

Combinatorial Approach to Generalization Bounds Tightening

Konstantin Vorontsov
(voron@ccas.ru, www.ccas.ru/voron)

Computing Centre of Russian Academy of Sciences,
Vavilova 40, 119991, Moscow, Russian Federation

7th Open German/Russian Workshop (OGRW-7)
on Pattern Recognition and Image Understanding
August 20–23, 2007
Ettlingen, Germany

Outline

- 1 Theory of Empirical Prediction**
 - Weak Probability Axioms and Empirical Prediction
 - Example: Transductive Form of the Law of Large Numbers
 - Discussion
- 2 Theory of Generalization Ability**
 - Classical Generalization Bounds
 - Data-Dependent Bounds
 - Measuring Effective Local Shatter
- 3 Experiments with Rule Induction System**
 - The Rule Induction Classifier
 - Experimental results: shatter coefficients
 - Experimental estimation of rules overfitting

└ Outline

Outline

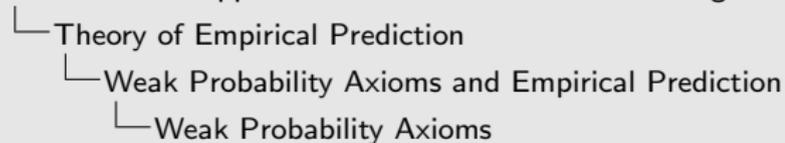
- 1 Theory of Empirical Prediction
 - Weak Probability Axioms and Empirical Prediction
 - Example: Transductive Form of the Law of Large Numbers
 - Discussion
- 2 Theory of Generalization Ability
 - Classical Generalization Bounds
 - Data-Dependent Bounds
 - Measuring Effective Local Shatter
- 3 Experiments with Rule Induction System
 - The Rule Induction Classifier
 - Experimental results: shatter coefficients
 - Experimental estimation of rules overfitting

Dear Colleagues, I shall speak about one of the most challenging problem in computational learning theory — the problem of generalization ability of learning algorithms. I shall start from a general theoretical framework, then consider a fundamental problem of learning theory — the looseness of generalization bounds. I shall finish by some empirical results.

Weak Probability Axioms

- 1 $X^L = \{x_i\}_{i=1}^L$ — a given finite subset from a set of objects \mathbb{X} .
- 2 All partitions $X^L = X_n^\ell \cup X_n^k$, $n = 1, \dots, N$,
where $N = C_L^k$, $L = \ell + k$, are equally probable.

Then...



- $X^L = \{x_i\}_{i=1}^L$ — a given finite subset from a set of objects X .
 - All partitions $X^L = X_1^L \cup X_2^L$, $n = 1, \dots, N$, where $N = C_L^k$, $L = \ell + k$, are equally probable.
- Then...

Let us start from two very simple axioms. First, [see. . .] we assume that in data analysis one can observe only a finite set of objects X^L . The set can be unknown, but it never can be infinite.

Second, [see. . .] we assume that objects appear at random, so that all partitions of our set into two subsets have equal chances to realize. The order of objects is the only source of randomness in our framework.

Weak Probability Axioms

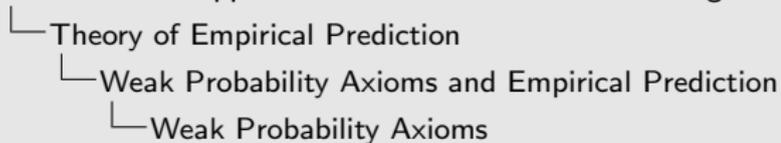
- 1 $X^L = \{x_i\}_{i=1}^L$ — a given finite subset from a set of objects \mathbb{X} .
- 2 All partitions $X^L = X_n^\ell \cup X_n^k$, $n = 1, \dots, N$, where $N = C_L^k$, $L = \ell + k$, are equally probable.

Then

- Consider an event A as a function $A: \{1, \dots, N\} \rightarrow \{0, 1\}$
- The fraction of partitions $n: A(n) = 1$ can be interpreted as probability or expectation:

$$P_n A(n) \equiv E_n A(n) = \frac{1}{N} \sum_{i=1}^N A(n).$$

Here “probability” P_n is simply averaging operator $\frac{1}{N} \sum_{i=1}^N$.



- $X^1 = \{x\}_{i=1}^N$ — a given finite subset from a set of objects X .
- All partitions $X^1 = X_1^1 \cup X_2^1, n = 1, \dots, N$, where $N = C_1^N, L = \ell + k$, are equally probable.

