

# Recent Advances on Generalization Bounds Part I

Konstantin Vorontsov

Computing Center RAS • Moscow Institute of Physics and Technology

4th International Conference on  
Pattern Recognition and Machine Intelligence (PReMI'11)  
Moscow, Russian Federation • June 27 – July 1, 2011

## Structure

- **Tutorial Part I. Overview of Generalization Bounds**

June 27, 13:40–14:40 PReMI tutorial 4 ( $\Gamma$ -313)

- VC, Occam Razor, Rademacher, and margin-based bounds
- How these bounds can be used for learning algorithm design?

- **Tutorial Part II. Combinatorial Generalization Bounds**

June 27, 15:00–16:00 PReMI tutorial 4 (continued) ( $\Gamma$ -313)

- Why complexity bounds are so loose (overestimated)?
- How to obtain tight or even exact bounds?
- Will they be useful?

- **Part III. Tight Combinatorial Generalization Bounds for Threshold Conjunction Rules**

June 29, 11:00–11:20 PReMI session 4 ( $\Gamma$ -408)

- A practical issue from Combinatorial Generalization Bounds

## Contents of the Part I

- 1 Generalization bounds: notations and definitions**
  - Classification problem
  - Generalization bounds
  - Techniques and tools
- 2 Complexity bounds**
  - Vapnik-Chervonenkis bounds (PAC learning)
  - Occam Razor bounds
  - Rademacher complexity bounds
- 3 Margin-based bounds**
  - Margin-based classifiers
  - Kernel Machine
  - Weighted voting of classifiers

## Classification problem

$X$  — a set of *objects*, usually  $\mathbb{R}^n$

$Y$  — a set of *class labels*, usually  $\{-1, +1\}$  or  $\{1, \dots, M\}$

$y: X \rightarrow Y$  — unknown *target function*

$X^\ell = \{(x_1, y_1), \dots, (x_\ell, y_\ell)\}$  — *training set*,  $y_i = y(x_i)$ ,  $i = 1, \dots, \ell$

Classification is a supervised learning problem:

find a classifier  $a: X \rightarrow Y$  from a given function set  $A$  that *generalizes well*, that is approximates well a target  $y$  not only on the training set  $X^\ell$  but everywhere on  $X$ .

One must specify accurately:

- what means “approximate well”?
- what means “approximate everywhere on  $X$ ”?

## Probabilistic model of data

Let  $X \times Y$  be a probability space with unknown distribution  $p(x, y)$  and observations  $(x_i, y_i)_{i=1}^{\ell}$  be drawn independently from  $p$

Define a *binary loss function*  $I(a, x)$ , usually

$$I(a, x) = [a(x) \neq y(x)] = \begin{cases} 1, & a(x) \neq y(x) \\ 0, & a(x) = y(x) \end{cases}$$

*Empirical error* (error rate, frequency of errors) of a classifier  $a$

$$\nu(a, X^{\ell}) = \frac{1}{\ell} \sum_{i=1}^{\ell} I(a, x_i)$$

*Probability of error* (generalization error) of a classifier  $a$

$$P(a) = P_x(I(a, x) = 1)$$

## Empirical risk minimization

*Learning algorithm*  $\mu$  is a function that takes a training sample  $X^\ell$  and gives a classifier  $a^*$  from  $A$ :

$$a^* = \mu(X^\ell)$$

*Empirical risk minimization* (ERM) is a classical example of the learning algorithm:

$$a^* = \arg \min_{a \in A} \nu(a, X^\ell)$$

Unfortunately, ERM can lead to *overfitting*, when

$$\nu(a^*, X^\ell) \ll P(a^*)$$

## How generalization bounds can help to reduce overfitting

There are two things to do:

- 1 to give an upper bound of the probability of error that holds for any  $A$ , any  $p$ , any  $\mu$  (and sometimes any  $X^\ell$ ):

$$P(a^*) \leq \hat{P}(a^*, X^\ell)$$

Two types of bounds exist:

