

On Consistency and Robustness Properties of SVMs for Heavy-Tailed Distributions

Arnout Van Messem

Andreas Christmann



Vrije
Universiteit
Brussel



UNIVERSITÄT
BAYREUTH

2nd Workshop of the ERCIM Working Group
on Computing and Statistics,
October 29-31, 2009, Limassol, Cyprus

Introduction

Known

Support Vector Machines (SVMs) are **consistent** and **robust**, if based on Lipschitz continuous loss and bounded kernel.

Christmann & Van Messem '08

Steinwart & Christmann '08

Christmann & Steinwart '07

Question

Can the assumptions $f \in L_1(P_X)$ and $\int |Y| dP < \infty$ be weakened?

(both for regression and classification problems)

Support Vector Machines

Definition

$$f_{L,P,\lambda} := \arg \inf_{f \in \mathcal{H}} \mathbb{E}_P L(X, Y, f(X)) + \lambda \|f\|_{\mathcal{H}}^2$$

- $\mathcal{X} \subseteq \mathbb{R}^d$ closed, $\mathcal{Y} \subseteq \mathbb{R}$ closed, $\mathcal{X} \neq \emptyset$, $\mathcal{Y} \neq \emptyset$
- (X_i, Y_i) i.i.d. $\sim P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$, **P (totally) unknown**
- $Y_i|x_i$ depends on an *unknown* function $f : \mathcal{X} \rightarrow \mathbb{R}$
- **RKHS** $\mathcal{H} \Leftrightarrow$ **kernel** $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, k **measurable**
- **Loss function** $L : \mathcal{X} \times \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$, $L(x, y, f(x))$
- $\lambda > 0$ regularization parameter
- $f_{L,D,\lambda}$, where D is empirical distribution for data set
 $D = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$

Support Vector Machines

Notions

- L is called **convex, continuous, Lipschitz continuous, differentiable**, if L has this property w.r.t. 3^{rd} argument
- k is called **bounded**, if $\|k\|_\infty := \sqrt{\sup_{x \in \mathcal{X}} k(x, x)} < \infty$
- $\Phi : \mathcal{X} \rightarrow \mathcal{H}$, $\Phi(x) := k(\cdot, x)$, is called **canonical feature map**
- Reproducing property:

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}} \quad \forall f \in \mathcal{H}, \forall x \in \mathcal{X}.$$

Risk

Definitions

Risk	$\mathcal{R}_{L,P}(f)$	$\mathbb{E}_P L(X, Y, f(X))$
Bayes risk	$\mathcal{R}_{L,P}^*$	$\inf_{f:\mathcal{X}\rightarrow\mathbb{R} \text{ measurable}} \mathcal{R}_{L,P}(f)$
Bayes function	$f_{L,P}^*$	$\arg \inf_{f:\mathcal{X}\rightarrow\mathbb{R} \text{ measurable}} \mathcal{R}_{L,P}(f)$

Questions

Under which conditions on \mathcal{X} , \mathcal{Y} , L , \mathcal{H} , and k do we have:

- 1 $f_{L,P,\lambda}$: existence, uniqueness, representation
- 2 Universal consistency to Bayes risk/function, i.e., $\forall P$

$$\mathcal{R}_{L,P}(f_{L,D,\lambda}) \xrightarrow{P} \mathcal{R}_{L,P}^* \text{ for } |D| = n \rightarrow \infty$$

$$f_{L,D,\lambda} \xrightarrow{P} f_{L,P}^* \text{ for } |D| = n \rightarrow \infty$$
- 3 Robustness of $f_{L,P,\lambda}$?

Shifted loss function

Loss function $L : \mathcal{X} \times \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$ measurable

Definition

$L^* : \mathcal{X} \times \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}$ with

$$L^*(x, y, t) := L(x, y, t) - L(x, y, 0).$$

Huber, 1967

L^* can be negative!

Properties

- L (strictly) convex, then L^* (strictly) convex.
- L Lipschitz continuous, then L^* Lipschitz continuous.

