

Empirical Bernstein Inequalities for U-Statistics

Authored by:

Liva Ralaivola
Thomas Peel
Sandrine Anthoine

Abstract

We present original empirical Bernstein inequalities for U-statistics with bounded symmetric kernels q . They are expressed with respect to empirical estimates of either the variance of q or the conditional variance that appears in the Bernstein-type inequality for U-statistics derived by Arcones [2]. Our result subsumes other existing empirical Bernstein inequalities, as it reduces to them when U-statistics of order 1 are considered. In addition, it is based on a rather direct argument using two applications of the same (non-empirical) Bernstein inequality for U-statistics. We discuss potential applications of our new inequalities, especially in the realm of learning ranking/scoring functions. In the process, we exhibit an efficient procedure to compute the variance estimates for the special case of bipartite ranking that rests on a sorting argument. We also argue that our results may provide test set bounds and particularly interesting empirical racing algorithms for the problem of online learning of scoring functions.

1 Paper Body

The motivation of the present work lies in the growing interest of the machine learning community for learning tasks that are richer than now well-studied classification and regression. Among those, we especially have in mind the task of ranking, where one is interested in learning a ranking function capable of predicting an accurate ordering of objects according to some attached relevance information. Tackling such problems generally implies the use of loss functions other than the 0-1 misclassification loss such as, for example, a misranking loss [6] or a surrogate thereof. For (x, y) and (x_0, y_0) two pairs from some space $Z := X \times Y$ (e.g., $X = \mathbb{R}^d$ and $Y = \mathbb{R}$) the misranking loss ‘rank and a surrogate convex loss ‘sur may be defined for a scoring function $f : Y \times X \rightarrow \mathbb{R}$ as: ‘rank $(f, (x, y), (x_0, y_0)) := 1\{(y - y_0)(f(x) - f(x_0)) \leq 0\}$, sur

0
0

0

(1) 0

2

$$\ell(f, (x, y), (x', y')) := (1 - (y - y'))(f(x) - f(x')) .$$

(2)

Given such losses or, more generally, a loss $\ell : Y \times X \times Z \times Z \rightarrow \mathbb{R}$, and a training sample $Z_n = \{(X_i, Y_i)\}_{i=1}^n$ of independent copies of some random variable $Z := (X, Y)$ distributed according to D , the learning task is to derive a function $f : X \rightarrow Y$ such that the expected risk $R^\ell(f)$ of f $R^\ell(f) := \mathbb{E}_{Z \sim D} \ell(f, Z, Z)$ is as small as possible. In practice, this naturally brings up the empirical estimate $R_n^\ell(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, (X_i, Y_i), (X_i, Y_i))$

$R_n^\ell(f)$ is as small as possible. In practice, this naturally brings up the empirical estimate $R_n^\ell(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, (X_i, Y_i), (X_i, Y_i))$

(3)

$i \leq j$

which is a U-statistic [6, 10]. $R_n^\ell(f)$ is related to $R^\ell(f)$ and, more An important question is to precisely characterize how $R_n^\ell(f)$ is related to $R^\ell(f)$ and, more specifically, one may want to derive an upper bound on $R_n^\ell(f)$ that is expressed in terms of $R^\ell(f)$ and other quantities such as a measure of the capacity of the class of functions f belongs to and the size n of Z_n . In other words, we may talk about generalization bounds [4]. Pivotal tools to perform such analysis are tail/concentration inequalities, which say how probable it is for a function of several independent variables to deviate from its expectation; of course, the sharper the concentration inequalities the more accurate the characterization of the relation between the empirical estimate and its expectation. It is therefore of the utmost importance to have at hand tail inequalities that are sharp; it is just as important that these inequalities rely as much as possible on empirical quantities. Here, we propose new empirical Bernstein inequalities for U-statistics. As indicated by the name (i) our results are Bernstein-type inequalities and therefore make use of information on the variance of the variables under consideration, (ii) instead of resting on some assumed knowledge about this variance, they only rely on empirical related quantities and (iii) they apply to U-statistics. Our new inequalities generalize those of [3] and [13], which also feature points (i) and (ii) (but not (iii)), while based on simple arguments. To the best of our knowledge, these are the first results that fulfill (i), (ii) and (iii); they may give rise to a few applications, of which we describe two in the sequel. The paper is organized as follows. Section 2 introduces the notations and briefly recalls the basics of U-statistics as well as tail inequalities our results are based upon. Our empirical Bernstein inequalities are presented in Section 3; we also provide an efficient way of computing the empirical variance when the U-statistics considered are based on the misranking loss ‘rank of (1). Section 4 discusses two applications of our new results: test set bounds for bipartite ranking and online ranking.