$\hat{P}(a)$  — data-independent bound (usually very loose)

$\hat{P}(a, X^\ell)$  — data-dependent bound (most recent bounds)

- 2 to construct the *learning algorithm*  $\mu$  that minimizes the generalization error bound:

$$a^* = \mu(X^\ell) \equiv \arg \min_{a \in A} \hat{P}(a, X^\ell)$$

## Probably Approximately Correct (PAC) learning

**Problem statement** [Vapnik & Chervonenkis, 1969; Valiant, 1984]:  
 Given only  $A$  and  $\ell$ , find a bound  $\eta(\varepsilon, \ell, A)$  on the *uniform deviation* of the error frequency from the error probability:

$$P_{X^\ell} \left( P(a^*) - \nu(a^*, X^\ell) \geq \varepsilon \right) \leq$$

$$P_{X^\ell} \left( \sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon \right) \leq \eta(\varepsilon, A, \ell) \quad - ?$$

**Plus:** it holds for any learning algorithm  $\mu$  and any distribution  $p$

**Minus:** it is a worst-case bound which can be very loose

The **uniform convergence principle** is an axiom in VC-theory, PAC-learning theory, and Rademacher Complexity theory



## The inversion technique

If the bound has been obtained

$$P_{X^\ell} \left( \sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon \right) \leq \eta(\varepsilon, A, \ell)$$

then, with probability at least  $1 - \eta$  for any classifier  $a \in A$

$$P(a) \leq \nu(a, X^\ell) + \varepsilon(\eta, A, \ell),$$

where  $\varepsilon(\eta, A, \ell)$  is the *inverse function* for  $\eta(\varepsilon, A, \ell)$ .

A new learning algorithm: try to minimize generalization error  $P(a)$

$$a^* = \arg \min_a \min_A \left( \nu(a, X^\ell) + \varepsilon(\eta, A, \ell) \right)$$

(differs from ERM by penalty  $\varepsilon(\eta, A, \ell)$  and extra optimization by  $A$ )

## The binomial tail bound for one-classifier case

The empirical error of a fixed classifier  $a$  is distributed binomially:

$$P(\nu(a, X^\ell) = \frac{s}{\ell}) = C_\ell^s p^s (1-p)^{\ell-s}, \quad \text{where } p = P(a)$$

*Binomial tail* — exact bound, tedious inversion:

$$P(P(a) - \nu(a, X^\ell) \geq \varepsilon) = \sum_{s=0}^{\ell p - \ell \varepsilon} C_\ell^s p^s (1-p)^{\ell-s}$$

*Chernoff's inequality* — inflated bound, easier inversion:

$$\leq \exp(-\ell \text{KL}(p - \varepsilon \| p))$$

where  $\text{KL}(q \| p) = q \ln \frac{q}{p} + (1-q) \ln \frac{1-q}{1-p}$  is Kullback–Leibler divergence.

*Hoeffding's inequality* — more inflated bound, trivial inversion:

$$\leq \exp(-2\ell \varepsilon^2)$$

## Vapnik-Chervonenkis bound for finite set

## Theorem (data-independent bound)

If  $A$  is a finite set of classifiers then for any  $\varepsilon \in (0, 1)$

$$P\left(\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon\right) \leq |A| \cdot \exp(-2l\varepsilon^2)$$

**Proof sketch:**

first, apply the *union bound*:

$$P\left(\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon\right) \leq \sum_{a \in A} P\left(P(a) - \nu(a, X^\ell) \geq \varepsilon\right)$$

second, apply the one-classifier bound from the previous slide:

$$\leq |A| \cdot \exp(-2l\varepsilon^2)$$

## Vapnik-Chervonenkis bound for infinite set

### Theorem (data-independent bound)

If  $A$  is an arbitrary set of classifiers then for any  $\varepsilon \in (0, 1)$

$$P\left(\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon\right) \leq \Delta^A(2\ell) \cdot \frac{3}{2} \exp(-\ell\varepsilon^2)$$

where  $\Delta^A(L)$  is the growth function of the set  $A$ .

