

# Computable Combinatorial Overfitting Bounds\*

*K. V. Vorontsov<sup>1</sup>, A. I. Frey<sup>2</sup>, E. A. Sokolov<sup>3</sup>*

voron@forecsys.ru, oleksandr.frei@gmail.com, sokolov.evg@gmail.com

<sup>1</sup>Dorodnitsyn Computing Centre, Russian Academy of Sciences; <sup>2</sup>Moscow Institute of Physics and Technology; <sup>3</sup>Moscow State University

Computable combinatorial data dependent on generalization bounds are studied. This approach is based on simplified probabilistic assumptions: it is assumed that the instance space is finite, the labeling function is deterministic, and the loss function is binary. A random walk across a set of linear classifiers with low error rate is used to compute the bound efficiently. The experimental evidence to confirm that this approach leads to practical overfitting bounds in classification tasks is provided.

**Keywords:** *probability of overfitting, empirical risk, combinatorial bounds, similarity between classifiers, cross-validation procedure.*

## Вычислимые комбинаторные оценки вероятности переобучения\*

*Воронцов К. В.<sup>1</sup>, Фрей А. И.<sup>2</sup>, Соколов Е. А.<sup>3</sup>*

<sup>1</sup>Москва, Вычислительный Центр РАН; <sup>2</sup>Московский Физико-Технический Институт; <sup>3</sup>Московский Государственный Университет

В данной статье изучаются комбинаторные оценки обобщающей способности, вычислимые по обучающей выборке. Эти оценки основаны на упрощенной вероятностной модели, в которой рассматривается лишь конечная генеральная совокупность объектов и бинарная функция потерь. Для линейных классификаторов предлагается новый эффективный метод вычисления комбинаторных оценок, использующий случайные блуждания по множеству классификаторов с низким числом ошибок. В заключении приводится экспериментальное обоснование предлагаемого метода.

**Ключевые слова:** *вероятность переобучения, эмпирический риск, оценка вероятности переобучения, сходство алгоритмов, скользящий контроль.*

### 1 Introduction

Accurate bounding of overfitting is an active area of research starting with the pioneer work [1]. A widely adapted approach is based on a probabilistic framework where an instance space  $\mathbb{X}$  (usually, of an infinite cardinality) is equipped with an unknown probability distribution. Consider an independent and identically distributed (i.i.d.) sample  $X = \{x_1, \dots, x_\ell\}$  from  $\mathbb{X}$ , a set  $A$  of all feasible classifiers (for example, all linear classifiers in the original feature space), and a learning algorithm  $\mu$  which selects a specific classifier  $a = \mu X$  from the set  $A$  based on the observed sample  $X$ . The goal of generalization bounds is to predict an average performance of the classifier  $a$  on the whole  $\mathbb{X}$ . The most generalization bounds are derived from various concentration inequalities [2] and take into account the dimensionality of  $A$  (such as

---

\*This work was supported by the Russian Foundation for Basic Research (projects Nos. 11-07-00480 and 12-07-33099-mol-a-ved) and by the program "Algebraic and Combinatorial Methods in Mathematical Cybernetics and New Generation Information Systems" of the Department of Mathematical Sciences of the Russian Academy of Sciences.

Vapnik–Chervonenko (VC) dimension, fat-shattering dimension, etc.), properties of the learning algorithm  $\mu$  (such as the local properties of empirical risk minimization [3]), and information drawn from the observed sample  $X$  (such as the normalized margin in margin-based bounds [4, 5]). Generalization bounds can be useful in structural risk minimization and in model selection, and, hope, some time in future they could replace costly cross-validation procedures.

Despite recent significant improvements [6], there is still a big gap between theory and practice. Even the latest PAC (Probably approximately correct) Bayesian bounds [7] vastly overestimate overfitting, especially when the cardinality of an instance space is small. Another difficulty is that intermediate bounds are usually expressed in terms of unobservable quantities and that makes impossible to measure and compare the factors of overestimation. Finally, many papers lack experimental evaluation. As a result, from practical perspective, the existing bounds are not well suitable for prediction and control of overfitting.

The present authors believe that a simplified probabilistic framework is quite sufficient for obtaining practical overfitting bounds. In this paper, it is assumed that the instance space  $\mathbb{X}$  is finite, and for each object  $x \in \mathbb{X}$  and classifier  $a \in A$ , there exists a deterministic binary loss  $I(a, x) \in \{0, 1\}$  associated with classification of  $x$  as  $a(x)$ . Let assume that all partitions of the set  $\mathbb{X}$  into an observed training sample  $X$  of size  $\ell$  and a hidden test sample  $\bar{X} = \mathbb{X} \setminus X$  of size  $k$  can occur with equal probability. Then, the overfitting probability can be defined by a purely combinatorial formula [8]:

$$Q_\varepsilon(\mu, \mathbb{X}) = \mathbf{P}[\nu(\mu(X), \bar{X}) - \nu(\mu(X), X) \geq \varepsilon] \tag{1}$$

where  $\mu$  is the learning algorithm;  $\nu(a, X)$  is the error rate of a classifier  $a \in A$  on a sample  $X$ ; and the square brackets denote a transformation of a logical value into a numerical one: [true] = 1, [false] = 0. A definition similar to (1) first appears in [9] for a specific case of  $k = 1$  and then in [10] (this time for an arbitrary  $k$ , but with significant notation differences). This definition closely resembles the procedure of complete cross-validation, which is known to provide sharp estimates of performance of a learning algorithm on data yet unknown during learning phase.

