

# A CONSISTENT REGULARIZATION APPROACH FOR STRUCTURED PREDICTION

Cédric Rommel

Journal Club - CMAP

February 23rd, 2017

**1** SUBJECT AND MOTIVATION**2** DERIVATION OF A GENERAL ALGORITHM**3** STATISTICAL PROPERTIES**4** EXPERIMENTS

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# SUBJECT AND MOTIVATION

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing,

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing, image reconstruction,

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing, image reconstruction, ranking problem, . . .

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

CÉDRIC  
ROMMEL

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing, image reconstruction, ranking problem, . . .

## MOTIVATION



# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing, image reconstruction, ranking problem, . . .

## MOTIVATION

- 1 Unifying theoretical framework for SPP,

# ARTICLE'S SUBJECT AND AUTHORS

## MOTIVATION

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

Article tackles *Structured Prediction Problems* (SPP) where we want to estimate

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

with  $\mathcal{Y}$  *structured*.

**Examples:** natural language parsing, image reconstruction, ranking problem, . . .

## MOTIVATION

- ① Unifying theoretical framework for SPP,
- ② Practical method for majority of these problems.

The article considers output spaces  $\mathcal{Y}$  whose structure is induced by some metric  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  and formalizes the problems to be solved as follows:

The article considers output spaces  $\mathcal{Y}$  whose structure is induced by some metric  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  and formalizes the problems to be solved as follows:

## MAIN PROBLEM (MP)

Estimate

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$

The article considers output spaces  $\mathcal{Y}$  whose structure is induced by some metric  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  and formalizes the problems to be solved as follows:

### MAIN PROBLEM (MP)

Estimate

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$

*Rem:*

- $\rho$  is unknown;

The article considers output spaces  $\mathcal{Y}$  whose structure is induced by some metric  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  and formalizes the problems to be solved as follows:

## MAIN PROBLEM (MP)

Estimate

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$

*Rem:*

- $\rho$  is unknown;
- using sample  $\{(x_i, y_i)\}_{i=1}^n$ .

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# DERIVATION OF A GENERAL ALGORITHM

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.



## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists! k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

- $\forall x \in \mathcal{X}, k(x, \cdot) \in H$  and

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

- $\forall x \in \mathcal{X}, k(x, \cdot) \in H$  and
- $\forall f \in H, f(x) = \langle f, k(x, \cdot) \rangle,$

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists! k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

- $\forall x \in \mathcal{X}, k(x, \cdot) \in H$  and
- $\forall f \in H, f(x) = \langle f, k(x, \cdot) \rangle,$

then  $H$  is an RKHS of kernel  $k$ .

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists! k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

- $\forall x \in \mathcal{X}, k(x, \cdot) \in H$  and
- $\forall f \in H, f(x) = \langle f, k(x, \cdot) \rangle$ ,

then  $H$  is an RKHS of kernel  $k$ .

*Rem:*

- $k(x, y) = \langle k(x, \cdot), k(y, \cdot) \rangle$ ,

## RKHS DEFINITION

Let  $\mathcal{X}$  be an arbitrary set and  $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  be a Hilbert space.

If  $\exists! k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric, such that

- $\forall x \in \mathcal{X}, k(x, \cdot) \in H$  and
- $\forall f \in H, f(x) = \langle f, k(x, \cdot) \rangle,$

then  $H$  is an RKHS of kernel  $k$ .

*Rem:*

- $k(x, y) = \langle k(x, \cdot), k(y, \cdot) \rangle,$
- A reproducing kernel  $k$  is *positive-definite*:

$$\sum_{i,j=1}^n c_i c_j k(x_i, x_j) \geq 0, \quad (1)$$

$$\forall n \in \mathbb{N}, x_1, \dots, x_n \in \mathcal{X}, c_1, \dots, c_n \in \mathbb{R}.$$

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$\forall x_i, x_j \in \mathcal{X}$ ,

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle$$



## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$$\forall x_i, x_j \in \mathcal{X},$$

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$$\forall x_i, x_j \in \mathcal{X},$$

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

where  $\psi : \mathcal{X} \rightarrow H$  is called *feature map*.

