

Applied Model Selection and Error Estimation: Some New Results and Open problems

Third Italian Workshop on Machine Learning and Data Mining
Workshop of the XIII AI*IA Symposium on Artificial Intelligence

Luca Oneto, Alessandro Ghio, Davide Anguita



University of Genoa

Dibris
DIBRIS

DITEN
DITEN



SmartLab

December 10, 2014

luca.oneto@unige.it - www.lucaoneto.com

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

In-Sample Methods

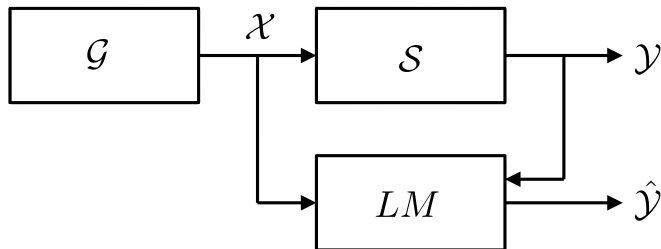
Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

Supervised Learning Framework¹

- \mathcal{G} Data Generator
- \mathcal{D} Supervisor
- LM Learning Machine
- Generalization (True) Error $L(f) = \mathbb{E}_{(X, Y)} \ell(f(X), Y)$
- Empirical Error $\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$

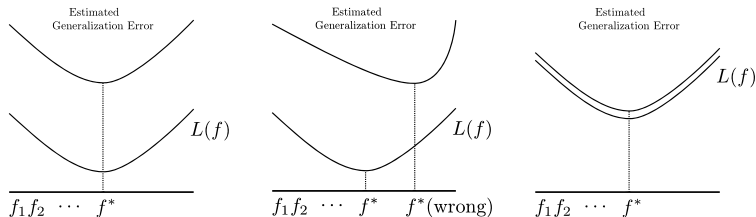


Model Selection and Error Estimation (I)

Model Selection is the effort to obtain a function f with the smallest possible $L(f)$ where $L(f)$ is obviously unknown since μ is unknown. (Selecting the number of hidden neurons in Multi Layer Perceptrons, the regularization parameter in Support Vector Machines or the depth of a Decision Tree)

Error Estimation is the effort to obtain the rigorously tightest possible estimate (with high probability) of $L(f)$

The two problems seem to be two faces of the same coin. The truth is that Model Selection is an easier task respect to Error Estimation.



Model Selection and Error Estimation (II)^{2,3}

- **Out-Of-Sample Methods (Hold-Out):** the idea behind this approach is to divide the data in different parts (Test Set, K-Fold Cross-Validation, Bootstrap, etc.). These methods are quite general and can be used for every possible Learning Algorithm.
- **In-Sample Methods:** all the data are used for learning the model f , for model selection (selecting the best hyperparameters of the algorithms or the right hypothesis space) and for error estimation.
 - **Hypotheses Space-Based Methods:** these methods require to know explicitly the hypotheses space from which we choose f
 - **Algorithm-Based Methods:** these methods are quite general and can be used for every possible Learning Algorithm

²Anguita D. Ghio A. Oneto L. Ridella S. In-sample and out-of-sample model selection and error estimation for support vector machines. IEEE Transactions on Neural Networks and Learning Systems. 2012

³Oneto L. Ghio A. Anguita D. Ridella S. Fully Empirical and Data-Dependent Stability-Based Bounds. IEEE Transactions on Cybernetics. –

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

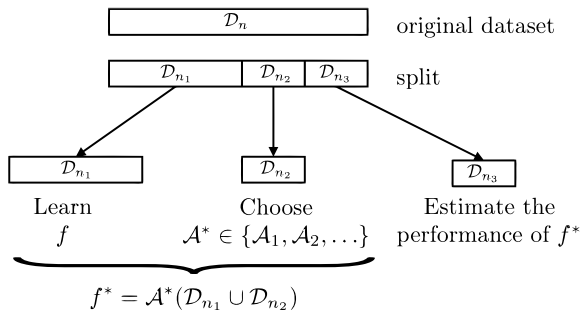
In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

Out-Of-Sample Methods



The key is the I.I.D. hypothesis

Model Selection: if the model is good, the error over previously unseen data (\mathcal{D}_{n_2}) should be small

Error Estimation: since the data in \mathcal{D}_{n_3} are I.I.D. with respect to $\mathcal{D}_{n_1} \cup \mathcal{D}_{n_2}$, the errors over the sample in \mathcal{D}_{n_3} are I.I.D. so we can rigorously estimate $L(f)$ (Clopper–Pearson, Hoeffding, etc.)