Then

- Consider an event A as a function $A: \{1, \dots, N\} \rightarrow \{0, 1\}$
- The fraction of partitions $n: A(n) = 1$ can be interpreted as probability or expectation:

$$P_n A(n) = E_n A(n) = \frac{1}{N} \sum_{i=1}^N A(n).$$

Here "probability" P_n is simply averaging operator $\frac{1}{N} \sum_{i=1}^N$.

Under these assumptions the probability of an event $A(n)$ [see. . .] is defined as the fraction of partitions n [see. . .] for which $A(n)$ is true. Note that there is no difference between expectation [see. . .] and probability [see. . .] in this framework.

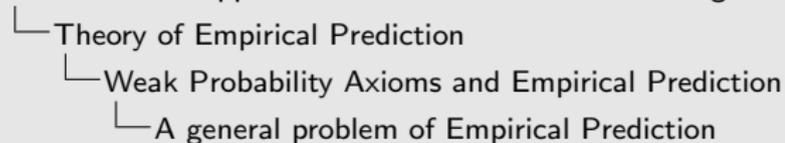
A general problem of Empirical Prediction

- Empirical framework:
 - some partition (X_n^ℓ, X_n^k) realizes, $n \in \{1, \dots, N\}$;
 - subsample X_n^ℓ is *observable*,
 - subsample X_n^k is *hidden*.
- Given a function $T: \mathbb{X}^k \times \mathbb{X}^\ell \rightarrow R$
 1. Chose a function $\hat{T}: \mathbb{X}^\ell \rightarrow R$ such that the value $\hat{T}_n = \hat{T}(X_n^\ell)$ predicts the value $T_n = T(X_n^k, X_n^\ell)$.
 2. Estimate prediction accuracy (obtain an upper bound):

$$P_n[d(\hat{T}_n, T_n) > \epsilon] \leq \eta(\epsilon),$$

where $d(\hat{r}, r)$ — discrepancy function, e.g. $d(\hat{r}, r) = |\hat{r} - r|$.

Combinatorial Approach to Generalization Bounds Tightening



- Empirical framework:
 - some partition $\{X_n^1, X_n^k\}$ realizes, $n \in \{1, \dots, N\}$;
 - subsample X_n^1 is observable;
 - subsample X_n^k is hidden.
- Given a function $T: X^k \times X^1 \rightarrow R$
 1. Choose a function $\hat{T}: X^1 \rightarrow R$ such that the value $\hat{T}_n = \hat{T}(X_n^1)$ predicts the value $T_n = T(X_n^1, X_n^k)$.
 2. Estimate prediction accuracy (obtain an upper bound):

$$P_n[d(\hat{T}_n, T_n) > \epsilon] \leq \eta(\epsilon),$$
 where $d(\hat{T}, T) = \text{discrepancy function, e.g. } d(\hat{T}, T) = |\hat{T} - T|$.

The problem of empirical prediction arise when some of the equiprobable partitions [see. . .] realizes, but one observe only a first subset X_n^1 [see. . .] whereas X_n^k [see. . .] remains hidden. We want to predict a value of a given function T [see. . .] that depends on both parts, having an information \hat{T} [see. . .] computed from the observed part only. Also we want to know in advance how accurate our prediction can be [see. . .] . Then, the problem is to upper bound the fraction of partitions for which our prediction fails. To formulate exactly what means “fails” we introduce a discrepancy function d [see. . .] (that can be simply a difference in most cases) and a threshold of exactness ϵ [see. . .] .

Why we call this framework “Weak Probability Axioms”?

Insight:

Weak PA leads to “simplified Probability Theory”

Indeed...

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening

- └ Theory of Empirical Prediction
 - └ Weak Probability Axioms and Empirical Prediction
 - └ Why we call this framework “Weak Probability Axioms”?

Why we call this framework “Weak Probability Axioms”?

Insight:

Weak PA leads to “Simplified Probability Theory”

Indeed...

Why we call this framework “the weak probability axioms” and even the weak Probability Theory?

Why we call this framework “Weak Probability Axioms”?

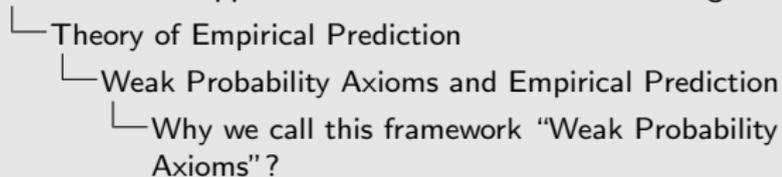
Insight:

Weak PA leads to “simplified Probability Theory”

Indeed. . .