**Definition.** The *growth function*  $\Delta^A(L)$  of the set  $A$  is the maximal number of distinct  $L$ -dimensional binary vectors  $\mathbf{a} = (I(a, x_1), \dots, I(a, x_L))$  induced by all classifiers  $a \in A$  on a sample  $X^L = (x_1, \dots, x_L)$

**Informally,**  $\Delta^A(L)$  is a *complexity measure* of the set  $A$ .

## Vapnik-Chervonenkis dimension

**Definition:** The VC-dimension of the set  $A$  is the maximal sample size  $h$  such that  $\Delta^A(h) = 2^h$ .

### Theorem

If such  $h$  exists then  $\Delta^A(L) \leq C_L^0 + \dots + C_L^h \leq \frac{3}{2} \frac{L^h}{h!}$

Consider a two-class classification problem  $Y = \{-1, +1\}$  and a set  $A$  of linear classifiers in  $n$ -dimensional object space  $X = \mathbb{R}^n$ :

$$a(x) = \text{sign}(w_1 x^1 + \dots + w_n x^n), \quad x = (x^1, \dots, x^n) \in X.$$

### Theorem

$\text{VCdim}(A) = n$

## The inversion technique

$$\text{VC-bound } P\left(\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon\right) \leq \Delta \cdot \exp(-\ell\varepsilon^2)$$

gives with probability at least  $1 - \eta$

$$P(a) \leq \underbrace{\nu(a, X^\ell)}_{\text{empirical risk}} + \underbrace{\sqrt{\frac{1}{\ell} \ln \Delta + \frac{1}{\ell} \ln \frac{1}{\eta}}}_{\text{complexity penalty}}$$

$$\text{VC-bound } P\left(\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon\right) \leq \frac{3}{2} \frac{L^h}{h!} \cdot \frac{3}{2} \exp(-\ell\varepsilon^2)$$

gives with probability at least  $1 - \eta$

$$P(a) \leq \underbrace{\nu(a, X^\ell)}_{\text{empirical risk}} + \underbrace{\sqrt{\frac{h}{\ell} \ln \left(\frac{2e\ell}{h}\right) + \frac{1}{\ell} \ln \frac{4}{9\eta}}}_{\text{complexity penalty}}$$

## Structural Risk Minimization (SRM)

Given a system of nested subsets of increasing dimensions

$$A_0 \subset A_1 \subset \dots \subset A_h \subset \dots$$

Find an optimal dimension  $h^*$ :

$$P(a) \leq \underbrace{\min_{a \in A_h} \nu(a, X^\ell)}_{\text{empirical risk minimization}} + \underbrace{\sqrt{\frac{h}{\ell} \ln \left( \frac{2e\ell}{h} \right) + \frac{1}{\ell} \ln \frac{4}{9\eta}}}_{\text{complexity penalty}} \rightarrow \min_h$$

**The main disadvantage of SRM approach:**

VC-bound is very loose (overestimated)

Then,  $h^*$  may be suboptimal (oversimplified)

Practitioners prefer to use Cross-Validation instead of the bound

## Two main reasons of the VC-bound looseness

- **The uniform deviation bound** is highly overestimated when most classifiers have a vanishing probability to be obtained by the learning algorithm.

In practice, the distribution over classifiers

$$q(a) = P(\mu(X^\ell) = a), \quad a \in A$$

is essentially nonuniform!

- **The union bound** is highly overestimated when there are a lot of similar classifiers.

In practice, this is usually the case!

Let us start with the first problem...



## Occam Razor bound

One can not know  $q(a) = P(\mu(X^\ell) = a)$ , but one can make a shot. Let  $p(a)$  be a normalized function — “prior” distribution over  $A$ .