Out-Of-Sample Methods: Open Problems (And Some Solutions)

- How to deal with small-sample setting ($n \leq 100$)?⁴
- Are the state-of-the-art bounds really the theoretical limit?⁵
- Does an Optimal splitting procedure exist?⁶
 - Size of the different sets
 - Randomized splitting
 - Stratified splitting
 - Nearly Homogeneous Splitting
- Why the methods work so well in practice?
 - No theoretical justification
 - The error of the Bootstrap procedure is not completely understood
 - Leave One Out works well in practice for Model Selection. For Error Estimation it is completely unreliable.
- Bootstrap or Cross Validation?⁷

⁴Anguita D. Ghio A. Greco N. Oneto L. Ridella S. Model selection for support vector machines: Advantages and disadvantages of the machine learning theory. IEEE International Joint Conference Neural Networks. 2010

⁵Anguita D. Ghelardoni L. Ghio A. Ridella S. Test error bounds for classifiers: A survey of old and new results. IEEE Symposium on Foundations of Computational Intelligence. 2011

⁶Anguita D. Ghelardoni L. Ghio A. Oneto L. Ridella S. The 'K' in K-fold Cross Validation. European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning. 2012

⁷Anguita D. Ghio A. Oneto L. Ridella S. In-sample and out-of-sample model selection and error estimation for support vector machines. IEEE Transactions on Neural Networks and Learning Systems. 2012

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

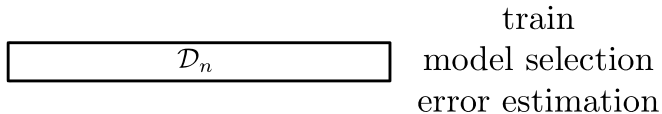
In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

In-Sample Methods^{8,9} (I)



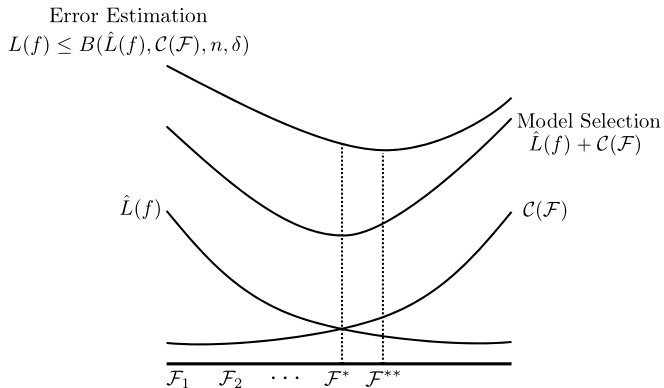
Problem:

Since we use the data for different purposes, the errors over the different samples are not I.I.D. anymore even if the the data in \mathcal{D}_n are I.I.D.

⁸Vapnik V. N. Statistical learning theory.1998

⁹Poggio T. Rifkin R. Mukherjee S. Niyogi P. General conditions for predictivity in learning theory. Nature. 2004

In-Sample Methods (II)



\mathcal{F} **can be:** the number of hidden neuron in MLP, the depth of a Decision Tree, a space of Functions, etc.

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

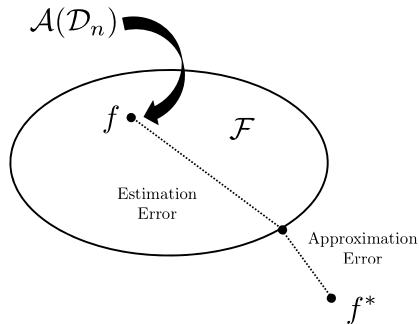
Final Remarks

Hypotheses Space–Based Methods¹⁰

**We need to define the hypotheses space
from which the algorithm will choose the model.**

NOTE: For many algorithms it may become computationally intractable:

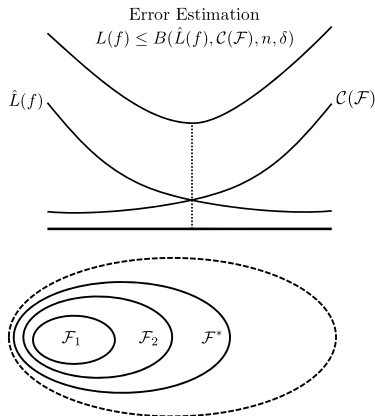
- K-NN
- Decision Tree
- SVM
- etc.