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonenkis generalization bounds (see later);
 - etc.

Combinatorial Approach to Generalization Bounds Tightening



Insight:

Weak PA leads to “simplified Probability Theory”

Indeed ...

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonenkis generalization bounds ([see later](#));
 - etc.

First, because it is very general. Empirical Prediction is one of the central problems in Probability Theory, Statistics and Learning Theory. In practice predictions are interesting only for finite sets of objects, because nobody can observe an infinite set of objects. Empirical Prediction is transductive by its nature. Many fundamental results in Statistics and Learning Theory can be reformulated in transductive form.

Why we call this framework “Weak Probability Axioms”?

Insight:

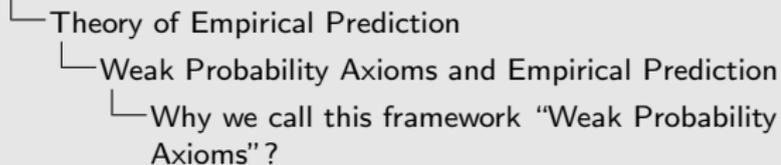
Weak PA leads to “simplified Probability Theory”

Indeed. . .

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonenkis generalization bounds (see later);
 - etc.
- Weak PA is based on less restrictive assumptions:
 - no need of probability measure on \mathbb{X} ;
 - no need of frequentist definition of probability via $L \rightarrow \infty$;

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening



Why we call this framework “Weak Probability Axioms”?

Insight:

Weak PA leads to “simplified Probability Theory”

Indeed...

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonemkis generalization bounds (see later);
 - etc.
- Weak PA is based on less restrictive assumptions:
 - no need of probability measure on \mathcal{X} ;
 - no need of frequentist definition of probability via $L \rightarrow \infty$;

On the other hand, Weak Probability Axioms are less restrictive if compared with classical Kolmogorov's Axioms. Here we don't need of probability measure on object space and we don't define a probability via passage to the limit.

Why we call this framework “Weak Probability Axioms”?

Insight:

Weak PA leads to “simplified Probability Theory”

Indeed. . .

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonenkis generalization bounds (see later);
 - etc.
- Weak PA is based on less restrictive assumptions:
 - no need of probability measure on \mathbb{X} ;
 - no need of frequentist definition of probability via $L \rightarrow \infty$;
 - **no need of “probability” at all!**

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening

- └ Theory of Empirical Prediction
 - └ Weak Probability Axioms and Empirical Prediction
 - └ Why we call this framework “Weak Probability Axioms”?

Why we call this framework “Weak Probability Axioms”?

Insight:
Weak PA leads to “simplified Probability Theory”

Indeed...

- Weak PA is sufficient to prove fundamental facts:
 - the Law of Large Numbers, with exact convergence rate;
 - Kolmogorov-Smirnov criterion, also exact;
 - many statistical hypothesis tests (order statistics etc.);
 - Vapnik-Chervonemkis generalization bounds ([see later](#));
 - etc.
- Weak PA is based on less restrictive assumptions:
 - no need of probability measure on \mathcal{X} ;
 - no need of frequentist definition of probability via $L \rightarrow \infty$;
 - **no need of “probability” at all!**

Frankly speaking, we don't need of the notion of probability at all. In our framework “probability” is no more than a synonym of “fraction of partitions”.

Example: transductive form of the Law of Large Numbers

Def. The *frequency* of $S \subset \mathbb{X}$ in a finite sample $U \subset \mathbb{X}$:

$$\nu_S(U) = \frac{1}{|U|} \sum_{u \in U} [u \in S].$$

Theorem (common knowledge)

The frequency $\nu_S(X_n^\ell)$ predicts the frequency $\nu_S(X_n^k)$.
Prediction accuracy is given by an **exact** bound

$$P_n[\nu_S(X_n^k) - \nu_S(X_n^\ell) \geq \epsilon] = H_{(L_m)}^{\ell s(\epsilon)},$$

where $H_{(L_m)}^{\ell s}$ is a tail of hypergeometric distribution,
 $s(\epsilon) = \lfloor \frac{\ell}{L}(m - \epsilon k) \rfloor$, $m = L\nu_S(X^L)$.

Remark. Here $\hat{T}(U) = T(U) = \nu_S(U)$.