Shifted loss function

Conditions for finite risk

For L Lipschitz continuous

- $\mathcal{R}_{L,P}(f) < \infty$ if $f \in L_1(P_X)$ and $\mathbb{E}_P|Y| < \infty$.
- $\mathcal{R}_{L^*,P}(f) < \infty$ if $f \in L_1(P_X)$.

Equality of SVMs

If $f_{L,P,\lambda}$ exists, then $f_{L^*,P,\lambda} = f_{L,P,\lambda}$.

Properties

- If L Lipschitz continuous, then

$$|\mathcal{R}_{L^*,P}(f)| \leq |L|_1 \mathbb{E}_{P_X} |f(X)|.$$

$$|\mathcal{R}_{L^*,P,\lambda}^{reg}(f)| \leq |L|_1 \mathbb{E}_{P_X} |f(X)| + \lambda \|f\|_{\mathcal{H}}^2.$$

- If L Lipschitz continuous and $f_{L^*,P,\lambda}$ exists, then

$$\|f_{L^*,P,\lambda}\|_{\mathcal{H}}^2 \leq \lambda^{-1} \min\{|L|_1 \mathbb{E}_{P_X} |f_{L^*,P,\lambda}(X)|, \mathcal{R}_{L,P}(0)\}.$$

If additionally k is bounded, then $\|f_{L^*,P,\lambda}\|_{\mathcal{H}} < \infty$.

Existence and Uniqueness of SVM solution

Uniqueness

- L convex and $\mathcal{R}_{L^*,P}(f) < \infty$ for some $f \in \mathcal{H}$ and $\mathcal{R}_{L^*,P}(f) > -\infty$ for all $f \in \mathcal{H}$

OR

- L is convex, Lipschitz continuous and $f \in L_1(P_X)$.

Then, for all $\lambda > 0$, there **exists at most one** SVM $f_{L^*,P,\lambda}$.

Existence

- L convex, Lipschitz continuous,
- \mathcal{H} RKHS of a bounded measurable kernel k .

Then, for all $P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ and for all $\lambda > 0$, there **exists** an SVM solution $f_{L^*,P,\lambda}$.

Representation

Theorem

- L convex, Lipschitz continuous loss function,
- k bounded, measurable kernel with separable RKHS \mathcal{H} .

Then, for all $P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ and for all $\lambda > 0$, there exists an $h \in \mathcal{L}_\infty(P)$ with

$$h(x, y) \in \partial L^*(x, y, f_{L^*, P, \lambda}(x)), \quad \forall (x, y)$$

$$\|h\|_\infty \leq |L^*|_1 = |L|_1$$

$$f_{L^*, P, \lambda} = -\frac{1}{2\lambda} \mathbb{E}_P(h\Phi)$$

$$\|f_{L^*, P, \lambda} - f_{L^*, Q, \lambda}\|_{\mathcal{H}} \leq \frac{1}{\lambda} \|\mathbb{E}_P(h\Phi) - \mathbb{E}_Q(h\Phi)\|_{\mathcal{H}}, \quad \forall Q.$$

Consistency

Theorem

- L convex, Lipschitz continuous loss function,
- \mathcal{H} separable RKHS of a bounded, measurable kernel k ,
- \mathcal{H} dense in $L_1(\mu)$ for all distributions μ on \mathcal{X} ,
- (λ_n) sequence of strictly positive numbers with $\lambda_n \rightarrow 0$.

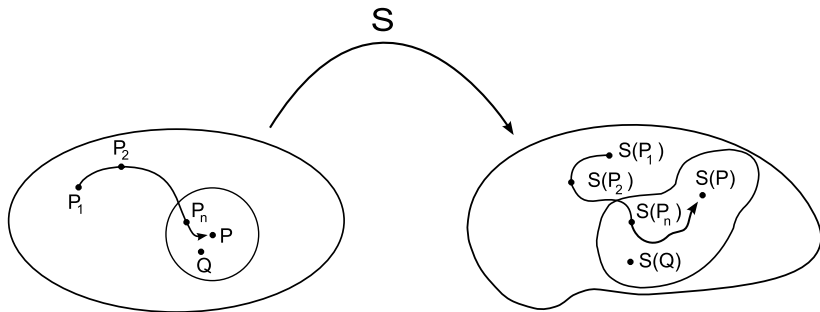
Then, for all $P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ and all D with $|D| = n$,

