)} R(f) \right]$$

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ \implies the target of the ERM over $\operatorname{Conv}(F)$ is $\min_{f \in \operatorname{Conv}(F)} R(f)$.

$$R(\hat{f}^{ERMconv}) \leq \min_{f \in \text{Conv}(F)} R(f) + c\phi(M, n)$$

$$\leq \min_{f \in F} R(f) + \left[c\phi(M, n) + \min_{f \in F} R(f) - \min_{f \in \text{Conv}(F)} R(f) \right]$$

 \implies gain in the approximation error : $\min_{f \in F} R(f) - \min_{f \in \text{Conv}(F)} R(f)$;

Second remark :

$$\hat{f}^{ERMconv} \in \operatorname{Arg}\min_{f\in\operatorname{Conv}(F)} R_n(f),$$

where $\operatorname{Conv}(F) = \{\sum_{j=1}^{M} \lambda_j f_j : \lambda_j \ge 0 \text{ and } \sum \lambda_j = 1\}.$ \implies the target of the ERM over $\operatorname{Conv}(F)$ is $\min_{f \in \operatorname{Conv}(F)} R(f)$.

$$R(\hat{f}^{ERMconv}) \leq \min_{f \in \operatorname{Conv}(F)} R(f) + c\phi(M, n)$$

$$\leq \min_{f \in F} R(f) + \left[c\phi(M, n) + \min_{f \in F} R(f) - \min_{f \in \operatorname{Conv}(F)} R(f) \right]$$

⇒ gain in the approximation error : $\min_{f \in F} R(f) - \min_{f \in \text{Conv}(F)} R(f)$; ⇒ the ERM over Conv(F) is an optimal aggregation procedure.

Applications

Suboptimality of ERM over Conv(F)

 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}.$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}.$$

 $(\phi_i)_{i \in \mathbb{N}}$ sequence of i.i.d Rademacher r.v.

 $\mathcal{F} := \{0, \pm \phi_1, \dots, \pm \phi_M\}$ and target $: \mathcal{Y} := \phi_{M+1}$

Applications

Suboptimality of ERM over Conv(F)

 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}.$$

 $(\phi_i)_{i \in \mathbb{N}}$ sequence of i.i.d Rademacher r.v.

 $F := \{0, \pm \phi_1, \dots, \pm \phi_M\}$ and target $: Y := \phi_{M+1}$

In the approximation error

$$\min_{f\in F} R(f) = \min_{f\in \operatorname{Conv}(F)} R(f)$$

Aggregation and Empirical Risk Minimization

 $\forall n, \exists F \text{ of cardinality } c\sqrt{n} \text{ such that, w.p.g. } 1/2$

$$R(\hat{f}^{ERMconv}) \geq \min_{f \in F} R(f) + \sqrt{\frac{c}{n}}.$$

 $(\phi_i)_{i \in \mathbb{N}}$ sequence of i.i.d Rademacher r.v.

$$F := \{0, \pm \phi_1, \dots, \pm \phi_M\}$$
 and target $: Y := \phi_{M+1}$

In the approximation error

$$\min_{f\in F} R(f) = \min_{f\in \operatorname{Conv}(F)} R(f)$$

2 maximization of the complexity of Conv(F)

< □ > < 同



Aggregation and Empirical Risk Minimization

3 CNRS, LATP Marseille

Sac

Э



Aggregation and Empirical Risk Minimization

э

Sac





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

(ロ)、(型)、(E)、(E)、 E) の(()



$$\tilde{f} \in \operatorname{Arg\,min}_{f \in \operatorname{Conv}(\hat{F}_1)} R_n(f)$$



<ロ> (日) (日) (日) (日) (日)

Data splitting

2*n* observations : $D_{2n} = ((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}), \dots, (X_{2n}, Y_{2n})).$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

3

Data splitting

2*n* observations : $D_{2n} = ((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}), \dots, (X_{2n}, Y_{2n})).$

 $D_1 = ((X_1, Y_1), \dots, (X_n, Y_n))$: construction of the set \hat{F}_1 of approximatively good elements of F.

CNRS, LATP Marseille

Data splitting

2*n* observations : $D_{2n} = ((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}), \dots, (X_{2n}, Y_{2n})).$

 $D_1 = ((X_1, Y_1), \dots, (X_n, Y_n))$: construction of the set \hat{F}_1 of approximatively good elements of F.