2 2.1

Background Notation

The following notation will hold from here on. Z is a random variable of distribution D taking values in $Z := X \times Y$; Z_0, Z_1, \dots, Z_n are independent copies of Z and $Z_n := \{Z_i = (X_i, Y_i)\}_{i=1}^n$ and $Z_{p:q} := \{Z_i\}_{i=p}^q$.

n denotes the set $A_n := \{(i_1, \dots, i_m) : 1 \leq i_1 \leq \dots \leq i_m \leq n\}$, with $0 \leq m \leq n$.

Finally, a function $q : Z^m \rightarrow \mathbb{R}$ is said to be symmetric if the value of $q(z)$ is independent of the order of the z_i 's in z .

U-statistics and Tail Inequalities

$U_q(Z)$ defined as Definition 1 (U-statistic, Hoeffding [10]). The random variable $U_n = U_q(Z_1, \dots, Z_n)$, $n \geq m$.

is a U-statistic of order m with kernel q , when $q : Z^m \rightarrow \mathbb{R}$ is a measurable function on Z^m . $U_q(Z)$; in addition, $E U_q(Z)$ is a lowest variance estimate of $E q(Z_1, \dots, Z_m)$ based on Z_n [10]. Also, reusing some notation from the ' (f, Z) of Eq. (3) is a U-statistic of order 2 with kernel $q_f(Z, Z_0) := (f(Z), Z_0)$. introduction, Remark 2. Two peculiarities of U-statistics that entail a special care are the following: (i) they are sums of identically distributed but dependent variables: special tools need be resorted to in order to deal with these dependencies to characterize the deviation of U_n algorithmic point of view, their direct computations may be expensive, as it scales as $O(nm)$; in Section 3, we show for the special case of bipartite ranking how this complexity can be reduced.

$m=2$ 0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0 101

$m=10$ Bernstein Hoeffding

102

103

0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0

$m=2$ 0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0 101

Bernstein Hoeffding

102

103

104

$m=10$ Bernstein Hoeffding

102

103

0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0 101

Bernstein Hoeffding

102

103

Figure 1: First two plots: values of the right-hand side of (5) and (6), for Duni and kernel q_m for $m = 2$ and $m = 10$ (see Example 1) as functions of n . Last two plots: same for DBer (0.15). We now recall three tail inequalities (Eq. (5), (6), (7)) that hold for U-statistics with symmetric and bounded kernels q . Normally, these inequalities make explicit use of the length $q_{\max} - q_{\min}$ of the range $[q_{\min}, q_{\max}]$ of q . To simplify the reading, we will consider without loss of generality that q has range $[0, 1]$ (an easy way of retrieving the results for bounded q is to consider q/k_q). One key quantity that appears in the original versions of tail inequalities (5) and (6) below is bn/mc , the integer part of the ratio n/m ; this quantity might be thought of as the effective number of

data. To simplify the notation, we will assume that n is a multiple of m and, therefore, $bn/mc = (n/m)$. $\varphi_q, [11]$). Hoeffding proved the following: Theorem 1 (First order tail inequality for U)

$\varphi_q(Z_0) \leq U \varphi_q(Z) \leq 2 \exp \left(-\frac{n}{m} \right)$, $\forall \delta \in (0, 1]$, $P(Z_0 \in U) \geq 1 - \delta$ over the random draw of Z : s

$$\begin{aligned} & 2 \ln \left(\frac{1}{\delta} \right) \leq \ln \left(\frac{EZ_n U_q(Z_n)}{U_n(n/m)} \right) \\ & (5) \end{aligned}$$

To go from the tail inequality (4) to the bound version (5), it suffices to make use of the elementary inequality reversal lemma (Lemma 1) provided in section 3, used also for the bounds given below. $\varphi_q, [2, 11]$). Hoeffding [11] and, later, Arcones [2] refined Theorem 2 (Bernstein Inequalities for U the previous result in the form of Bernstein-type inequalities of the form

$\varphi_q(Z_0) \leq U \varphi_q(Z) \leq a \exp \left(-\frac{n}{m} \right)$, $\forall \delta \in (0, 1]$, $P(Z_0 \in U) \geq 1 - \delta$ where, φ_q is the variance of $q(Z_1, \dots, Z_m)$ and $bm = 2/3$. Hence, $\varphi_q \in (0, 1]$, with probability at least $1 - \delta$: s

