

Introduction

Learning problem in the least-squares case

Seeing a machine learning problem as a linear inverse problem

Tikhonov regularization

A finer analysis

Conclusion

# Linear Systems and Inverse problems

Ulysse Marteau-Ferey

- 1 Introduction
- 2 Learning problem in the least-squares case
- 3 Seeing a machine learning problem as a linear inverse problem
- 4 Tikhonov regularization
  - For a generic inverse problem with approximation
  - For a machine learning problem
- 5 A finer analysis
  - Tikhonov regularization : version 2
  - Learning problems, version 2

## Learning problem and regularized empirical risk minimization

- Input: random variable  $X \in \mathcal{X}$ ;
- Output: random variable  $Y \in \mathcal{Y}$ ;
- Law of  $(X, Y)$  :  $\rho(x, y)$
- Objective : find a predictor  $f : \mathcal{X} \rightarrow \mathcal{Y}$ .
- Loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

### Learning problem

$$\inf_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} L(f) := \mathbb{E}[\ell(f(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(f(x), y) d\rho(x, y)$$

## Classical way of tackling the learning problem

### Learning problem

$$\inf_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} L(f) := \mathbb{E} [\ell(f(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(f(x), y) d\rho(x, y)$$

Access to  $\rho$  only through  $\mathbf{z} = (z_i)_{1 \leq i \leq n}$  where  $z_i = (x_i, y_i)$  are training samples

### Regularized empirical risk minimization

$$\inf_{f \in \mathcal{H}} \hat{L}_\lambda(f) := \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \frac{\lambda}{2} \Omega(f)$$

- $\mathcal{H}$  is a space of functions;
- $\Omega$  is a regularizer.

## Least squares case

- We take  $\ell(y, y') = \frac{1}{2} \|y - y'\|^2$ ;
- we assume that  $Y \in L^2(\mathcal{Y}, \rho_Y)$ .

### The minimizer exists

$$g_* = \mathbb{E}[Y|X], \quad g_*(x) = \int_{\mathcal{Y}} y \, d\rho(y|x), \quad \in L^2(\mathcal{X}, \rho_X)$$

Consequence: one needs only to solve

$$\inf_{f \in L^2(\mathcal{X}, \rho_X)} \mathbb{E} \left[ \|f(X) - Y\|^2 \right]$$

## ERM

## Regularized empirical risk minimization

$$\inf_{f \in \mathcal{H}} \widehat{L}_\lambda(f) := \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|^2 + \lambda \Omega(f)$$

## Question

What space  $\mathcal{H}$  of functions and what penalty  $\Omega$  ?

Requirements:

- $\mathcal{H} \hookrightarrow L^2(\mathcal{X}, \rho_X)$
- Solvable : finite dimensional ?

Classical parameterized version :

$$\mathcal{H} = \{f_\theta : \theta \in \Theta\}, \Theta \subset \mathbb{R}^d$$

Alternative ?

## RKHS

RKHS  $\mathcal{H}_K$ 

- $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  measurable
- $K$  positive definite, i.e.

$$\forall n \in \mathbb{N}, \mathbf{x} = (x_i)_{1 \leq i \leq n}, \alpha = (\alpha_i)_{1 \leq i \leq n}, \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j K(x_i, x_j) \geq 0$$

- $\forall x \in \mathcal{X}, K_x := K(x, \cdot) \in \mathcal{H}_K$
- $\mathcal{H}_K := \overline{\text{Span}(\{K_x, x \in \mathcal{X}\})}$  endowed with  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  so that  $\langle K_x, K_{x'} \rangle_{\mathcal{H}} = K(x, x')$

## Properties/Assumptions

- $\forall f \in \mathcal{H}, f(x) = \langle f, K_x \rangle_{\mathcal{H}}$
- Assume  $K(x, x) \leq \kappa^2$  such that

$$\forall f \in \mathcal{H}, \int_{\mathcal{X}} \|f(x)\|^2 d\rho_X(x) = \int_{\mathcal{X}} \|f \cdot K_x\|^2 d\rho_X(x) \leq \kappa^2 \|f\|_{\mathcal{H}}^2$$

## Formulation of the problem

### Regularized empirical risk minimization

$$\hat{f}_\lambda = \arg \min_{f \in \mathcal{H}} \hat{L}_\lambda(f) := \frac{1}{n} \sum_{i=1}^n \underbrace{\|f(x_i) - y_i\|^2}_{\|f \cdot K_{x_i} - y_i\|^2} + \lambda \|f\|_{\mathcal{H}}^2$$

Representer theorem :

$$\hat{f}_\lambda = \sum_{i=1}^n \alpha_i K_{x_i}$$

Equivalent problem, linear and finite dimensional !