 $D_2 = ((X_{n+1}, Y_{n+1}), \dots, (X_{2n}, Y_{2n})) :$ construction of the ERM over $\operatorname{Conv}(\hat{F}_1)$

$$\widetilde{f} \in \operatorname{Arg} \min_{f \in \operatorname{Conv}(\widehat{F}_1)} \frac{1}{n} \sum_{i=n+1}^{2n} (f(X_i) - Y_i)^2.$$

Construction of \hat{F}_1 (preselection step)

```
1st step : ERM over F (over D_1) :
```

 $\hat{f} \in \operatorname{Arg}\min_{f \in F} R_n(f)$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

Lower bound for the ERM

Applications

Construction of \hat{F}_1 (preselection step)

```
1st step : ERM over F (over D_1) :
```

```
\hat{f} \in \operatorname{Arg\,min}_{f \in F} R_n(f)
```

2nd step : set of almost minimizers of the ER (over D_1) :

$$\hat{F}_1 = \Big\{ f \in F : R_n(f) \le R_n(\hat{f}) + C_1 \max(\alpha \| \hat{f} - f \|_{L_2^n}, \alpha^2) \Big\},\$$

where $\alpha = \sqrt{(x + \log M)/n}$, x > 0 is the confidence and

$$\|\hat{f} - f\|_{L_2^n}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - f(X_i))^2.$$

< □ > < 同

ERM over $\operatorname{Conv}(\hat{F}_1)$

$$\operatorname{Conv}(\hat{F}_1) = \Big\{ \sum_{f \in \hat{F}_1} \lambda_f f : \lambda_f \ge 0 \text{ and } \sum_{f \in \hat{F}_1} \lambda_f = 1 \Big\},$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Э

н

ERM over $\operatorname{Conv}(\hat{F}_1)$

$$\operatorname{Conv}(\hat{F}_1) = \Big\{ \sum_{f \in \hat{F}_1} \lambda_f f : \lambda_f \ge 0 \text{ and } \sum_{f \in \hat{F}_1} \lambda_f = 1 \Big\},$$

Construction of the ERM over $Conv(F_1)$ with ER based on D_2 :

$$\tilde{f} \in \operatorname{Arg} \min_{f \in \operatorname{Conv}(\tilde{F}_1)} \frac{1}{n} \sum_{i=n+1}^{2n} (f(X_i) - Y_i)^2.$$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

Exact Oracle Inequality

Theorem

 $\forall x > 0, \forall F \text{ of } M \text{ functions, any target } Y \text{ (all bounded by b), with } \nu^{2n}\text{-probability at least } 1 - 2 \exp(-x),$

$$R(\tilde{f}) \leq \min_{f \in F} R(f) + c_1(1+x) \frac{\log M}{n}$$

Exact Oracle Inequality

Theorem

 $\forall x > 0, \forall F \text{ of } M \text{ functions, any target } Y \text{ (all bounded by b), with } \nu^{2n}\text{-probability at least } 1 - 2 \exp(-x),$

$$R(\tilde{f}) \leq \min_{f \in F} R(f) + c_1(1+x) \frac{\log M}{n}$$

Conclusion : $(\log M)/n$ is the optimal rate of aggregation and \tilde{f} is an optimal aggregation procedure with confidence $2 \exp(-x)$.

Exact Oracle Inequality

Theorem

 $\forall x > 0, \forall F \text{ of } M \text{ functions, any target } Y \text{ (all bounded by b), with } \nu^{2n}\text{-probability at least } 1 - 2 \exp(-x),$

$$R(\tilde{f}) \leq \min_{f \in F} R(f) + c_1(1+x) \frac{\log M}{n}$$

Conclusion : $(\log M)/n$ is the optimal rate of aggregation and \hat{f} is an optimal aggregation procedure with confidence $2 \exp(-x)$. Remark : This aggregation procedure is sparse in the sense that non-relevant elements in F have a zero-coefficient.

Model : agnostic functional learning

• Observations $(X_i, T(X_i))_{i=1,...,n}$ (Y = T(X)); T : Target function



CNRS, LATP Marseille

Model : agnostic functional learning

- Observations $(X_i, T(X_i))_{i=1,...,n}$ (Y = T(X)); T : Target function
- $T \in \tau$; τ is convex;

Model : agnostic functional learning

- Observations $(X_i, T(X_i))_{i=1,...,n}$ (Y = T(X)); T : Target function
- $T \in \tau$; τ is convex;
- *F* (set of candidates) $\subset \tau$;
Model : agnostic functional learning

- Observations $(X_i, T(X_i))_{i=1,...,n}$ (Y = T(X)); T : Target function
- $T \in \tau$; τ is convex;
- *F* (set of candidates) $\subset \tau$;

agnostic functional learning problem : Find a procedure which has a risk as close as possible to