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$\forall x_i, x_j \in \mathcal{X}$ ,

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

where  $\psi : \mathcal{X} \rightarrow H$  is called *feature map*.

*Kernel trick*

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$\forall x_i, x_j \in \mathcal{X}$ ,

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

where  $\psi : \mathcal{X} \rightarrow H$  is called *feature map*.

*Kernel trick*

$$f(x) = w^\top \psi(x) + b$$

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$\forall x_i, x_j \in \mathcal{X}$ ,

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

where  $\psi : \mathcal{X} \rightarrow H$  is called *feature map*.

*Kernel trick*

$$f(x) = w^\top \psi(x) + b \quad \rightsquigarrow \quad \text{SVM, KRR, ...}$$

## MOORE-ARONSZAJN THEOREM

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  symmetric positive-definite  $\Rightarrow \exists!$  RKHS  
 $H \subset \mathcal{F}(\mathcal{X}, \mathbb{R})$  of kernel  $k$ .

*Rem:*

$\forall x_i, x_j \in \mathcal{X}$ ,

$$k(x_i, x_j) = \langle k(x_i, \cdot), k(x_j, \cdot) \rangle = \langle \psi(x_i), \psi(x_j) \rangle,$$

where  $\psi : \mathcal{X} \rightarrow H$  is called *feature map*.

*Kernel trick*

$$f(x) = w^\top \psi(x) + b \quad \underset{\text{SVM, KRR, ...}}{\rightsquigarrow} \quad \hat{f}(x) = \sum_{i=1}^n c_i k(x_i, x)$$

Let

- $\lambda \in \mathbb{R}$ ,
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  *reproducing kernel*,
- $\mathbf{K} \in \mathbb{R}^{n \times n}$  with  $\mathbf{K}_{ij} = k(x_i, x_j)$ ,
- $\mathbf{K}_x \in \mathbb{R}^n$  with  $(\mathbf{K}_x)_i = k(x, x_i)$ .

Let

- $\lambda \in \mathbb{R}$ ,
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  *reproducing kernel*,
- $\mathbf{K} \in \mathbb{R}^{n \times n}$  with  $\mathbf{K}_{ij} = k(x_i, x_j)$ ,
- $\mathbf{K}_x \in \mathbb{R}^n$  with  $(\mathbf{K}_x)_i = k(x, x_i)$ .

## ALGORITHM 1



Let

- $\lambda \in \mathbb{R}$ ,
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  *reproducing kernel*,
- $\mathbf{K} \in \mathbb{R}^{n \times n}$  with  $\mathbf{K}_{ij} = k(x_i, x_j)$ ,
- $\mathbf{K}_x \in \mathbb{R}^n$  with  $(\mathbf{K}_x)_i = k(x, x_i)$ .

## ALGORITHM 1

LEARNING:

$$\alpha(x) = (\mathbf{K} + n\lambda\mathbf{I})^{-1}\mathbf{K}_x \in \mathbb{R}^n$$

Let

- $\lambda \in \mathbb{R}$ ,
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  *reproducing kernel*,
- $\mathbf{K} \in \mathbb{R}^{n \times n}$  with  $\mathbf{K}_{ij} = k(x_i, x_j)$ ,
- $\mathbf{K}_x \in \mathbb{R}^n$  with  $(\mathbf{K}_x)_i = k(x, x_i)$ .

## ALGORITHM 1

LEARNING:

$$\alpha(x) = (\mathbf{K} + n\lambda\mathbf{I})^{-1}\mathbf{K}_x \in \mathbb{R}^n$$

PREDICTION:

$$\hat{f}(x) = \arg \min_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i)$$

Special case of *Kernel Dependency Estimation*:

$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y')$$

Special case of *Kernel Dependency Estimation*:

$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y')$$

where  $h$  is a reproducing kernel and  $\psi : \mathcal{Y} \rightarrow \mathcal{H}_y$  its nonlinear feature map.