Definition (1) is not guaranteed to be an upper bound for new objects beyond  $\mathbb{X}$  (not even up to some probability). In this paper, the authors do not discuss how to mitigate this problem. It is just assumed that the lack of guaranties is acceptable in the same way as people normally accept results of a fare 10-fold cross-validation in experimental evaluations.

Within (1), an empirical risk minimization  $\mu X = \arg \min_{a \in A} \nu(a, X)$  is used as a learning algorithm. This requires an explicit representation of the set of classifiers  $A$ , which might be of enormous cardinality ( $10^9$  and higher) in real applications. In this paper, a special case of linear classifiers is studied and an efficient algorithm that samples a small set of classifiers (up to  $10^4$ ) to recover accurate overfitting bound (1) is presented. Also note that a direct computation of (1) is intractable because it involves a sum across all  $\mathbb{C}_{\ell+k}^\ell$  subsets  $X \subset \mathbb{X}$ . For empirical risk minimization, efficient upper bounds on (1), obtained in [11] are used, and compared with Monte-Carlo estimate of (1).

Overfitting of logistic regression in experiments is studied on 15 datasets from the UCI repository [12]. The results confirm that the present approach provides sharp estimates of overfitting, which correlate with the actual overfitting, recovers a correct shape of the learning curves, and outperform the state-of-the-art of PAC-Bayesian bounds.

The rest of this paper is organized as follows. Section 2 contains a brief review of combinatorial bounds on (1). Section 3 describes the algorithm of sampling a representative set of linear classifiers. Section 4 provides experimental results and Section 5 concludes the paper.

## 2 Background

Let  $\mathbb{X} = \{x_1, \dots, x_L\}$  be a finite instance space and  $A$  be a set of classifiers. By  $I: A \times X \rightarrow \{0, 1\}$  denote a binary loss function such that  $I(a, x) = 1$  if a classifier  $a$  produces an error on an object  $x$ . For further consideration, there is no need to specify what is “classifier.” Particularly, a regression function can also be a “classifier” if a binary loss function is used.

The binary vector  $(a_i) \equiv (I(a, x_i))_{i=1}^L$  of size  $L$  is called an *error vector* of the classifier  $a$ . Assume that all classifiers from  $A$  have pairwise different error vectors. The number of errors of a classifier  $a$  on a sample  $X \subseteq \mathbb{X}$  is defined as  $n(a, X) = \sum_{x \in X} I(a, X)$ . The error rate is defined as  $\nu(a, X) = (1/|X|)n(a, X)$ . The subset  $A_m = \{a \in A: n(a, \mathbb{X}) = m\}$  is called *m-layer* of classifiers.

A *learning algorithm* is a mapping  $\mu: 2^{\mathbb{X}} \rightarrow A$  that takes a training sample  $X \subseteq \mathbb{X}$  and gives a classifier  $\mu X \in A$ . The learning algorithm  $\mu$  is called *empirical risk minimization* (ERM) whenever for all  $X \in \mathbb{X}$  it satisfies  $\mu X \in A(X)$  where

$$A(X) \equiv \operatorname{Arg} \min_{a \in A} n(a, X).$$

The choice of a classifier that minimizes empirical risk may be ambiguous because of discreteness of the function  $n(a, X)$ . An ERM algorithm  $\mu$  is said to be *pessimistic* if

$$\mu X \in \operatorname{Arg} \max_{a \in A(X)} n(a, \bar{X}).$$

The pessimistic ERM cannot be implemented in practice because it looks into a hidden part of data  $\bar{X}$  unknown at the learning stage. Nevertheless, pessimistic ERM is a very useful theoretical concept because it gives tight upper bounds of overfitting probability for any ERM.

**Permutational probability.** By  $[\mathbb{X}]^\ell$  denote a set of all  $\mathbb{C}_L^\ell = L!/(\ell!(L-\ell)!)$  samples  $X \subset \mathbb{X}$  of size  $\ell$ . Define a probability operator  $\mathbf{P}$  and an expectation operator  $\mathbf{E}$  for a predicate  $\varphi: [\mathbb{X}]^\ell \rightarrow \{0, 1\}$  and a real function  $\psi: [\mathbb{X}]^\ell \rightarrow \mathbb{R}$ :

$$\mathbf{P} \varphi = \mathbb{C}_L^{\ell-1} \sum_{X \in [\mathbb{X}]^\ell} \varphi(X); \quad \mathbf{E} \psi = \mathbb{C}_L^{\ell-1} \sum_{X \in [\mathbb{X}]^\ell} \psi(X).$$

If the *discrepancy*  $\delta(a, X) = \nu(a, \bar{X}) - \nu(a, X)$  is greater than a given nonnegative threshold  $\varepsilon$ , then the classifier  $a = \mu X$  is said to be *overfitted*. The goal is to estimate the *probability of overfitting*:

$$Q_\varepsilon(\mu, \mathbb{X}) = \mathbf{P}[\delta(\mu, X) \geq \varepsilon].$$

where  $\delta(\mu, X) = \delta(\mu X, X)$  for short.