Combinatorial Approach to Generalization Bounds Tightening

- Theory of Empirical Prediction

- Example: Transductive Form of the Law of Large Numbers

- Example: transductive form of the Law of Large Numbers

Example: transductive form of the Law of Large Numbers

Def. The frequency of $S \subset \mathcal{X}$ in a finite sample $U \subset \mathcal{X}$:

$$\nu_S(U) = \frac{1}{|U|} \sum_{u \in U} \mathbb{1}_{\{u \in S\}}$$

Theorem (common knowledge)The frequency $\nu_S(X_i^*)$ predicts the frequency $\nu_S(X_i^*)$. Prediction accuracy is given by an **exact** bound

$$P_{\mathcal{X}}[\nu_S(X_i^*) - \nu_S(X_i) \geq \epsilon] = H\left(\frac{4\epsilon^2}{m}\right),$$

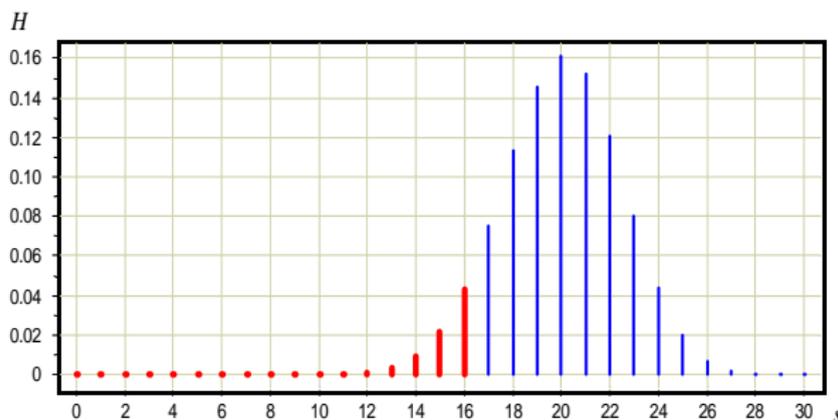
where $H\left(\frac{4\epsilon^2}{m}\right)$ is a tail of hypergeometric distribution, $s(\epsilon) = \lfloor \frac{4}{\epsilon}(m - \epsilon k) \rfloor$, $m = L_{\nu_S}(X^*)$.**Remark.** Here $\tilde{T}(U) = T(U) = \nu_S(U)$.

Let us consider the Law of Large Numbers as an example. It is commonly known that the frequency of an event S in the hidden sample [see. . .] can be predicted by its frequency in the observed sample [see. . .] . In classical Probability Theory one predict not frequency but a probability of the event S and one use inequalities of Hoeffding, Bernstein and Chernoff to give asymptotical bounds of the prediction accuracy. In our transductive framework the bound is given by hypergeometric distribution [see. . .] . Remarkable than it is an exact [see. . .] non-asymptotic bound.

Example: transductive form of the Law of Large Numbers

The tail of hypergeometric distribution:

$$H(L, m, \ell, s(\epsilon)) = \sum_{t=s_0}^{s(\epsilon)} \frac{C_m^t C_{L-m}^{\ell-t}}{C_L^\ell}, \quad s_0 = \max\{0, m - k\}$$



Here $\epsilon = 0.05$, $L = 300$, $\ell = 200$, $m = 30$.

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening

└ Theory of Empirical Prediction

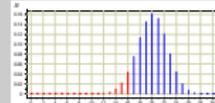
└ Example: Transductive Form of the Law of Large Numbers

└ Example: transductive form of the Law of Large Numbers

Example: transductive form of the Law of Large Numbers

The tail of hypergeometric distribution:

$$H(L, m) = \sum_{i=0}^{k-L} \frac{C_i^k C_{L-i}^{m-k}}{C_L^m}, \quad k_0 = \max\{0, m-k\}$$



Here $\epsilon = 0.05$, $L = 300$, $\ell = 200$, $m = 30$.

In this case convergence is a trivial consequence of the fact that the relative width [see . . .] of hypergeometric peak tends to zero when the sample size L tends to infinity.

