

Stability Analysis and Learning Bounds for Transductive Regression Algorithms

Corinna Cortes

CORINNA@GOOGLE.COM

Google Research

76 Ninth Avenue,

New York, NY 10011, USA.

Mehryar Mohri

MOHRI@CS.NYU.EDU

Courant Institute of Mathematical Sciences and Google Research

251 Mercer Street,

New York, NY 10012, USA.

Dmitry Pechyony*

PECHYONY@NEC-LABS.COM

NEC Laboratories America,

4 Independence Way, Suite 200,

Princeton, NJ 08648, USA.

Ashish Rastogi

RASTOGI@CS.NYU.EDU

Google Research

76 Ninth Avenue,

New York, NY 10011, USA.

Editor: X

Abstract

This paper uses the notion of algorithmic stability to derive novel generalization bounds for several families of transductive regression algorithms, both by using convexity and closed-form solutions. Our analysis helps compare the stability of these algorithms. It also shows that a number of widely used transductive regression algorithms are in fact unstable. Finally, it reports the results of experiments with local transductive regression demonstrating the benefit of our stability bounds for model selection, for one of the algorithms, in particular for determining the radius of the local neighborhood used by the algorithm.

Contents

1	Introduction	2
2	Definitions	4
2.1	Transductive learning set-up	4
2.2	Notions of stability	4
3	General transduction stability bounds	5
3.1	Concentration bound for sampling without replacement	6
3.2	Transductive stability bound	8

*. This author's work was done at the Computer Science Department, Technion - Israel Institute of Technology.

4	Stability of local transductive regression algorithms	9
4.1	Local transductive regression algorithms	10
4.2	Generalization bounds	10
5	Stability of unconstrained regularization algorithms	14
5.1	Unconstrained regularization algorithms	14
5.2	Score-based stability analysis	14
5.3	Application	15
5.3.1	Consistency method (CM)	15
5.3.2	Local learning regularization (LL – Reg)	15
5.3.3	Gaussian Mean Fields algorithm	16
5.4	Lower bound on stability coefficient	16
6	Stability of constrained regularization algorithms	17
6.1	Constrained graph regularization algorithms	17
6.2	Score stability of graph Laplacian regularization algorithm	18
6.3	Cost stability of graph Laplacian regularization algorithm	20
6.4	General case	22
7	Experiments	22
8	Conclusion	25

1. Introduction

The problem of *transductive inference* was originally introduced by Vapnik (1982). Many learning problems in information extraction, computational biology, natural language processing and other domains can be formulated as a transductive inference problem. In the transductive setting, the learning algorithm receives both a labeled training set, as in the standard induction setting, and a set of unlabeled test points. The objective is to predict the labels of the test points. No other test points will ever be considered. This setting arises in a variety of applications. Often, there are orders of magnitude more unlabeled points than labeled ones and they have not been assigned a label due to the prohibitive cost of labeling. This motivates the use of transductive algorithms which leverage the unlabeled data during training to improve learning performance.

This paper deals with transductive regression, which arises in problems such as predicting the real-valued labels of the nodes of a fixed (known) graph in computational biology, or the scores associated with known documents in information extraction or search engine tasks. Several algorithms have been devised for the specific setting of transductive regression (Belkin et al., 2004b; Chapelle et al., 1999; Schuurmans and Southey, 2002; Cortes and Mohri, 2007). Several other algorithms introduced for transductive classification can be viewed in fact as transductive regression ones as their objective function is based on the square loss, for example, in Belkin et al. (2004a,b). Cortes and Mohri (2007) gave explicit VC-dimension generalization bounds for transductive regression that hold for all bounded loss functions and coincide with the tight classification bounds of Vapnik (1998) when applied to classification.

We present novel algorithm-dependent generalization bounds for transductive regression. Since they are algorithm-specific, these bounds can often be tighter than bounds based on general com-

plexity measures such as the VC-dimension. Our analysis is based on the notion of algorithmic stability and our learning bounds generalize to the transduction scenario the stability bounds given by Bousquet and Elisseeff (2002) for the inductive setting and extend to regression the stability-based transductive classification bounds of El-Yaniv and Pechyony (2006).

In Section 2 we give a formal definition of the transductive inference learning set-up, including a precise description and discussion of two related transductive settings. We also introduce the notions of cost and score stability used in the following sections.

Standard concentration bounds such as McDiarmid’s bound (McDiarmid, 1989) cannot be readily applied to the transductive regression setting since the points are not drawn independently but uniformly without replacement from a finite set. Instead, Section 3.1 proves a concentration bound generalizing McDiarmid’s bound to the case of random variables sampled without replacement. This bound is slightly stronger than that of El-Yaniv and Pechyony (2006, 2007) and the proof much simpler and more concise. This concentration bound is used to derive a general transductive regression stability bound in Section 3.2. Figure 1 shows the outline of the paper.

Section 4 introduces and examines a very general family of transductive algorithms, that of local transductive regression (LTR) algorithms, a generalization of the algorithm of Cortes and Mohri (2007). It gives general bounds for the stability coefficients of LTR algorithms and uses them to derive stability-based learning bounds for these algorithms. The stability analysis in this section is based on the notion of cost stability and based on convexity arguments.

In Section 5, we analyze a general class of unconstrained optimization algorithms that includes a number of recent algorithms (Wu and Schölkopf, 2007; Zhou et al., 2004; Zhu et al., 2003). The optimization problems for these algorithms admit a closed-form solution. We use that to give a score-based stability analysis of these algorithms. Our analysis shows that in general these algorithms may not be stable. In fact, in Section 5.4 we prove a lower bound on the stability coefficient of these algorithms under some assumptions.

Section 6 examines a class of constrained regularization optimization algorithms for graphs that enjoy better stability properties than the unconstrained ones just mentioned. This includes the graph Laplacian algorithm of Belkin et al. (2004a). In Section 6.2, we give a score stability analysis with novel generalization bounds for this algorithm, simpler and more general than those given by Belkin et al. (2004a). Section 6.3 shows that algorithms based on constrained graph regularizations are in fact special instances of the LTR algorithms by showing that the regularization term can be written in terms of a norm in a reproducing kernel Hilbert space. This is used to derive a cost stability analysis and novel learning bounds for the graph Laplacian algorithm of Belkin et al. (2004a) in terms of the second smallest eigenvalue of the Laplacian and the diameter of the graph. Much of the results of these sections generalize to other constrained regularization optimization algorithms. These generalizations are briefly discussed in Section 6.4 where it is indicated, in particular, how similar constraints can be imposed to the algorithms of Wu and Schölkopf (2007); Zhou et al. (2004); Zhu et al. (2003) to derive new and stable versions of these algorithms.

Finally, Section 7 shows the results of experiments with local transductive regression demonstrating the benefit of our stability bounds for model selection, in particular for determining the radius of the local neighborhood used by the algorithm, which provides a partial validation of our bounds and analysis.

2. Definitions

Let \mathcal{X} denote the input space and \mathcal{Y} a measurable subset of \mathbb{R} .

2.1 Transductive learning set-up

In transductive learning settings, the algorithm receives a labeled training set S of size m , $S = ((x_1, y_1), \dots, (x_m, y_m)) \in \mathcal{X} \times \mathcal{Y}$, and an unlabeled test set T of size u , $x_{m+1}, \dots, x_{m+u} \in \mathcal{X}$. The transductive learning problem consists of predicting accurately the labels y_{m+1}, \dots, y_{m+u} of the test examples, no other test example is ever considered. Two different settings can be distinguished to formalize this problem, see (Vapnik, 1998).

Setting 1 In this setting, a full sample X of $m + u$ examples is given. The learning algorithm further receives the labels of a training sample S of size m selected from X uniformly at random without replacement. The remaining u unlabeled examples serve as a test sample T . We denote by $X = (S, T)$ a partitioning of X into a training set S and test set T .

Setting 2 Here, the training sample S and test sample T are both drawn i.i.d. according to some distribution D . The labeled sample S and the test points T , without their labels, are made available to the learning algorithm.

As in previous theoretical studies of the transduction problem, e.g., (Vapnik, 1998; Derbeko et al., 2004; Cortes and Mohri, 2007; El-Yaniv and Pechyony, 2006), we analyze setting 1 and derive generalization bounds for this specific setting. However, as pointed out by Vapnik (1998), any generalization bound in the setting we analyze directly yields a bound for setting 2 by taking the expectation.

The specific problem where the labels are real-valued numbers, as in the case studied in this paper, is that of *transductive regression*. It differs from the standard *inductive regression* since the learning algorithm is given the unlabeled test examples beforehand and can thus possibly exploit that information to improve its performance.

2.2 Notions of stability

We denote by $c(h, x)$ the cost of an error of a hypothesis h on a point x labeled with $y(x)$. The cost function commonly used in regression is the square loss $c(h, x) = [h(x) - y(x)]^2$. We shall assume a square loss for the remaining of this paper, but many of our results generalize to other convex cost functions. The training error $\hat{R}(h)$ and test error $R(h)$ of a hypothesis h are defined as follows:

$$\hat{R}(h) = \frac{1}{m} \sum_{k=1}^m c(h, x_k) \quad R(h) = \frac{1}{u} \sum_{k=1}^u c(h, x_{m+k}). \quad (1)$$

The generalization bounds we derive are based on the notion of algorithmic stability. We shall use the following two notions of stability in our analysis.

Definition 1 (Cost stability) Let L be a transductive learning algorithm and let h denote the hypothesis returned by L for $X = (S, T)$ and h' the hypothesis returned for $X = (S', T')$, where S and S' differ in exactly one point. L is said to be uniformly β -stable with respect to the cost function c if there exists $\beta \geq 0$ such that for all $x \in X$,

$$|c(h', x) - c(h, x)| \leq \beta. \quad (2)$$

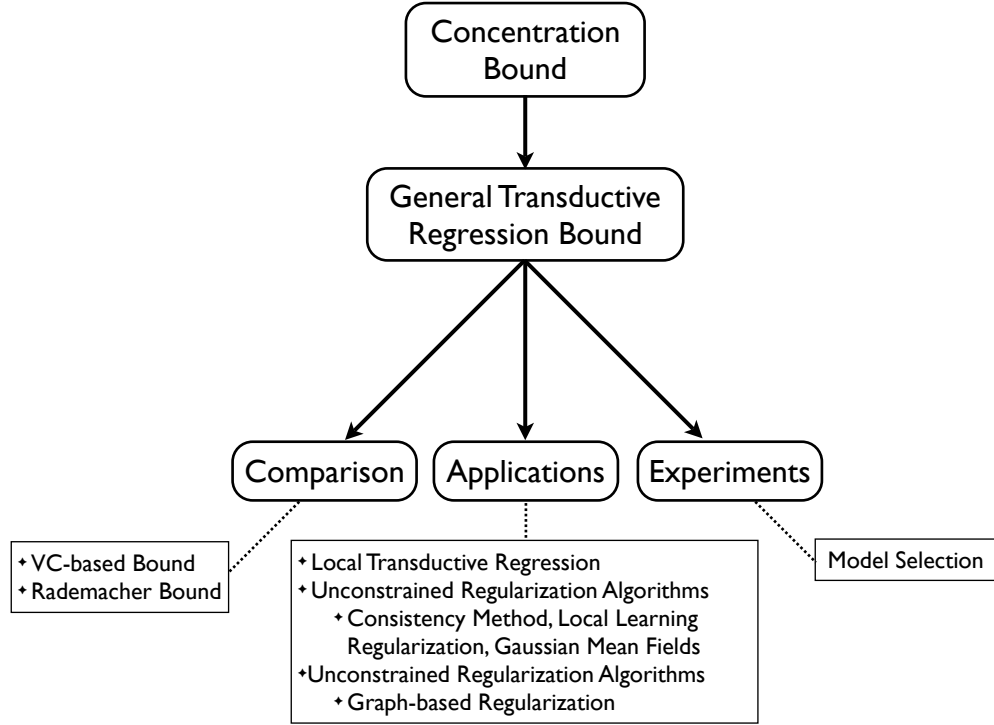


Figure 1: A high-level outline of the paper.