 $\inf_{f\in F} \mathbb{E}(f(X) - T(X))^2$







Empirical risk minimization algorithm

Exact Oracle Inequality : Construction of a procedure $\hat{f}(\cdot, D)$ such that $\mathbb{E}(\hat{f}(X) - T(X))^2 \leq \inf_{f \in F} \mathbb{E}(f(X) - T(X))^2 + r_n(F)$

CNRS, LATP Marseille

Empirical risk minimization algorithm

Exact Oracle Inequality : Construction of a procedure $\hat{f}(\cdot, D)$ such that $\mathbb{E}(\hat{f}(X) - T(X))^2 \leq \inf_{f \in F} \mathbb{E}(f(X) - T(X))^2 + r_n(F)$

Empirical risk minimization algorithm :

$$\hat{f} \in \operatorname{Arg}\min_{f \in F} \frac{1}{n} \sum_{i=1}^{n} (f(X_i) - T(X_i))^2$$

< □ > < 同

Upper bound

Symmetrization

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

э

н

Symmetrization+majoration of Rademacher processes by Gaussian processes

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

$$\label{eq:symmetrization} \begin{split} & \mathsf{Symmetrization} + \mathsf{majoration} \text{ of Rademacher processes by Gaussian} \\ & \mathsf{processes} + \mathsf{Dudley's entropy integral}: \end{split}$$

$$\label{eq:symmetrization} \begin{split} & \mathsf{Symmetrization} + \mathsf{majoration} \text{ of Rademacher processes by Gaussian} \\ & \mathsf{processes} + \mathsf{Dudley's entropy integral}: \end{split}$$

$$\mathbb{E}(\widehat{f}(X) - \mathcal{T}(X))^2 - \inf_{f \in \mathcal{F}} \mathbb{E}(f(X) - \mathcal{T}(X))^2$$

 $\lesssim \frac{1}{\sqrt{n}} \mathbb{E}_X \int_0^{\operatorname{diam}(P_\sigma \mathcal{F}, |\cdot|_{2,n})} \sqrt{\mathcal{N}(P_\sigma \mathcal{F}, |\cdot|_{2,n}, \epsilon)} d\epsilon$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

$$\label{eq:symmetrization} \begin{split} & \mathsf{Symmetrization} + \mathsf{majoration} \ \text{of Rademacher processes by Gaussian} \\ & \mathsf{processes} + \mathsf{Dudley's \ entropy \ integral} : \end{split}$$

$$\mathbb{E}(\widehat{f}(X) - \mathcal{T}(X))^2 - \inf_{f \in \mathcal{F}} \mathbb{E}(f(X) - \mathcal{T}(X))^2$$

 $\lesssim rac{1}{\sqrt{n}} \mathbb{E}_X \int_0^{\operatorname{diam}(P_\sigma \mathcal{F}, |\cdot|_{2,n})} \sqrt{\mathcal{N}(P_\sigma \mathcal{F}, |\cdot|_{2,n}, \epsilon)} d\epsilon$

 $\mathcal{F} = \{(f - T)^2 - (f^* - T)^2 : f \in F\} \quad (\text{excess loss class})$ $P_{\sigma}\mathcal{F} = \{(g(X_1), \dots, g(X_n)) : g \in \mathcal{F}\} \quad (\text{coordinate projection})$ $|u|_{2,n} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} u_i^2} \quad (\text{normalized } l_2^n \text{-norm})$

Aim

Lower bound for the excess risk of the Empirical risk minimization algorithm

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

Aim

Lower bound for the excess risk of the Empirical risk minimization algorithm

 $\exists T \in \tau \text{ s.t.}$

$$\mathbb{E}(\widehat{f}(X) - T(X))^2 - \inf_{f \in F} \mathbb{E}(f(X) - T(X))^2 \ge r_n(F')$$

where $F' \subset F$.



CNRS, LATP Marseille

Sac

Aim

Lower bound for the excess risk of the Empirical risk minimization algorithm

 $\exists T \in \tau \text{ s.t.}$

$$\mathbb{E}(\widehat{f}(X) - T(X))^2 - \inf_{f \in F} \mathbb{E}(f(X) - T(X))^2 \ge r_n(F')$$

where $F' \subset F$. Assumption : $\exists T_0 \in \tau \text{ s.t. } card(\mathcal{M}(T_0)) \ge 2$ where

 $\mathcal{M}(T) = \{f \in F : \mathbb{E}(f(X) - T(X))^2 = \inf_{f \in F} \mathbb{E}(f(X) - T(X))^2\}$











Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Outline of the proof

The core of the proof is to find a set Q that can "compete" with $B_r = \{f \in F : \mathbb{E}\mathcal{L}_{\lambda}(f) \leq r\}$ in the sense that the empirical excess risk function

$$\mathcal{E}_n: f \in F \longmapsto rac{1}{n} \sum_{i=1}^n \mathcal{L}_\lambda(f)(X_i) = P_n \mathcal{L}_\lambda(f)$$

will be more negative on Q than on it can possibly be on B_r $(\mathcal{L}_{\lambda}(f) := (f - T_{\lambda})^2 - (f^* - T_{\lambda})^2).$

Outline of the proof

The core of the proof is to find a set Q that can "compete" with $B_r = \{f \in F : \mathbb{E}\mathcal{L}_{\lambda}(f) \leq r\}$ in the sense that the empirical excess risk function

$$\mathcal{E}_n: f \in F \longmapsto rac{1}{n} \sum_{i=1}^n \mathcal{L}_\lambda(f)(X_i) = P_n \mathcal{L}_\lambda(f)$$

will be more negative on Q than on it can possibly be on B_r $(\mathcal{L}_{\lambda}(f) := (f - T_{\lambda})^2 - (f^* - T_{\lambda})^2).$ Thus, the ERM $\hat{f}_{\lambda} \notin B_r$, and thus, with a certain probability,

 $\mathbb{E}\big[\mathcal{L}_{\lambda}(\hat{f}_{\lambda})|D\big] > r.$

Proof in two parts :

- \mathcal{E}_n is likely to be very negative on $\{f \in F : \mathbb{E}\mathcal{L}(f) = \min_{f \in F} \mathbb{E}\mathcal{L}(f)\}$;
- find some r on which the oscillations of \mathcal{E}_n in B_r are small.

Gaussian process : Let $Q \subset L_2(\mu)$.

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Sac

Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N}).$



Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N})$. A common measure of the "complexity" of Q is

 $H(Q) := \mathbb{E} \sup_{q \in Q} G_q$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N})$. A common measure of the "complexity" of Q is

$$H(Q) := \mathbb{E} \sup_{q \in Q} G_q$$

Multidimensional CLT : $(V_i)_{i \in \mathbb{N}}$: sequence of *d*-dimensional i.i.d.r.v. with zero mean and finite L_2 -norm.

Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N})$. A common measure of the "complexity" of Q is

$$H(Q) := \mathbb{E} \sup_{q \in Q} G_q$$

Multidimensional CLT : $(V_i)_{i \in \mathbb{N}}$: sequence of *d*-dimensional i.i.d.r.v. with zero mean and finite L_2 -norm.

$$\left|\mathbb{P}ig(rac{1}{\sqrt{n}}\sum_{i=1}^n V_i\in A(t_1,\ldots,t_d)ig)-\mathbb{P}(G\in A(t_1,\ldots,t_d))
ight|\longrightarrow 0 ext{ as } n\longrightarrow +\infty,$$

Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N})$. A common measure of the "complexity" of Q is

$$H(Q) := \mathbb{E} \sup_{q \in Q} G_q$$

Multidimensional CLT : $(V_i)_{i \in \mathbb{N}}$: sequence of *d*-dimensional i.i.d.r.v. with zero mean and finite L_2 -norm.

$$\left|\mathbb{P}ig(rac{1}{\sqrt{n}}\sum_{i=1}^n V_i\in A(t_1,\ldots,t_d)ig)-\mathbb{P}(G\in A(t_1,\ldots,t_d))
ight|\longrightarrow 0 ext{ as } n\longrightarrow +\infty,$$

where $A(t_1, \ldots, t_d) = \{v = (v_1, \ldots, v_d) \in \mathbb{R}^d : x_j \leq t_j, \forall j\}$

Gaussian process : Let $Q \subset L_2(\mu)$. $(G_q)_{q \in Q}$ is a Canonical Gaussian process associated with Q when $\forall N \in \mathbb{N}, \forall q_1, \ldots, q_N \in Q$, $(G_{q_1}, \ldots, G_{q_N}) \sim \mathcal{N}_N(0, (\langle q_i, q_j \rangle)_{1 \leq i, j \leq N})$. A common measure of the "complexity" of Q is

$$H(Q) := \mathbb{E} \sup_{q \in Q} G_q$$

Multidimensional CLT : $(V_i)_{i \in \mathbb{N}}$: sequence of *d*-dimensional i.i.d.r.v. with zero mean and finite L_2 -norm.

$$\left|\mathbb{P}ig(rac{1}{\sqrt{n}}\sum_{i=1}^n V_i\in A(t_1,\ldots,t_d)ig)-\mathbb{P}(G\in A(t_1,\ldots,t_d))
ight|\longrightarrow 0 ext{ as } n\longrightarrow +\infty,$$

where $A(t_1, \ldots, t_d) = \{v = (v_1, \ldots, v_d) \in \mathbb{R}^d : x_j \leq t_j, \forall j\}$ and G is a d-dimensional Gaussian process with zero mean and covariance matrix $(\mathbb{E}V^{(i)}V^{(j)})_{1 \leq i,j \leq d}$.