$$\inf_{\alpha \in \mathbb{R}^n} \alpha^T (K_{nn}^2 + \lambda n K_{nn}) \alpha - 2(K_{nn} \mathbf{y}) \cdot \alpha \Leftrightarrow (K_{nn} + \lambda n I) \alpha = \mathbf{y}$$



- 1 Introduction
- 2 Learning problem in the least-squares case
- 3 Seeing a machine learning problem as a linear inverse problem
- 4 Tikhonov regularization
  - For a generic inverse problem with approximation
  - For a machine learning problem
- 5 A finer analysis
  - Tikhonov regularization : version 2
  - Learning problems, version 2

# Linear inverse problem

## Inverse problem setting

- $\mathcal{H}, \mathcal{K}$  two Hilbert spaces
- An operator  $\mathbf{A} \in \mathcal{L}(\mathcal{H}, \mathcal{K})$
- An objective  $g \in \mathcal{K}$ .

The aim is to reconstruct a solution to

$$\mathbf{A}f = g$$

## Machine Learning problem

### Ideal problem

$f(X) = Y$ ,  $f \in \mathcal{H}$  which can be reformulated as  $Sf = Y$ .

- $S : \mathcal{H} \rightarrow L^2(\mathcal{X} \times \mathcal{Y}, \rho)$  such that  $(Sf)(x, y) = f \cdot K_x = f(x)$
- $S^* : L^2(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{H}$ , such that  $S^*g = \int_{\mathcal{X} \times \mathcal{Y}} g(x, y)K_x d\rho(x, y)$
- $C = S^*S : \mathcal{H} \rightarrow \mathcal{H}$  such that  $\langle f, Cf \rangle_{\mathcal{H}} = \int_{\mathcal{X}} \|f(x)\|^2 d\rho_X(x)$

## Ill-posed problems

Possibly no solution to  $\mathbf{A}f = g$  (or obviously  $Sf = Y$ )!

### Alternative problem

$$\inf_{f \in \mathcal{H}} \|\mathbf{A}f - g\|_{\mathcal{K}}^2, \quad \inf_{f \in \mathcal{H}} \|Sf - Y\|_{L^2(\mathcal{X} \times \mathcal{Y}, \rho)}$$

$\mathbf{P}_A$  orthogonal projector on  $\overline{\text{range}(\mathbf{A})}$

$$\|\mathbf{A}f - g\|_{\mathcal{K}}^2 = \|\mathbf{A}f - \mathbf{P}_A g\|_{\mathcal{K}}^2 + \underbrace{\|(\mathbf{I} - \mathbf{P}_A)g\|_{\mathcal{K}}^2}_{\text{inevitable error}}$$

$$\text{range}(S) \subset L^2(\mathcal{X}, \rho_X) \implies \mathbf{P}_S Y = \mathbf{P}_S g_*, \quad g_*(x) = \int_{\mathcal{Y}} y \, d\rho(y|x)$$

Solution if it exists :

$$\mathbf{A}^* \mathbf{A} f = \mathbf{A}^* g, \quad C f = S^* Y = S^* \mathbf{P}_S g_*$$

## Approximations

### Inverse problem setting

Let  $\delta = (\delta_1, \delta_2) \in \mathbb{R}_+^2$ .

- An approximation space  $\tilde{\mathcal{K}}$
- An approximation of  $\mathbf{A}$  :  $\mathbf{A}_{\delta_1} : \mathcal{H} \rightarrow \tilde{\mathcal{K}}$  such that  $\|\mathbf{A}^* \mathbf{A} - \mathbf{A}_{\delta_1}^* \mathbf{A}_{\delta_1}\| \leq \delta_1$
- An approximation of  $g$  :  $g_{\delta_2} \in \tilde{\mathcal{K}}$  such that  $\|\mathbf{A}_{\delta_1}^* g_{\delta_2} - \mathbf{A}^* g\| \leq \delta_2$

### Machine Learning : approximation space limited by the data

- An approximation space  $\mathbb{R}^n$
- An approximation of  $Y$  :  $\mathbf{y} = \frac{1}{\sqrt{n}}(y_i)_{1 \leq i \leq n} \in \mathbb{R}^n$
- An approximation of  $S$  :  $S_n : \mathcal{H} \rightarrow \mathbb{R}^n$ ,  
 $S_n f = \frac{1}{\sqrt{n}}(f \cdot K_{x_i})_{1 \leq i \leq n} = \frac{1}{\sqrt{n}}(f(x_i))_{1 \leq i \leq n}$ .