Definition 2 (Score stability) Let L be a transductive learning algorithm and let h denote the hypothesis returned by L for $X = (S, T)$ and h' the hypothesis returned for $X = (S', T')$. L is said to be uniformly β -stable with respect to its output scores if there exists $\beta \geq 0$ such that for all $x \in X$,

$$|h'(x) - h(x)| \leq \beta. \quad (3)$$

We will say that a hypothesis set H is bounded by $B > 0$ when $|h(x) - y(x)| \leq B$ for all $x \in X$ and $h \in H$. For such a hypothesis set and the square loss, for any two hypotheses $h, h' \in H$ and $x \in \mathcal{X}$, the following inequality holds:

$$|c(h', x) - c(h, x)| = |[h'(x) - y(x)]^2 - [h(x) - y(x)]^2| \quad (4)$$

$$= |h'(x) - h(x)| |h'(x) - y(x) + h(x) - y(x)| \quad (5)$$

$$\leq 2B |h'(x) - h(x)|. \quad (6)$$

Thus, for H bounded by B and the square loss, β -score-stability implies $2B\beta$ -cost-stability.

For the remainder of this paper, unless otherwise specified, stability is meant as cost-based stability.

3. General transduction stability bounds

Stability-based generalization bounds in the inductive setting are derived using McDiarmid's inequality (McDiarmid, 1989). The main technique used is to show that under suitable conditions on

the stability of the algorithm, the difference of the test error and the training error is sharply concentrated around its expected value, and that this expected value itself is small. Roughly speaking, this implies that with high probability, the test error is close to the training error. Since the points in the training and test sample are drawn in an i.i.d. fashion, McDiarmid's inequality can be applied.

However, in the transductive setting, the sampling random variables are not drawn independently. Thus, McDiarmid's concentration bound cannot be readily used in this case. Instead, a generalization of McDiarmid's bound that holds for random variables sampled without replacement is needed. We present such a generalization in this section with a concise proof. A slightly weaker version of this bound with a somewhat more complex proof was derived by El-Yaniv and Pechyony (2006, 2007).

3.1 Concentration bound for sampling without replacement

To derive this concentration bound, we use the method of averaged bounded differences and the following theorem due to Azuma (1967) and McDiarmid (1989). where we denote by \mathbf{S}_i^j the subsequence of random variables S_i, \dots, S_j and write $\mathbf{S}_i^j = \mathbf{x}_i^j$ as a shorthand for the event $S_i = x_i, \dots, S_j = x_j$.

Theorem 3 (McDiarmid (1989), Th. 6.10) *Let \mathbf{S}_1^m be a sequence of random variables with each S_i taking values in \mathcal{X} . Let $\phi: \mathcal{X}^m \rightarrow \mathbb{R}$ be a measurable function satisfying the following conditions:*

$$\forall i \in [1, m], \forall x_i, x'_i \in \mathcal{X}, \left| \mathbb{E}_{\mathbf{S}_{i+1}^m} [\phi | \mathbf{S}_1^{i-1}, S_i = x_i] - \mathbb{E}_{\mathbf{S}_{i+1}^m} [\phi | \mathbf{S}_1^{i-1}, S_i = x'_i] \right| \leq c_i.$$

Then, for all $\epsilon > 0$,

$$\Pr \left[\phi - \mathbb{E} [\phi] \geq \epsilon \right] \leq \exp \left[\frac{-2\epsilon^2}{\sum_{i=1}^m c_i^2} \right]. \quad (7)$$

The following definition is needed for the presentation of our concentration bound.

Definition 4 (Symmetric Functions) *A function $\phi: \mathcal{X}^m \rightarrow \mathbb{R}$ is said to be symmetric if its value does not depend on the order of its arguments, that is for any two permutations σ and σ' over $[1, m]$ and any m points $x_1, \dots, x_m \in \mathcal{X}$, $\Phi(x_{\sigma(1)}, \dots, x_{\sigma(m)}) = \Phi(x_{\sigma'(1)}, \dots, x_{\sigma'(m)})$.*

Theorem 5 (Concentration bound for sampling without replacement) *Let \mathbf{x}_1^m be a sequence of random variables, sampled uniformly without replacement from a fixed set X of $m + u$ elements, and let $\phi: \mathcal{X}^m \rightarrow \mathbb{R}$ be a symmetric function such that for all $i \in [1, m]$ and for all $x_1, \dots, x_m \in X$ and $x'_1, \dots, x'_m \in X$,*

$$\left| \phi(x_1, \dots, x_m) - \phi(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_m) \right| \leq c.$$

Then, for all $\epsilon > 0$,

$$\Pr \left[\phi - \mathbb{E} [\phi] \geq \epsilon \right] \leq \exp \left[\frac{-2\epsilon^2}{\alpha(m, u)c^2} \right], \quad (8)$$

where $\alpha(m, u) = \frac{mu}{m+u-1/2} \frac{1}{1-1/(2 \max\{m, u\})}$.

Proof Fix $i \in [1, m]$ and define $g(\mathbf{S}_1^{i-1}, x_i, x'_i)$ as follow:

$$g(\mathbf{S}_1^{i-1}, x_i, x'_i) = \mathbb{E}_{\mathbf{S}_{i+1}^m} [\phi | \mathbf{S}_1^{i-1}, S_i = x_i] - \mathbb{E}_{\mathbf{S}_{i+1}^m} [\phi | \mathbf{S}_1^{i-1}, S_i = x'_i]. \quad (9)$$

Then,

$$\begin{aligned} g(\mathbf{x}_1^{i-1}, x_i, x'_i) &= \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x_i, \mathbf{x}_{i+1}^m) \Pr[\mathbf{S}_{i+1}^m = \mathbf{x}_{i+1}^m | \mathbf{S}_1^{i-1} = \mathbf{x}_1^{i-1}, S_i = x_i] \\ &\quad - \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x'_i, \mathbf{x}_{i+1}^m) \Pr[\mathbf{S}_{i+1}^m = \mathbf{x}_{i+1}^m | \mathbf{S}_1^{i-1} = \mathbf{x}_1^{i-1}, S_i = x'_i]. \end{aligned}$$

We show that $g(\mathbf{x}_1^{i-1}, x_i, x'_i)$ can be bounded by $c_i = \frac{uc}{m+u-i}$ and apply Theorem 3 to obtain the bound claimed. For uniform sampling without replacement, the probability terms can be written explicitly:

$$\Pr[\mathbf{S}_{i+1}^m = \mathbf{x}_{i+1}^m | \mathbf{S}_1^{i-1} = \mathbf{x}_1^{i-1}, S_i = x_i] = \prod_{k=i}^{m-1} \frac{1}{m+u-k} = \frac{u!}{(m+u-i)!}.$$

Thus,

$$g(\mathbf{x}_1^{i-1}, x_i, x'_i) = \frac{u!}{(m+u-i)!} \left[\sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x_i, \mathbf{x}_{i+1}^m) - \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x'_i, \mathbf{x}_{i+1}^m) \right].$$

To compute $\sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x_i, \mathbf{x}_{i+1}^m) - \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x'_i, \mathbf{x}_{i+1}^m)$, we divide the set of permutations $\{\mathbf{x}_{i+1}^m\}$ into two sets, those that contain the element x_i and those that do not. If a permutation \mathbf{x}_{i+1}^m contains x_i we can write it as $\mathbf{x}_{i+1}^{k-1} x_i \mathbf{x}_{k+1}^m$, where k is such that $x'_k = x_i$. We then match it up with the permutation $x_i \mathbf{x}_{i+1}^{k-1} \mathbf{x}_{k+1}^m$ from the set $\{x_i \mathbf{x}_{i+1}^m\}$. These two permutations contain exactly the same elements, and since the function ϕ is symmetric in its arguments, the difference in the value of the function on the permutations is zero.

In the other case, if a permutation \mathbf{x}_{i+1}^m does not contain the element x_i , then we simply match it up with the same permutation in $\{\mathbf{x}_{i+1}^m\}$. The matching permutations appearing in the summation are then $x_i \mathbf{x}_{i+1}^m$ and $x'_i \mathbf{x}_{i+1}^m$ which clearly only differ with respect to x_i . The difference in the value of the function ϕ in this case can be bounded by c . The number of such permutations can be counted as follows: it is the number of permutations of length $m-i$ from the set X of $m+u$ elements that do not contain any of the elements of \mathbf{x}_1^{i-1} , x_i and x'_i , which is equal to $\frac{(m+u-i-1)!}{(u-1)!}$. This leads us to the following upper bound on $\sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x_i, \mathbf{x}_{i+1}^m) - \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x'_i, \mathbf{x}_{i+1}^m)$:

$$\sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x_i, \mathbf{x}_{i+1}^m) - \sum_{\mathbf{x}_{i+1}^m} \phi(\mathbf{x}_1^{i-1}, x'_i, \mathbf{x}_{i+1}^m) \leq \frac{(m+u-i-1)!}{(u-1)!} c, \quad (10)$$

which implies that $|g(\mathbf{x}_1^{i-1}, x_i, x'_i)| \leq \frac{u!}{(m+u-i)!} \cdot \frac{(m+u-i-1)!c}{(u-1)!} \leq \frac{uc}{m+u-i}$. To apply Theorem 3, we need to bound $\sum_{i=1}^m \left(\frac{uc}{m+u-i} \right)^2$. To this end, note that

$$\sum_{i=1}^m \frac{1}{(m+u-i)^2} = \sum_{j=u}^{m+u-1} \frac{1}{j^2} \leq \int_{u-1/2}^{m+u-1/2} \frac{dx}{x^2} = \frac{m}{m+u-1/2} \frac{1}{u-1/2}. \quad (11)$$

The application of Theorem 3 then yields:

$$\Pr \left[\phi - \mathbb{E} [\phi] \geq \epsilon \right] \leq \exp \left[\frac{-2\epsilon^2}{\alpha_u(m, u)c^2} \right], \quad (12)$$

where $\alpha_u(m, u) = \frac{mu}{m+u-1/2} \frac{1}{1-1/(2u)}$. Function ϕ is symmetric in m and u in the sense that selecting one of the sets uniquely determines the other. The statement of the theorem then follows by obtaining a similar bound with $\alpha_m(m, u) = \frac{mu}{m+u-1/2} \frac{1}{1-1/(2m)}$ and taking the tighter of the two bounds. \blacksquare

3.2 Transductive stability bound

Observe that, since the full sample X is given, the average error of a hypothesis $h \in H$ over X defined by $R_X(h) = \frac{1}{m+u} \sum_{i=1}^{m+u} h(x_i)$ is not a random variable. Also, for any training sample S , the test error $R(h)$ can be expressed in terms of $R_X(h)$ and the empirical error $\hat{R}(h)$ as follows:

$$R(h) = \frac{1}{u} \sum_{i=1}^u h(x_{m+i}) = \frac{1}{u} \left((m+u)R_X(h) - \sum_{i=1}^m h(x_i) \right) = \frac{m+u}{u} R_X(h) - \frac{m}{u} \hat{R}(h). \quad (13)$$

Thus, for a fixed h , the quantity $R(h) - \hat{R}(h) = \frac{m+u}{u} R_X(h) - \frac{m+u}{u} \hat{R}(h)$ only varies with $\hat{R}(h)$ and is only a function of the training sample S . Let Φ be defined by $\phi(S) = R(h) - \hat{R}(h)$. Since permuting the points of S does not affect $\hat{R}(h)$, Φ is symmetric.