Links with classical Kolmogorov's Probability Axioms

Theorem (Correspondence principle)

If we have obtained a bound under Weak PA

$$Q_\epsilon(X^L) = P_n[d(\hat{T}(X_n^\ell), T(X_n^\ell, X_n^k)) > \epsilon] \leq \eta(\epsilon, X^L)$$

*then analogous bound will be true under Strong PA
(classical Kolmogorov's PA)*

$$\mathbf{E}_{X^L} Q_\epsilon(X^L) = P_{X^L}[d(\hat{T}(X^\ell), T(X^\ell, X^k)) > \epsilon] \leq \mathbf{E}_{X^L} \eta(\epsilon, X^L)$$

- Theory of Empirical Prediction
 - Discussion
 - Links with classical Kolmogorov's Probability Axioms

Theorem (Correspondence principle)*If we have obtained a bound under Weak PA*

$$Q_n(X^i) = P_n[d(\hat{T}(X_n^i), T(X_n^i, X_n^i)) > \epsilon] \leq \eta(\epsilon, X^i)$$

*then analogous bound will be true under Strong PA
(classical Kolmogorov's PA)*

$$E_{X^i} Q_n(X^i) = P_{X^i}[d(\hat{T}(X^i), T(X^i, X^i)) > \epsilon] \leq E_{X^i} \eta(\epsilon, X^i)$$

Each bound obtained under Weak Axioms can be easily restated under classical Kolmogorov's Axioms. To do this one may take expectation [see. . .] of both sides of the inequality.

Links with classical Kolmogorov's Probability Axioms

Theorem (Correspondence principle)

If we have obtained a bound under Weak PA

$$Q_\epsilon(X^L) = P_n[d(\hat{T}(X_n^\ell), T(X_n^\ell, X_n^k)) > \epsilon] \leq \eta(\epsilon, X^L)$$

*then analogous bound will be true under Strong PA
(classical Kolmogorov's PA)*

$$E_{X^L} Q_\epsilon(X^L) = P_{X^L}[d(\hat{T}(X^\ell), T(X^\ell, X^k)) > \epsilon] \leq E_{X^L} \eta(\epsilon, X^L)$$

When a transduction becomes an induction

If $\eta(\epsilon, X^L) \equiv \eta(\epsilon)$ then the same bound is true for any sample.

- Theory of Empirical Prediction
 - Discussion
 - Links with classical Kolmogorov's Probability Axioms

Theorem (Correspondence principle)

If we have obtained a bound under Weak PA

$$Q_\epsilon(X^L) = P_{\text{in}}[d(\hat{T}(X^L), T(X^L, X^L)) > \epsilon] \leq \eta(\epsilon, X^L)$$

then analogous bound will be true under Strong PA
(classical Kolmogorov's PA)

$$E_{X^L} Q_\epsilon(X^L) = P_{\text{in}}[d(\hat{T}(X^L), T(X^L, X^L)) > \epsilon] \leq E_{X^L} \eta(\epsilon, X^L)$$

When a transduction becomes an inductionIf $\eta(\epsilon, X^L) = \eta(\epsilon)$ then the same bound is true for any sample

Philosophic remark. It is commonly to think that transduction is more restrictive than induction. This is not the case when one writes a functional in transductive form and obtains its bound that is true for any sample X^L [see...]. Really this means that transduction transforms into induction.

Links with empirical techniques

Cross-Validation

If we did not succeed to obtain a bound theoretically, or if we obtained very loose (overestimated) bound:

$$Q_\epsilon(X^L) = P_n[d(\hat{T}_n, T_n) > \epsilon] \leq \boxed{???},$$

then we can measure it empirically:

$$Q_\epsilon(X^L) \approx \frac{1}{|N'|} \sum_{n \in N'} [d(\hat{T}_n, T_n) > \epsilon]$$

where $N' \subset \{1, \dots, N\}$ is small enough to compute the sum and big enough to estimate be accurate.

- Theory of Empirical Prediction

- Discussion

- Links with empirical techniques

Cross-Validation

If we did not succeed to obtain a bound theoretically, or if we obtained very loose (overestimated) bound:

$$Q_n(X^k) = P_n[d(\bar{T}_n, T_n) > \epsilon] \leq \frac{1}{N}$$

then we can measure it empirically:

$$Q_n(X^k) \approx \frac{1}{N'} \sum_{i \in N'} [d(\bar{T}_n, T_n) > \epsilon]$$

where $N' \subset \{1, \dots, N\}$ is small enough to compute the sum and big enough to estimate be accurate.

The keystone advantage of Weak Probability Axioms. It provides the unique starting point for both theoretical and empirical consideration. If we did not succeed to obtain a bound theoretically, we can measure the prediction functional empirically, replacing the average of all partitions [see. . .] by the average of some partitions [see. . .]. This leads to the well known empirical technique — Cross Validation.

The main idea of the further part of presentation is that Cross Validation can help to understand the causes of bounds looseness.