Special case of *Kernel Dependency Estimation*:

$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y')$$

where  $h$  is a reproducing kernel and  $\psi : \mathcal{Y} \rightarrow \mathcal{H}_y$  its nonlinear feature map.

$$\Delta(f(x), y) = \|\psi(f(x)) - \psi(y)\|_{\mathcal{H}_y}^2$$

Special case of *Kernel Dependency Estimation*:

$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y')$$

where  $h$  is a reproducing kernel and  $\psi : \mathcal{Y} \rightarrow \mathcal{H}_y$  its nonlinear feature map.

$$\Delta(f(x), y) = \|\psi(f(x)) - \psi(y)\|_{\mathcal{H}_y}^2$$

As  $\psi(f(x))$  is hard to minimize wrt to  $f$ , we replace it by  $g \in \mathcal{G} \subset \mathcal{F}(\mathcal{X}, \mathcal{H}_y)$  easier to optimize:

Special case of *Kernel Dependency Estimation*:

$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y')$$

where  $h$  is a reproducing kernel and  $\psi : \mathcal{Y} \rightarrow \mathcal{H}_\mathcal{Y}$  its nonlinear feature map.

$$\Delta(f(x), y) = \|\psi(f(x)) - \psi(y)\|_{\mathcal{H}_\mathcal{Y}}^2$$

As  $\psi(f(x))$  is hard to minimize wrt to  $f$ , we replace it by  $g \in \mathcal{G} \subset \mathcal{F}(\mathcal{X}, \mathcal{H}_\mathcal{Y})$  easier to optimize:

$$\min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|g(x_i) - \psi(y_i)\|_{\mathcal{H}_\mathcal{Y}}^2 + \lambda \|g\|_{\mathcal{G}}^2$$

If we chose

$$\mathcal{G} = \left\{ g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_y) : g(x) = \sum_{i=1}^n c_i k(x, x_i), c_i \in \mathcal{H}_y \right\}$$



If we chose

$$\mathcal{G} = \left\{ g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_y) : g(x) = \sum_{i=1}^n c_i k(x, x_i), c_i \in \mathcal{H}_y \right\}$$

with  $k$  a reproducing kernel, the last problem becomes a *Kernel Ridge Regression* problem and

If we chose

$$\mathcal{G} = \left\{ g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_y) : g(x) = \sum_{i=1}^n c_i k(x, x_i), c_i \in \mathcal{H}_y \right\}$$

with  $k$  a reproducing kernel, the last problem becomes a *Kernel Ridge Regression* problem and

$$\hat{g}(x) = \sum_{i=1}^n \alpha_i(x) \psi(y_i),$$

with

$$\alpha(x) = (\mathbf{K} + n\lambda\mathbf{I})^{-1} \mathbf{K}_x \in \mathbb{R}^n.$$

In this case, some references advise to choose  $\hat{f}$  as follows:

$$\hat{f}(x) = \arg \min_{y \in \mathcal{Y}} \|\psi(y) - \hat{g}(x)\|_{\mathcal{H}_y}^2$$

In this case, some references advise to choose  $\hat{f}$  as follows:

$$\begin{aligned}\hat{f}(x) &= \arg \min_{y \in \mathcal{Y}} \|\psi(y) - \hat{g}(x)\|_{\mathcal{H}_y}^2 \\ &= \arg \min_{y \in \mathcal{Y}} h(y, y) - 2 \sum_{i=1}^n \alpha_i(x) h(y, y_i)\end{aligned}$$

In this case, some references advise to choose  $\hat{f}$  as follows:

$$\begin{aligned}\hat{f}(x) &= \arg \min_{y \in \mathcal{Y}} \|\psi(y) - \hat{g}(x)\|_{\mathcal{H}_y}^2 \\ &= \arg \min_{y \in \mathcal{Y}} h(y, y) - 2 \sum_{i=1}^n \alpha_i(x) h(y, y_i)\end{aligned}$$

which is equal to the *prediction* step of *Algorithm 1*

$$\hat{f}(x) = \arg \min_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i)$$

if  $h$  is *normalized*.