To obtain a general transductive regression stability bound, we apply the concentration bound of Theorem 5 to the random variable $\phi(S)$. To do so, we need to bound $\mathbb{E}_S [\phi(S)]$, where S is a random subset of X of size m , and $|\phi(S) - \phi(S')|$ where S and S' are samples differing by exactly one point. The following lemma proves a Lipschitz condition for Φ .

Lemma 6 *Let H be a hypothesis set bounded by B . Let L be a β -cost-stable algorithm and let S and S' be two training sets of size m that differ in exactly one point. Let $h \in H$ be the hypothesis returned by L when trained on S and $h' \in H$ the one returned when L is trained on S' . Then,*

$$|\phi(S) - \phi(S')| \leq 2\beta + \frac{m+u}{mu} B^2. \quad (14)$$

Proof By the definition of S' , there exist $i \in [1, m]$ and $j \in [1, u]$ such that $S' = S \setminus \{x_i\} \cup \{x_{m+j}\}$. $\phi(S) - \phi(S')$ can be written as follows:

$$\begin{aligned} \phi(S) - \phi(S') &= \frac{1}{u} \sum_{k=1, k \neq j}^u [c(h, x_{m+k}) - c(h', x_{m+k})] + \frac{1}{m} \sum_{k=1, k \neq i}^m [c(h', x_k) - c(h, x_k)] \\ &\quad + \frac{1}{u} [c(h, x_{m+j}) - c(h', x_i)] + \frac{1}{m} [c(h', x_{m+j}) - c(h, x_i)]. \end{aligned}$$

Since the hypothesis set H is bounded by B , the square loss c is bounded by B^2 , $c(h, x) \leq B^2$ for all $x \in X, h \in H$. Thus,

$$|\phi(S) - \phi(S')| \leq \frac{(u-1)\beta}{u} + \frac{(m-1)\beta}{m} + \frac{B^2}{u} + \frac{B^2}{m} \leq 2\beta + B^2 \left(\frac{1}{u} + \frac{1}{m} \right). \quad (15)$$

■

The next lemma bounds the expectation of Φ .

Lemma 7 *Let h be the hypothesis returned by a β -cost-stable algorithm L . Then, the following inequality holds for the expectation of Φ :*

$$|\mathbb{E}_S [\phi(S)]| \leq \beta. \quad (16)$$

Proof By the definition of $\phi(S)$, we can write

$$\mathbb{E}_S [\phi(S)] = \mathbb{E}_S [R(h)] - \mathbb{E}_S [\widehat{R}(h)] = \frac{1}{u} \sum_{k=1}^u \mathbb{E}_S [c(h, x_{m+k})] - \frac{1}{m} \sum_{k=1}^m \mathbb{E}_S [c(h, x_k)]. \quad (17)$$

$\mathbb{E}_S [c(h, x_{m+k})]$ is the same for all $1 \leq k \leq u$, and similarly, $\mathbb{E}_S [c(h, x_k)]$ is the same for all $1 \leq k \leq m$. Let $i \in [1, m]$ and $j \in [1, u]$, and let S' be defined as in the previous lemma: $S' = S \setminus \{x_i\} \cup \{x_{m+j}\}$, and let h' denote a hypothesis trained on S' , then the following holds:

$$\mathbb{E}_S [\phi(S)] = \mathbb{E}_S [c(h, x_{m+j})] - \mathbb{E}_S [c(h, x_i)] \quad (18)$$

$$= \mathbb{E}_{S'} [c(h', x_i)] - \mathbb{E}_S [c(h, x_i)] \quad (19)$$

$$= \mathbb{E}_{S, S'} [c(h', x_i) - c(h, x_i)] \leq \beta, \quad (20)$$

by the cost β -stability of the algorithm. ■

Theorem 8 *Let H be a hypothesis set bounded by B and L a β -cost-stable algorithm. Let h be the hypothesis returned by L when trained on $X = (S, T)$. Then, for any $\delta > 0$, with probability at least $1 - \delta$,*

$$R(h) \leq \widehat{R}(h) + \beta + \left(2\beta + \frac{B^2(m+u)}{mu}\right) \sqrt{\frac{\alpha(m, u) \ln \frac{1}{\delta}}{2}}. \quad (21)$$

Proof The result follows directly from Theorem 5 and Lemmas 6 and 7. ■

The bound of Theorem 8 is a general bound that applies to *any* transductive algorithm. To apply it, the stability coefficient β , which depends on m and u , needs to be determined. In the subsequent sections, we derive bounds on β for a number of transductive regression algorithms (Cortes and Mohri, 2007; Belkin et al., 2004a; Wu and Schölkopf, 2007; Zhou et al., 2004; Zhu et al., 2003). Note that when $\beta = O(1/\min(m, u))$, the slack term of this bound is in $O(1/\sqrt{\min(m, u)})$.

4. Stability of local transductive regression algorithms

This section describes and analyzes a general family of local transductive regression algorithms (LTR) generalizing the algorithm of Cortes and Mohri (2007).

4.1 Local transductive regression algorithms

LTR algorithms can be viewed as a generalization of the so-called kernel regularization-based learning algorithms to the transductive setting. The objective function that is minimized is of the form:

$$F(f, S) = \|f\|_K^2 + \frac{C}{m} \sum_{k=1}^m c(f, x_k) + \frac{C'}{u} \sum_{k=1}^u \tilde{c}(f, x_{m+k}), \quad (22)$$

where $\|\cdot\|_K$ is the norm in the reproducing kernel Hilbert space (RKHS) with associated kernel K , $C \geq 0$ and $C' \geq 0$ are trade-off parameters, f is the hypothesis and $\tilde{c}(f, x) = (f(x) - \tilde{y}(x))^2$ is the error of f on the unlabeled point x with respect to a pseudo-target \tilde{y} .

Pseudo-targets are obtained from neighborhood labels $y(x)$ by a local weighted average or other regression algorithms applied locally. Neighborhoods can be defined as a ball of radius r around each point in the feature space. We will denote by β_{loc} the score-stability coefficient (Definition 2).

4.2 Generalization bounds

In this section, we use the bounded-labels assumption, that is we shall assume that for all $x \in S$, $|y(x)| \leq M$ for some $M > 0$. We also assume that for any $x \in X$, $K(x, x) \leq \kappa^2$. We will use the following bound based on the reproducing property and the Cauchy-Schwarz inequality valid for any hypothesis $h \in H$, and for all $x \in X$,

$$|h(x)| = |\langle h, K(x, \cdot) \rangle| \leq \|h\|_K \sqrt{K(x, x)} \leq \kappa \|h\|_K. \quad (23)$$

Lemma 9 *Let h be the hypothesis minimizing (22). Assume that for any $x \in X$, $K(x, x) \leq \kappa^2$. Then, for any $x \in X$, $|h(x)| \leq \kappa M \sqrt{C + C'}$.*

Proof The proof is an adaptation of the technique of Bousquet and Elisseeff (2002) to LTR algorithms. By Equation 23, $|h(x)| \leq \kappa \|h\|_K$. Let $\mathbf{0} \in \mathbb{R}^{m+u}$ be the hypothesis assigning label zero to all examples. By the definition of h ,

$$F(h, S) \leq F(\mathbf{0}, S) \leq (C + C')M^2. \quad (24)$$

Using the fact that $\|h\|_K \leq \sqrt{F(h, S)}$ yields the statement of the lemma. \blacksquare

Since $|h(x)| \leq \kappa M \sqrt{C + C'}$, this immediately gives us a bound on $|h(x) - y(x)|$:

$$|h(x) - y(x)| \leq M(1 + \kappa \sqrt{C + C'}), \quad (25)$$

and we are in a position to apply Theorem 8 with $B = AM$, $A = 1 + \kappa \sqrt{C + C'}$.

Let h be a hypothesis obtained by training on S and h' by training on S' . To determine the cost-stability coefficient β , we must upper-bound $|c(h, x) - c(h', x)|$. Let $\Delta h = h - h'$. Then, for all $x \in X$,

$$|c(h, x) - c(h', x)| = \left| \Delta h(x) [(h(x) - y(x)) + (h'(x) - y(x))] \right| \quad (26)$$

$$\leq 2M(1 + \kappa \sqrt{C + C'}) |\Delta h(x)|. \quad (27)$$

As in Inequality 23, for all $x \in X$, $|\Delta h(x)| \leq \kappa \|\Delta h\|_K$, thus for all $x \in X$,

$$|c(h, x) - c(h', x)| \leq 2M(1 + \kappa \sqrt{C + C'}) \kappa \|\Delta h\|_K. \quad (28)$$

It remains to bound $\|\Delta h\|_K$. Our approach towards bounding $\|\Delta h\|_K$ is similar to the one used by Bousquet and Elisseeff (2000), and relies on the convexity of $h \mapsto c(h, x)$. Note however, that in the case of \tilde{c} , the pseudo-targets may depend on the training set S . This dependency matters when we wish to apply convexity of two hypotheses h and h' obtained by training on different samples S and S' . For convenience, for any two such fixed hypotheses h and h' , we extend the definition of \tilde{c} as follows. For all $t \in [0, 1]$,

$$\tilde{c}(th + (1-t)h', x) = ((th + (1-t)h')(x) - (t\tilde{y} + (1-t)\tilde{y}'))^2. \quad (29)$$

This allows us to use the same convexity property for \tilde{c} as for c for any two fixed hypotheses h and h' as verified by the following lemma.

Lemma 10 *Let h be a hypothesis obtained by training on S and h' by training on S' . Then, for all $t \in [0, 1]$,*

$$t\tilde{c}(h, x) + (1-t)\tilde{c}(h', x) \geq \tilde{c}(th + (1-t)h', x). \quad (30)$$

Proof Let $\tilde{y} = \tilde{y}(x)$ be the pseudo-target value at x when the training set is S and $\tilde{y}' = \tilde{y}'(x)$ when the training set is S' . For all $t \in [0, 1]$,

$$\begin{aligned} & tc(h, x) + (1-t)c(h', x) - c(th + (1-t)h', x) \\ &= t(h(x) - \tilde{y})^2 + (1-t)(h'(x) - \tilde{y}')^2 - [(th(x) + (1-t)h'(x) - (t\tilde{y} + (1-t)\tilde{y}'))]^2 \\ &= t(h(x) - \tilde{y})^2 + (1-t)(h'(x) - \tilde{y}')^2 - [t(h(x) - \tilde{y}) + (1-t)(h'(x) - \tilde{y}')]^2. \end{aligned}$$

The statement of the lemma follows directly by the convexity of the function $x \mapsto x^2$ defined over \mathbb{R} . ■

Recall that β_{loc} denotes the score-stability of the algorithm that produces the pseudo-targets. In Lemma 12 we present an upper-bound $\|\Delta h\|_K$, which can then be plugged into Equation 28 to determine the stability of LTR.

Lemma 11 *Assume that for all $x \in X$, $|y(x)| \leq M$. Let S and S' be two samples differing by exactly one point. Let h be the hypothesis returned by the algorithm minimizing the objective function $F(f, S)$, h' be the hypothesis obtained by minimization of $F(f, S')$ and let \tilde{y} and \tilde{y}' be the corresponding pseudo-targets. Then for all $i \in [1, m + u]$,*

$$\begin{aligned} & \frac{C}{m} [c(h', x_i) - c(h, x_i)] + \frac{C'}{u} [\tilde{c}(h', x_i) - \tilde{c}(h, x_i)] \\ & \leq 2AM \left(\kappa \|\Delta h\|_K \left(\frac{C}{m} + \frac{C'}{u} \right) + \beta_{loc} \frac{C'}{u} \right), \end{aligned} \quad (31)$$

where $\Delta h = h' - h$ and $A = 1 + \kappa\sqrt{C + C'}$.