$2 \ln \left(\frac{1}{\delta} \right) \leq \ln \left(\frac{EZ_n U_q(Z_n)}{U_n(n/m)} \right) \leq 3 \ln \left(\frac{1}{\delta} \right)$ For Arcones, $a = 4$, $\varphi_q = m \varphi_q^2$ where φ_q^2 is the variance of $q(Z_1, Z_2, \dots, Z_m)$ (this is a function of Z_1) and $bm = 2m + 3 \ln m + (2/3)m^2$. $\varphi_q \in (0, 1]$, with probability at least $1 - \delta$: s

$2 \ln \left(\frac{1}{\delta} \right) \leq \ln \left(\frac{EZ_n U_q(Z_n)}{U_n(n/m)} \right) \leq \ln \left(\frac{1}{\delta} \right)$ With a slight abuse, we will now refer to Eq. (5), (6) and (7) as tail inequalities. In essence, these $\varphi_q(Z)$ are confidence intervals at level $1 - \delta$ for $EZ_m q(Z_m) = EZ_n U_n$ Remark 3. Eq. (7) is based on the so-called Hoeffding decomposition of U -statistics [11]. It provides a more accurate Bernstein-type inequality than that of Eq. (6), as $m \varphi_q^2$ is known to be smaller than φ_q (see [16]). However, for moderate values of n/m (e.g. $n/m \leq 105$) and reasonable values of δ (e.g. $\delta = 0.05$), the influence of the log terms might be such that the advantage of (7) over (6) goes unnoticed. Thus, we detail our results focusing on an empirical version of (6). Example 1. To illustrate how the of the variance information provides smaller confidence intervals, consider the kernel $q_m := \frac{1}{m} \sum_{i=1}^m z_i$ and two distributions Duni and DBer (p). Duni is the uniform distribution on $[0, 1]$, for which $\varphi^2 = 3/12$. DBer (p) is the Bernoulli distribution with parameter $p \in [0, 1]$, for which $\varphi^2 = p(1 - p)$. Figure 1 shows the behaviors of (6) and (5) for various values of m as functions of n . Observe that the variance information renders the bound smaller. 3

3 Main Results

This section presents the main results of the paper. We first introduce the inequality reversal lemma, which allows to transform tail inequalities into upper

bounds (or confidence intervals), as in (5)-(7). Lemma 1 (Inequality Reversal lemma). Let X be a random variable and $a, b \geq 0, c, d \geq 0$ such that

$b \geq 0, P(X \leq -X) \leq a \exp(-b)$, (8) $c + d \geq 0$ then, with probability at least $1 - e^{-r}$

$$\frac{c}{a} \leq \frac{d}{a} \ln \left(\frac{b}{b} \right) + \ln \left(\frac{b}{b} \right) \quad (9)$$

Proof. Solving for r such that the right hand side of (8) is equal to r gives:

$r \leq \frac{1}{a} \ln \left(\frac{a}{a} \right) = \frac{1}{a} \ln \left(\frac{d}{d} \ln 2 + 4bc \ln \left(\frac{2b}{b} \right) \right)$ Using $a + b \geq a + b$ gives an upper bound on r and provides the result. 3.1

Empirical Bernstein Inequalities

Let us now define the empirical variances we will use in our main result. \hat{q}_2 be the U-statistic of order $2m$ defined as: Definition 2. Let \hat{q}_2

$$\hat{q}_2(Z) := \frac{1}{n} \sum_{i=1}^n q(Z_{i1}, \dots, Z_{im}) q(Z_{im+1}, \dots, Z_{i2m}) \quad (10)$$

and \hat{q}_2 be the U-statistic of order $2m$ defined as:

$$\hat{q}_2(Z) := \frac{1}{n} \sum_{i=1}^n q(Z_{i1}, Z_{i2}, \dots, Z_{im}) q(Z_{i1}, Z_{im+1}, \dots, Z_{i2m}), \quad (11)$$

It is straightforward to see that (cf. the definitions of \hat{q}_2 in (6) and \hat{q}_2 in (7)) $\hat{q}_2(Z) = \hat{q}_2, E \hat{q}_2 = n$

and

$$E \hat{q}_2(Z) = \hat{q}_2 + E \hat{q}_2(Z_1, \dots, Z_m).$$

We have the following main result. Theorem 3 (Empirical Bernstein Inequalities/Bounds). With probability at least $1 - e^{-s}$ over Z, n, s

$$\hat{q}_2 \leq \frac{2}{4} \ln \left(\frac{4}{5} \right) + \ln \left(\frac{4}{5} \right)$$