- ① And if we don't have
- $$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y') ?$$

- 1 And if we don't have  
$$\Delta(y, y') = h(y, y) - 2h(y, y') + h(y', y') ?$$
- 2 Is the transfer between  $\hat{g}$  and  $\hat{f}$  theoretically justified ?

## MAIN PROBLEM (MP)

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$



## MAIN PROBLEM (MP)

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$

We want to find  $\mathcal{L}(g(x), y)$  a relaxation of  $\Delta(f(x), y)$  on some space  $\mathcal{H}_y$  easy to optimize in order to replace (MP) by

## MAIN PROBLEM (MP)

$$f^* \in \arg \min_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \Delta(f(x), y) d\rho(x, y).$$

We want to find  $\mathcal{L}(g(x), y)$  a relaxation of  $\Delta(f(x), y)$  on some space  $\mathcal{H}_y$  easy to optimize in order to replace (MP) by

## SURROGATE PROBLEM (SP)

$$g^* \in \arg \min_{g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_y)} \mathcal{R}(g) = \int_{\mathcal{X} \times \mathcal{Y}} \mathcal{L}(g(x), y) d\rho(x, y)$$

(MP) and (SP) will be equivalent if we find an appropriate *decoding* function  $d : \mathcal{H}_y \rightarrow \mathcal{Y}$  satisfying

(MP) and (SP) will be equivalent if we find an appropriate *decoding* function  $d : \mathcal{H}_Y \rightarrow \mathcal{Y}$  satisfying

FISHER CONSISTENCY:

$$\mathcal{E}(d \circ g^*) = \mathcal{E}(f^*)$$

(MP) and (SP) will be equivalent if we find an appropriate *decoding* function  $d : \mathcal{H}_Y \rightarrow \mathcal{Y}$  satisfying

FISHER CONSISTENCY:

$$\mathcal{E}(d \circ g^*) = \mathcal{E}(f^*)$$

COMPARISON INEQUALITY:

$$\mathcal{E}(d \circ g) - \mathcal{E}(f^*) \leq \varphi(\mathcal{R}(g) - \mathcal{R}(g^*))$$

for all  $g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_Y)$ , with  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  such that  $\lim_{s \rightarrow 0} \varphi(s) = 0$ .

## ASSUMPTION 1

There is

- $\mathcal{H}_y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_y}$ ;

## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;

## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;
- $V : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$  bounded linear operator;



## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;
- $V : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$  bounded linear operator;

such that

$$\forall y, y' \in \mathcal{Y}, \quad \Delta(y, y') = \langle \psi(y), V\psi(y') \rangle_{\mathcal{H}_Y}.$$

## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;
- $V : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$  bounded linear operator;

such that

$$\forall y, y' \in \mathcal{Y}, \quad \Delta(y, y') = \langle \psi(y), V\psi(y') \rangle_{\mathcal{H}_Y}.$$

*Rem:*

- $V$  is not required to be positive definite, nor even symmetric;

# SURROGATE PROBLEM EQUIVALENCE (1/2)

## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;
- $V : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$  bounded linear operator;

such that

$$\forall y, y' \in \mathcal{Y}, \quad \Delta(y, y') = \langle \psi(y), V\psi(y') \rangle_{\mathcal{H}_Y}.$$

*Rem:*

- $V$  is not required to be positive definite, nor even symmetric;
- looks like reproducing kernel, but its not;

# SURROGATE PROBLEM EQUIVALENCE (1/2)

## ASSUMPTION 1

There is

- $\mathcal{H}_Y$  separable Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_Y}$ ;
- $\psi : \mathcal{Y} \rightarrow \mathcal{H}_Y$  continuous embedding;
- $V : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$  bounded linear operator;

such that

$$\forall y, y' \in \mathcal{Y}, \quad \Delta(y, y') = \langle \psi(y), V\psi(y') \rangle_{\mathcal{H}_Y}.$$

*Rem:*

- $V$  is not required to be positive definite, nor even symmetric;
- looks like reproducing kernel, but its not;
- "wide range of functions"

# EXAMPLES OF $\Delta$ VERIFYING ASSUMPTION 1

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

- Loss functions on  $\mathcal{Y}$  of finite cardinality: Multi-class classification, ranking, ...