Proof From Equation 28, we know that:

$$|c(h', x_i) - c(h, x_i)| \leq 2M(1 + \kappa\sqrt{C + C'})\kappa\|\Delta h\|_K. \quad (32)$$

It remains to bound $|\tilde{c}(h', x_i) - \tilde{c}(h, x_i)|$.

$$\begin{aligned}\tilde{c}(h', x_i) - \tilde{c}(h, x_i) &= (h'(x) - \tilde{y}'(x))^2 - (h(x) - \tilde{y}(x))^2 \\ &= ((h'(x) - \tilde{y}'(x)) + (h(x) - \tilde{y}(x))) (\Delta h(x) - (\tilde{y}'(x) - \tilde{y}(x))) \\ &\leq 2M(1 + \kappa\sqrt{C + C'}) (\kappa\|\Delta h\|_K + \beta_{loc})\end{aligned}$$

Here, we are using score-stability β_{loc} of the local algorithm in $|\tilde{y}'(x) - \tilde{y}(x)| \leq \beta_{loc}$ and that $|h(x) - \tilde{y}(x)| \leq M(1 + \kappa\sqrt{C + C'})$ when $|\tilde{y}(x)| \leq M$ (by Lemma 9).

Plugging the bounds for $|c(h', x_i) - c(h, x_i)|$ and $|\tilde{c}(h', x_i) - \tilde{c}(h, x_i)|$ into the left hand side of Equation 31 yields the statement of the lemma. \blacksquare

Lemma 12 *Assume that for all $x \in X$, $|y(x)| \leq M$. Let S and S' be two samples differing by exactly one point. Let h be the hypothesis returned by the algorithm minimizing the objective function $F(f, S)$, h' the hypothesis obtained by minimization of $F(f, S')$ and let \tilde{y} and \tilde{y}' be the corresponding pseudo-targets. Then*

$$\|\Delta h\|_K^2 \leq 2AM \left(\kappa\|\Delta h\|_K \left(\frac{C}{m} + \frac{C'}{u} \right) + \beta_{loc} \frac{C'}{u} \right), \quad (33)$$

where $\Delta h = h' - h$ and $A = 1 + \kappa\sqrt{C + C'}$.

Proof By the definition of h and h' , we have

$$h = \operatorname{argmin}_{f \in H} F(f, S) \quad \text{and} \quad h' = \operatorname{argmin}_{f \in H} F(f, S').$$

Let $t \in [0, 1]$. Then $h + t\Delta h$ and $h' - t\Delta h$ satisfy:

$$F(h, S) - F(h + t\Delta h, S) \leq 0 \quad (34)$$

$$F(h', S') - F(h' - t\Delta h, S') \leq 0 \quad (35)$$

For notational ease, let $h_{t\Delta}$ denote $h + t\Delta h$ and $h'_{t\Delta}$ denote $h' - t\Delta h$. Summing the two inequalities in Equations 34 and 35 yields:

$$\begin{aligned}& \frac{C}{m} \sum_{k=1}^m [c(h, x_k) - c(h_{t\Delta}, x_k)] + \frac{C'}{u} \sum_{k=1}^u [\tilde{c}(h, x_{m+k}) - \tilde{c}(h_{t\Delta}, x_{m+k})] + \\ & \frac{C}{m} \sum_{k=1, k \neq i}^m [c(h', x_k) - c(h'_{t\Delta}, x_k)] + \frac{C'}{u} \sum_{k=1, k \neq j}^u [\tilde{c}(h', x_{m+k}) - \tilde{c}(h'_{t\Delta}, x_{m+k})] + \\ & \frac{C}{m} [c(h', x_{m+j}) - c(h'_{t\Delta}, x_{m+j})] + \frac{C'}{u} [\tilde{c}(h', x_i) - \tilde{c}(h'_{t\Delta}, x_i)] + \\ & \|h\|_K^2 - \|h_{t\Delta}\|_K^2 + \|h'\|_K^2 - \|h'_{t\Delta}\|_K^2 \leq 0.\end{aligned}$$

By the convexity of $c(h, \cdot)$ in h , it follows that for all $k \in [1, m + u]$

$$c(h, x_k) - c(h_{t\Delta}, x_k) \geq t [c(h, x_k) - c(h + \Delta h, x_k)], \quad (36)$$

and

$$c(h', x_k) - c(h'_{t\Delta}, x_k) \geq t [c(h', x_k) - c(h' - \Delta h, x_k)]. \quad (37)$$

By Lemma 10, similar inequalities hold for \tilde{c} . These observations lead to:

$$\begin{aligned} & \frac{Ct}{m} \sum_{k=1}^m [c(h, x_k) - c(h', x_k)] + \frac{C't}{u} \sum_{k=1}^u [\tilde{c}(h, x_{m+k}) - \tilde{c}(h', x_{m+k})] + \\ & \frac{Ct}{m} \sum_{k=1, k \neq i}^m [c(h', x_k) - c(h, x_k)] + \frac{C't}{u} \sum_{k=1, k \neq j}^u [\tilde{c}(h', x_{m+k}) - \tilde{c}(h, x_{m+k})] + \\ & \frac{Ct}{m} [c(h', x_{m+j}) - c(h, x_{m+j})] + \frac{C't}{u} [\tilde{c}(h', x_i) - \tilde{c}(h, x_i)] + \\ & \|h\|_K^2 - \|h_{t\Delta}\|_K^2 + \|h'\|_K^2 - \|h'_{t\Delta}\|_K^2 \leq 0. \end{aligned}$$

Let E denote $\|h\|_K^2 - \|h_{t\Delta}\|_K^2 + \|h'\|_K^2 - \|h'_{t\Delta}\|_K^2$. Simplifying the previous inequality leads to:

$$\begin{aligned} E & \leq \frac{Ct}{m} [c(h', x_i) - c(h, x_i) + c(h, x_{m+j}) - c(h', x_{m+j})] - \\ & \frac{C't}{u} [\tilde{c}(h', x_i) - \tilde{c}(h, x_i) + \tilde{c}(h, x_{m+j}) - \tilde{c}(h', x_{m+j})]. \end{aligned}$$

Let $A = 1 + \kappa\sqrt{C + C'}$. Using Lemma 11 twice (with x_i and x_{m+j}), the expression above can be bounded by

$$E \leq 4AMt \left(\kappa \|\Delta h\|_K \left(\frac{C}{m} + \frac{C'}{u} \right) + \beta_{loc} \frac{C'}{u} \right). \quad (38)$$

Finally, since $\|h\|_K^2 = \langle h, h \rangle_K$ for any $h \in H$, it is not hard to show that:

$$\|h\|_K^2 - \|h + t\Delta h\|_K^2 + \|h'\|_K^2 - \|h' - t\Delta h\|_K^2 = 2t\|\Delta h\|_K^2(1 - t). \quad (39)$$

Using Equation 39 in Equation 38, it follows that:

$$\|\Delta h\|_K^2(1 - t) \leq 2AM \left(\kappa \|\Delta h\|_K \left(\frac{C}{m} + \frac{C'}{u} \right) + \beta_{loc} \frac{C'}{u} \right). \quad (40)$$

Taking the limit as $t \rightarrow 0$ yields the statement of the lemma. ■

The following is the main result of this section, a stability-based generalization bound for LTR.

Theorem 13 *Assume that for all $x \in X$, $|y(x)| \leq M$ and there exists κ such that for all $x \in X$, $K(x, x) \leq \kappa^2$. Further, assume that the local estimator has score-stability β_{loc} . Let $A = 1 + \kappa\sqrt{C + C'}$. Then, LTR is uniformly β -cost-stable with*

$$\beta \leq 2(AM)^2 \kappa^2 \left[\frac{C}{m} + \frac{C'}{u} + \sqrt{\left(\frac{C}{m} + \frac{C'}{u} \right)^2 + \frac{2C'\beta_{loc}}{AM\kappa^2 u}} \right].$$

Proof From Lemma 12, we know that

$$\|\Delta h\|_K^2 \leq 2AM \left(\kappa \|\Delta h\|_K \left(\frac{C}{m} + \frac{C'}{u} \right) + \beta_{loc} \frac{C'}{u} \right), \quad (41)$$

where $\Delta h = h' - h$ and $A = 1 + \kappa \sqrt{C + C'}$. This implies that $\|\Delta h\|$ is bounded by the non-negative root of the second-degree polynomial which gives

$$\|\Delta h\|_K \leq AM\kappa \left[\left(\frac{C}{m} + \frac{C'}{u} \right) + \sqrt{\left(\frac{C}{m} + \frac{C'}{u} \right)^2 + \frac{2C'\beta_{loc}}{AM\kappa^2 u}} \right]. \quad (42)$$

Using the above bound on $\|\Delta h\|_K$ in Equation 28 yields the desired bound on the stability coefficient of LTR and completes the proof. \blacksquare

Our experiments with local transductive regression in Section 7 will show the benefit of this bound for model selection.

5. Stability of unconstrained regularization algorithms

5.1 Unconstrained regularization algorithms

In this section, we consider a family of transductive regression algorithms that can be formulated as the following optimization problem:

$$\min_{\mathbf{h}} \mathbf{h}^\top \mathbf{Q} \mathbf{h} + (\mathbf{h} - \mathbf{y})^\top \mathbf{C} (\mathbf{h} - \mathbf{y}), \quad (43)$$

where $\mathbf{Q} \in \mathbb{R}^{(m+u) \times (m+u)}$ is a symmetric regularization matrix, $\mathbf{C} \in \mathbb{R}^{(m+u) \times (m+u)}$ a symmetric matrix of empirical weights (in practice it is often a diagonal matrix), $\mathbf{y} \in \mathbb{R}^{(m+u) \times 1}$ the target values of the m labeled points together with the pseudo-target values of the u unlabeled points (in some formulations, the pseudo-target value is 0), and $\mathbf{h} \in \mathbb{R}^{(m+u) \times 1}$ a column matrix whose i th row is the predicted target value for the x_i . The closed-form solution of (43) is given by

$$\mathbf{h} = (\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1} \mathbf{y}. \quad (44)$$

The formulation (43) is quite general and includes as special cases the algorithms of Belkin et al. (2004a); Wu and Schölkopf (2007); Zhou et al. (2004); Zhu et al. (2003). We present a general framework for bounding the stability coefficient of these algorithms and then examine the stability coefficient of each of these algorithms in turn.

5.2 Score-based stability analysis

For a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ we denote by $\lambda_M(\mathbf{A})$ its largest and by $\lambda_m(\mathbf{A})$ its smallest eigenvalue. Thus, for any $\mathbf{v} \in \mathbb{R}^{n \times 1}$, $\lambda_m(\mathbf{A}) \|\mathbf{v}\|_2 \leq \|\mathbf{A} \mathbf{v}\|_2 \leq \lambda_M(\mathbf{A}) \|\mathbf{v}\|_2$. We will also use, in the proof of the following proposition, the fact that for symmetric matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\lambda_M(\mathbf{AB}) \leq \lambda_M(\mathbf{A}) \lambda_M(\mathbf{B})$.