## Naive method

### Potential solutions:

$$\mathbf{A}^* \mathbf{A} f = \mathbf{A}^* g, \quad \mathbf{S}^* \mathbf{S} f = \mathbf{S}^* Y$$

$$\mathbf{A}_{\delta_1}^* \mathbf{A}_{\delta_1} f_{\delta} = \mathbf{A}_{\delta_1}^* g_{\delta_2}, \quad \mathbf{S}_n^* \mathbf{S}_n \hat{f} = \mathbf{S}_n^* \mathbf{y}$$

Define  $C_n = \mathbf{S}_n^* \mathbf{S}_n : \mathcal{H} \rightarrow \mathcal{H}$  such that  $C_n = \sum_{i=1}^n K_{x_i} \otimes K_{x_i}$

### Equivalence with empirical risk minimization for learning

$$C_n \hat{f} = \mathbf{S}_n^* \mathbf{y} \Leftrightarrow \hat{f} \in \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|^2$$

### Problems :

- The second problem is not necessarily solvable (it is in finite dimension)
- We want the reconstruction error to be small, i.e.

$$\|A f_{\delta} - \mathbf{P}_A g\|^2 \xrightarrow{\delta \rightarrow 0} 0$$

- Generalization error

- 1 Introduction
- 2 Learning problem in the least-squares case
- 3 Seeing a machine learning problem as a linear inverse problem
- 4 Tikhonov regularization**
  - For a generic inverse problem with approximation
  - For a machine learning problem
- 5 A finer analysis
  - Tikhonov regularization : version 2
  - Learning problems, version 2

# Tikhonov regularization

Approximations :  $\|\mathbf{A}^* \mathbf{A} - \mathbf{A}_{\delta_1}^* \mathbf{A}_{\delta_1}\| \leq \delta_1, \|\mathbf{A}_{\delta_1}^* \mathbf{g}_{\delta_2} - \mathbf{A}^* \mathbf{g}\| \leq \delta_2$

## Regularization method

For a given  $\lambda > 0$ , choose

$$f_{\lambda, \delta} = (\mathbf{A}_{\delta_1}^* \mathbf{A}_{\delta_1} + \lambda \mathbf{I})^{-1} \mathbf{A}_{\delta_1}^* \mathbf{g}_{\delta_2} = \arg \min_{f \in \mathcal{H}} \|\mathbf{A}_{\delta_1} f - \mathbf{g}_{\delta_2}\|^2 + \lambda \|f\|^2$$

Aim : choose  $\lambda(\delta)$  well to have

$$\|\mathbf{A} f_{\lambda(\delta), \delta} - \mathbf{P}_A \mathbf{g}\|^2 \xrightarrow{\delta \rightarrow 0} 0$$

## Analysis of the error for Tikhonov regularization

Define  $f_\lambda = (\mathbf{A}^* \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^* g = (\mathbf{A}^* \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^* \mathbf{P}_A g = \arg \min_{\|Af - g\|^2 + \lambda \|f\|^2}$ .

### Decomposition

If  $\lambda \geq 2\delta_1$ , then the following holds:

$$\|\mathbf{A}f_{\lambda,\delta} - \mathbf{P}_A g\| \leq \underbrace{\|\mathbf{A}f_\lambda - \mathbf{P}_A g\|}_{S(\lambda)} + \frac{\delta_1}{\lambda} + \frac{\delta_2}{\sqrt{\lambda}}.$$

$S(\lambda)$  characterizes the between  $\mathbf{P}_A g$  and the range of  $\mathbf{A}$ .

If  $\mathbf{P}_A g \in \text{range}(\mathbf{A})$ ,  $S(\lambda) \leq C\lambda^{1/2}$



## Equivalent to solving the regularized ERM

### Equivalence

$$\hat{f}_\lambda = (\mathbf{S}_n^* \mathbf{S}_n + \lambda I)^{-1} \mathbf{S}_n^* \mathbf{y} \Leftrightarrow \hat{f}_\lambda \in \arg \min_{f \in \mathcal{H}} \|\mathbf{S}_n f - \mathbf{y}\|_{\mathbb{R}^n}^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\mathbf{S}_n f = \frac{1}{\sqrt{n}} (f(x_i))_{1 \leq i \leq n} \quad \Rightarrow \quad \hat{f}_\lambda \in \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|^2 + \lambda \|f\|_{\mathcal{H}}^2$$

**Tikhonov regularization amounts to solving the regularized ERM problem**

$$\delta_2 = \|\mathbf{S}_n^* \mathbf{y} - \mathbf{S}^* y\| = \left\| \frac{1}{n} \sum_{i=1}^n y_i K_{x_i} - \mathbb{E}[YK_X] \right\|$$

$$\delta_1 = \|\mathbf{C}_n - \mathbf{C}\| = \left\| \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i} - \mathbb{E}[K_X \otimes K_X] \right\|$$

## Bounds in high probability

Very naive bounds: with probability at least  $1 - \delta$

$$\delta_1 \leq \frac{\kappa^2 \log \frac{1}{\delta}}{\sqrt{n}}, \quad \delta_2 \leq \frac{\|y\|_\infty \kappa \log \frac{1}{\delta}}{\sqrt{n}}$$