# EXAMPLES OF $\Delta$ VERIFYING ASSUMPTION 1

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

- Loss functions on  $\mathcal{Y}$  of finite cardinality: Multi-class classification, ranking, ...
- Least-squares, Logistic, Hinge,  $\epsilon$ -sensitivity, ...

# EXAMPLES OF $\Delta$ VERIFYING ASSUMPTION 1

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

- Loss functions on  $\mathcal{Y}$  of finite cardinality: Multi-class classification, ranking, ...
- Least-squares, Logistic, Hinge,  $\epsilon$ -sensitivity, ...
- Robust loss functions: Huber,  $L_2 - L_1$ , ...

# EXAMPLES OF $\Delta$ VERIFYING ASSUMPTION 1

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

- Loss functions on  $\mathcal{Y}$  of finite cardinality: Multi-class classification, ranking, ...
- Least-squares, Logistic, Hinge,  $\epsilon$ -sensitivity, ...
- Robust loss functions: Huber,  $L2 - L1$ , ...
- Kernel Dependency Estimation (loss function from intuition)



# EXAMPLES OF $\Delta$ VERIFYING ASSUMPTION 1

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

- Loss functions on  $\mathcal{Y}$  of finite cardinality: Multi-class classification, ranking, ...
- Least-squares, Logistic, Hinge,  $\epsilon$ -sensitivity, ...
- Robust loss functions: Huber,  $L2 - L1$ , ...
- Kernel Dependency Estimation (loss function from intuition)
- Distances on histograms and probabilities

# SURROGATE PROBLEM EQUIVALENCE (2/2)

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

## LEMMA 1

If  $\Delta$  satisfies *Assumption 1* with  $\psi$  bounded, then we have

## LEMMA 1

If  $\Delta$  satisfies *Assumption 1* with  $\psi$  bounded, then we have

$$\mathcal{E}(f) = \int_{\mathcal{X}} \langle \psi(f(x), Vg^*(x)), \mathcal{H}_y \rangle d\rho_{\mathcal{X}}(x)$$

for all  $f : \mathcal{X} \rightarrow \mathcal{Y}$ ,

## LEMMA 1

If  $\Delta$  satisfies *Assumption 1* with  $\psi$  bounded, then we have

$$\mathcal{E}(f) = \int_{\mathcal{X}} \langle \psi(f(x)), \nabla g^*(x) \rangle_{\mathcal{H}_Y} d\rho_{\mathcal{X}}(x)$$

for all  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $g^* : \mathcal{X} \rightarrow \mathcal{H}_Y$  solves

$$\min_{g \in \mathcal{F}(\mathcal{X}, \mathcal{H}_Y)} \mathcal{R}(g) = \int_{\mathcal{X} \times \mathcal{Y}} \|g(x) - \psi(y)\|_{\mathcal{H}_Y}^2 d\rho(x, y).$$

## THEOREM 2

If  $\Delta$  satisfies *Assumption 1* and  $\mathcal{Y}$  is compact, then for any  $g : \mathcal{X} \rightarrow \mathcal{H}_{\mathcal{Y}}$  measurable and for  $d : \mathcal{H}_{\mathcal{Y}} \rightarrow \mathcal{Y}$  such that