Theorem 14 *Let \mathbf{h}^* and \mathbf{h}'^* solve (43), under test and training sets that differ exactly in one point and let $\mathbf{C}, \mathbf{C}', \mathbf{y}, \mathbf{y}'$ be the corresponding empirical weight and the target value matrices. Then,*

$$\|\mathbf{h}^* - \mathbf{h}'^*\|_\infty \leq \|\mathbf{h}^* - \mathbf{h}'^*\|_2 \leq \frac{\|\mathbf{y} - \mathbf{y}'\|_2}{\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1} + \frac{\lambda_M(\mathbf{Q}) \|\mathbf{C}'^{-1} - \mathbf{C}^{-1}\|_2 \|\mathbf{y}'\|_2}{\left(\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C}')} + 1 \right) \left(\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1 \right)} \quad (45)$$

Proof The first inequality holds as a result of the general relation between norm-infinity and norm-2. Let $\Delta \mathbf{h}^* = \mathbf{h}^* - \mathbf{h}'^*$ and $\Delta \mathbf{y} = \mathbf{y} - \mathbf{y}'$. By definition,

$$\Delta \mathbf{h}^* = (\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1} \mathbf{y} - (\mathbf{C}'^{-1} \mathbf{Q} + \mathbf{I})^{-1} \mathbf{y}' \quad (46)$$

$$= (\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1} \Delta \mathbf{y} + ((\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1} - (\mathbf{C}'^{-1} \mathbf{Q} + \mathbf{I})^{-1}) \mathbf{y}' \quad (47)$$

$$= (\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1} \Delta \mathbf{y} + [(\mathbf{C}'^{-1} \mathbf{Q} + \mathbf{I})^{-1} - (\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1}] \mathbf{y}'. \quad (48)$$

Since $\|(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1}\|_2 = \lambda_{\max}[(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})^{-1}] = \lambda_{\min}(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})$, and $\lambda_m(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I}) \geq \frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1$, $\|\Delta \mathbf{h}^*\|_2$ can be bounded as follows:

$$\|\Delta \mathbf{h}^*\|_2 \leq \frac{\|\Delta \mathbf{y}\|_2}{\lambda_m(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})} + \frac{\lambda_M(\mathbf{Q}) \|\mathbf{C}'^{-1} - \mathbf{C}^{-1}\|_2 \|\mathbf{y}'\|_2}{\lambda_m(\mathbf{C}'^{-1} \mathbf{Q} + \mathbf{I}) \lambda_m(\mathbf{C}^{-1} \mathbf{Q} + \mathbf{I})} \quad (49)$$

$$\leq \frac{\|\Delta \mathbf{y}\|_2}{\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1} + \frac{\lambda_M(\mathbf{Q}) \|\mathbf{C}'^{-1} - \mathbf{C}^{-1}\|_2 \|\mathbf{y}'\|_2}{\left(\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1\right) \left(\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1\right)}. \quad (50)$$

This proves the second inequality. ■

The theorem helps derive score-stability bounds for various transductive regression algorithms (Zhou et al., 2004; Wu and Schölkopf, 2007; Zhu et al., 2003) based on the closed-form solution for the hypothesis. Recall that score-stability (Definition 2) is the maximum change in the hypothesis score on any point x as the learning algorithm is trained on two training sets that differ in exactly one point, that is precisely an upper-bound on $\|\mathbf{h}^* - \mathbf{h}'^*\|_\infty$.

5.3 Application

For each of the algorithms in (Zhou et al., 2004; Wu and Schölkopf, 2007; Zhu et al., 2003), an estimate of 0 is used for unlabeled points. Thus, the vector \mathbf{y} has the following structure: the entries corresponding to training examples are their true labels and those corresponding to the unlabeled examples are 0.

For each one of the three algorithms, we make the bounded labels assumption (for all $x \in X$, $|y(x)| \leq M$ for some $M > 0$). It is then not difficult to show that $\|\mathbf{y} - \mathbf{y}'\|_2 \leq \sqrt{2}M$ and $\|\mathbf{y}'\|_2 \leq \sqrt{m}M$. Furthermore, all the stability bounds derived are based on the notion of score-stability (Definition 2).

5.3.1 CONSISTENCY METHOD (CM)

In the CM algorithm (Zhou et al., 2004), the matrix \mathbf{Q} is a normalized Laplacian of a weight matrix $\mathbf{W} \in \mathbb{R}^{(m+u) \times (m+u)}$ that captures affinity between pairs of points in the full sample X . Thus, $\mathbf{Q} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$, where $\mathbf{D} \in \mathbb{R}^{(m+u) \times (m+u)}$ is a diagonal matrix, with $[\mathbf{D}]_{i,i} = \sum_j [\mathbf{W}]_{i,j}$. Note that $\lambda_m(\mathbf{Q}) = 0$. Furthermore, matrices \mathbf{C} and \mathbf{C}' are identical in CM, both diagonal matrices with (i, i) th entry equal to a positive constant $\mu > 0$. Thus $\mathbf{C}^{-1} = \mathbf{C}'^{-1}$ and using Proposition 14, we obtain the following bound on the score-stability of the CM algorithm: $\beta_{\text{CM}} \leq \sqrt{2}M$.

5.3.2 LOCAL LEARNING REGULARIZATION (LL - Reg)

In the LL - Reg algorithm (Wu and Schölkopf, 2007), the regularization matrix \mathbf{Q} is $(\mathbf{I} - \mathbf{A})^\top (\mathbf{I} - \mathbf{A})$, where $\mathbf{I} \in \mathbb{R}^{(m+u) \times (m+u)}$ is an identity matrix and $\mathbf{A} \in \mathbb{R}^{(m+u) \times (m+u)}$ is a non-negative

weight matrix that captures the local similarity between all pairs of points in X . \mathbf{A} is normalized, i.e. each of its rows sum to 1. Let $C_l, C_u > 0$ be two positive constants. The matrix \mathbf{C} is a diagonal matrix with $[\mathbf{C}]_{i,i} = C_l$ if $x_i \in S$ and C_u otherwise. Let $C_{\max} = \max\{C_l, C_u\}$ and $C_{\min} = \min\{C_l, C_u\}$. Thus, $\|\mathbf{C}'^{-1} - \mathbf{C}^{-1}\|_2 = \sqrt{2} \left(\frac{1}{C_{\min}} - \frac{1}{C_{\max}} \right)$. By the Perron-Frobenius theorem, its eigenvalues lie in the interval $(-1, 1]$ and $\lambda_M(\mathbf{A}) \leq 1$. Thus, $\lambda_m(\mathbf{Q}) \geq 0$ and $\lambda_M(\mathbf{Q}) \leq 4$ and we have the following bound on the score-stability of the LL-Reg algorithm: $\beta_{\text{LL-Reg}} \leq \sqrt{2}M + 4\sqrt{2m}M \left(\frac{1}{C_{\min}} - \frac{1}{C_{\max}} \right) \leq \sqrt{2}M + \frac{4\sqrt{2m}M}{C_{\min}}$.

5.3.3 GAUSSIAN MEAN FIELDS ALGORITHM

GMF (Zhu et al., 2003) is very similar to the LL-Reg, and admits exactly the same stability coefficient.

Thus, using our bounding technique, the stability coefficients of the algorithms of CM, LL-Reg, and GMF can be large. Without additional constraints on the matrix \mathbf{Q} , these algorithms do not seem to be stable enough for the generalization bound of Theorem 8 to converge. The next section in fact demonstrates that by presenting a constant lower bound on their score-stability.

5.4 Lower bound on stability coefficient

The stability coefficient is a function of the sample size. For stability learning bounds to converge, it must go to zero as a function of the sample size. The following theorem proves that the stability coefficient of the CM algorithm is lower-bounded by a constant for some problems. A similar lower bound can be given for the other two algorithms examined.

Theorem 15 *There exists a transductive regression problem with $m \geq 2$ labeled samples and m unlabeled samples and a diagonal matrix \mathbf{C} for which the score-stability β of the CM algorithm admits the following lower bound:*

$$\beta \geq \frac{1}{2} \frac{C}{C+1}. \quad (51)$$

Proof Consider a transductive regression problem with $2m$ instances where m instances have a target value of 0 and the other m instances a target value of 1. Let the labeled sample S include exactly the instances x_1, \dots, x_m with target value 0 and U be defined by the complement x_{m+1}, \dots, x_{2m} . Let \mathbf{L} denote an $m \times m$ normalized graph Laplacian matrix, with 1s along the diagonal and all off-diagonal terms equal to $-\frac{1}{m-1}$. Then the matrix \mathbf{Q} is defined with the following block structure:

$$\mathbf{Q} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{0} & \mathbf{L} \end{bmatrix} \quad (52)$$

In our example, we set \mathbf{C} to be a diagonal matrix with all its entries equal to the constant C . The matrix $\mathbf{M} = \mathbf{C}^{-1}\mathbf{Q} + \mathbf{I}$ has the following block structure:

$$\mathbf{M} = \begin{bmatrix} \mathbf{N} & \mathbf{0} \\ \mathbf{0} & \mathbf{N} \end{bmatrix}, \quad (53)$$

where \mathbf{N} is the $m \times m$ matrix whose diagonal entries are all equal to $1 + \frac{1}{C}$ and whose off-diagonal entries all equal to $-\frac{1}{C(m-1)}$.

Now, consider the training sample S' obtained from S by swapping a labeled point with an unlabeled point. For the sake of convenience, let the index of this point be m . The \mathbf{y} vector changes (to \mathbf{y}') only in the m th position. Thus, all the entries of $\Delta\mathbf{y} = \mathbf{y} - \mathbf{y}'$ are zero except from its m th entry which equals 1. By Equation 48, $\Delta\mathbf{h}^* = \mathbf{M}^{-1}\Delta\mathbf{y}$, thus, $\Delta\mathbf{h}^*$ is exactly the m th column of \mathbf{M}^{-1} . Let $[\mathbf{M}^{-1}]_{m,m} = [\mathbf{N}^{-1}]_{m,m}$ denote the (m, m) entry of \mathbf{M}^{-1} which coincides with the (m, m) entry of \mathbf{N}^{-1} . Since $\|\Delta\mathbf{h}^*\|_2 \geq |[\mathbf{M}^{-1}]_{m,m}|$, to give a lower bound on $\|\Delta\mathbf{h}^*\|_2$, it suffices to lower bound $|[\mathbf{N}^{-1}]_{m,m}|$. To do so, we can compute $[\mathbf{N}^{-1}]_{m,m}$.

By symmetry, the diagonal entries of \mathbf{N}^{-1} are all equal to some value a , thus $ma = \text{Tr}(\mathbf{N}^{-1})$, which can be computed from the inverses of the eigenvalues of \mathbf{N} . Observe that $\mathbf{N}_0 = \mathbf{N} - (1 + \frac{1}{C} + \frac{1}{C(m-1)})\mathbf{I}$ is a matrix with all entries equal to $-\frac{1}{C(m-1)}$. Thus, it is a rank one matrix and its only non-zero eigenvalue coincides with its trace: $\text{Tr}(\mathbf{N}_0) = -\frac{m}{C(m-1)}$. Since $1 + \frac{1}{C} + \frac{1}{C(m-1)} = \frac{(m-1)C+m}{(m-1)C}$, this shows that the eigenvalues of \mathbf{N} are $\frac{m}{C(m-1)} + \frac{(m-1)C+m}{(m-1)C} = 1$ with multiplicity 1, and $\frac{(m-1)C+m}{(m-1)C}$ with multiplicity $m-1$. Thus, $ma = \text{Tr}(\mathbf{N}^{-1}) = 1 + \frac{(m-1)^2C}{(m-1)C+m}$, which gives

$$a = \frac{1}{m} + \frac{\frac{m-1}{m}C}{C + \frac{m}{m-1}} \geq \frac{\frac{1}{2}C}{C+1}, \quad (54)$$

since for $m \geq 2$, $\frac{1}{2} \leq \frac{m-1}{m} \leq 1$. ■

An example of constraint that can help guarantee stability is the condition $\sum_{i=1}^{m+u} h(x_i) = 0$ used in the algorithm of Belkin et al. (2004a). In the next section, we give a generalization bound for a family of algorithms based on this constraint and then describe a general method for making the algorithms just examined stable.