$0 \leq \hat{q}_2(Z) \leq \ln \left(\frac{4}{5} \right) + \ln \left(\frac{4}{5} \right)$ (12) $E \hat{q}_2(Z) \leq U \hat{q}_2(Z) \leq (n/m) \hat{q}_2(Z) \leq (n/m) \hat{q}_2$ And, also, with probability at least $1 - e^{-s}$, (bm is the same as in (7)) $s \geq$

$$\frac{2m}{8} \hat{q}_2 \leq \frac{5}{8} m + \frac{bm}{8} \leq 0$$

$$0 \leq E \hat{q}_2(Z) \leq \ln \left(\frac{4}{5} \right) + \ln \left(\frac{4}{5} \right) \leq U \hat{q}_2(Z) \leq (n/m) \hat{q}_2(Z) \leq (n/m) \hat{q}_2 \quad (13)$$

Proof. We provide the proof of (12) for the upper bound of the confidence interval; the same reasoning carries over to prove the lower bound. The proof of (13) is very similar. \hat{q}_2 : First, let us call Q the kernel of \hat{q}_2

$$Q(Z_1, \dots, Z_{2m}) := (q(Z_1, \dots, Z_m) q(Z_{m+1}, \dots, Z_{2m})) \cdot 4$$

Q is of order $2m$, has range $[0, 1]$ but it is not necessarily symmetric. An equivalent symmetric \hat{q}_2 is Q_{sym} : kernel for \hat{q}_2 $Q_{\text{sym}}(Z_1, \dots, Z_{2m}) := q(Z_1, \dots, Z_m) q(Z_{m+1}, \dots, Z_{2m}) (2m)! P_m$

where P_m is the set of all the permutations over $\{1, \dots, m\}$. This kernel is symmetric (and has range $[0, 1]$) and Theorem 2 can be applied to bound \hat{q}_2 as follows: with prob. at least $1 - e^{-s}$ $\hat{q}_2 \leq \frac{2}{4} \ln \left(\frac{4}{5} \right) + \ln \left(\frac{4}{5} \right) + 2V(Q_{\text{sym}})$

$\ln 2 + \frac{1}{2} = \mathbb{E} Z_{02m} Q_{\text{sym}}(Z_{2m}) = \mathbb{E} Z_{0n} \ln \frac{1}{q} \frac{q}{n} \ln \frac{(n/2m)}{3(n/2m)}$
where $V(Q_{\text{sym}})$ is the variance of Q_{sym} . As Q_{sym} has range $[0, 1]$,

$$\begin{aligned} V(Q_{\text{sym}}) &= \mathbb{E} Q_{2\text{sym}}^2 - \mathbb{E}^2 Q_{\text{sym}} = \mathbb{E} Q_{2\text{sym}}^2 - \mathbb{E}^2 Q_{\text{sym}} = \frac{1}{2}, \text{ and therefore} \\ &\leq \frac{1}{2} \\ &\leq 2q(Z) \leq n \\ &+ \\ &\leq 2 \ln 2 + \ln \frac{(n/m)}{3(n/m)} \end{aligned}$$

(To establish (13) we additionally use $\frac{1}{2} q^2(Z_n) \leq \frac{1}{2} q^2$). Following the approach of [13], we introduce $\frac{1}{2} \ln \frac{(m/n)}{3(n/m)}$ and we get $\frac{1}{2} \ln \frac{(m/n)}{3(n/m)} + \ln \frac{1}{2} \ln \frac{(m/n)}{3(n/m)} \leq \frac{1}{2} \ln \frac{(m/n)}{3(n/m)} + \frac{1}{2} \ln \frac{(m/n)}{3(n/m)}$ and taking the square root of both side, using $1 + \frac{7}{3} \leq 3$ and $a + b \leq a + b$ again gives $\frac{1}{2} \ln \frac{(m/n)}{3(n/m)} + \frac{1}{2} \ln \frac{(m/n)}{3(n/m)} \leq \frac{1}{2} \ln \frac{(m/n)}{3(n/m)} + \frac{1}{2} \ln \frac{(m/n)}{3(n/m)}$ —, and plug in the latter equation, adjusting We now apply Theorem 2 to bound $\mathbb{E} Z_{0n} U_{nn} \leq \frac{1}{2}$ so the obtained inequality still holds with probability $1 - \frac{1}{2}$. Bounding appropriate constants gives the desired result. Remark 4. In addition to providing an empirical Bernstein bound for U-statistics based on arbitrary bounded kernels, our result differs from that of Maurer and Pontil [13] by the way we derive it. Here, we apply the same tail inequality twice, taking advantage of the fact that estimates for the variances we are interested in are also U-statistics. Maurer and Pontil use a tail inequality on self bounded random variables and do not explicitly take advantage of the estimates they use being U-statistics. 3.2