### Statistical bound on ERM

With probability at least  $1 - \delta$ ,

$$\|\hat{f}_\lambda(x) - \mathbf{P}_S g_*(x)\|_{L^2(\mathcal{X}, \rho_X)} \leq \mathcal{S}(\lambda) + \left( \frac{\kappa^2}{\sqrt{n\lambda}} + \frac{\kappa \|y\|_\infty}{\sqrt{\lambda n}} \right).$$

- 1 Introduction
- 2 Learning problem in the least-squares case
- 3 Seeing a machine learning problem as a linear inverse problem
- 4 Tikhonov regularization
  - For a generic inverse problem with approximation
  - For a machine learning problem
- 5 **A finer analysis**
  - Tikhonov regularization : version 2
  - Learning problems, version 2

## Refinement of the bound

Consider  $\mathbf{A}_\delta^* \mathbf{A}_\delta \approx \mathbf{A}^* \mathbf{A}$ ,  $\mathbf{A}_\delta^* \mathbf{g}_\delta \approx \mathbf{A}^* \mathbf{g}$ . Define

$$\epsilon_1 = \|(\mathbf{A}^* \mathbf{A} + \lambda \mathbf{I})^{-1/2} (\mathbf{A}_\delta^* \mathbf{A}_\delta - \mathbf{A}^* \mathbf{A}) (\mathbf{A}^* \mathbf{A} + \lambda \mathbf{I})^{-1/2}\|$$

$$\epsilon_2 = \|(\mathbf{A}^* \mathbf{A} + \lambda \mathbf{I})^{-1/2} (\mathbf{A}_\delta^* \mathbf{g}_\delta - \mathbf{A}^* \mathbf{g})\|$$

$$f_{\lambda, \delta} = \arg \min_f \|\mathbf{A}_\delta f - \mathbf{g}_\delta\|^2 + \lambda \|f\|^2$$

### Tikhonov, refinement

If  $\epsilon_1 < \frac{1}{2}$ ,

$$\|\mathbf{A} f_{\lambda, \delta} - \mathbf{P}_A \mathbf{g}\| \leq \mathcal{S}(\lambda) + \epsilon_1 + \epsilon_2$$

$$\mathcal{S}(\lambda) = \|\mathbf{A} f_\lambda - \mathbf{P}_A \mathbf{g}\|$$

$$\begin{aligned}\epsilon_1 &= \|(\mathbf{C}_n + \lambda \mathbf{I})^{-1/2}(\mathbf{C}_n - \mathbf{C})(\mathbf{C}_n + \lambda \mathbf{I})^{-1/2}\| \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{C}_\lambda^{-1/2} \mathbf{K}_{x_i} \otimes \mathbf{K}_{x_i} \mathbf{C}_\lambda^{-1/2} - \mathbb{E} \left[ \mathbf{C}_\lambda^{-1/2} \mathbf{K}_X \otimes \mathbf{K}_X \mathbf{C}_\lambda^{-1/2} \right] \right\|.\end{aligned}$$

$$\epsilon_2 = \|(\mathbf{C} + \lambda \mathbf{I})^{-1/2}(\mathbf{S}_n^* \mathbf{y} - \mathbf{S}^* y)\| = \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{C}_\lambda^{-1/2} y_i \mathbf{K}_{x_i} - \mathbb{E} \left[ \mathbf{C}_\lambda^{-1/2} \mathbf{Y} \mathbf{K}_X \right] \right\|$$

## High probability results

Let  $d_\infty(\lambda) = \sup_{x \in \text{supp} \rho_X} \|\mathbf{C}_\lambda^{-1/2} \mathbf{K}_x\|^2$ . Then with probability at least  $1 - \delta$ , for any  $\lambda \leq \|\mathbf{C}\|$ ,  $\epsilon_1 \leq \sqrt{\frac{d_\infty(\lambda)}{n}} \log \frac{1}{\delta}$  and  $\epsilon_2 = \|\mathbf{Y}\|_\infty \sqrt{\frac{d_\infty(\lambda)}{n}} \log \frac{1}{\delta}$

## Final bound

$$\|\hat{f}_\lambda(x) - \mathbf{P}_S g_*(x)\|_{L^2(\mathcal{X}, \rho_X)} \leq \mathcal{S}(\lambda) + (\|\mathbf{Y}\|_\infty \vee 1) \log \frac{1}{\delta} \sqrt{\frac{d_\infty(\lambda)}{n}}.$$

Introduction

Learning problem in the least-squares case

Seeing a machine learning problem as a linear inverse problem

Tikhonov regularization

A finer analysis

Conclusion

## Conclusion

**Thank you for your attention !**