6. Stability of constrained regularization algorithms

6.1 Constrained graph regularization algorithms

Here, we examine constrained regularization algorithms such as the graph Laplacian regularization algorithm of Belkin et al. (2004a). Given a weighted graph $G = (X, E)$ in which edge weights can be interpreted as similarities between vertices, the task consists of predicting the vertex labels. The input space X is thus reduced to the set of vertices, and a hypothesis $h: X \rightarrow \mathbb{R}$ can be identified with the finite-dimensional vector \mathbf{h} of its predictions $\mathbf{h} = [h(x_1), \dots, h(x_{m+u})]^\top$. The hypothesis set H can thus be identified with \mathbb{R}^{m+u} here. Let \mathbf{h}_S denote the restriction of \mathbf{h} to the training points, $[h(x_1), \dots, h(x_m)]^\top \in \mathbb{R}^m$, and similarly let \mathbf{y}_S denote $[y_1, \dots, y_m]^\top \in \mathbb{R}^m$.

The general family of constrained graph regularization algorithms can then be defined by the following optimization problem:

$$\begin{aligned} \min_{\mathbf{h} \in H} \quad & \mathbf{h}^\top \mathbf{L} \mathbf{h} + \frac{C}{m} (\mathbf{h}_S - \mathbf{y}_S)^\top (\mathbf{h}_S - \mathbf{y}_S) \\ \text{subject to: } & \mathbf{h}^\top \mathbf{u} = 0, \end{aligned} \quad (55)$$

where $\mathbf{L} \in \mathbb{R}^{(m+u) \times (m+u)}$ is a positive semi-definite symmetric matrix, y_i , $i \in [1, m]$, the target values of the m labeled nodes, and $\mathbf{u} \in \mathbb{R}^{m+u}$ a fixed vector. The constraint of the optimization thus restricts the space of solutions to be in H_1 , the hyperplane in H of the vectors orthogonal to \mathbf{u} . We denote by \mathbf{P} the projection matrix over the hyperplane H_1 . As further discussed later, for stability

reasons, \mathbf{u} is typically selected to be orthogonal to the range of \mathbf{L} , $\text{range}(\mathbf{L})$. More generally, the optimizations constraint can be generalized to orthogonality with respect to a subspace U such that the space of solutions H_1 be a subset of $\text{range}(\mathbf{L})$.

In the case of the regularization algorithm of Belkin et al. (2004a), \mathbf{L} is the graph Laplacian. Thus, $\mathbf{h}^\top \mathbf{L} \mathbf{h} = \sum_{i,j=1}^m w_{ij} (h(x_i) - h(x_j))^2$, for some weight matrix (w_{ij}) . The vector \mathbf{u} is defined to be $\mathbf{1}$, that is all its entries equal 1. For this algorithm, the authors further assume the label vector \mathbf{y} to be centered, which implies that $\mathbf{u}^\top \mathbf{y} = 0$, and also that the graph G is connected. This last assumption implies that the zero eigenvalue of the Laplacian has multiplicity one and that H_1 coincides with $\text{range}(\mathbf{L})$.

For a sample S drawn without replacement from X , define $\mathbf{I}_S \in \mathbb{R}^{(m+u) \times (m+u)}$ as the diagonal matrix with $[\mathbf{I}_S]_{i,i} = 1$ if $x_i \in S$ and 0 otherwise. Similarly, let $\mathbf{y}_S \in \mathbb{R}^{(m+u) \times 1}$ be the vector with $[\mathbf{y}_S]_{i,1} = y_i$ if $x_i \in S$ and 0 otherwise. Then, the Lagrangian associated to the problem (55) is $\mathcal{L} = \mathbf{h}^\top \mathbf{L} \mathbf{h} + \frac{C}{m} (\mathbf{h}_S - \mathbf{y}_S)^\top (\mathbf{h}_S - \mathbf{y}_S) + \beta \mathbf{h}^\top \mathbf{u}$, where $\beta \in \mathbb{R}$ is a Lagrange variable. Setting its gradient with respect to \mathbf{h} to zero gives

$$\mathbf{L} \mathbf{h} + \frac{C}{m} (\mathbf{h}_S - \mathbf{y}_S) + \beta \mathbf{u} = 0. \quad (56)$$

Multiplying by the projection matrix \mathbf{P} gives

$$\mathbf{P}(\mathbf{L} + \frac{C}{m} \mathbf{I}_S) \mathbf{h} = \frac{C}{m} \mathbf{P} \mathbf{y}_S - \beta \mathbf{P} \mathbf{u} = \frac{C}{m} \mathbf{P} \mathbf{y}_S. \quad (57)$$

6.2 Score stability of graph Laplacian regularization algorithm

This section gives a simple generalization bound for the graph Laplacian regularization algorithm using a closed-form solution of (57) and a score-stability analysis.

In the case of the graph Laplacian regularization algorithm of Belkin et al. (2004a) with the assumptions already indicated, matrix $\mathbf{P}(\frac{m}{C} \mathbf{L} + \mathbf{I}_S)$ is invertible. Then, Equation (57) gives the closed-form solution:

$$\mathbf{h} = [\mathbf{P}(\frac{m}{C} \mathbf{L} + \mathbf{I}_S)]^{-1} \mathbf{P} \mathbf{y}_S, \quad (58)$$

which clearly verifies the constraint of the optimization problem.

Theorem 16 *Assume that the graph $G = (X, E)$ is connected and that its vertex labels are bounded: for all x , $|y(x)| \leq M$ for some $M > 0$. Let h denote the solution of the optimization problem (55) where \mathbf{L} is the graph Laplacian and $\mathbf{u} = \mathbf{1}$, and let $A = 1 + \kappa \sqrt{C}$. Then, for any $\delta > 0$, with probability at least $1 - \delta$,*

$$R(h) \leq \widehat{R}(h) + \beta + \left(2\beta + \frac{(AM)^2(m+u)}{mu} \right) \sqrt{\frac{\alpha(m, u) \ln \frac{1}{\delta}}{2}}, \quad (59)$$

where

$$\alpha(m, u) = \frac{mu}{m+u-1/2} \frac{1}{1-1/(2 \max\{m, u\})} \quad \text{and} \quad \beta \leq \frac{4\sqrt{2}M^2}{m\lambda_2/C-1} + \frac{4\sqrt{2m}M^2}{(m\lambda_2/C-1)^2},$$

λ_2 being the second smallest eigenvalue of the Laplacian \mathbf{L} .

Proof Our proof is similar to that of Theorem 5 in (Belkin et al., 2004a), with the important exception that we no longer need to cope with vertex multiplicity in sampling since S is sampled from X without replacement. This makes our proof and the resulting bound considerably simpler and more concise.

By Lemma 9 and Equation 25, since the labels are bounded by M , for any x , the following inequality holds: $|h(x) - y(x)| \leq M(1 + \kappa\sqrt{C}) = AM$. To determine the stability coefficient, it suffices to bound $\max_{S, S'} \|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty$, where S and S' are two training sets that differ only in one vertex. Let $\mathbf{M}_S = \mathbf{P} \left(\frac{m}{C} \mathbf{L} + \mathbf{I}_S \right)$ and $\mathbf{M}_{S'} = \mathbf{P} \left(\frac{m}{C} \mathbf{L} + \mathbf{I}_{S'} \right)$. Then,

$$\|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty \leq \|\mathbf{h}_S - \mathbf{h}_{S'}\| \quad (60)$$

$$= \|\mathbf{M}_S^{-1} \mathbf{P} \mathbf{y}_S - \mathbf{M}_{S'}^{-1} \mathbf{P} \mathbf{y}_{S'}\| \quad (61)$$

$$= \|\mathbf{M}_S^{-1} \mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'}) + (\mathbf{M}_S^{-1} - \mathbf{M}_{S'}^{-1}) \mathbf{P} \mathbf{y}_{S'}\| \quad (62)$$

$$\leq \|\mathbf{M}_S^{-1} \mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'})\| + \|(\mathbf{M}_S^{-1} - \mathbf{M}_{S'}^{-1}) \mathbf{P} \mathbf{y}_{S'}\|. \quad (63)$$

For any column matrix $\mathbf{v} \in \mathbb{R}^{(m+u) \times 1}$, by the triangle inequality and the projection property $\|\mathbf{P} \mathbf{v}\| \leq \|\mathbf{v}\|$, the following inequalities hold:

$$\left\| \frac{m}{C} \mathbf{P} \mathbf{L} \right\| = \left\| \frac{m}{C} \mathbf{P} \mathbf{L} + \mathbf{P} \mathbf{I}_S \mathbf{v} - \mathbf{P} \mathbf{I}_S \mathbf{v} \right\| \quad (64)$$

$$\leq \left\| \frac{m}{C} \mathbf{P} \mathbf{L} + \mathbf{P} \mathbf{I}_S \mathbf{v} \right\| + \|\mathbf{P} \mathbf{I}_S \mathbf{v}\| \quad (65)$$

$$\leq \left\| \mathbf{P} \left(\frac{m}{C} \mathbf{L} + \mathbf{I}_S \right) \mathbf{v} \right\| + \|\mathbf{I}_S \mathbf{v}\|. \quad (66)$$

This yields the lower bound:

$$\|\mathbf{M}_S \mathbf{v}\| = \left\| \mathbf{P} \left(\frac{m}{C} \mathbf{L} + \mathbf{I}_S \right) \mathbf{v} \right\| \geq \frac{m}{C} \|\mathbf{P} \mathbf{L}\| - \|\mathbf{I}_S \mathbf{v}\| \geq \left(\frac{m}{C} \lambda_2 - 1 \right) \|\mathbf{v}\|, \quad (67)$$

which gives the following upper bound on $\|\mathbf{M}_S^{-1}\|, \|\mathbf{M}_{S'}^{-1}\|$:

$$\|\mathbf{M}_S^{-1}\| \leq \frac{1}{\frac{m}{C} \lambda_2 - 1} \quad \text{and} \quad \|\mathbf{M}_{S'}^{-1}\| \leq \frac{1}{\frac{m}{C} \lambda_2 - 1}. \quad (68)$$

We bound each of the two terms, $\|\mathbf{M}_S^{-1} \mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'})\|$ and $\|(\mathbf{M}_S^{-1} - \mathbf{M}_{S'}^{-1}) \mathbf{P} \mathbf{y}_{S'}\|$ separately. $\|\mathbf{M}_S^{-1} \mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'})\|$ can be bounded straightforwardly:

$$\|\mathbf{M}_S^{-1} \mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'})\| \leq \|\mathbf{M}_S^{-1}\| \|\mathbf{P} (\mathbf{y}_S - \mathbf{y}_{S'})\| \leq \|\mathbf{M}_S^{-1}\| \|\mathbf{y}_S - \mathbf{y}_{S'}\| \leq \frac{\sqrt{2}M}{\frac{m}{C} \lambda_2 - 1}. \quad (69)$$

$\|(\mathbf{M}_S^{-1} - \mathbf{M}_{S'}^{-1}) \mathbf{P} \mathbf{y}_{S'}\|$ is bounded as follows:

$$\|(\mathbf{M}_S^{-1} - \mathbf{M}_{S'}^{-1}) \mathbf{P} \mathbf{y}_{S'}\| = \|\mathbf{M}_{S'}^{-1} (\mathbf{M}_{S'} - \mathbf{M}_S) \mathbf{M}_S^{-1} \mathbf{P} \mathbf{y}_{S'}\| \quad (70)$$

$$= \|\mathbf{M}_{S'}^{-1} \mathbf{P} (\mathbf{I}_{S'} - \mathbf{I}_S) \mathbf{M}_S^{-1} \mathbf{P} \mathbf{y}_{S'}\| \quad (71)$$

$$\leq \frac{\sqrt{2}mM}{\left(\frac{m}{C} \lambda_2 - 1 \right)^2}. \quad (72)$$

This leads to the following bound on $\|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty$:

$$\|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty \leq \frac{\sqrt{2}M}{\frac{m}{C}\lambda_2 - 1} + \frac{\sqrt{2m}M}{\left(\frac{m}{C}\lambda_2 - 1\right)^2} \quad (73)$$

Note that this is the hypothesis stability of the algorithm. Let $\mathbf{h}_S(x_i)$ denote the predicted target value of the i th vertex under \mathbf{h}_S (i.e. the i th coordinate of \mathbf{h}_S). The cost-stability is given by:

$$|(\mathbf{h}_S(x_i) - y_i)^2 - (\mathbf{h}_{S'}(x_i) - y_i)^2| \leq 4M\|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty. \quad (74)$$

Substituting the upper bound on $\|\mathbf{h}_S - \mathbf{h}_{S'}\|_\infty$ derived in Equation 73 into the above expression yields the statement of the theorem. \blacksquare

The generalization bound we just presented differs in several respects from that of Belkin et al. (2004a). Our bound explicitly depends on both m and u while theirs shows only a dependency on m . Also, our bound does not depend on the number of times a point is sampled in the training set (parameter t), thanks to our analysis based on sampling without replacement.