### Theorem (Occam Razor bound)

For any “prior”  $p(a)$  over  $A$ , for any  $\eta \in (0, 1)$ , for all  $a \in A$

$$P(a) \leq \nu(a, X^\ell) + \sqrt{\frac{1}{\ell} \ln \frac{1}{p(a)} + \frac{1}{\ell} \ln \frac{1}{\eta}}.$$

with probability at least  $1 - \eta$ .

**Statement 1.** If one guess well,  $p(a) = q(a)$ , then this bound is most tight.

**Statement 2.** Although it is still very overestimated because only first of two reasons of looseness has been treated...

## Occam Razor bound: how to specify $p(a)$ ?

### Example 1.

The *uniform prior*  $p(a) = \frac{1}{|A|}$  gives the VC-bound:

$$P(a) \leq \nu(a, X^\ell) + \sqrt{\frac{1}{\ell} \ln |A| + \frac{1}{\ell} \ln \frac{1}{\eta}}.$$

with probability at least  $1 - \eta$ . Nothing new...

## Occam Razor bound: how to specify $p(a)$ ?

### Example 2.

Consider a two-class classification problem  $Y = \{-1, +1\}$  and a set  $A$  of linear classifiers in  $n$ -dimensional object space  $X = \mathbb{R}^n$ :

$$a(x) = \text{sign}(w_1 x^1 + \dots + w_n x^n), \quad x = (x^1, \dots, x^n) \in X.$$

The *Gaussian prior*: weights  $w \in \mathbb{R}^n$  are independent, with zero expectation, and equal variance  $\sigma^2$ :

$$p(a) = \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left(-\frac{1}{2\sigma^2} \|w\|^2\right)$$

Substituting this prior into Occam Razor bound gives

$$P(a) \leq \nu(a, X^\ell) + \sqrt{\frac{n}{\ell} \ln \sigma\sqrt{2\pi} + \frac{\|w\|^2}{2\ell\sigma^2} + \frac{1}{\ell} \ln \frac{1}{\eta}}.$$

## Regularization

The minimization of the obtained bound

$$P(a) \leq \nu(a, X^\ell) + \sqrt{\frac{n}{\ell} \ln \sigma \sqrt{2\pi} + \frac{\|w\|^2}{2\ell\sigma^2} + \frac{1}{\ell} \ln \frac{1}{\eta}} \rightarrow \min_w$$

can be considered as a nontrivial mixture of  $L_0$ - and  $L_2$ -regularization:

$$\nu(a, X^\ell) \rightarrow \min_w \quad \text{ERM}$$

$$\nu(a, X^\ell) + C_0 n \rightarrow \min_w \quad L_0\text{-regularization}$$

$$\nu(a, X^\ell) + C_1 \sum_{j=1}^n |w_j| \rightarrow \min_w \quad L_1\text{-regularization}$$

$$\nu(a, X^\ell) + C_2 \sum_{j=1}^n w_j^2 \rightarrow \min_w \quad L_2\text{-regularization}$$

Bound minimization leads to shrinkage and features selection.

## Two main reasons of the VC-bound looseness (revisited)

- **The uniform deviation bound** is loose when  $q(a)$  is essentially nonuniform distribution.
- **The union bound** is loose when there are a lot of similar classifiers.

Occam Razor bound treats only the first problem.  
Its main difficulty is to guess the prior  $p(a)$  well.

Below we consider approaches which treat the second problem...

### Further readings on Occam Razor bounds:

- [1] *Langford J.* Quantitatively Tight Sample Complexity Bounds: Ph.D. thesis. Carnegie Mellon Thesis. 2002.
- [2] *Langford J.* Tutorial on practical prediction theory for classification. Journal of Machine Learning Research. 2005. Vol. 6. Pp. 273–306.