- 1 if $\lambda_n^2 n \rightarrow \infty$, then $\mathcal{R}_{L^*,P}(f_{L^*,D,\lambda_n}) \xrightarrow{P} \mathcal{R}_{L^*,P}^*$.
- 2 if $\lambda_n^{2+\delta} n \rightarrow \infty$ for some $\delta \in (0, \infty)$, then

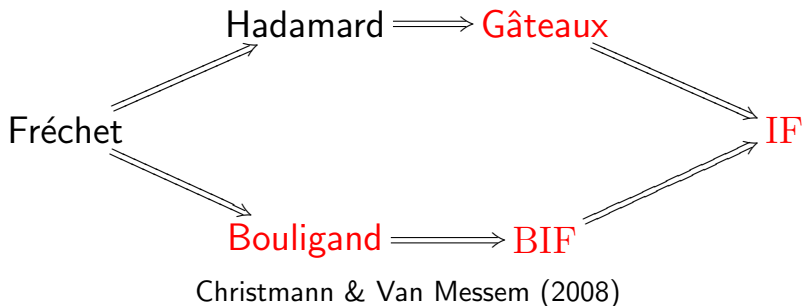
$$\mathcal{R}_{L^*,P}(f_{L^*,D,\lambda_n}) \xrightarrow{\text{a.s.}} \mathcal{R}_{L^*,P}^*.$$
- 3 if $L = L_\tau$ pinball loss: $d(f_{L^*,D,\lambda_n}, f_{L_\tau,P}^*) \rightarrow 0$.
 d is a metric describing convergence in probability.

Robustness

- ① What if (X_i, Y_i) i.i.d. $\sim P$, $P \in \mathcal{M}_1$ unknown is invalid?
- ② What is the impact on $S : P \mapsto f_{L^*, P, \lambda}$?



Derivatives and Influence Functions



Notation: ∇^F , ∇^G , ∇^B , ∇_3^B , etc.

Property: $\nabla_3^F L^* = \nabla_3^F L$, $\nabla_3^B L^* = \nabla_3^B L$

Influence Function

Definition (Hampel, '68, Hampel et al. '86)

The **influence function** (IF) of a function $S : \mathcal{M}_1 \rightarrow \mathcal{H}$ for a distribution P is given by

$$\text{IF}(z; S, P) := \lim_{\varepsilon \downarrow 0} \frac{S((1 - \varepsilon)P + \varepsilon\delta_z) - S(P)}{\varepsilon},$$

in those $z := (x, y) \in \mathcal{X} \times \mathcal{Y}$ where this limit exists.

If $\nabla^G(z; S, P)$ exists: $\nabla^G = \text{IF}$ and IF is linear and continuous

Goal: **Bounded IF**

Problem: **Loss function L often not Fréchet-differentiable**

Bouligand Influence Function

Definition (C&VM '08)

The **Bouligand influence function** (BIF) of a function $S : \mathcal{M}_1 \rightarrow \mathcal{H}$ for a distribution P in the direction of a distribution $Q \neq P$ is the special Bouligand-derivative

$$\lim_{\varepsilon \downarrow 0} \frac{\|S((1 - \varepsilon)P + \varepsilon Q) - S(P) - \text{BIF}(Q; S, P)\|_{\mathcal{H}}}{\varepsilon} = 0$$

(if it exists).