Structural Risk Minimization (SRM):

Different Approach to estimate $\mathcal{C}(\mathcal{F})$

- Vapnik–Chervonenkis dimension
 - Data-independent
 - Takes into account functions that will be never piked up during the learning procedure
- Rademacher Complexity
 - Data-dependent (on $\mathbb{P}(\mathcal{X})$ not on $\mathbb{P}(\mathcal{Y}|\mathcal{X})$)
 - Takes into account functions that will be never piked up during the learning procedure
- Local Rademacher Complexity
 - Data-dependent (on $\mathbb{P}(\mathcal{X})$ and on $\mathbb{P}(\mathcal{Y}|\mathcal{X})$)
 - Takes into account functions with small empirical error
- PAC Bayes
 - Data-dependent (on $\mathbb{P}(\mathcal{X})$ and on $\mathbb{P}(\mathcal{Y}|\mathcal{X})$)
 - Takes into account functions with small empirical error



Open Problems (And Some Solutions)

- Does not work for all the algorithms (Decision Tree, k -NN, etc.¹¹)
- We need to improve the Concentration Inequalities¹²
- Hard to compute in practice¹³
- Still loose in practice

Can we exploit these methods for optimizing the amount of resources needed for learning?

LINK - VIDEO

¹¹Bousquet O. Andre E. Stability and generalization. The Journal of Machine Learning Research. 2002

¹²Boucheron S. Gabor L. Pascal M.. Concentration Inequalities: A nonasymptotic theory of independence. Oxford University Press. 2013

¹³Anguita D. Ghio A. Oneto L. Ridella S. A Deep Connection Between the Vapnik–Chervonenkis Entropy and the Rademacher Complexity. IEEE Transaction on Neural Networks and Learning Systems. 2014

Resource Aware Learning (I)¹⁴

- Reproducing kernel Hilbert space

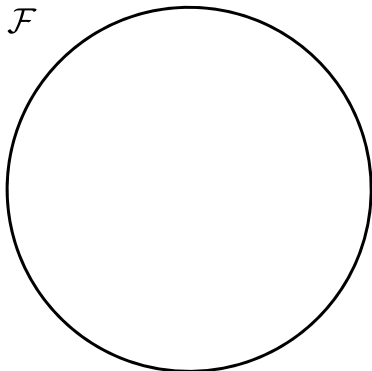
$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \sum_{i=1}^n \alpha K(\mathbf{x}_i, \mathbf{x})$$

- Regularization (limiting the space of functions)

$$\|\mathbf{w}\|_2^2 \leq w_{\text{MAX}}$$

- Empirical Risk Minimization

$$\inf_{f \in \mathcal{F}} \hat{L}_n(f)$$



¹⁴L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Learning resource-aware models for mobile devices: from regularization to energy efficiency. Neurocomputing, 2014.

Resource Aware Learning (II)¹⁵

- Bound on the True Error

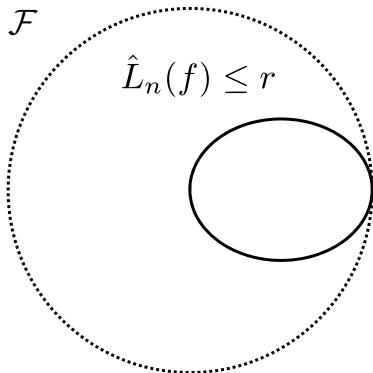
$$L(f)$$

- Local Rademacher Complexity

$$\hat{L}_n(f) \leq \rho$$

- SRM with Local Rademacher Complexity

$$\mathcal{F}^* \in \{\mathcal{F}_1, \mathcal{F}_2, \dots\}$$

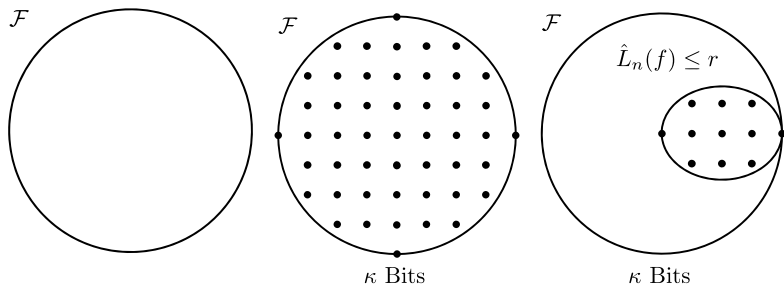


¹⁵L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Learning resource-aware models for mobile devices: from regularization to energy efficiency. Neurocomputing, 2014.