## The notion of Rademacher Complexity

$\mathcal{L}: A \times X \rightarrow [-1, +1]$  — the real-valued bounded loss function

$\mathcal{L}(a, x)$  — the loss of a classifier  $a$  at the object  $x$

$\mathbf{a} = (\mathcal{L}(a, x_1), \dots, \mathcal{L}(a, x_\ell))$  — loss vector of a classifier  $a$ .

**Definition 1.** *Local Rademacher complexity* of the set  $A$  on  $X^\ell$

$$\mathcal{R}(A, X^\ell) = \mathbb{E}_\sigma \sup_{a \in A} \left| \frac{1}{\ell} \sum_{i=1}^{\ell} \sigma_i \mathcal{L}(a, x_i) \right|,$$

where  $\sigma_1, \dots, \sigma_\ell$  are independent Rademacher random variables, i.e.  $P(\sigma_i = -1) = P(\sigma_i = +1) = \frac{1}{2}$ .

**Interpretation:** If for any noise vector  $(\sigma_1, \dots, \sigma_\ell)$  one can find in  $A$  a highly covariated loss vector, then the set  $A$  is complex.

**Definition 2.** *Rademacher complexity* of the set  $A$ :

$$\mathcal{R}(A) = \mathbb{E}_X \mathcal{R}(A, X^\ell)$$

## Generalization bound via Rademacher Complexity

$\tilde{P}(a) = E\mathcal{L}(a, x)$  — expected loss

$\tilde{v}(a, X^\ell) = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}(a, x_i)$  — empirical loss

### Theorem

*With probability at least  $1 - \eta$  for all  $a \in A$*

$$\begin{aligned} \tilde{P}(a) &\leq \tilde{v}(a, X) + 2\mathcal{R}(A) + \sqrt{\frac{1}{2\ell} \ln \frac{2}{\eta}} \\ &\leq \tilde{v}(a, X) + 2\mathcal{R}(A, X) + 3\sqrt{\frac{1}{2\ell} \ln \frac{2}{\eta}} \end{aligned}$$

## Most important properties of Rademacher Complexity

- 1 Relationship with the growth function:

$$\mathcal{R}(A) \leq \sqrt{\frac{2}{\ell} \ln \Delta^A(\ell)}$$

- 2 For any sets of classifiers  $A$ ,  $B$  and any constant  $c \in \mathbb{R}$

$$\mathcal{R}(A \cup B) \leq \mathcal{R}(A) + \mathcal{R}(B);$$

$$\mathcal{R}(c \cdot A) = |c| \cdot \mathcal{R}(A), \quad c \cdot A = \{c\mathbf{a} : \mathbf{a} \in A\};$$

$$\mathcal{R}(A \oplus B) \leq \mathcal{R}(A) + \mathcal{R}(B), \quad A \oplus B = \{\mathbf{a} + \mathbf{b} : \mathbf{a} \in A, \mathbf{b} \in B\};$$

- 3 The convex hull of the set of loss vectors  $A$  has the same Rademacher complexity as  $A$ :

$$\mathcal{R}\left\{\sum_{\mathbf{a} \in A} c_{\mathbf{a}} \mathbf{a} : \sum_{\mathbf{a} \in A} |c_{\mathbf{a}}| \leq 1\right\} = \mathcal{R}(A).$$



## Properties of Rademacher Complexity

Rademacher Complexity being defined via covariance has a lot of convenient algebraical properties.

Due to this fact Rademacher Complexity can be estimated for nontrivial and practically useful sets of classifiers.

Below we consider two of them: kernel machines and boosting.

### Further readings on Rademacher Complexity

[1] *Bartlett P. L., Mendelson S.* Rademacher and Gaussian complexities: risk bounds and structural results. JMLR. 2002 No 3. Pp. 463–482.

[2] *Bartlett P., Bousquet O., Mendelson S.* Local rademacher complexities. Vol. 33. Institute of Mathematical Statistics, 2005. Pp. 1497–1537.