Efficient Computation of the Variance Estimate for Bipartite Ranking

We have just showed how empirical Bernstein inequalities can be derived for U-statistics. The estimates that enter into play in the presented results are U-statistics with kernels of order $2m$ (or $2m - 1$), meaning that a direct approach to practically compute them would scale as $O(n^{2m})$ (or $O(n^{2m-1})$). This scaling might be prohibitive as soon as n gets large. $\frac{1}{2} q^2$ (a similar reasoning carries over for Here, we propose an efficient way of evaluating the estimate $\frac{1}{2} \ln \frac{(m/n)}{3(n/m)}$ in the special case where $Y = \{-1, +1\}$ and the kernel q_f induces the misranking loss (1): $q_f((x, y), (x_0, y_0)) := \frac{1}{2} \{ (y - y_0)(f(x) - f(x_0)) \}_+$, $\frac{1}{2} f^T R X$, which is a symmetric kernel of order $m = 2$ with range $[0, 1]$. In other words, we address the bipartite ranking problem. We have the following result. $\frac{1}{2} q^2$). $\frac{1}{2} n$, the computation of Proposition 1 (Efficient computation of $\frac{1}{2} f^T X$ $\frac{1}{2} q(z_n) = \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \{ (y_{i1} - y_{i2})(f(x_{i1}) - f(x_{i2})) \}_+ \{ (y_{i3} - y_{i4})(f(x_{i3}) - f(x_{i4})) \}_+ - \frac{A}{4n} - \frac{4}{n} \frac{1}{A} n$

can be performed in $O(n \ln n)$. 5

$\frac{1}{2} q^2(z_n)$. To simplify the reading, we Proof. We simply provide an algorithmic way to compute $\frac{1}{2} f^T$ replace i_1, i_2, i_3, i_4 by i, j, k, l , respectively. We also drop the normalization factor $-\frac{A}{4n} - \frac{1}{n}$ (hence the use of $\frac{1}{2}$ instead of $\frac{1}{4}$ in the first line below). We have $\frac{1}{2} q^2(z_n) \leq \frac{1}{2} f^T$

$$\begin{aligned} &\sum_{i,j,k,l} \\ &\sum_{i,j,k,l} \\ &(\frac{1}{2} q_f(z_i, z_j) - \frac{1}{2} q_f(z_k, z_l))^2 = \\ &\sum_{i,j,k,l} \frac{1}{4} (q_f(z_i, z_j) - q_f(z_k, z_l))^2 \end{aligned}$$

The first term of the last line is proportional to the well-known Wilcoxon-Mann-Whitney statistic [9]. There exist efficient ways ($O(n \ln n)$) to compute based on sorting the values of the $f(x_i)$'s. We show how to deal with the second term, using sorting arguments as well. Note that $\sum_{i=1}^n X_i$

P We have subtracted from the square of $\sum_{i,j} qf(z_i, z_j)$ all the products $qf(z_i, z_j)qf(z_k, z_l)$ such that exactly one of the variables appears both in $qf(z_i, z_j)$ and $qf(z_k, z_l)$, which happens when $i = k, i = l, j = k, j = l$; using the symmetry of qf then provides the second term (together with the factor 4). We also have subtracted all the products $qf(z_i, z_j)qf(z_k, z_l)$ where $i = k$ and $j = l$ or $i = l$ and $j = k$, in which case the product reduces to $qf^2(z_i, z_j)$ (hence the factor 2) ? this gives the last term. P Thus (using $qf^2 = qf$), defining $R(z_n) = \sum_{i,j} qf(z_i, z_j)$ and doing some simple calculations: ? ? $X_{1,2,2}^2 = q(z_n) = 2R(z_n) + 2(n^2 - 5n + 8)R(z_n) + 8 \sum_{i,j} qf(z_i, z_j)qf(z_i, z_k) = f - 4n -$
(14)