Multivariate CLT outside B_r

Fix a finite set $Q' \subset Q := \{\mathcal{L}(f) : \mathbb{E}\mathcal{L}(f) = \min_{f \in F} \mathbb{E}\mathcal{L}(f)\}$ for which $H(Q') \ge H(Q)/2$ and $0 \in Q'$.



CNRS, LATP Marseille

Multivariate CLT outside B_r

Fix a finite set $Q' \subset Q := \{\mathcal{L}(f) : \mathbb{E}\mathcal{L}(f) = \min_{f \in F} \mathbb{E}\mathcal{L}(f)\}$ for which $H(Q') \ge H(Q)/2$ and $0 \in Q'$.



Uniform Central Limit Theorem

Recall that a bounded class of functions F is μ -Donsker if and only if for $\forall u > 0, \exists \delta > 0, \exists n_0 \text{ s.t. } \forall n \ge n_0, \operatorname{osc}_n(F, \delta) \le u$ where

$$\operatorname{osc}_n(F,\delta) = \frac{1}{\sqrt{n}} \mathbb{E} \sup_{\{f,h\in F: \|f-h\|\leq \delta\}} \left| \sum_{i=1}^n g_i(f-h)(X_i) \right|,$$

where g_1, \ldots, g_n are *n* i.i.d. standard Gaussian variables.

Uniform Central Limit Theorem

Recall that a bounded class of functions F is μ -Donsker if and only if for $\forall u > 0, \exists \delta > 0, \exists n_0 \text{ s.t. } \forall n \ge n_0, \operatorname{osc}_n(F, \delta) \le u$ where

$$\operatorname{osc}_n(F,\delta) = \frac{1}{\sqrt{n}} \mathbb{E} \sup_{\{f,h\in F: \|f-h\|\leq \delta\}} \left| \sum_{i=1}^n g_i(f-h)(X_i) \right|,$$

where g_1, \ldots, g_n are *n* i.i.d. standard Gaussian variables.

δ s.t. $\forall n \geq N(F)$, $\operatorname{osc}_n(F, f^*, \delta) \leq C_2 H(Q)/\sqrt{n}$.

UCLT around f^*

Now we are ready to control the oscillation of the empirical excess risk function uniformly over the set $B_r = \{f \in F : \mathbb{E}\mathcal{L}_\lambda \leq r\}$.

UCLT around f^*

Now we are ready to control the oscillation of the empirical excess risk function uniformly over the set $B_r = \{f \in F : \mathbb{E}\mathcal{L}_\lambda \leq r\}$.

Theorem

 $\exists c_3 \text{ s.t. } \forall n \geq n_1 \text{, with } \mu^n \text{-probability at least } 1 - c_1/2,$

$$\inf_{\{f\in F:\mathbb{E}\mathcal{L}_{\lambda_n}(f)\leq r_n\}}P_n\mathcal{L}_{\lambda_n}(f)\geq -\frac{c_2H(Q)}{2\sqrt{n}}$$

where

$$r_n = c_3 \frac{H(Q)}{\sqrt{n}} \delta^2 ||T - f^*||^2$$

Main Theorem

Theorem

Let $F \subset B(L_{\infty})$ which is μ -Donsker and assume that $\exists T_0 \in \tau$ s.t. $\mathcal{M}(T_0) \geq 2$.

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

-

Sac

Main Theorem

Theorem

Let $F \subset B(L_{\infty})$ which is μ -Donsker and assume that $\exists T_0 \in \tau$ s.t. $\mathcal{M}(T_0) \geq 2$. Set $Q = \{\mathcal{L}(f) : f \in F, \mathbb{E}\mathcal{L}(f) = 0\}$.

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

-
Main Theorem

Theorem

Let $F \subset B(L_{\infty})$ which is μ -Donsker and assume that $\exists T_0 \in \tau$ s.t. $\mathcal{M}(T_0) \geq 2$. Set $Q = \{\mathcal{L}(f) : f \in F, \mathbb{E}\mathcal{L}(f) = 0\}$. $\exists C_1, C_2, N(F)$ s.t. $\forall n \geq N(F)$, with μ^n -probability at least C_1 ,

$$\mathbb{E}\mathcal{L}_{\lambda_n}(\hat{f}_{\lambda_n}) \geq C_2 \frac{H(Q)}{\sqrt{n}} \delta^2 \|T - f^*\|$$

where δ is s.t. $\forall n \ge N(F)$, $\operatorname{osc}_n(F, f^*, \delta) \le C_2 H(Q)/\sqrt{n}$ and $\lambda_n = C_2 H(Q)/\sqrt{n}$.