Resource Aware Learning (III)¹⁶

Reducing the number of bits for representing \mathbf{w}

$$w_j \in \frac{w_{\text{MAX}}}{2^\kappa - 1} \{-2^\kappa + 1, \dots, 2^\kappa - 1\}$$

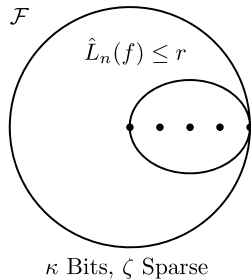
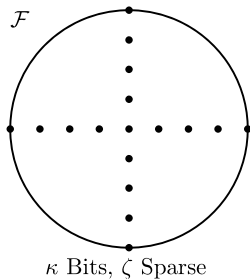
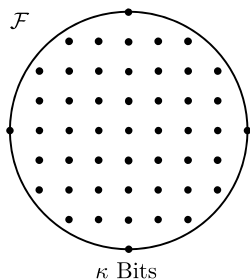


¹⁶L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Learning resource-aware models for mobile devices: from regularization to energy efficiency. Neurocomputing, 2014.

Resource Aware Learning (IV)¹⁷

Increasing the sparsity of w

$$\sum_{i=1}^D [w_i \neq 0] \leq D\zeta$$



¹⁷L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Learning resource-aware models for mobile devices: from regularization to energy efficiency. Neurocomputing, 2014.



Human Activity Recognition Using Smartphones Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Human Activity Recognition database built from the recordings of 30 subjects performing activities of daily living (ADL) while carrying

Data Set Characteristics:	Multivariate, Time-Series	Number of Instances:	10299	Area:	Computer
Attribute Characteristics:	N/A	Number of Attributes:	561	Date Donated	2012-12-10
Associated Tasks:	Classification, Clustering	Missing Values?	N/A	Number of Web Hits:	131283

¹⁸D. Anguita, A. Ghio, L. Oneto, X. Parra, J. L. Reyes-Ortiz, A Public Domain Dataset for Human Activity Recognition using Smartphones, European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN), pp. 437-442, Bruges, Belgium, 24-26 Apr. 2013.

Standing vs Sitting

n	ζ	κ	<i>Error</i>
25	0.5	1	21.3
100	0.5	1	11.7
225	0.5	4	7.7
400	0.5	8	7.2
625	0.5	8	7.1
900	0.5	8	7.0

¹⁹D. Anguita, A. Ghio, L. Oneto, X. Parra, J. L. Reyes-Ortiz, A Public Domain Dataset for Human Activity Recognition using Smartphones, European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN), pp. 437-442, Bruges, Belgium, 24-26 Apr. 2013.

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

Algorithm–Based Methods

These methods do not require the knowledge of the Space of Function \mathcal{F} from which the Algorithm \mathcal{A} will choose.

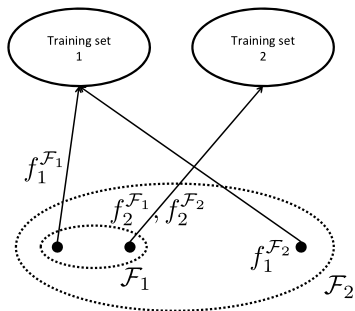
$$\mathcal{A} : \mathcal{D}_n \rightarrow f$$

Different Approach to estimate $\mathcal{C}(\mathcal{A})$

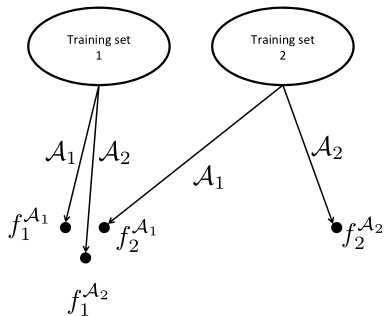
- Algorithmic Stability: how much the solution found by the algorithm is affected by changing of just one sample?
- Compression Bound: how much data can be removed from the training set without affecting the performance of the learning algorithms?