The *inversion* of an upper bound  $Q_\varepsilon \leq \eta(\varepsilon)$  is an inequality  $\nu(\mu X, \bar{X}) - \nu(\mu X, X) \leq \varepsilon(\eta)$  that holds with probability at least  $1 - \eta$  where  $\varepsilon(\eta)$  is the inverse function for  $\eta(\varepsilon)$ . The *median of an upper bound*  $Q_\varepsilon \leq \eta(\varepsilon)$  is the inversion at  $\eta = 1/2$ .

Average train and test errors are defined as follows:

$$\nu_\ell(\mu, \mathbb{X}) = \mathbf{E} \nu(\mu X, X); \tag{2}$$

$$\bar{\nu}_\ell(\mu, \mathbb{X}) = \mathbf{E} \nu(\mu X, \bar{X}). \tag{3}$$

**Hypergeometric distribution.** For a classifier  $a$  such that  $m = n(a, \mathbb{X})$ , the probability to have  $s$  errors on a sample  $X$  is given by a hypergeometric function:

$$\mathbf{P}[n(a, X) = s] = \mathbb{C}_m^s \mathbb{C}_{L-m}^{\ell-s} \mathbb{C}_L^{\ell-1} \equiv h_L^{\ell, m}(s)$$

where argument  $s$  runs from  $s_0 = \max\{0, m - k\}$  to  $s_1 = \min\{m, \ell\}$ , and parameter  $m$  takes values  $0, \dots, L$ . It is assumed that  $\mathbb{C}_m^s = h_L^{\ell, m}(s) = 0$  for all other integers  $m, s$ .

Define the hypergeometric cumulative distribution function (left tail of the distribution):

$$H_L^{\ell, m}(z) = \sum_{s=s_0}^{\lfloor z \rfloor} h_L^{\ell, m}(s).$$

Consider a set  $A = \{a\}$  containing a fixed classifier so that  $\mu X = a$  for any  $X$ . Then the probability of overfitting  $Q_\varepsilon$  transforms into the probability of large deviation between error rates on two samples  $X$  and  $\bar{X}$ . If the number of errors  $n(a, \mathbb{X})$  is known, then an exact  $Q_\varepsilon$  bound can be obtained.

**Theorem 1 (FC-bound [11]).** *For a fixed classifier (FC)  $a$  such that  $m = n(a, \mathbb{X})$ , any set  $\mathbb{X}$ , and any  $\varepsilon \in [0, 1]$ , the probability of overfitting is given by the left tail of the hypergeometric distribution:*

$$Q_\varepsilon(a, \mathbb{X}) = H_L^{\ell, m}\left(\frac{\ell}{L}(m - \varepsilon k)\right). \tag{4}$$

The hypergeometric distribution plays a fundamental role in combinatorial bounds. Together with union bound, Eq. (4) provides an upper estimate of  $Q_\varepsilon(\mu, \mathbb{X})$  that holds for any learning algorithm  $\mu$ .

**Theorem 2 (VC-type bound [11]).** *For any set  $\mathbb{X}$ , any learning algorithm  $\mu$ , and any  $\varepsilon \in [0, 1]$ , the probability of overfitting is bounded by the sum of FC-bounds over the set  $A$ :*

$$Q_\varepsilon(\mu, \mathbb{X}) \leq \mathbb{P}\left[\max_{a \in A} \delta(a, X) \geq \varepsilon\right] \leq \sum_{a \in A} H_L^{\ell, m}\left(\frac{\ell}{L}(m - \varepsilon k)\right), \quad m = n(a, \mathbb{X}). \tag{5}$$

There are two reasons for looseness of (5). First, most classifiers in  $A$  are bad and should have vanishing probability to be obtained as a result of learning. Nevertheless, the uniform deviation bound ignores the learning algorithm  $\mu$ . Second, similar classifiers share their contribution, which is ignored by union bound. Better bound should account for actual learning algorithm and similarity between classifiers.

**Splitting and connectivity bounds.** Define an order relation on classifiers  $a \leq b$  as a natural order over their error vectors:  $a_i \leq b_i$  for all  $i = 1, \dots, L$ . Define a metric on classifiers as a Hamming distance between error vectors:  $\rho(a, b) = \sum_{i=1}^L |a_i - b_i|$ .

**Theorem 3 (Splitting and connectivity (SC) bound [11]).** *If learning algorithm  $\mu$  is pessimistic ERM, then for any  $\varepsilon \in [0, 1]$ , the probability of overfitting is bounded by the weighted sum of FC-bounds over the set  $A$ :*

$$Q_\varepsilon(\mu, \mathbb{X}) \leq \sum_{a \in A} \mathbb{C}_{L-q-r}^{\ell-q} \mathbb{C}_L^{\ell-1} H_{L-q-r}^{\ell-q, m-r}\left(\frac{\ell}{L}(m - \varepsilon k)\right). \tag{6}$$

Here,  $m = n(a, \mathbb{X})$ ;  $q = q(a)$  is the upper connectivity; and  $r = r(a)$  is the inferiority of a classifier  $a$ :

$$\begin{aligned} q(a) &= \#\{b \in A : a < b \text{ and } \rho(a, b) = 1\}; \\ r(a) &= \#\{x \in \mathbb{X} : I(a, x) = 1 \text{ and } \exists b \in A \text{ such that } b \leq a \text{ and } I(x, b) = 0\} \end{aligned}$$

where for any set  $S$  notation  $\#S$  stands for cardinality of  $S$ .