If BIF exists and $Q = \delta_z$: IF exists and $\text{BIF} = \text{IF}$

Goal: **Bounded BIF**

Result for IF

Assumptions

- \mathcal{H} is RKHS with **bounded**, continuous kernel k
- L **convex** and **Lipschitz continuous**
- $\nabla_3^F L(x, y, \cdot)$ and $\nabla_{3,3}^F L(x, y, \cdot)$ continuous with

$$\kappa_1 := \sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left\| \nabla_3^F L(x, y, \cdot) \right\|_\infty \in (0, \infty),$$

$$\kappa_2 := \sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left\| \nabla_{3,3}^F L(x, y, \cdot) \right\|_\infty < \infty$$

Theorem IF

Then $\text{IF}(z; S, P)$ with $S(P) := f_{L^*, P, \lambda}$ and $z := (x, y)$

- 1 exists,
- 2 equals

$$\mathbb{E}_P \nabla_3^F L^*(X, Y, f_{L^*, P, \lambda}(X)) T^{-1} \Phi(X) \\ - \nabla_3^F L^*(x, y, f_{L^*, P, \lambda}(x)) T^{-1} \Phi(x),$$

where $T : \mathcal{H} \rightarrow \mathcal{H}$ with $T(\cdot) :=$

$$2\lambda \text{id}_{\mathcal{H}}(\cdot) + \mathbb{E}_P \nabla_{3,3}^F L^*(X, Y, f_{L^*, P, \lambda}(X)) \langle \Phi(X), \cdot \rangle_{\mathcal{H}} \Phi(X),$$

- 3 is bounded.

Bounds for bias

Maxbias and IF

- \mathcal{H} is separable RKHS with **bounded**, measurable kernel k
- L **convex** and **Lipschitz continuous**

Then, for all $\lambda > 0$, all $\varepsilon \in [0, 1]$ and **all** $P, Q \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$

$$\|f_{L^*, (1-\varepsilon)P + \varepsilon Q} - f_{L^*, P, \lambda}\|_{\mathcal{H}} \leq c_{P, Q} \varepsilon,$$

where $c_{P, Q} = \lambda^{-1} \|k\|_{\infty} |L|_1 \|P - Q\|_{\mathcal{M}}$.

- $Q = \delta_z$ with $z := (x, y)$
- IF($z; S, P$) with $S(P) := f_{L^*, P, \lambda}$ exists

Then $\|IF(z; S, P)\|_{\mathcal{H}} \leq c_{P, \delta_z}$.

Result for BIF

Assumptions

- \mathcal{H} is RKHS with **bounded**, continuous kernel k
- L **convex** and **Lipschitz continuous** with $|L|_1 \in (0, \infty)$
- $\nabla_3^B L(x, y, \cdot)$ and $\nabla_{3,3}^B L(x, y, \cdot)$ measurable with

$$\kappa_1 := \sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left\| \nabla_3^B L(x, y, \cdot) \right\|_\infty \in (0, \infty),$$

$$\kappa_2 := \sup_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left\| \nabla_{3,3}^B L(x, y, \cdot) \right\|_\infty < \infty$$

Assumptions

- $\delta_1 > 0, \delta_2 > 0$
- $\mathcal{N}_{\delta_1}(f_{L^*,P,\lambda}) := \{f \in \mathcal{H} : \|f - f_{L^*,P,\lambda}\|_{\mathcal{H}} < \delta_1\}$
- $\lambda > \frac{1}{2}\kappa_2\|k\|_{\infty}^3$ ($\kappa_2 = 0$ for eps-insensitive and pinball)
- $P \neq Q$, probability measures on $\mathcal{X} \times \mathcal{Y}$
- Define $G : (-\delta_2, \delta_2) \times \mathcal{N}_{\delta_1}(f_{L^*,P,\lambda}) \rightarrow \mathcal{H}$,

$$G(\varepsilon, f) := 2\lambda f + \mathbb{E}_{(1-\varepsilon)P+\varepsilon Q} \nabla_3^B L^*(X, Y, f(X)) \Phi(X)$$

- $G(0, f_{L^*,P,\lambda}) = 0$ and $\nabla_2^B G(0, f_{L^*,P,\lambda})$ is **strong**

Theorem BIF

Then $\text{BIF}(Q; S, P)$ with $S(P) := f_{L^*, P, \lambda}$ and $Q \neq P \in \mathcal{M}_1$