The only term that now requires special care is the last one (which is proportional to $\frac{1}{2} \sum_{i,j} \frac{1}{\|z_i - z_j\|} \frac{1}{\|z_i - z_j\|} \frac{1}{\|z_i - z_j\|}$). Recalling that $q_f(z_i, z_j) = \frac{1}{2} \{ (y_i - y_j)(f(z_i) - f(z_j)) \}$, we observe that

Let us define $E_+(i)$ and $E_-(i)$ as $E_+(i) := \{j : y_j = +1, f(x_j) \leq f(x_i)\}$, and $E_-(i) := \{j : y_j = -1, f(x_j) \leq f(x_i)\}$. $n_+ = |E_+|$ and $n_- = |E_-|$ and their sizes $n_+ := |E_+|$, and $n_- := |E_-|$.

For i such that $y_i = 1$, n_+^i is the number of negative instances that have been scored higher than x_i by f . From (15), we see that the contribution of i to the last term of (14) corresponds to the number $n_+^i + n_-^i (n_+^i + 1)$ of ordered pairs of indices in $E_-(i)$ (similarly for n_-^i , with $y_i = -1$). Henceforth: X

$$\sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i \neq y_j} \mathbb{1}_{f(x_i) \leq f(x_j)} =$$

$$\sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i \neq y_j} \mathbb{1}_{f(x_i) \leq f(x_j)} =$$

$$X$$

$$+ \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i = +1} \mathbb{1}_{f(x_i) \leq f(x_j)} +$$

$$\sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i = -1} \mathbb{1}_{f(x_i) \leq f(x_j)} =$$

$$X$$

$$+ \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i = -1} \mathbb{1}_{f(x_i) \leq f(x_j)}.$$

$$\sum_{i=1}^n \sum_{j=1}^n \mathbb{1}_{y_i = -1} \mathbb{1}_{f(x_i) \leq f(x_j)} =$$

A simple way to compute the first sum (on i such that $y_i = +1$) is to sort and visit the data by descending order of scores and then to incrementally compute the n_+^i 's and the corresponding sum: when a negative instance is encountered, n_+^i is incremented by 1 and when a positive instance is $i + 1$ visited, $n_+^i (n_+^i + 1)$ is added to the current sum. An identical reasoning works for the second sum. $n_+^i + 1$ is therefore that of sorting the scores, which has cost $O(n \ln n)$. The cost of computing $\sum f$

4

Applications and Discussion

Here, we mention potential applications of the new empirical inequalities we have just presented. 6

Bernstein vs Hoeffding

Banana dataset 0.5 2

Hoeffding Bernstein

0.4

1

0.3

0

0.2

-1

0.1 -2

-1

0

1

0 102

2

103

104

Figure 2: Left: UCI banana dataset, data labelled +1 (-1) in red (green). Right: half the confidence interval of the Hoeffding bound and that of the empirical Bernstein bound as functions of n_{test} . 4.1

Test Set Bounds

A direct use of the empirical Bernstein inequalities is to draw test set bounds. In this scenario, a sample Z_n is split into a training set $Z_{\text{train}} := Z_{1:n_{\text{train}}}$ of n_{train} data and a hold-out set $Z_{\text{test}} := Z_{n_{\text{train}}+1:n}$ of size n_{test} . Z_{train} is used to train a model f that minimizes an empirical risk based on a U-statistic inducing loss (such as in (1) or (2)) and Z_{test} is used to compute a confidence interval on the expected risk of f . For instance, if we consider the bipartite ranking problem, the loss is ‘rank’, the corresponding kernel is $qf(Z, Z_0) = \text{‘rank’}(f, Z, Z_0)$, and, with probability at least $1 - \delta$, $\text{‘rank’}(f, Z) + R \sqrt{\frac{\text{var}(qf(Z, Z_0))}{n_{\text{test}}}} \leq \mathbb{E}[\text{‘rank’}(f, Z, Z_0)] \leq \text{‘rank’}(f, Z) + R \sqrt{\frac{\text{var}(qf(Z, Z_0))}{n_{\text{test}}}}$, (16) where $\text{var}(qf(Z, Z_0))$ is naturally the empirical variance of qf computed on Z . where $\text{‘rank’}(f, Z) = \sum_{i=1}^{n_{\text{train}}} \sum_{j=1}^{n_{\text{test}}} f(X_i, X_j) - \sum_{i=1}^{n_{\text{train}}} f(X_i, X_i)$. Figure 2 displays the behavior of such test set bounds as n_{test} grows for the UCI banana dataset. To produce this plot, we have learned a linear scoring function $f(x) = \langle w, x \rangle$ by minimizing $\sum_{i=1}^n \|X_i - Y_i\|^2 + \lambda \sum_{i=1}^n \|X_i\|^2$.