Subgaussian regression

 $Y = f_0(X) + \sigma \varepsilon$

 $f_0: \mathbb{R}^d \to \mathbb{R}$; ε : noise (noise level σ is known)

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Subgaussian regression

$$Y = t_0(X) + \sigma\varepsilon$$
$$f_0 : \mathbb{R}^d \to \mathbb{R}; \varepsilon : \text{noise (noise level } \sigma \text{ is known)}$$

n observations :

$$D_n := [(X_i, Y_i); 1 \le i \le n],$$

10 000

Notation : $P_n = P[\cdot|X_1, \ldots, X_n]$ and E_n : expectation w.r.t. P_n .

글 🕨 🖌 글 🕨

Assumptions

• ε is centered and subgaussian $(\exists b > 0; E_n[\exp(t\varepsilon)] \le \exp(b^2t^2/2), \forall t > 0).$

Assumptions

• ε is centered and subgaussian $(\exists b > 0; E_n[\exp(t\varepsilon)] \le \exp(b^2t^2/2), \forall t > 0).$

② X has a compact support, (we do not need $P_X << \lambda_d$). Take $\operatorname{Supp}(P_X) \subset [0,1]^d$.

Assumptions

• ε is centered and subgaussian $(\exists b > 0; E_n[\exp(t\varepsilon)] \le \exp(b^2t^2/2), \forall t > 0).$

② X has a compact support, (we do not need $P_X << \lambda_d$). Take $\operatorname{Supp}(P_X) \subset [0,1]^d$.

Complexity assumption

Aggregation and Empirical Risk Minimization

Complexity assumption

$\ \, {\cal S} \ \, {\cal F} \subset C([0,1]^d)$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Complexity assumption

- $\ \, \bullet \ \, {\mathcal F} \ \, {\rm is \ endowed \ \, with \ \, a \ \, semi-norm \ \, } |\cdot|_{{\mathcal F}}.$
- $\begin{array}{l} \bullet \quad (C_{\beta}) \ \exists \beta \in (0,2), D > 0 \ \text{s.t.} \ \forall \delta, R > 0 \\ \\ H_{\infty} \left(\delta, \mathcal{F}(R) \right) \leq D(R/\delta)^{\beta} \\ \mathcal{F}(R) := \{ f \in \mathcal{F} : |f|_{\mathcal{F}} \leq R \} \\ \\ H_{\infty} (\delta, \mathcal{F}(R)) = \log \min \left(N \in \mathbb{N} : \mathcal{F}(R) \subset \cup_{j=1}^{N} \mathcal{B}_{\infty}(f_{j}, \delta) \right). \end{array}$

s = (s₁,..., s_d) with s_i > 0 : vector of directional smoothness (s_i smoothness in direction e_i)

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

- s = (s₁,..., s_d) with s_i > 0 : vector of directional smoothness (s_i smoothness in direction e_i)
- *f* ∈ *L^p*(\mathbb{R}^d) belongs to the anisotropic Besov space $B^s_{p,q}(\mathbb{R}^d)$ if the semi-norm

$$|f|_{B^{s}_{p,q}(\mathbb{R}^{d})} := \sum_{i=1}^{d} \Big(\int_{0}^{1} (t^{-s_{i}} \|\Delta_{te_{i}}^{k_{i}}f\|_{p})^{q} \frac{dt}{t} \Big)^{1/q} < +\infty$$

 $\Delta_h^1 f(x) = f(x+h) - f(x)$ and $\Delta_h^k f = \Delta_h^1 (\Delta_h^{k-1} f)(x)$ and $k_i > s_i$.

- s = (s₁,..., s_d) with s_i > 0 : vector of directional smoothness (s_i smoothness in direction e_i)
- *f* ∈ *L^p*(\mathbb{R}^d) belongs to the anisotropic Besov space $B^s_{p,q}(\mathbb{R}^d)$ if the semi-norm

$$|f|_{B^{s}_{p,q}(\mathbb{R}^{d})} := \sum_{i=1}^{d} \Big(\int_{0}^{1} (t^{-s_{i}} \|\Delta_{te_{i}}^{k_{i}}f\|_{p})^{q} \frac{dt}{t} \Big)^{1/q} < +\infty$$

 $\Delta_h^1 f(x) = f(x+h) - f(x) \text{ and } \Delta_h^k f = \Delta_h^1 (\Delta_h^{k-1} f)(x) \text{ and } k_i > s_i.$ **a** $\mathbf{s} = (s, \dots, s) \Longrightarrow B_{\rho,q}^{\mathbf{s}}$ is the standard isotropic Besov space.