Contrasting the stability coefficient of Belkin's algorithm with the stability coefficient of LTR (Theorem 13), we note that it does not depend on C' and β_{loc} . This is because unlabeled points do not enter the objective function, and thus $C' = 0$ and $\tilde{y}(x) = 0$ for all $x \in X$. However, the stability does depend on the second smallest eigenvalue λ_2 and the bound diverges as λ_2 approaches $\frac{C}{m}$. Actually, the bound in Theorem 16 will converge so long as $\lambda_2 = \Omega(1/m)$. As observed empirically by Cortes and Mohri (2007), this algorithm does not perform as well in comparison with LTR.

6.3 Cost stability of graph Laplacian regularization algorithm

Here we give a cost-stability analysis of the graph Laplacian regularization algorithm of Belkin et al. (2004a). To do so, we show that the algorithm can in fact be viewed as a special instance of the family of LTR algorithms. Theorem 13 can then be applied in this instance with a bound on the cost stability coefficient.

To show that the graph Laplacian algorithm is a specific LTR algorithm, we need to prove that the regularization term $\mathbf{h}^\top \mathbf{L} \mathbf{h}$ corresponds to the square of a norm in some reproducing kernel Hilbert space (RKHS). We show a more general result valid for all positive semi-definite symmetric matrices \mathbf{L} . We denote by \mathbf{L}^+ the pseudo-inverse of a matrix \mathbf{L} .

Theorem 17 *Let H_1 be a vector space such that $H_1 \subseteq \text{range}(\mathbf{L})$, then the regularization term $\mathbf{h}^\top \mathbf{L} \mathbf{h}$ coincides with the square of the norm in the RKHS defined by the kernel matrix \mathbf{L}^+ .*

Proof We need to show that there exists a kernel K such that $\mathbf{h}^\top \mathbf{L} \mathbf{h} = \|\mathbf{h}\|_K^2$ for all $\mathbf{h} \in H_1$, where $\|\cdot\|_K$ is the norm in the RKHS associated to K . This condition can be rewritten as $\mathbf{h}^\top \mathbf{L} \mathbf{h} = \langle \mathbf{h}, \mathbf{h} \rangle_K$, and more generally in terms of the inner product of $\mathbf{h}, \mathbf{h}' \in H_1$ as

$$\mathbf{h}'^\top \mathbf{L} \mathbf{h} = \langle \mathbf{h}', \mathbf{h} \rangle_K. \quad (75)$$

Let \mathbf{K} denote the Gram matrix of K for the sample S . Select \mathbf{h}' to be $\mathbf{K} \mathbf{e}_i$, where \mathbf{e}_i the i th unit vector of H . Then, the equality is equivalent to

$$\forall i \in [1, m + u], \mathbf{e}_i^\top \mathbf{K} \mathbf{L} \mathbf{h} = \langle \mathbf{K} \mathbf{e}_i, \mathbf{h} \rangle_K = \langle K(x_i, \cdot), \mathbf{h} \rangle_K = h(x_i) = \mathbf{e}_i^\top \mathbf{h}, \quad (76)$$

where we used the reproducing property of the inner product. Since the equality $\mathbf{e}_i^\top \mathbf{K} \mathbf{L} \mathbf{h} = \mathbf{e}_i^\top \mathbf{h}$ holds for all $i \in [1, m + u]$, this is equivalent to the following,

$$\forall \mathbf{h} \in H_1, \mathbf{K} \mathbf{L} \mathbf{h} = \mathbf{h}. \quad (77)$$

$\mathbf{K} = \mathbf{L}^+$ verifies this equality. Indeed, by the properties of the pseudo-inverse, $\mathbf{L}^+ \mathbf{L}$ is the projection over $\text{range}(\mathbf{L})$. Since by assumption $H_1 \subseteq \text{range}(\mathbf{L})$, we can write $\mathbf{L}^+ \mathbf{L} \mathbf{h} = \mathbf{h}$. ■

In the particular case of the graph Laplacian, when the graph is connected and the space H_1 orthogonal to $\mathbf{1}$ coincides with $\text{range}(\mathbf{L})$ and the result of the theorem holds.

Corollary 18 *Any constraint optimization algorithm of the form (55) with $H_1 \subseteq \text{range}(\mathbf{L})$ is a special instance of the LTR algorithms. In particular, the graph Laplacian regularization algorithm of Belkin et al. (2004a) is a specific instance of the LTR algorithms.*

The following theorem gives a bound on the cost stability of the graph Laplacian algorithm.

Theorem 19 *Assume that the hypothesis set H is bounded; that is, for all $h \in H$, and $x \in X$, $|h(x) - y(x)| \leq M$. Then, the graph Laplacian regularization algorithm of Belkin et al. (2004a) has uniform stability β with*

$$\beta \leq \frac{4CM^2}{m} \min \left\{ \frac{1}{\lambda_2}, \rho_G \right\}, \quad (78)$$

where λ_2 is the second smallest eigenvalue of the Laplacian matrix and ρ_G the diameter of the graph G .

Proof By Corollary 18, the graph Laplacian regularization algorithm of Belkin et al. (2004a) is a special case of the LTR algorithms. Thus, Theorem 13 can be applied to determine its stability coefficient, with the term AM bounding $|h(x) - y(x)|$ in that theorem replaced by M here:

$$\beta \leq \frac{4CM^2\kappa^2}{m}. \quad (79)$$

Furthermore, using the same techniques as (Herbster et al., 2005), we can bound $\mathbf{h}^\top \mathbf{L}^+ \mathbf{h}$ and thus κ^2 in terms of the second smallest eigenvalue of the Laplacian matrix λ_2 and the diameter of the graph ρ_G as: $\kappa^2 \leq \min \left\{ \frac{1}{\lambda_2}, \rho_G \right\}$. Substituting this upper bound in Equation 79 yields the statement of the theorem. ■

The following theorem gives a novel stability generalization bound for the algorithm of Belkin et al. (2004a) in terms of the second eigenvalue of the Laplacian and the diameter of the graph.

Theorem 20 *Let H be a bounded hypothesis set. Let G be a connected graph with diameter ρ_G and L be the associated Laplacian kernel with second smallest eigenvalue λ_2 . Let S be a random subset of labeled points of size m drawn from the vertex set X . Let \mathbf{h} be the hypothesis returned by Equation 55 when trained on $X = (S, T)$. Then, for any $\epsilon > 0$,*

$$R(\mathbf{h}) \leq \widehat{R}(\mathbf{h}) + \frac{4CM^2\kappa^2}{m} + \left(\frac{8CM^2\kappa^2}{m} + \frac{M^2(m+u)}{mu} \right) \sqrt{\frac{\ln(1/\delta)\alpha(m, u)}{2}}, \quad (80)$$

where $\alpha(m, u) = \frac{mu}{m+u-1/2}$ and $\kappa^2 = \min\{1/\lambda_2, \rho_G\}$.

Proof The result follows directly from Theorem 8 and the stability coefficient β derived in Theorem 19. ■

6.4 General case

The previous sections demonstrated the stability benefits of constraints of the type $\mathbf{u}^\top \mathbf{h} = 0$, which helped us bound the stability of the graph Laplacian regularization algorithm and derive stability-based generalization bounds.

This idea and in fact much of the results presented for this particular algorithm can be generalized. To ensure stability, it suffices that the optimization constraint restricts the hypothesis set H_1 to be a subset of $\text{range}(\mathbf{L})$. By Theorem 17, the regularization term then corresponds to an RKHS norm regularization. Orthogonality with respect to a single vector may not be sufficient to ensure $H_1 \subseteq \text{range}(\mathbf{L})$. In fact, this does not hold even for the graph Laplacian regularization algorithm if the graph G is not connected since the dimension of the null space of \mathbf{L} is then more than one. But, the constraints can be augmented to guarantee this property by imposing orthogonality with respect to the null space. More generally, one might wish to impose orthogonality with respect to some space that guarantees that the smallest non-zero eigenvalue over H_1 is not too small, for example by excluding eigenvalue λ_2 if it is too small.

In particular, “stable” versions of the algorithms presented in Section 5 CM, LL – Reg, and GMF can be derived by augmenting their optimization problems with such constraints. Recall that the stability bound in Proposition 14 is inversely proportional to the smallest eigenvalue $\lambda_m(\mathbf{Q})$. The main difficulty with using the proposition for these algorithms is that $\lambda_m(\mathbf{Q}) = 0$ in each case. Let \mathbf{v}_m denote the eigenvector corresponding to $\lambda_m(\mathbf{Q})$ and let λ_2 be the second smallest eigenvalue of \mathbf{Q} . One can modify (43) and constrain the solution to be orthogonal to \mathbf{v}_m by imposing $\mathbf{h} \cdot \mathbf{v}_m = 0$. In the case of Belkin et al. (2004a), $\mathbf{v}_m = \mathbf{1}$. This modification, motivated by the algorithm of Belkin et al. (2004a), is equivalent to increasing the smallest eigenvalue to be λ_2 .

As an example, by imposing the additional constraint, we can show that the stability coefficient of CM becomes bounded by $O(C/\lambda_2)$, instead of $\Theta(1)$. Thus, if $C = O(1/m)$ and $\lambda_2 = \Omega(1)$, it is bounded by $O(1/m)$ and the generalization bound converges as $O(1/m)$.

7. Experiments

This section reports the results of experiments using our stability-based generalization bound for model selection for the LTR algorithm. A crucial parameter of this algorithm is the stability coefficient $\beta_{loc}(r)$ of the local algorithm, which computes pseudo-targets \tilde{y}_x based on a ball of radius r around each point. We derive an expression for $\beta_{loc}(r)$ and show, using extensive experiments with multiple data sets, that the value r^* minimizing the bound is a remarkably good estimate of the best r for the test error. This demonstrates the benefit of our generalization bound for model selection, avoiding the need for a held-out validation set.