for $\lambda = 1.0$. Of course, a purely linear scoring function would not make it possible to achieve good ranking accuracy so we in fact work in the reproducing kernel hilbert space associated with the Gaussian kernel $k(x, x_0) = \exp(-\|x - x_0\|^2 / 2)$. We train our scoring function on $n_{\text{train}} = 1000$ data points and evaluate the test set bound on $n_{\text{test}} = 100, 500, 1000, 5000, 10000$ data points. Figure 2 (right) reports the size of half the confidence interval of the Hoeffding bound (5) and that of the empirical Bernstein bound given in (16). Just as in the situation described in Example 1, the use of variance information gives rise to smaller confidence intervals, even for moderate sizes of test sets. 4.2

Online Ranking and Empirical Racing Algorithms

Another application that we would like to describe is online bipartite ranking. Due to space limitation, we only provide the main ideas on how we think our empirical tail inequalities and the efficient computation of the variance estimates we propose might be particularly useful in this scenario. First, let us precise what we mean by online bipartite ranking. Obviously, this means that $Y = \{-1, +1\}$ and that the loss of interest is ‘rank’. In addition, it means that given a training set $Z = \{Z_i := (X_i, Y_i)\}_{i=1}^n$ the learning procedure will process the data of Z incrementally to give rise to hypotheses f_1, f_2, \dots, f_T . As ‘rank’ entails a kernel of order $m = 2$, we assume that $n = 2T$ and that we process the data from Z pairs by pairs, i.e. (Z_1, Z_2) are used to learn f_1 , (Z_3, Z_4) and f_1 are used to learn f_2 and, more generally, (Z_{2t-1}, Z_{2t}) and f_{t-1} are used to produce f_t (there exist more clever ways to handle the data but this goes out of the scope of the present paper). We do not specify any learning algorithm but we may imagine trying to minimize a penalized empirical risk based on the surrogate loss ‘sur’: if linear functions $f(x) = \langle w, x \rangle$ are considered and a penalization γ

like $\|w\|^2$ is used then the optimization problem to solve is of the same form as in the batch case: $\sum_{i=1}^n \|X_i - Y_i\|^2 + \lambda \sum_{i=1}^n \|X_i\|^2$, $i=1, \dots, n$

but is solved incrementally here. Rank-1 update formulas for inverses of matrices easily provide means to incrementally solve this problem as new data arrive (this is the main reason why we have mentioned this surrogate function).

As evoked by [5], a nice feature of online learning is that the expected risk of hypothesis f_t can be estimated on the $n + 2t$ examples of Z it was not trained on. Namely, when $2t$ data have been processed, there exist $2t$ hypotheses f_1, \dots, f_{2t} and, for $t \leq n$, with probability at least $1 - \delta$:

$$\mathbb{E} \sum_{t=1}^{2t} (Z(f_t) - \min_{f \in \mathcal{H}} Z(f)) \leq \ln(4/\delta)$$

If one wants to have these confidence intervals to simultaneously hold for all t and all δ with probability $1 - \delta$, basic computations to calculate the number of pairs (t, δ) , with $1 \leq t \leq n$ show that it suffices to adjust δ to $4/(n+1)^2$. Hence, with probability at least $1 - \delta$:

$$\mathbb{E} \sum_{t=1}^{2t} (Z(f_t) - \min_{f \in \mathcal{H}} Z(f)) \leq \ln((n+1)/\delta)$$

If f_t is $\text{rank}(f_t, Z + \ln \cdot \text{R}\text{rank}(f_t) - \text{R}\text{rank}(f_t))$ then, if one wants to have these confidence intervals to simultaneously hold for all t and all δ with probability $1 - \delta$, basic computations to calculate the number of pairs (t, δ) , with $1 \leq t \leq n$ show that it suffices to adjust δ to $4/(n+1)^2$. Hence, with probability at least $1 - \delta$:

We would like to draw the attention of the reader on two features: one has to do with statistical considerations and the other with algorithmic ones. First, if the confidence intervals simultaneously hold for all t and all δ as in (17), it is possible, as the online learning process goes through, to discard the hypotheses f_t which have their lower bound (according to (17)) on $\text{R}\text{rank}(f_t)$ that is higher than the upper bound (according to (17) as well) on $\text{R}\text{rank}(f_t)$ for some other hypothesis f_{t_0} . This corresponds to a racing algorithm as described in [12]. Theoretically analyzing the relevance of such a race can be easily done with the results of [14], which deal with empirical Bernstein racing, but for non-U-statistics. This full analysis will be provided in a long version of the present paper. Second, it is algorithmically possible to preserve some efficiency in computing the various variance estimates through the online learning process: these computations rely on sorting arguments, and it is possible to take advantage of structures like binary search trees such as AVL trees, that are precisely designed to efficiently maintain and update sorted lists of numbers. The remaining question is whether it is possible to have shared such structures to summarize the sorted lists of scores for various hypotheses (recall that the scores are computed on the same data). This will be the subject of further research.

5

Conclusion

We have proposed new empirical Bernstein inequalities designed for U-statistics. They generalize the empirical inequalities of [13] and [3] while they merely result from two applications of the same non-empirical tail inequality for U-statistics. We also show how, in the bipartite ranking situation, the empirical variance can be efficiently computed. We mention potential applications, with illustrative results for the case of test set bounds in the realm of bipartite ranking. In addition to the possible extensions discussed in the previous section, we wonder whether it is possible to draw similar empirical inequalities for other types of rich statistics such as, e.g., linear rank statistics [8]. Obviously, we plan to work on establishing generalization bounds derived from the new concentration inequalities presented. This would require to carefully define a sound notion of capacity for U-statistic-based classes of functions (inspired, for example, from localized Rademacher complexities). Such new bounds would be compared with those proposed in [1, 6, 7, 15] for the bipartite ranking and/or pairwise classi-

fication problems. Finally, we also plan to carry out intensive simulations ?in particular for the task of online ranking? to get even more insights on the relevance of our contribution. Acknowledgments This work is partially supported by the IST Program of the EC, under the FP7 Pascal 2 Network of Excellence, ICT-216886-NOE. LR is partially supported by the ANR project ASAP. 8

2 References

- [1] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth. Generalization Bounds for the Area under the ROC Curve. *Journal of Machine Learning Research*, 6:393?425, 2005.
- [2] M. A. Arcones. A bernstein-type inequality for u-statistics and u-processes. *Statistics & probability letters*, 22(3):239?247, 1995.
- [3] J.-Y. Audibert, R. Munos, and C. Szepesv?ari. Tuning bandit algorithms in stochastic environments. In *ALT '07: Proceedings of the 18th international conference on Algorithmic Learning Theory*, pages 150?165, Berlin, Heidelberg, 2007. Springer-Verlag.
- [4] S. Boucheron, O. Bousquet, and G. Lugosi. Theory of classification : A survey of some recent advances. *ESAIM. P&S*, 9:323?375, 2005.
- [5] N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of online learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050?2057, 2004.
- [6] S. Cl?emenc?on, G. Lugosi, and N. Vayatis. Ranking and empirical minimization of u -statistics. *The Annals of Statistics*, 36(2):844?874, April 2008.
- [7] Y. Freund, R. Iyer, R.E. Schapire, and Y. Singer. An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research*, 4:933?969, 2003.
- [8] J. H?ajek and Z. Sid?ak. *Theory of Rank Tests*. Academic Press, 1967.
- [9] J. A. Hanley and B. J. Mcneil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29?36, April 1982.
- [10] W. Hoeffding. A Class of Statistics with Asymptotically Normal Distribution. *Annals of Mathematical Statistics*, 19(3):293?325, 1948.
- [11] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13?30, 1963.
- [12] O. Maron and A. Moore. Hoeffding races: Accelerating model selection search for classification and function approximation. In *Adv. in Neural Information Processing Systems NIPS 93*, pages 59?66, 1993.
- [13] A. Maurer and M. Pontil. Empirical bernstein bounds and sample-variance penalization. In *COLT 09: Proc. of The 22nd Annual Conference on Learning Theory*, 2009.
- [14] V. Mnih, C. Szepesv?ari, and J.-Y. Audibert. Empirical bernstein stopping. In *ICML '08: Proceedings of the 25th international conference on Machine learning*, pages 672?679, New York, NY, USA, 2008. ACM.
- [15] C. Rudin and R. E. Schapire. Margin-based ranking and an equivalence between AdaBoost and RankBoost. *Journal of Machine Learning Research*, 10:2193?2232, Oct 2009.
- [16] R. J. Serfling. *Approximation theorems of mathematical statistics*. J. Wiley & Sons, 1980.