Splitting and connectivity bound (6) turns into VC-type bound (5) when all  $q(a)$  and  $r(a)$  are set to zeros.

The weight  $P_a = \mathbb{C}_{L-q-r}^{\ell-q} \mathbb{C}_L^{\ell-1}$  in sum (6) is an upper bound on the probability  $\mathbb{P}[\mu X = a]$  to get a given classifier  $a$  as a result of learning. This quantity decreases exponentially as the connectivity  $q(a)$  or the inferiority  $r(a)$  increase. This implies that approximate calculation of  $Q_\varepsilon(\mu, \mathbb{X})$  requires knowledge not about the full set  $A$ , but only about few bottom layers of  $A$ . This fact motivates an algorithm presented in the next section.

### 3 Sampling Linear Classifiers

One has to deal with the set of all classifiers  $A$  to calculate bounds (5), (6), or to estimate  $Q_\varepsilon(\mu, \mathbb{X})$  directly from definition (1). In this section, an efficient algorithm which samples a small set of classifiers (about  $10^4$ ) sufficient to recover accurate overfitting bound is described.

Consider binary classification problem with labels  $y_i \in \{-1, +1\}$ ,  $i = 1, \dots, L$ , assigned to objects  $x_i \in \mathbb{X} \subset \mathbb{R}^d$ , respectively. Consider a set of unbiased linear classifiers  $a(x; w) = \text{sign}\langle w, x \rangle$  where  $w \in \mathbb{R}^d$  is a real vector of weights. A pair of classifiers  $(w_1, w_2)$  is called *neighbors* if their classification differs only by one object:  $x \in \mathbb{X}$  such that  $\text{sign}(\langle w_1, x \rangle) \text{sign}(\langle w_2, x \rangle) = -1$ .

The immediate goal is to find all or some of neighbors of a given classifier  $w_0$ . Then, this procedure will be used to organize random walk on the graph  $G = (A, E)$  where vertices correspond to classifiers in  $A$  and edges connect neighbor classifiers.

**Finding neighbor classifiers along specific direction.** *Dual transformation*  $D$  maps a point  $x \in \mathbb{R}^d$  into hyperplane  $D(x) = \{w \in \mathbb{R}^d: \langle w, x \rangle = 0\}$  and a hyperplane  $h = \{x \in \mathbb{R}^d: \langle w, x \rangle = 0\}$  into point  $D(h) = w$ . Applying dual transformation  $D$  to finite set of points  $\mathbb{X} \subset \mathbb{R}^d$  produces a set of hyperplanes  $\mathbb{H} \equiv \{D(x_i)\}_{i=1}^L$ . Each hyperplane  $h_i \in \mathbb{H}$  divides  $\mathbb{R}^d$  into two half-spaces:

$$\begin{aligned} h_i^+ &= \{w \in \mathbb{R}^d: \text{sign}\langle w, x_i \rangle = y_i\}; \\ h_i^- &= \{w \in \mathbb{R}^d: \text{sign}\langle w, x_i \rangle = -y_i\}. \end{aligned}$$

These half-spaces  $h_i^+$  and  $h_i^-$  correspond to linear classifiers giving correct and incorrect answer on  $x_i$ , respectively. So, to find all classifiers with given error vector  $I = (I_i)_{i=1}^L$ ,  $I_i \in \{+, -\}$  where “+” corresponds to correct answer and “-” corresponds to incorrect, let just find the intersection of half-spaces  $\bigcap_{i=1}^L h_i^{I_i}$ . This intersection contains all linear classifiers with error vector  $I$  (and only them). So, a set of hyperplanes  $\mathbb{H}$  dissects  $\mathbb{R}^d$  into convex polytopes called *cells*, and the partition itself is called *an arrangement of hyperplanes* [13]. It can be shown that finding neighbors of classifier  $w_0 \in \mathbb{R}^d$  is equivalent to finding cells adjacent to the cell of  $w_0$  in arrangement  $\mathbb{H}$ .

In order to find a neighbor of the classifier  $w_0$ , let select an arbitrary vector  $u \in \mathbb{R}^d$  and consider a parametric set of classifiers  $\{w_0 + tu: t \geq 0\}$ . This set corresponds to a ray in the space of classifiers which starts from  $w_0$  and goes along the direction of  $u$ . An intersection of this ray with hyperplane  $h_i \in \mathbb{H}$  is defined by condition  $\langle w_0 + tu, x_i \rangle = 0$ , e. g., for  $t_i = -\langle w_0, x_i \rangle / \langle u, x_i \rangle$ . Let  $t_{(1)}$  and  $t_{(2)}$  be the first and the second smallest positive values from  $\{t_i\}$ ,  $i = 1, \dots, L$ . Whenever  $t_{(1)} \neq t_{(2)}$ , one can conclude that  $w' = w_0 + (t_{(1)} + t_{(2)})u/2$  defines an adjacent classifier along direction  $u$ .