 $\Omega \subset \mathbb{R}^d$: arbitrary domain. $B^s_{p,q}(\Omega)$: set of all $f \in L^p(\Omega)$ such that $\exists g \in B^s_{p,q}(\mathbb{R}^d)$ with restriction $g \mid \Omega$ to Ω equal to f in $L^p(\Omega)$. Moreover,

$$\|f\|_{B^{\mathbf{s}}_{p,q}(\Omega)} = \inf_{g:g\mid \Omega=f} \|g\|_{B^{\mathbf{s}}_{p,q}(\mathbb{R}^d)},$$

where the infimum is taken over all $g \in B_{p,q}^{s}(\mathbb{R}^{d})$ such that $g|\Omega = f$.

Applications

example : Anisotropic Besov space

 $\Omega \subset \mathbb{R}^d$: arbitrary domain. $B^s_{p,q}(\Omega)$: set of all $f \in L^p(\Omega)$ such that $\exists g \in B^s_{p,q}(\mathbb{R}^d)$ with restriction $g \mid \Omega$ to Ω equal to f in $L^p(\Omega)$. Moreover,

$$\|f\|_{B^{\mathsf{s}}_{p,q}(\Omega)} = \inf_{g:g|\Omega=f} \|g\|_{B^{\mathsf{s}}_{p,q}(\mathbb{R}^d)},$$

where the infimum is taken over all $g \in B_{p,q}^{s}(\mathbb{R}^{d})$ such that $g|\Omega = f$.

Theorem (Birman and Solomjak (67))

Let $1 \le p, q \le \infty$ and $\mathbf{s} = (s_1, \dots, s_d)$ where $s_i > 0$, and let s be the harmonic mean of \mathbf{s}

$$\frac{1}{s} = \frac{1}{d} \sum_{j=1}^d \frac{1}{s_j}$$

If s > d/p, then

 $egin{aligned} B^{f s}_{p,q}(\Omega) \subset C(\Omega), \ H_{\infty}(\delta,B^{f s}_{p,q}(R)) &\leq D(R/\delta)^{f s/d}, orall \delta,R>0 \end{aligned}$

Estimation : PERM

Definition (PERM)

Let $\lambda = (h, \mathcal{F})$ be fixed. We say that \overline{f}_{λ} is a penalized empirical risk minimizer if it minimizes

 $R_n(f) + \operatorname{pen}(f)$

over \mathcal{F} , where $pen(f) := h^2 |f|_{\mathcal{F}}^{\alpha}$ for some $\alpha > 0$ and where

$$R_n(f) := \|Y - f\|_n^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$

is the empirical risk of f over the sample D_n .

Non-adaptive rate of convergence

Theorem

 \mathcal{F} satisfying (C_{β}) ,

$$h = an^{-1/(2+\beta)}$$

and $\alpha > 2\beta/(\beta+2)$.

Non-adaptive rate of convergence

Theorem



Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

< E

Applications

Non-adaptive rate of convergence

Theorem

 $\mathcal{F} \text{ satisfying } (C_{\beta}),$ $h = an^{-1/(2+\beta)}$ and $\alpha > 2\beta/(\beta+2).$ $\mathbf{E}_{n} \|\overline{f}_{\lambda} - f_{0}\|_{n}^{2} \leq C_{1}(1+|f_{0}|_{\mathcal{F}}^{\alpha})n^{-2/(2+\beta)}$ $\mathbf{I}f \|\overline{f}_{\lambda} - f_{0}\|_{\infty} \leq Q \text{ and } \|f_{0}\|_{\infty} \leq Q \text{ then}$ $E^{n} \|\overline{f}_{\lambda} - f_{0}\|^{2} \leq C_{2}(1+|f_{0}|_{\mathcal{F}}^{\alpha})n^{-2/(2+\beta)}$

Rate of convergence for anisotropic Besov space

Corollary

 $\mathcal{F} = B_{p,\infty}^{s}$ and $h = an^{-s/(2s+d)}$ where s is the harmonic mean of s such that s > d/p and $\|\bar{f}_{\lambda} - f_{0}\|_{\infty}, \|f_{0}\|_{\infty} \leq Q$. Then, we have :

Rate of convergence for anisotropic Besov space

Corollary

$$\begin{split} \mathcal{F} &= B_{p,\infty}^{s} \text{ and } h = an^{-s/(2s+d)} \text{ where } s \text{ is the harmonic mean of } s \text{ such } \\ & \text{that } s > d/p \text{ and } \|\bar{f}_{\lambda} - f_{0}\|_{\infty}, \|f_{0}\|_{\infty} \leq Q. \text{ Then, we have } : \\ & E\|\bar{f}_{\lambda} - f_{0}\|^{2} \leq C_{3}(1 + |f_{0}|_{B_{p,\infty}^{s}}^{2})n^{-2s/(2s+d)}. \end{split}$$

Lower bound for the ERM

< □ > < 同

Applications

Problem of adaptation

Case
$$\mathcal{F} = B^{\mathbf{s}}_{p,\infty}(\Omega)$$
.