[3] *Boucheron S., Bousquet O., Lugosi G.* Theory of classification: A survey of some recent advances. ESAIM: Probability and Statistics. 2005. No. 9. Pp. 323–375.

## Continuous approximations of threshold loss function

Consider a two-class classification problem  $Y = \{-1, +1\}$  and a set  $A$  of classifiers  $a(x, w) = \text{sign } f(x, w)$ .

For linear classifier in  $n$ -dimensional object space  $X = \mathbb{R}^n$ :

$$f(x, w) = w_1 x^1 + \dots + w_n x^n = \langle x_i, w \rangle.$$

**Definition.** *Margin* of the object  $x_i$  with a class label  $y_i \in \{-1, 1\}$

$$M_i(w) = y_i f(x_i, w)$$

$M_i(w) < 0 \iff$  classifier  $a(x, w)$  makes an error on  $x_i$ .

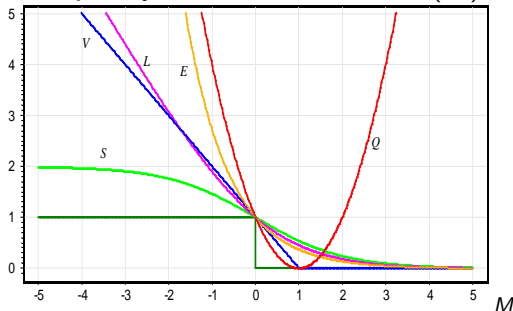
From Empirical Risk Minimization to *Approximated ERM*:

$$\nu(w, X^\ell) = \sum_{i=1}^{\ell} [M_i(w) < 0] \leq \tilde{\nu}(w, X^\ell) = \sum_{i=1}^{\ell} \mathcal{L}(M_i(w)) \rightarrow \min_w$$

*loss function*  $\mathcal{L}(M)$  is continuous, nonincreasing, nonnegative.

## Continuous approximations of threshold loss function

Frequently used loss functions  $\mathcal{L}(M)$ :



- |                             |                          |
|-----------------------------|--------------------------|
| $Q(M) = (1 - M)^2$          | — Fisher's Discriminant; |
| $V(M) = (1 - M)_+$          | — SVM;                   |
| $S(M) = 2(1 + e^M)^{-1}$    | — sigmoidal ANN;         |
| $L(M) = \log_2(1 + e^{-M})$ | — Logistic Regression;   |
| $E(M) = e^{-M}$             | — AdaBoost.              |

## Approximation and Regularization of the Empirical Risk

Many practical learning algorithms are based on both Approximation and Regularization of the Empirical Risk, e. g.

$$\tilde{v}(w, X^\ell) + C\|w\|^2 \rightarrow \min_w;$$

This can be justified from generalization bounds.

### Theorem

Let  $A$  be a set of linear classifiers, loss function is bounded  $[M < 0] \leq \mathcal{L}(M) \leq \mathcal{L}_{\max}$  and has a Lipschitz constant  $\lambda$ . Then with probability at least  $1 - \eta$  for all  $a \in A$

$$P(a) \leq \tilde{v}(w, X^\ell) + 2\lambda\mathcal{R}(A, X^\ell) + \mathcal{L}_{\max} \sqrt{\frac{2}{\ell} \ln \frac{1}{\eta}}.$$

## Rademacher Complexity bound for Kernel Machines

Consider a kernel based linear classifier

$$a(x, w) = \text{sign} \left( \sum_{i=1}^{\ell} w_i K(x_i, x) - w_0 \right),$$

### Theorem

If  $w$  is bounded in a sense of the norm

$$\|w\|_K^2 = \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} w_i w_j K(x_i, x_j) \leq B^2, \text{ then}$$

$$\mathcal{R}(A, X^{\ell}) \leq \frac{2B}{\ell} \sqrt{\sum_{i=1}^{\ell} K(x_i, x_i)}.$$

**Interpretation:**  $\|w\|_K \leq B$  is a data-dependent regularization.  
Learning algorithm: minimize  $\|w\|_K$  until  $\tilde{v}(a, X^{\ell})$  grow.