**Random walk on classifiers graph.** Techniques of random walk [14, 15, 16] provide common approach to sample vertices from huge graphs. They are based on stationary distributions of Markov chains and have nice properties when the sample is large. The goal is to get

---

**Algorithm 1** Random walk on classifiers graph

---

**Require:** starting point  $w_0$ ; sample  $\mathbb{X} \subset \mathbb{R}^d$ ; integer parameters  $N, m, n$ ; float parameter  $p \in (0, 1]$ ;

**Ensure:** set of classifiers  $A$  with unique error vectors.

- 1: Initialize concurrent random walk:  $v_i = w_0, i = 1, \dots, N$ ;
  - 2: Create set  $A := \emptyset$ ;
  - 3: **while**  $A.size() < n$
  - 4:   **for all**  $i \in 1, \dots, N$
  - 5:     Find neighbor  $v'_i$  of  $v_i$  along random direction  $u \in \mathbb{R}^d$ ;
  - 6:     **if**  $n(v'_i, \mathbb{X}) > n(v_i, \mathbb{X})$  **then**
  - 7:       with probability  $(1 - p)$  **continue**;
  - 8:       **else if**  $n(v'_i, \mathbb{X}) > n(w_0, \mathbb{X}) + m$  **then**
  - 9:        **continue**;
  - 10:      $v_i = v'_i$ ;
  - 11:      $A.add(v_i)$ ;
  - 12: **return**  $A$
- 

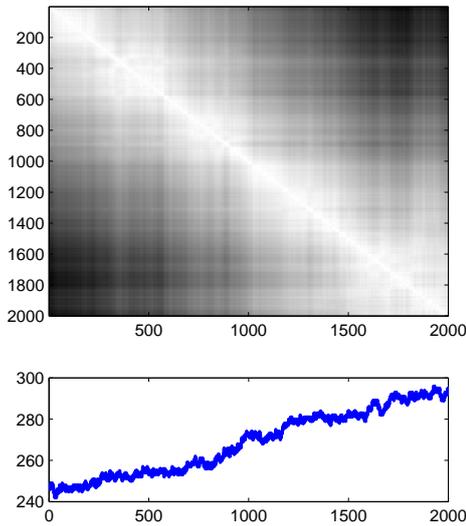


Fig. 1: Map of hamming distances between classifiers (top chart) and error profile (bottom chart) produced by a simple random walk

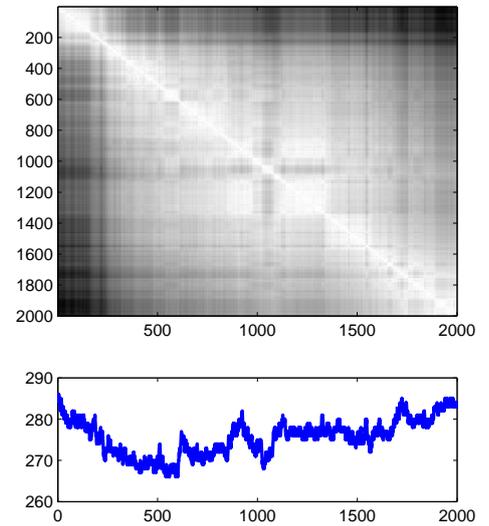


Fig. 2: Map of hamming distances between classifiers (top chart) and error profile (bottom chart) produced by random walk where a step to upper vertex is made with probability 0.5

a small sample sufficient to estimate overfitting from (1), (5), or (6). In this paragraph, it will be discussed how to organize such random walk on  $A$  based on procedure that finds a random neighbor for  $w \in A$ .

The algorithm is given in listing 1. It is controlled by the desired number of classifiers  $n$ , maximal number of layer  $m$ , the number of concurrent walks  $N$ , and the probability  $p$  of transition towards classifier with higher number of errors. The computational complexity of this algorithm is  $O(Ldn)$ .

To explain the necessity of the parameter  $p$ , the results of the simplest random walk with  $n = 2000$  iterations are presented in Fig. 1. The bottom chart displays the number of errors  $n(v_i, \mathbb{X})$  as a function of step. The upper chart displays a color map of pairwise hamming distances  $\rho(v_i, v_j)$  between sampled vertices  $v_i$  and  $v_j$ . As a starting point, a classifier learned by logistic regression is used. It is natural to expect that it has relatively small number of errors which drifts upwards along random walk. This effect is undesired, because classifiers with high number of errors have too small chance to be selected by learning algorithm.

Figure 2 presents similar result for updated random walk where a step to upper vertex is made with probability  $p = 0.5$ . This enforces random walk to stay within the lower layers of the graph.

## 4 Experiment

The goal of the experiment on the benchmark datasets is twofold. First, it was checked whether combinatorial functionals  $Q_\varepsilon(\mu, \mathbb{X})$  (1) and  $\bar{\nu}_\ell(\mu, \mathbb{X})$  (3) together with algorithm 1 provide an accurate estimates of the overfitting on the holdout testing sample. Second, direct Monte-Carlo estimates of overfitting based on functional (1) were compared with VC-type bound (5), SC-bound (6), and with recent PAC-Bayesian dimension dependent (DD) and dimension independent (DI) margin bounds proposed in [7].

Table 1: Description of datasets

Dataset	#Examples	#Features	Dataset	#Examples	#Features
Sonar	208	60	Statlog	2310	19
Glass	214	9	Wine	4898	11
Liver dis.	345	6	Waveform	5000	21
Ionosphere	351	34	Pageblocks	5473	10
Wdbc	569	30	Optdigits	5620	64
Australian	690	6	Pendigits	10992	16
Pima	768	8	Letter	20000	16
Faults	1941	27			

Fifteen datasets from the UCI repository [12] have been used. If the dataset is a multiclass problem, the data were manually grouped into two classes since the binary classification problem has been studied. For preprocessing, objects with one or more missing features have been eliminated and all features have been normalized into  $[0, 1]$  interval. A description of the datasets is given in Table 1 with number of examples after elimination.