## THEOREM 2

If  $\Delta$  satisfies *Assumption 1* and  $\mathcal{Y}$  is compact, then for any  $g : \mathcal{X} \rightarrow \mathcal{H}_{\mathcal{Y}}$  measurable and for  $d : \mathcal{H}_{\mathcal{Y}} \rightarrow \mathcal{Y}$  such that

$$d(h) = \arg \min_{y \in \mathcal{Y}} \langle \psi(y), Vh \rangle_{\mathcal{H}_{\mathcal{Y}}},$$

## THEOREM 2

If  $\Delta$  satisfies *Assumption 1* and  $\mathcal{Y}$  is compact, then for any  $g : \mathcal{X} \rightarrow \mathcal{H}_Y$  measurable and for  $d : \mathcal{H}_Y \rightarrow \mathcal{Y}$  such that

$$d(h) = \arg \min_{y \in \mathcal{Y}} \langle \psi(y), Vh \rangle_{\mathcal{H}_Y},$$

we have

$$\begin{aligned} \mathcal{E}(d \circ g^*) &= \mathcal{E}(f^*) \\ \mathcal{E}(d \circ g) - \mathcal{E}(f^*) &\leq c_\Delta \sqrt{\mathcal{R}(g) - \mathcal{R}(g^*)}, \end{aligned}$$

with  $c_\Delta = \|V\| \max_{y \in \mathcal{Y}} \|\psi(y)\|_{\mathcal{H}_Y}$ .

## THEOREM 2

If  $\Delta$  satisfies *Assumption 1* and  $\mathcal{Y}$  is compact, then for any  $g : \mathcal{X} \rightarrow \mathcal{H}_Y$  measurable and for  $d : \mathcal{H}_Y \rightarrow \mathcal{Y}$  such that

$$d(h) = \arg \min_{y \in \mathcal{Y}} \langle \psi(y), Vh \rangle_{\mathcal{H}_Y},$$

we have

$$\begin{aligned} \mathcal{E}(d \circ g^*) &= \mathcal{E}(f^*) \\ \mathcal{E}(d \circ g) - \mathcal{E}(f^*) &\leq c_\Delta \sqrt{\mathcal{R}(g) - \mathcal{R}(g^*)}, \end{aligned}$$

with  $c_\Delta = \|V\| \max_{y \in \mathcal{Y}} \|\psi(y)\|_{\mathcal{H}_Y}$ .

*Rem:* Any SPP on  $\mathcal{Y}$  of finite cardinality satisfies *Th. 2*



## LEMMA 3

If  $\Delta$  verify *Assumption 1*,  $\mathcal{Y}$  is compact and  $\hat{g} \in \mathcal{G} = \mathcal{F}(\mathcal{X}, \mathcal{H}_{\mathcal{Y}})$  solves the *Empirical surrogate problem*

## LEMMA 3

If  $\Delta$  verify *Assumption 1*,  $\mathcal{Y}$  is compact and  $\hat{g} \in \mathcal{G} = \mathcal{F}(\mathcal{X}, \mathcal{H}_{\mathcal{Y}})$  solves the *Empirical surrogate problem*

$$\min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|g(x_i) - \psi(y_i)\|_{\mathcal{H}_{\mathcal{Y}}}^2 + \lambda \|g\|_{\mathcal{G}}^2$$

## LEMMA 3

If  $\Delta$  verify *Assumption 1*,  $\mathcal{Y}$  is compact and  $\hat{g} \in \mathcal{G} = \mathcal{F}(\mathcal{X}, \mathcal{H}_{\mathcal{Y}})$  solves the *Empirical surrogate problem*

$$\min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|g(x_i) - \psi(y_i)\|_{\mathcal{H}_{\mathcal{Y}}}^2 + \lambda \|g\|_{\mathcal{G}}^2$$

then

$$\forall x \in \mathcal{X}, \quad \hat{f} = d \circ \hat{g} = \arg \min_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i)$$

with

$$\alpha(x) = (\mathbf{K} + n\lambda\mathbf{I})^{-1} \mathbf{K}_x.$$

## LEMMA 3

If  $\Delta$  verify *Assumption 1*,  $\mathcal{Y}$  is compact and  $\hat{g} \in \mathcal{G} = \mathcal{F}(\mathcal{X}, \mathcal{H}_{\mathcal{Y}})$  solves the *Empirical surrogate problem*