## Rademacher Complexity bound for Kernel Machines

Learning algorithm: minimize  $\|w\|_K$  until  $\tilde{v}(a, X^\ell)$  grow.

**Learning algorithm as an optimization problem:**  
 data-dependent and kernel-dependent regularization:

$$\underbrace{\sum_{i=1}^{\ell} \mathcal{L}(M_i(w))}_{\text{approximated ERM}} + \underbrace{\frac{2}{\ell} \sqrt{\sum_{i=1}^{\ell} \sum_{j=1}^{\ell} w_i w_j K(x_i, x_j)} \sqrt{\sum_{i=1}^{\ell} K(x_i, x_i)}}_{\text{regularization term (complexity penalty)}} \rightarrow \min_w$$

**Note:** regularization term may be nontrivial...  
 in contrast with usual  $\|w\|^2$  or  $L_p$ -norms

## Weighted voting of classifiers

Weighted voting of classifiers (*boosting*, *bagging*, etc.):

$$a(x) = \text{sign} \sum_{t=1}^T w_t b_t(x), \quad w_t \geq 0,$$

where  $b_t(x)$  are *base classifiers* of VC-dimension  $h$ .

$b_t(x)$  can be learned independently (bagging) or subsequently (boosting). It is no matter for generalization!

From the property  $\mathcal{R}(\text{conv}A) = \mathcal{R}(A)$  one obtain

$$\mathcal{R}(\text{conv}A) \leq \sqrt{\frac{2h}{\ell} \ln \frac{\ell e}{h}},$$

where  $h$  is VC-dimension of the set of base classifiers.

## Rademacher Complexity bound for weighted voting

### Theorem

For any  $a$  with probability at least  $1 - \eta$

$$P(a) \leq \sum_{i=1}^{\ell} \mathcal{L}(M_i(w)) + 2\lambda \sqrt{\frac{2h}{\ell} \ln \frac{\ell e}{h}} + \mathcal{L}_{\max} \sqrt{\frac{2}{\ell} \ln \frac{1}{\eta}}.$$

Conclusions for weighted voting learning algorithms:

- a great variety of loss functions  $\mathcal{L}(M)$  can be used;
- generalization of weighted voting does not depend on  $T$ ;
- boosting maximizes margins  $M_i$  effectively, then minimizing the first term of the bound;
- one can use very simple base classifiers



## Further reading

### Further readings on margin-based generalization bounds

[1] *Koltchinskii V., Panchenko D.* Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*. 2002. Vol. 30, No. 1. Pp 1–50.

[2] *Boucheron S., Bousquet O., Lugosi G.* Theory of classification: A survey of some recent advances. *ESAIM: Probability and Statistics*. 2005. No. 9. Pp. 323–375.

## Conclusion (Part I)

- Generalization bounds give optimization problems to construct learning algorithms with better performance.
- Typically, this is ERM Approximation + Regularization.
- The better performance is not always successfully attained because of the looseness of the bounds.
- There are two reasons for the looseness:
  - nonuniform splitting of the set of classifiers;
  - similarity of classifiers.
- None of recent generalization bounds can treat both problems.

**To be continued in 20 minutes...**

We will consider a *combinatorial approach*, the first in Learning Theory that takes into account **both splitting and similarity** of a classifier set and sometimes gives **exact generalization bounds**.

## Questions?

Konstantin Vorontsov  
[vokov@forecsys.ru](mailto:vokov@forecsys.ru)  
<http://www.ccas.ru/voron>

[www.MachineLearning.ru/wiki](http://www.MachineLearning.ru/wiki) (in Russian):

- Участник:Vokov
- Слабая вероятностная аксиоматика
- Расслоение и сходство алгоритмов (виртуальный семинар)