In all experiments, the original dataset  $\mathbb{X}$  was splitted into a training sample  $\mathbb{X}_L$  and a testing sample  $\mathbb{X}_K$ . The training sample  $\mathbb{X}_L$  is used to train a logistic regression and to calculate overfitting bounds. Then, the predictions of the bounds were compared with the actual error rate on  $\mathbb{X}_K$ .

In the first experiment, the learning curves of logistic regression, where  $L$  runs from 5% to 95% of the original dataset size with 5 percent steps were built. For each  $L$ ,  $M = 100$  splits  $\mathbb{X} = \mathbb{X}_L^i \cup \mathbb{X}_K^i$ ,  $i = 1, \dots, M$ , were generated and used to get Monte-Carlo estimates of train error rate  $\nu_L(\mu_{LR}, \mathbb{X})$  from (2) and test error rate  $\bar{\nu}_L(\mu_{LR}, \mathbb{X})$  from (3) for logistic regression learning algorithm  $\mu_{LR}$ :

$$\hat{\nu}_L(\mu_{LR}, \mathbb{X}) = \frac{1}{M} \sum_{i=1}^M \nu(\mu_{LR} \mathbb{X}_L^i, \mathbb{X}_L^i), \quad \hat{\bar{\nu}}_L(\mu_{LR}, \mathbb{X}) = \frac{1}{M} \sum_{i=1}^M \nu(\mu_{LR} \mathbb{X}_L^i, \mathbb{X}_K^i).$$

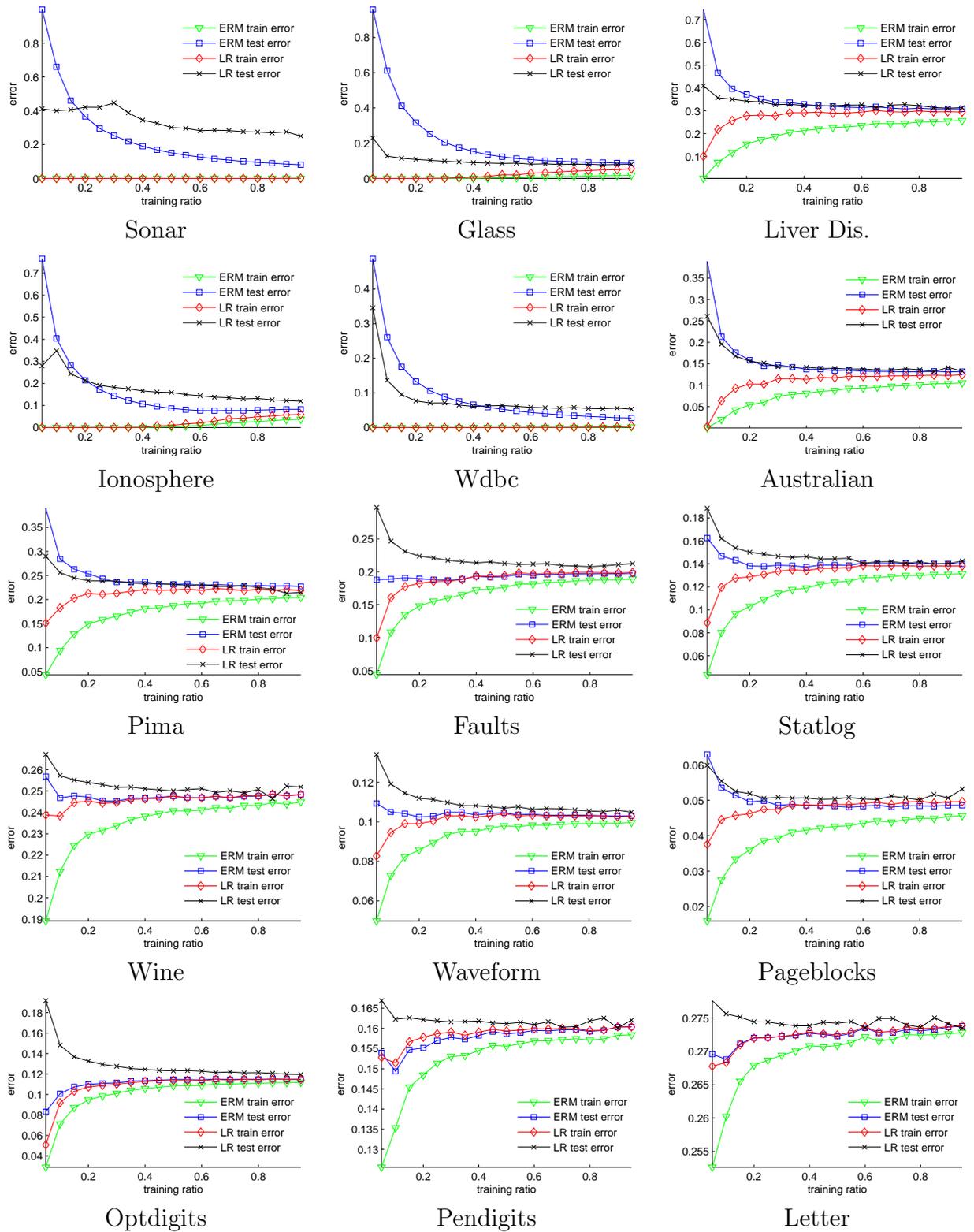


Fig. 3: Learning curves of logistic regression and ERM. The error ratio of logistic regression is estimated by Monte-Carlo method on splits of the original dataset  $\mathbb{X} = \mathbb{X}_L \cup \mathbb{X}_K$ . The error ratio of ERM is estimated on splits of the training set  $\mathbb{X}_L = X_\ell \cup X_k$