$$\min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \|g(x_i) - \psi(y_i)\|_{\mathcal{H}_{\mathcal{Y}}}^2 + \lambda \|g\|_{\mathcal{G}}^2$$

then

$$\forall x \in \mathcal{X}, \quad \hat{f} = d \circ \hat{g} = \arg \min_{y \in \mathcal{Y}} \sum_{i=1}^n \alpha_i(x) \Delta(y, y_i)$$

with

$$\alpha(x) = (\mathbf{K} + n\lambda\mathbf{I})^{-1} \mathbf{K}_x.$$

*Rem:*  $\hat{f}$  does not depend on  $\psi$  nor on  $V$ : *loss trick*  $\simeq$  *kernel trick*.

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# STATISTICAL PROPERTIES

## THEOREM 4

Let  $n \in \mathbb{N}$  and  $\rho$  an arbitrary distribution on  $\mathcal{X} \times \mathcal{Y}$ .

## THEOREM 4

Let  $n \in \mathbb{N}$  and  $\rho$  an arbitrary distribution on  $\mathcal{X} \times \mathcal{Y}$ . If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{X}$  and  $\mathcal{Y}$  are compact,  $k$  is a *continuous universal kernel*

## THEOREM 4

Let  $n \in \mathbb{N}$  and  $\rho$  an arbitrary distribution on  $\mathcal{X} \times \mathcal{Y}$ . If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{X}$  and  $\mathcal{Y}$  are compact,  $k$  is a *continuous universal kernel*, and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/4}$ , then



## THEOREM 4

Let  $n \in \mathbb{N}$  and  $\rho$  an arbitrary distribution on  $\mathcal{X} \times \mathcal{Y}$ . If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{X}$  and  $\mathcal{Y}$  are compact,  $k$  is a *continuous universal kernel*, and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/4}$ , then

$$\lim_{n \rightarrow +\infty} \mathcal{E}(\hat{f}_n) = \mathcal{E}(f^*) \quad \text{with probability 1.}$$

## THEOREM 4

Let  $n \in \mathbb{N}$  and  $\rho$  an arbitrary distribution on  $\mathcal{X} \times \mathcal{Y}$ . If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{X}$  and  $\mathcal{Y}$  are compact,  $k$  is a *continuous universal kernel*, and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/4}$ , then

$$\lim_{n \rightarrow +\infty} \mathcal{E}(\hat{f}_n) = \mathcal{E}(f^*) \quad \text{with probability 1.}$$

*Rem:* First universal consistency result for general form of SPP, including  $\mathcal{Y}$  of infinite cardinality.

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then if  $\mathcal{R}$  admits a minimizer  $g^*$ ...

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then if  $\mathcal{R}$  admits a minimizer  $g^*$ ...

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq c\tau^2 n^{-1/4}$$

holds with probability  $1 - 8e^{-\tau}$  for any  $\tau > 0$ , with  $c$  a constant independant of  $n$  and  $\tau$ .

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then if  $\mathcal{R}$  admits a minimizer  $g^*$ ...

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq c\tau^2 n^{-1/4}$$

holds with probability  $1 - 8e^{-\tau}$  for any  $\tau > 0$ , with  $c$  a constant independant of  $n$  and  $\tau$ .

*Rem:*

- Not possible to prove uniform convergence rates according to ref [25];

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then if  $\mathcal{R}$  admits a minimizer  $g^*$ ...

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq c\tau^2 n^{-1/4}$$

holds with probability  $1 - 8e^{-\tau}$  for any  $\tau > 0$ , with  $c$  a constant independant of  $n$  and  $\tau$ .