The experiments were carried out on several publicly available regression data sets: *Boston Housing*, *Elevators* and *Ailerons*¹. For each of these data sets, we used $m = u$, inspired by the observation that, all other parameters being fixed, the bound of Theorem 8 is tightest when $m = u$. The value of the input variables were normalized to have mean zero and variance one. For the Boston Housing data set, the total number of examples was 506. For the Elevators and the Ailerons data set, a random subset of 2000 examples was used. For both of these data sets, other random subsets of 2000 samples led to similar results. The Boston Housing experiments were repeated for 50 random partitions, while for the Elevators and the Ailerons data set, the experiments were

1. www.liaad.up.pt/~ltorgo/Regression/DataSets.html.

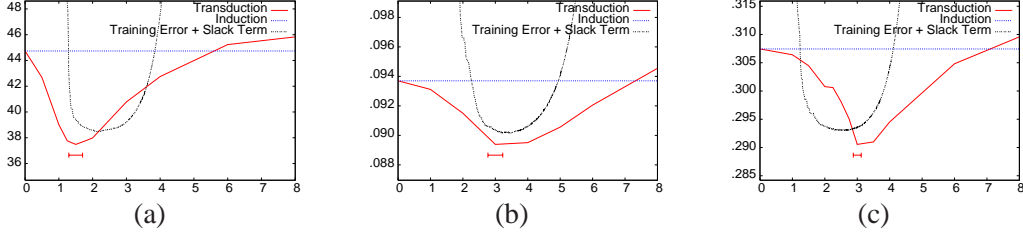


Figure 2: MSE against the radius r of LTR for three data sets: (a) Boston Housing. (b) Ailerons. (c) Elevators. The small horizontal bar indicates the location (mean \pm one standard deviation) of the minimum of the empirically determined r .

repeated for 20 random partitions each. Since the target values for the Elevators and the Ailerons data set were extremely small, they were scaled by a factor 1000 and 100 respectively in a pre-processing step.

In our experiments, we estimated the pseudo-target of a point $x' \in T$ as a weighted average of the labeled points $x \in N(x')$ in a neighborhood of x' . Thus, $\tilde{y}_{x'} = \sum_{x \in N(x')} \alpha_x y(x) / \sum_{x \in N(x')} \alpha_x$. We considered two weighting approaches, as discussed in (Cortes and Mohri, 2007), defining them in terms of the inverse of the distance between $\Phi(x)$ and $\Phi(x')$ (i.e. $\alpha_x = (1 + \|\Phi(x) - \Phi(x')\|)^{-1}$), and in terms of a similarity measure $K(x, x')$ captured by a kernel K (i.e. $\alpha_x = K(x, x')$). In our experiments, the two approaches produced similar results. We report the results of kernelized weighted average with a Gaussian kernel.

Lemma 21 *Let $r \geq 0$ be the radius of the ball around an unlabeled point $x' \in X$ that determines the neighborhood $N(x')$ of x' and let $m(r)$ be the number of labeled points in $N(x')$. Furthermore, assume that the values of the labels are bounded (i.e. for all $x \in X$, $|y(x)| \leq M$ for some $M > 0$) and that all the weights in (7) are non-negative (i.e. for all x , $\alpha_x \geq 0$). Then, the stability coefficient of the weighted average algorithm for determining the estimate of the unlabeled point x' is bounded by:*

$$\beta_{loc} \leq \frac{4\alpha_{\max}M}{\alpha_{\min}m(r)}, \quad (81)$$

where $\alpha_{\max} = \max_{x \in N(x')} \alpha_x$ and $\alpha_{\min} = \min_{x \in N(x')} \alpha_x$.

Proof We consider the change in the estimate as a point is removed from $N(x')$ and show that this is at most $\frac{2\alpha_{\max}M}{\alpha_{\min}m(r)}$. The statement of the lemma then follows straightforwardly from the observation that changing one point is equivalent to removing one point and adding another point.

Let $N(x') = \{x_1, \dots, x_{m(r)}\}$. For ease of notation, assume that $n = m(r)$. Consider the effect of removing x_n from the neighborhood $N(x')$. The estimate changes by:

$$\frac{\sum_{i=1}^n \alpha_i y_i}{\sum_{i=1}^n \alpha_i} - \frac{\sum_{i=1}^{n-1} \alpha_i y_i}{\sum_{i=1}^{n-1} \alpha_i}.$$

Thus, the stability β_{loc} can be bounded as follows:

$$\begin{aligned}
\beta_{loc} &\leq \left| \frac{\sum_{i=1}^n \alpha_i y_i}{\sum_{i=1}^n \alpha_i} - \frac{\sum_{i=1}^{n-1} \alpha_i y_i}{\sum_{i=1}^{n-1} \alpha_i} \right| \\
&\leq \frac{\alpha_n |y_n|}{\sum_{i=1}^n \alpha_i} + \sum_{i=1}^{n-1} \alpha_i |y_i| \left(\frac{1}{\sum_{i=1}^{n-1} \alpha_i} - \frac{1}{\sum_{i=1}^n \alpha_i} \right) \\
&= \frac{\alpha_n |y_n|}{\sum_{i=1}^n \alpha_i} + \frac{\sum_{i=1}^{n-1} \alpha_i |y_i|}{\sum_{i=1}^{n-1} \alpha_i} \cdot \frac{\alpha_n}{\sum_{i=1}^n \alpha_i} \\
&\leq \frac{2\alpha_n M}{\sum_{i=1}^n \alpha_i} \leq \frac{2\alpha_{\max} M}{\alpha_{\min} n} \leq \frac{2\alpha_{\max} M}{\alpha_{\min} m(r)},
\end{aligned}$$

which proves the statement of the lemma. \blacksquare

Corollary 22 *Using the notation of Lemma 21, the stability coefficient of the kernelized weighted average algorithm with a Gaussian kernel K with parameter σ is bounded by:*

$$\beta_{loc} \leq \frac{4M}{m(r)e^{-2r^2/\sigma^2}}.$$

Proof This follows directly from Lemma 21 using the observation that for a Gaussian kernel K , $K(x, x') \leq 1$, and for x, x' such that $\|x\| \leq r$ and $\|x'\| \leq r$, $\|x - x'\| \leq 2r$. Thus, $K(x, x') \geq e^{-2r^2/\sigma^2}$. \blacksquare

Corollary 23 *Using the notation of Lemma 21, the stability coefficient of the weighted average algorithm, where weights are determined by the inverse of the distance in the feature space, i.e. $\alpha_x = (1 + \|\Phi(x) - \Phi(x')\|)^{-1}$ is bounded by:²*

$$\beta_{loc} \leq \frac{(2r + 1)2M}{m(r)}.$$

Proof This follows directly from Lemma 21 using the observation that for all $x \in N(x')$,

$$0 \leq \|\Phi(x) - \Phi(x')\| \leq 2r.$$

To estimate β_{loc} , one needs an estimate of $m(r)$, the number of examples in a ball of radius r from an unlabeled point x' . In our experiments, we estimated $m(r)$ as the number of labeled examples in a ball of radius r from the origin. Since all features are normalized to mean zero and variance one, the origin is also the centroid of the set X .

We implemented a dual solution of LTR and used Gaussian kernels, for which, the parameter σ was selected using cross-validation on the training set. Experiments were repeated across 36

2. 1 is added to the weight to make the weights between 0 and 1.

different pairs of values of (C, C') . For each pair, we varied the radius r of the neighborhood used to determine estimates from zero to the radius of the ball containing all points.

Figure 2(a) shows the mean values of the test MSE of our experiments on the Boston Housing data set for typical values of C and C' . Figures 2(b)-(c) show similar results for the Ailerons and Elevators data sets. For the sake of comparison, we also report results for induction. The induction algorithm we chose was Kernel Ridge Regression (since it is analogous to LTR with the choice of $C' = 0$). The relative standard deviations on the MSE are not indicated, but were typically of the order of 10%. LTR generally achieves a significant improvement over induction.

The generalization bound we derived in Equation 21 consists of the training error and a complexity term that depends on the parameters of the LTR algorithm $(C, C', M, m, u, \kappa, \beta_{loc}, \delta)$. Only two terms depend upon the choice of the radius r : $\hat{R}(h)$ and β_{loc} . Thus, keeping all other parameters fixed, the theoretically optimal radius r^* is the one that minimizes the training error plus the slack term. The figures also include plots of the training error combined with the complexity term, appropriately scaled. The empirical minimization of the radius r coincides with or is close to r^* . The optimal r based on test MSE is indicated with error bars.

8. Conclusion

We presented a comprehensive analysis of the stability of transductive regression algorithms with novel generalization bounds for a number of algorithms. Since they are algorithm-dependent, our bounds are often tighter than those based on complexity measures such as the VC-dimension. Our experiments also show the effectiveness of our bounds for model selection and the good performance of LTR algorithms in practice. Our analysis can also guide the design of algorithms with better stability properties and thus generalization guarantees, as discussed in Section 6.4. The general concentration bound for uniform sampling without replacement proved here can be of independent interest in a variety of other machine learning and algorithmic analyses.

References

- Kazuoki Azuma. Weighted sums of certain dependent random variables. In *Tohoku Mathematical Journal*, volume 19, pages 357–367, 1967.
- Mikhail Belkin, Irina Matveeva, and Partha Niyogi. Regularization and semi-supervised learning on large graphs. In *Conference on Learning Theory (COLT 2004)*, pages 624–638. Springer, 2004a.
- Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization; a geometric framework for learning from examples. Technical Report TR-2004-06, University of Chicago, 2004b.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research (JMLR)*, 2:499–526, 2002. ISSN 1533-7928.
- Olivier Bousquet and André Elisseeff. Algorithmic stability and generalization performance. In *Advances in Neural Information Processing Systems (NIPS 2000)*, pages 196–202. MIT Press, 2000.
- Olivier Chapelle, Vladimir Vapnik, and Jason Weston. Transductive inference for estimating values of functions. In *Neural Information Processing Systems (NIPS 1999)*, pages 421–427. MIT Press, 1999.

- Corinna Cortes and Mehryar Mohri. On transductive regression. In *Advances in Neural Information Processing Systems (NIPS 2006)*, pages 305–312. MIT Press, 2007.
- Philip Derbeko, Ran El-Yaniv, and Ron Meir. Explicit learning curves for transduction and application to clustering and compression algorithms. *J. Artif. Intell. Res. (JAIR)*, 22:117–142, 2004.
- Ran El-Yaniv and Dmitry Pechyony. Stable transductive learning. In *Conference on Learning Theory (COLT 2006)*, pages 35–49. Springer, 2006.
- Ran El-Yaniv and Dmitry Pechyony. Transductive rademacher complexity and its applications. *Conference on Learning Theory (COLT 2007)*, 2007.
- Mark Herbster, Massimiliano Pontil, and Lisa Wainer. Online learning over graphs. In *International Conference on Machine learning*, pages 305–312, New York, NY, USA, 2005. ACM Press.
- Colin McDiarmid. On the method of bounded differences. In *Surveys in Combinatorics*, pages 148–188. Cambridge University Press, Cambridge, 1989.
- Dale Schuurmans and Finnegan Southey. Metric-based methods for adaptive model selection and regularization. *Machine Learning*, 48:51–84, 2002.
- Vladimir N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer, Berlin, 1982.
- Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley-Interscience, New York, 1998.
- Mingrui Wu and Bernhard Schölkopf. Transductive classification via local learning regularization. In *Artificial Intelligence and Statistics (AISTATS 2007)*, 2007.
- Dengyong (Denny) Zhou, Olivier Bousquet, Thomas N. Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. In *Advances in Neural Information Processing Systems (NIPS 2004)*. MIT Press, 2004.
- Xiaojin Zhu, Zoubin Ghahramani, and John D. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *International Conference on Machine Learning (ICML 2003)*, pages 912–919. ACM Press, 2003.