- 1 exists,
- 2 equals

$$T^{-1} \left(\mathbb{E}_P \nabla_3^B L^*(X, Y, f_{L^*, P, \lambda}(X)) \Phi(X) - \mathbb{E}_Q \nabla_3^B L^*(X, Y, f_{L^*, P, \lambda}(X)) \Phi(X) \right),$$

where $T : \mathcal{H} \rightarrow \mathcal{H}$ with $T(\cdot) :=$

$$2\lambda \text{id}_{\mathcal{H}}(\cdot) + \mathbb{E}_P \nabla_{3,3}^B L^*(X, Y, f_{L^*, P, \lambda}(X)) \langle \Phi(X), \cdot \rangle_{\mathcal{H}} \Phi(X),$$

- 3 is bounded.

Conclusions

SVMs based on $L^*(x, y, t) := L(x, y, t) - L(x, y, 0)$

- ① Weaker assumption on P : only $f \in L_1(P_X)$ is needed
e.g. f bounded and $\mathcal{X} \subset \mathbb{R}^d$ bounded
- ② Existence and uniqueness of $f_{L^*, P, \lambda}$
- ③ Representation of SVM solution
- ④ Consistency of risk and SVM solution
- ⑤ Robustness
 - Existence of IF and BIF
 - IF(Q; S, P) bounded if $\nabla_3^F L$, $\nabla_{3,3}^F L$ and k continuous and bounded
 - BIF(Q; S, P) bounded if $\nabla_3^B L$, $\nabla_{3,3}^B L$ measurable and bounded as well as k continuous and bounded
 - Bounds for bias

References

- Christmann, Van Messem & Steinwart (2009). *Statistics and Its Interface*, **2**, 311-327.
- Christmann & Van Messem (2008). *Journal of Machine Learning Research*, **9**, 915-936.
- Steinwart & Christmann (2008). *Support Vector Machines*. Springer, New York.
- Christmann & Steinwart (2007). *Bernoulli*, **13**, 799-819.
- Hampel (1974). *J. Amer. Statist. Assoc.*, **69**, 383-393.
- Huber (1967). *Proceedings of the 5th Berkeley Symposium*.
- Koenker (2005). *Quantile regression*. Cambridge University Press.
- Schölkopf & Smola (2002). *Learning with kernels*. MIT Press.
- Vapnik (1998). *Statistical learning theory*. Wiley.

Reason

Conditions for finite risk

For L Lipschitz continuous

- $\mathbb{E}_{\mathbb{P}} L(X, Y, f(X)) < \infty$ if $f \in L_1(\mathbb{P}_X)$ and $Y \in L_1(\mathbb{P}_{Y|x})$.

$$\mathcal{R}_{L,\mathbb{P}}(f) \leq |L|_1 \left(\int_{\mathcal{X}} |f(x)| d\mathbb{P}_X(x) + \int_{\mathcal{X}} \int_{\mathcal{Y}} |y| d\mathbb{P}(y|x) d\mathbb{P}_X(x) \right)$$

- $\mathbb{E}_{\mathbb{P}} L^*(X, Y, f(X)) < \infty$ if $f \in L_1(\mathbb{P}_X)$.

$$\mathcal{R}_{L,\mathbb{P}}(f) \leq |L|_1 \int_{\mathcal{X}} |f(x)| d\mathbb{P}_X(x)$$

Sketch of proof for IF

- $G(\varepsilon, f) := 2\lambda f + \mathbb{E}_{(1-\varepsilon)P+\varepsilon\Delta_z} \nabla_2^F L^*(Y, f(X))\Phi(X)$
- $G(\varepsilon, f) = \nabla_2^F \mathcal{R}_{L^*, (1-\varepsilon)P+\varepsilon\Delta_z, \lambda}^{reg}(f), \quad \varepsilon \in [0, 1]$
- $G(\varepsilon, f)$ fulfills conditions of a standard implicit function theorem on Banach spaces

Sketch: Proof for IF

For the proof of the theorem about the IF we showed:

- i.** $G(0, f) = 0 \Leftrightarrow f = f_{L^*, P, \lambda}$.
- ii.** G continuously F-differentiable.
- iii.** $\frac{\partial G}{\partial \mathcal{H}}(0, f_{L^*, P, \lambda})$ invertible.
- iv.** Then there exist $\delta > 0$, a neighborhood $\mathcal{N}_\delta(f_{L^*, P, \lambda}) := \{f \in \mathcal{H}; \|f - f_{L^*, P, \lambda}\|_{\mathcal{H}} < \delta\}$, and a function $f^* : (-\delta, \delta) \rightarrow \mathcal{N}_\delta(f_{L^*, P, \lambda})$ satisfying

iv.1) $f^*(0) = f_{L^*, P, \lambda}$.

iv.2) It holds

$$\nabla^F f^*(0) = -(\nabla_2^F G(0, f_{L^*, P, \lambda}))^{-1} - \nabla_1^B G(0, f_{L^*, P, \lambda}).$$

Sketch of proof for BIF

- $\nabla_2^B L^*(Y, f(X)) = \nabla_2^B L(Y, f(X))$ hence

$$G(\varepsilon, f) = 2\lambda f + \mathbb{E}_{(1-\varepsilon)P+\varepsilon Q} \nabla_2^B L(Y, f(X)) \Phi(X)$$

- $G(\varepsilon, f) = \nabla_2^B \mathcal{R}_{L^*, (1-\varepsilon)P+\varepsilon Q, \lambda}^{reg}(f), \quad \varepsilon \in [0, 1]$
 - $G(\varepsilon, f)$ fulfills the conditions of Robinson's (1991) implicit function theorem on Bouligand-derivatives for non-smooth functions in Banach or normed linear spaces
- ⇒ Rest of proof uses same arguments as Christmann & Van Messem (2008).

Sketch: Proof for BIF

For the proof of the theorem about the BIF we showed:

- i.** For some χ and each $f \in \mathcal{N}_{\delta_1}(f_{L^*,P,\lambda})$, $G(\cdot, f)$ is Lipschitz continuous on $(-\delta_2, \delta_2)$ with Lipschitz constant χ .
- ii.** G has partial B-derivatives with respect to ε and f at $(0, f_{L^*,P,\lambda})$.
- iii.** $\nabla_2^B G(0, f_{L^*,P,\lambda})(\mathcal{N}_{\delta_1}(f_{L^*,P,\lambda}) - f_{L^*,P,\lambda})$ is a neighborhood of $0 \in \mathcal{H}$.
- iv.** $\delta(\nabla_2^B G(0, f_{L^*,P,\lambda}), \mathcal{N}_{\delta_1}(f_{L^*,P,\lambda}) - f_{L^*,P,\lambda}) =: d_0 > 0$.

- v.** For each $\xi > d_0^{-1}\chi$ there exist $\delta_3, \delta_4 > 0$, a neighborhood $\mathcal{N}_{\delta_3}(f_{L^*,P,\lambda}) := \{f \in \mathcal{H}; \|f - f_{L^*,P,\lambda}\|_{\mathcal{H}} < \delta_3\}$, and a function $f^* : (-\delta_4, \delta_4) \rightarrow \mathcal{N}_{\delta_3}(f_{L^*,P,\lambda})$ satisfying
- v.1)** $f^*(0) = f_{L^*,P,\lambda}$.
 - v.2)** $f^*(\cdot)$ is Lipschitz continuous on $(-\delta_4, \delta_4)$ with Lipschitz constant $|f^*|_1 = \xi$.
 - v.3)** For each $\varepsilon \in (-\delta_4, \delta_4)$ is $f^*(\varepsilon)$ the unique solution of $G(\varepsilon, f) = 0$ in $(-\delta_4, \delta_4)$.
 - v.4)** It holds $\nabla^B f^*(0)(u) = (\nabla_2^B G(0, f_{L^*,P,\lambda}))^{-1} (-\nabla_1^B G(0, f_{L^*,P,\lambda})(u))$.