After that, for each training sample  $\mathbb{X}_L$ , the classifiers were used and an average ERM errors were estimated: train error  $\nu_\ell(\mu, \mathbb{X}_L)$  and test error  $\bar{\nu}_\ell(\mu, \mathbb{X}_L)$  where  $\mu$  is ERM learning

Table 2: Comparison between real overfitting and various overfitting bounds: TrainErr stands for  $\nu_L(\mu_{LR}, \mathbb{X})$ , TestErr for  $\bar{\nu}_L(\mu_{LR}, \mathbb{X})$ , Overfit is their difference,  $\delta_\ell(\mu) \equiv \bar{\nu}_\ell(\mu, \mathbb{X}_L) - \nu_\ell(\mu, \mathbb{X}_L)$

Task	Monte-Carlo estimates				Generalization bounds			
	TrainErr	TestErr	Overfit	$\delta_\ell(\mu)$	VC	SC	PAC DI	PAC DD
Sonar	0.000	0.271	0.271	0.095	0.185	0.119	1.287	1.287
Glass	0.046	0.075	0.029	0.078	0.211	0.140	1.126	0.738
Liver dis.	0.299	0.314	0.015	0.060	0.261	0.209	1.207	1.067
Ionosphere	0.049	0.125	0.077	0.052	0.150	0.112	1.219	1.153
Wdbc	0.001	0.056	0.055	0.032	0.071	0.043	1.174	0.705
Australian	0.122	0.136	0.013	0.030	0.137	0.110	1.146	0.678
Pima	0.220	0.227	0.007	0.028	0.159	0.127	0.971	0.749
Faults	0.198	0.210	0.012	0.010	0.108	0.087	1.110	1.061
Statlog	0.138	0.142	0.005	0.010	0.096	0.082	1.102	0.747
Wine	0.248	0.250	0.002	0.004	0.134	0.109	0.776	0.637
Waveform	0.103	0.105	0.002	0.004	0.099	0.079	0.561	0.354
Pageblocks	0.050	0.050	0.001	0.004	0.073	0.057	0.737	0.186
Optdigits	0.115	0.121	0.006	0.004	0.102	0.084	1.068	0.604
Pendigits	0.160	0.161	0.001	0.002	0.127	0.103	0.774	0.432
Letter	0.274	0.274	0.001	0.001	0.165	0.137	0.818	0.636

algorithm. To sample classifiers on  $\mathbb{X}_L$ , algorithm 1 was launched with parameters  $n = 8192$ ,  $N = 64$ ,  $m = 15$ ,  $p = 0.8$ , and classifier  $\mu_{LR}\mathbb{X}_L$  was used as a starting point. To estimate  $\nu_\ell(\mu, \mathbb{X}_L)$  and  $\bar{\nu}_\ell(\mu, \mathbb{X}_L)$ , let again compute Monte-Carlo type estimates of definitions (2) and (3) by randomly generating  $M' = 4096$  splits  $\mathbb{X}_L = X_\ell^j \cup X_k^j$ ,  $j = 1, \dots, M'$ , at constant ratio  $\frac{\ell}{L} = 0.8$ :

$$\hat{\nu}_\ell(\mu, \mathbb{X}_L) = \frac{1}{M'} \sum_{j=1}^M \nu(\mu X_\ell^j, X_k^j), \quad \hat{\bar{\nu}}_\ell(\mu, \mathbb{X}_L) = \frac{1}{M'} \sum_{j=1}^M \bar{\nu}(\mu X_\ell^j, X_k^j).$$

These estimates are then averaged over all partitions  $\mathbb{X} = \mathbb{X}_L^i \cup \mathbb{X}_K^i$ .

The four values (the actual train and test errors of logistic regression  $\nu_L(\mu_{LR}, \mathbb{X})$  and  $\bar{\nu}_L(\mu_{LR}, \mathbb{X})$ , ERM train error  $\nu_\ell(\mu, \mathbb{X}_L)$ , and ERM test error  $\bar{\nu}_\ell(\mu, \mathbb{X}_L)$ ) are charted as a functions of a training sample size ratio (Fig. 3) sorted according to sizes of datasets, from the smallest to the largest. Note that ERM test error might be either below or above actual test error rate of logistic regression because  $\mu$  and  $\mu_{LR}$  are quite different learning algorithms. However, from charts, one may conclude that  $\bar{\nu}_\ell(\mu, \mathbb{X}_L)$ , estimated only based on  $\mathbb{X}_L$  provides reasonably good estimate of actual test error rate  $\bar{\nu}_L(\mu_{LR}, \mathbb{X})$  and of the learning curve on test sample  $\mathbb{X}_K$ .