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

э

Case
$$\mathcal{F} = B^{\mathsf{s}}_{p,\infty}(\Omega).$$

 $h = an^{-s/(2s+d)}$

depends on \boldsymbol{s} the harmonic mean of \boldsymbol{s}



Case
$$\mathcal{F} = B^{s}_{p,\infty}(\Omega).$$

 $h = an^{-s/(2s+d)}$

depends on s the harmonic mean of s which is unknown

₩

Problem of adaptation

Case
$$\mathcal{F} = B^{\mathbf{s}}_{p,\infty}(\Omega).$$

 $h = an^{-s/(2s+d)}$

depends on s the harmonic mean of \mathbf{s} which is unknown

₩

Problem of adaptation

We want to construct a procedure \tilde{f}

• independent of s

Case
$$\mathcal{F} = B^{s}_{p,\infty}(\Omega).$$

 $h = an^{-s/(2s+d)}$

depends on s the harmonic mean of s which is unknown

 \Downarrow

Problem of adaptation

We want to construct a procedure \tilde{f}

- independent of s
- if $f \in B^{\mathbf{s}}_{p,\infty}$ with s the harmonic mean of \mathbf{s} :

 $E\|\tilde{f}-f_0\|^2 \leq C_3(1+|f_0|^2_{B^s_{p,\infty}})n^{-2s/(2s+d)}.$

Aggregation of PERM

Split of the data

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

Aggregation of PERM

Split of the data

② $D_1 = ((X_1, Y_1), ..., (X_m, Y_m))$: construction of \overline{f}_s the PERM with $h = an^{-s/(2s+d)}$ for different $s \in S$:

$$S = \left\{ b + \frac{k}{\log n} : k = 0, \dots, B \log n \right\}.$$

Aggregation of PERM

Split of the data

• $D_1 = ((X_1, Y_1), \dots, (X_m, Y_m))$: construction of \overline{f}_s the PERM with $h = an^{-s/(2s+d)}$ for different $s \in S$:

$$\mathcal{S} = \left\{ b + \frac{k}{\log n} : k = 0, \dots, B \log n \right\}.$$

• $D_2 = ((X_{m+1}, Y_{m+1}), \dots, (X_n, Y_n))$: construction of the aggregation method

$$\tilde{f} \in \operatorname{Arg} \min_{f \in \operatorname{Conv}(\hat{F}_1)} R^{(2)}_{(n-m)/2}(f),$$

where the dictionary is $F = \{\overline{f}_{s} 1\!\!1_{|\overline{f}_{s}| \leq Q}, s \in S\}.$

Result

Theorem

 $\forall s \in \mathbb{R}^d \text{ s.t. the harmonic mean s is s.t. } s > d/p \text{ and } s \in [b, B], \text{ then } \forall f_0 \in B_{p,\infty}^s \cap \mathcal{B}_{\infty}(Q),$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

-

Result

Theorem

 $\forall s \in \mathbb{R}^d \text{ s.t. the harmonic mean } s \text{ is s.t. } s > d/p \text{ and } s \in [b, B], \text{ then } \forall f_0 \in B^s_{p,\infty} \cap \mathcal{B}_{\infty}(Q),$

 $|E^n\|\tilde{f}-f_0\|^2 \leq C_3(1+|f_0|^2_{B^s_{n,\infty}})n^{-2s/(2s+d)}.$

Aggregation and Empirical Risk Minimization

CNRS, LATP Marseille

< E

Result

Theorem

 $\forall \mathbf{s} \in \mathbb{R}^d \text{ s.t. the harmonic mean s is s.t. } s > d/p \text{ and } s \in [b, B], \text{ then } \forall f_0 \in B^{\mathbf{s}}_{p,\infty} \cap \mathcal{B}_{\infty}(Q),$

$$E^n \|\tilde{f} - f_0\|^2 \le C_3 (1 + |f_0|^2_{B^s_{n,\infty}}) n^{-2s/(2s+d)}.$$

\widetilde{f} adapts automatically to the "regularity parameter" s