*Rem:*

- Not possible to prove uniform convergence rates according to ref [25];
- Better bounds possible in the case of least-squares classifiers by using Tsybakov condition to regularize  $\rho$ ;

## THEOREM 5

If  $\Delta$  satisfies *Assumption 1*,  $\mathcal{Y}$  is compact,  $k$  is a *bounded* continuous universal kernel and if  $\hat{f}_n$  is the estimator obtained through *Alg.1* with  $n$  i.i.d. training points with  $\lambda_n = n^{-1/2}$ , then if  $\mathcal{R}$  admits a minimizer  $g^*$ ...

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq c\tau^2 n^{-1/4}$$

holds with probability  $1 - 8e^{-\tau}$  for any  $\tau > 0$ , with  $c$  a constant independant of  $n$  and  $\tau$ .

*Rem:*

- Not possible to prove uniform convergence rates according to ref [25];
- Better bounds possible in the case of least-squares classifiers by using Tsybakov condition to regularize  $\rho$ ;
- Not sure if the same strategy would work on infinite setting.



A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

# EXPERIMENTS

$$\Delta_{rank}(y, y') = \frac{1}{2} \sum_{i,j=1}^M \gamma(y')_{ij} (1 - \text{sign}(y_i - y_j))$$

	Rank Loss
<b>Linear</b> [7]	0.430 ± 0.004
<b>Hinge</b> [27]	0.432 ± 0.008
<b>Logistic</b> [28]	0.432 ± 0.012
<b>SVM Struct</b> [4]	0.451 ± 0.008
<b>Alg. 1</b>	<b>0.396 ± 0.003</b>

Table 1: Normalized  $\Delta_{rank}$  for ranking methods on the MovieLens dataset [29].

# EXPERIEMENTS - DIGIT RECONSTRUCTION

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS

$$\Delta_G(y, y') = 1 - k_G(y, y')$$

$$\Delta_H(y, y') = \sum_{i=1}^M \left| (y_i)^{1/2} - (y'_i)^{1/2} \right|, \quad \text{for } y = (y_i)_{i=1}^M$$

$\Delta_R(y, y')$  = "Recongnition" loss wrt SVM classifier

	KDE [18]	Alg. 1
Loss	(Gaussian)	(Hellinger)
$\Delta_G$	<b>0.149 ± 0.013</b>	0.172 ± 0.011
$\Delta_H$	0.736 ± 0.032	<b>0.647 ± 0.017</b>
$\Delta_R$	0.294 ± 0.012	<b>0.193 ± 0.015</b>

Table 2: Digit reconstruction using Gaussian (KDE [18]) and Hellinger loss.

# EXPERIEMENTS - ROBUST ESTIMATION

A  
CONSISTENT  
REGULARIZA-  
TION  
APPROACH  
FOR  
STRUCTURED  
PREDICTION

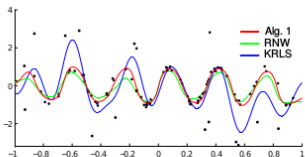
CÉDRIC  
ROMMEL

SUBJECT AND  
MOTIVATION

DERIVATION  
OF A GENERAL  
ALGORITHM

STATISTICAL  
PROPERTIES

EXPERIMENTS



n	Alg. 1	RNW	KRR
50	$0.39 \pm 0.17$	$0.45 \pm 0.18$	$0.62 \pm 0.13$
100	$0.21 \pm 0.04$	$0.29 \pm 0.04$	$0.47 \pm 0.09$
200	$0.12 \pm 0.02$	$0.24 \pm 0.03$	$0.33 \pm 0.04$
500	$0.08 \pm 0.01$	$0.22 \pm 0.02$	$0.31 \pm 0.03$
1000	$0.07 \pm 0.01$	$0.21 \pm 0.02$	$0.19 \pm 0.02$

Figure 1: Robust estimation on the regression problem in Sec. 6 by minimizing the Cauchy loss with Alg. 1 (Ours) or Nadaraya-Watson (Nad). KRLS as a baseline predictor. **Left.** Example of one run of the algorithms. **Right.** Average distance of the predictors to the actual function (without noise and outliers) over 100 runs with respect to training sets of increasing dimension.