Algorithmic Stability^{20,21} (I)

Structural Risk Minimization



Stability



²⁰Bousquet O. Andre E. Stability and generalization. The Journal of Machine Learning Research. 2002

²¹Poggio T. Rifkin R. Mukherjee S. Niyogi P. General conditions for predictivity in learning theory. Nature. 2004

Algorithmic Stability (II)

Approach 1: Hypothesis Stability

$$H_{\text{emp}}(\mathcal{A}_{\mathcal{F}}, n) = \mathbb{E}_{\mathcal{D}_n, Z_i'} \left| \ell(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, Z_i) - \ell(\mathcal{A}_{(\mathcal{D}_n^i, \mathcal{F})}, Z_i) \right| \leq \beta_{\text{emp}}$$

$$H_{\text{loo}}(\mathcal{A}_{\mathcal{F}}, n) = \mathbb{E}_{\mathcal{D}_n, Z} \left| \ell(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, Z) - \ell(\mathcal{A}_{(\mathcal{D}_n^i, \mathcal{F})}, Z) \right| \leq \beta_{\text{loo}}$$

Bound on the Generalization Error

$$L(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}) \leq \hat{L}_{\text{emp}}(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \mathcal{D}_n) + \sqrt{\frac{1}{2n\delta} + \frac{3\beta_{\text{emp}}}{\delta}}$$

$$L(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}) \leq \hat{L}_{\text{loo}}(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \mathcal{D}_n) + \sqrt{\frac{1}{2n\delta} + \frac{3\beta_{\text{loo}}}{\delta}}$$

Algorithmic Stability (III)

Approach 2: Uniform Stability

$$U^i(\mathcal{A}_{\mathcal{F}}, n) = \left| \ell(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \cdot) - \ell(\mathcal{A}_{(\mathcal{D}_n^i, \mathcal{F})}, \cdot) \right|_{\infty} \leq \beta^i$$

$$U^{\setminus i}(\mathcal{A}_{\mathcal{F}}, n) = \left| \ell(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \cdot) - \ell(\mathcal{A}_{(\mathcal{D}_n^{\setminus i}, \mathcal{F})}, \cdot) \right|_{\infty} \leq \beta^{\setminus i}$$

Bound on the Generalization Error

$$L(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}) \leq \hat{L}_{\text{emp}}(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \mathcal{D}_n) + 2\beta^i + (4n\beta^i + 1) \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{2n}},$$

$$L(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}) \leq \hat{L}_{\text{loo}}(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}, \mathcal{D}_n) + \beta^{\setminus i} + (4n\beta^{\setminus i} + 1) \sqrt{\frac{\log\left(\frac{1}{\delta}\right)}{2n}}.$$

Open Problems (And Some Solutions)

- We need to improve and Concentration Inequalities^{22,23}
- Loose in practice

**Uniform Stability is data independent and (almost) algorithmic independent.
Hypothesis Stability is unknown for kernel-based algorithms.**

²²Hoeffding W. Probability inequalities for sums of bounded random variables. Journal of the American statistical association. 1963

²³Boucheron S. Gabor L. Pascal M.. Concentration Inequalities: A nonasymptotic theory of independence. Oxford University Press. 2013

Hypothesis

$$H^{\text{loo}}(\mathcal{A}_{\mathcal{F}}, n) \leq H^{\text{loo}}(\mathcal{A}_{\mathcal{F}}, \sqrt{n}/2)$$

We point out that the above property is a desirable requirement for any learning algorithm: in fact, the impact on the learning procedure of removing samples from \mathcal{D}_n should decrease, on average, as n grows.

²⁴L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Fully empirical and data-dependent stability-based bounds. IEEE Transactions on Cybernetics, 2014.

Fully Empirical Algorithmic Stability Based Bound (II)²⁵

Thanks to this hypothesis we can estimate the Hypothesis Stability based on the empirical data.

$$\hat{H}^{\text{loo}}(\mathcal{A}_{(\mathcal{D}_{\sqrt{n}/2}, \mathcal{F})}, \mathcal{D}_{\sqrt{n}/2}) = \frac{8}{n\sqrt{n}} \sum_{k=1}^{\sqrt{n}/2} \sum_{j=1}^{\sqrt{n}/2} \sum_{i=1}^{\sqrt{n}/2} \left| \ell(\mathcal{A}_{(\check{\mathcal{D}}_{\sqrt{n}/2}^k, \mathcal{F})}, \check{\mathbf{z}}_j^k) - \ell(\mathcal{A}_{((\check{\mathcal{D}}_{\sqrt{n}/2}^k)^{\setminus i}, \mathcal{F})}, \check{\mathbf{z}}_j^k) \right|$$