Now, let turn to comparison of different overfitting bounds. For each dataset, 5-fold cross validation were used and the results were averaged over 20 runs (for a total 100 runs). As before, training sample  $\mathbb{X}_L$  was used to learn logistic regression, run algorithm 1, and estimate  $\nu_\ell(\mu, \mathbb{X}_L)$  and  $\bar{\nu}_\ell(\mu, \mathbb{X}_L)$  based on 4096 randomly generated splits  $\mathbb{X}_L = X_\ell^j \cup X_k^j$ . In addition,  $\mathbb{X}_L$  was used to estimate overfitting  $\bar{\nu}_\ell(\mu, \mathbb{X}_L) - \nu_\ell(\mu, \mathbb{X}_L)$  by medians of VC-type bound (5) and SC-bound (6) and to calculate DD-margin and DI-margin bounds from [7]. The results are presented in Table 2. Note that while all combinatorial bounds estimate overfitting, PAC DI and PAC DD are the upper bounds on the test error.

The key observation is that  $\delta_\ell(\mu) \equiv \bar{\nu}_\ell(\mu, \mathbb{X}_L) - \nu_\ell(\mu, \mathbb{X}_L)$  is in the order of magnitude sharper than any of the other bounds. It works well for all datasets except Sonar (which is the smallest dataset in the selection described). Across combinatorial bounds, the SC-bound outperform VC-type bound, but still vastly overestimates the target quantity  $\delta_\ell(\mu)$ . All combinatorial bounds provide tighter estimates on overfitting and test error rate than PAC-Bayesian bounds.

Note that the VC-bound is estimated by a small subset of  $A$  obtained from a random walk. This is a “localized” VC-bound. The usual VC-bound estimated from VC-dimension  $d$  of a whole set  $A$  should be greater than 1 on all datasets.

## 5 Concluding Remarks

In this paper, new random walk technique is presented for efficient calculation of combinatorial data-dependent generalization bounds. Although combinatorial bounds are obtained for

empirical risk minimization under binary loss, it is shown that they provide sharp overfitting estimates for logistic regression. The bounds recover correct shape of the learning curves for logistic regression, correlate well with the actual overfitting, and outperform both classical VC-bound and recent state-of-the-art PAC-Bayesian bounds in experiments on 15 datasets from the UCI repository.

## References

- [1] *Vapnik V. N., Chervonenkis A. Y.* On the uniform convergence of relative frequencies of events to their probabilities // *Theory of Probability and Its Applications*. 1971. Vol. 16(2). Pp. 264–280.
- [2] *Lugosi G.* On concentration-of-measure inequalities // *Machine Learning Summer School*, Australian National University, Canberra, 2003.
- [3] *Koltchinskii V.* Local Rademacher complexities and oracle inequalities in risk minimization (with discussion) // *The Annals of Statistics*, 2006. Vol. 34. Pp. 2593–2706.
- [4] *Koltchinskii V., Panchenko D.* Empirical margin distributions and bounding the generalization error of combined classifiers // *The Annals of Statistics*, 2002. Vol. 30(1). Pp. 1–50.
- [5] *Koltchinskii V., Panchenko D.* Bounding the generalization error of convex combinations of classifiers: balancing the dimensionality and the margins // *The Annals of Applied Probability*, 2003. Vol. 13(1). Pp. 213–252.
- [6] *Boucheron S., Bousquet O., Lugosi G.* Theory of classification: A survey of some recent advances // *ESAIM: probability and statistics*. 2005. Vol. 9(1). Pp. 323–375.
- [7] *Jin C., Wang L.* Dimensionality dependent PAC-Bayes margin bound // *Advances in Neural Information Processing Systems*. 2012. Vol. 25. Pp. 1043–1051.
- [8] *Vorontsov K. V.* Splitting and similarity phenomena in the sets of classifiers and their effect on the probability of overfitting // *Pattern Recognition and Image Analysis*. 2009. Vol. 19(3). Pp. 412–420.
- [9] *Haussler D., Littlestone N., Warmuth M. K.* Predicting  $\{0, 1\}$ -functions on randomly drawn points // *Information and Computation*, 1994. Vol. 115(2). Pp. 248–292, 1994.
- [10] *Bax E.* Similar classifiers and VC error bounds, 1997.
- [11] *Vorontsov K. V., Ivahnenko A. A.* Tight combinatorial generalization bounds for threshold conjunction rules. // *4-th International Conference on Pattern Recognition and Machine Intelligence (PReMI'11)*. Lecture Notes in Computer Science. Springer-Verlag, 2011. Pp. 66–73.
- [12] *Asuncion A, Newman D. J.* UCI Machine Learning Repository // University of California, Irvine, School of Information and Computer Sciences, 2007.
- [13] *Agarwal P. K., Sharir P.* Arrangements and their applications // *Handbook of Computational Geometry*, 1998. Pp. 49–119.
- [14] *Avrachenkov K., Ribeiro B., Towsley D.* (2010) Improving random walk estimation accuracy with uniform restarts // *Proc. of WAW 2010*.
- [15] *Ribeiro B., Towsley D.* Estimating and sampling graphs with multidimensional random walks // *10th Conference on Internet Measurement*, 2010. Pp. 390–403.
- [16] *Lee C., Xu X., Eun D.* Beyond random walk and Metropolis-Hastings samplers: Why you should not backtrack for unbiased graph sampling // *ACM SIGMETRICS Performance Evaluation Review*, 2012. Vol. 40(1). Pp. 319–330.