where:

$$\check{\mathcal{D}}_{\sqrt{n}/2}^k : \left\{ \mathbf{z}_{(k-1)\sqrt{n}+1}, \dots, \mathbf{z}_{(k-1)\sqrt{n}+\sqrt{n}/2} \right\}, \quad \check{\mathbf{z}}_j^k : \mathbf{z}_{(k-1)\sqrt{n}+\sqrt{n}/2+j}, \quad k \in \{1, \dots, \sqrt{n}/2\}$$

then:

$$H^{\text{loo}}(\mathcal{A}_{\mathcal{F}}, \sqrt{n}/2) \leq \hat{H}^{\text{loo}}(\mathcal{A}_{(\mathcal{D}_{\sqrt{n}/2}, \mathcal{F})}, \mathcal{D}_{\sqrt{n}/2}) + \sqrt{\log(1/\delta)/\sqrt{n}}$$

²⁵L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Fully empirical and data-dependent stability-based bounds. IEEE Transactions on Cybernetics, 2014.

Fully Empirical Algorithmic Stability Based Bound (III)²⁶

Bound on the true error

$$L(f) \leq \hat{L}_n^{\text{loo}}(f) + \sqrt{\frac{2}{\delta} \left[\frac{1}{2n} + 3 \left(\hat{H}^{\text{loo}}(\mathcal{A}_{(\mathcal{D}_{\sqrt{n}/2}, \mathcal{F})}, \mathcal{D}_{\sqrt{n}/2}) + \sqrt{\frac{\log(\frac{2}{\delta})}{\sqrt{n}}} \right) \right]}.$$

Method for comparing the performance of an algorithm or hyperparameters

$$\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}^* = \arg \min_{\mathcal{A}_{\mathcal{F}} \in \{\mathcal{A}_{\mathcal{F}_1^1}, \mathcal{A}_{\mathcal{F}_2^1}, \dots, \mathcal{A}_{\mathcal{F}_1^2}, \mathcal{A}_{\mathcal{F}_2^2}, \dots\}} E(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}),$$

$$E(\mathcal{A}_{(\mathcal{D}_n, \mathcal{F})}) = \hat{L}_n^{\text{loo}}(f) + \sqrt{\frac{2}{\delta} \left[\frac{1}{2n} + 3 \left(\hat{H}^{\text{loo}}(\mathcal{A}_{(\mathcal{D}_{\sqrt{n}/2}, \mathcal{F})}, \mathcal{D}_{\sqrt{n}/2}) + \sqrt{\frac{\log(\frac{2}{\delta})}{\sqrt{n}}} \right) \right]}.$$

²⁶L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Fully empirical and data-dependent stability-based bounds. IEEE Transactions on Cybernetics, 2014.

Fully Empirical Algorithmic Stability Based Bound (IV)²⁷

Coverttype

$n = 25$				$n = 50$				$n = 100$				$n = 200$				$n = 400$			
KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG
7	7	21	4	9	7	20	2	8	8	17	6	9	7	21	5	12	8	17	6

Mnist OVO

$n = 25$				$n = 50$				$n = 100$				$n = 200$				$n = 400$			
KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG
16	13	45	9	8	11	43	8	10	11	41	5	19	10	43	0	18	11	35	4

Not-Mnist OVO

$n = 25$				$n = 50$				$n = 100$				$n = 200$				$n = 400$			
KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG	KCV	LOO	STAB	MARG
16	14	45	15	13	13	45	11	12	12	42	11	15	14	42	11	21	16	36	11

²⁷L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Fully empirical and data-dependent stability-based bounds. IEEE Transactions on on Cybernetics, 2014.

Outline

Introduction: Model Selection and Error Estimation

Out-Of-Sample Methods (Hold-Out)

In-Sample Methods

Hypotheses Space-Based Methods

Algorithm-Based Methods

Final Remarks

Final Remarks^{28,29}

New results:

- SRM → Energy Aware Learning
- Fully Empirical Algorithmic Stability

Open Problems:

- A long list...

²⁸L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Learning resource-aware models for mobile devices: from regularization to energy efficiency. *Neurocomputing*, 2014.

²⁹L. Oneto, A. Ghio, D. Anguita, and S. Ridella. Fully empirical and data-dependent stability-based bounds. *IEEE Transactions on Cybernetics*, 2014.