Weak Probability Axioms: pro and con

+ Weak PA is suitable for data analysis, statistics, COLT

- └ Theory of Empirical Prediction

- └ Discussion

- └ Weak Probability Axioms: pro and con

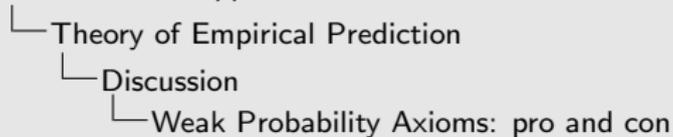
Intermediate summary. Our framework is suitable for data analysis, statistics and learning theory, where all samples are finite and all variables can be calculated from data. In classical Probability Theory one operates with hypothetical asymptotic values such that probabilities, expectations, distribution functions, etc. We intend to manage without this hypothetical values. For what reason? Because asymptotic considerations are often the cause of bound looseness and can lead to numerous misunderstandings that are very difficult to reveal.

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ...but not suitable for physics

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening



Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics

The same time I think that Weak PA will not be suitable for physics and other applications where continuity is crucial. . . so every theory may have its own limitations.

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ...but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ...represented by complicated combinatorial formula

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening

└ Theory of Empirical Prediction

└ Discussion

└ Weak Probability Axioms: pro and con

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ... represented by complicated combinatorial formula

It can give exact bounds but we are to elaborate fast algorithms to calculate them effectively.

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ... represented by complicated combinatorial formula
- + Weak PA satisfies a “correspondence principle”
- ... but not all theorems in classical Probability Theory have analogs in Weak PA

└ Theory of Empirical Prediction

└ Discussion

└ Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ... represented by complicated combinatorial formula
- + Weak PA satisfies a "correspondence principle"
- ... but not all theorems in classical Probability Theory have analogs in Weak PA

Each bound obtained under Weak PA can be restated under classical PA.
But many measure-specific theorems can not be transferred to the weak form.

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ... represented by complicated combinatorial formula
- + Weak PA satisfies a “correspondence principle”
- ... but not all theorems in classical Probability Theory have analogs in Weak PA

Open problem:

What part of mathematical statistics can be restated in Weak PA?

2007-08-31

Combinatorial Approach to Generalization Bounds Tightening

└ Theory of Empirical Prediction

└ Discussion

└ Weak Probability Axioms: pro and con

Weak Probability Axioms: pro and con

- + Weak PA is suitable for data analysis, statistics, COLT
- ... but not suitable for physics
- + Weak PA gives non-asymptotic, exact bounds
- ... represented by complicated combinatorial formula
- + Weak PA satisfies a "correspondence principle"
- ... but not all theorems in classical Probability Theory have analogs in Weak PA

Open problem:

What part of mathematical statistics can be restated in Weak PA?

A big open problem is "what part of mathematical statistics can be restated in Weak PA?" My opinion is that a sufficiently big part.

The learning problem (classification, regression, etc.)

- \mathbb{X} — set of objects, \mathbb{Y} — set of outputs
- Binary *loss function* $\mathcal{L}: \mathbb{X} \times \mathbb{Y} \rightarrow \{0, 1\}$
 - In classification: $\mathcal{L}(x, y) = [y \neq y^*(x)]$,
 - In regression: $\mathcal{L}(x, y) = [|y - y^*(x)| > \delta]$
where $y^*(x)$ — unknown *target function*
- *Training set* $X^\ell = \{x_i\}_{i=1}^\ell \subset \mathbb{X}$ with known losses $\mathcal{L}(x_i, y)$
- *Learning algorithm*
 $\mu: X^\ell \mapsto f$, where $f: \mathbb{X} \rightarrow \mathbb{Y}$ — some function
- *Average error* of a function $f: \mathbb{X} \rightarrow \mathbb{Y}$ on a set $U \subset \mathbb{X}$
$$\nu(f, U) = \frac{1}{|U|} \sum_{u \in U} \mathcal{L}(u, f(u))$$
- *Generalization ability*:
 $\nu(\mu X^\ell, U)$ must be **sufficiently** small for **most** $U \in \mathbb{X}^*$

Combinatorial Approach to Generalization Bounds Tightening

Theory of Generalization Ability

Classical Generalization Bounds

The learning problem (classification, regression, etc.)

The learning problem (classification, regression, etc.)

- X — set of objects, Y — set of outputs
- Binary loss function $\mathcal{L}: X \times Y \rightarrow \{0, 1\}$
 - in classification: $\mathcal{L}(x, y) = [y \neq y^*(x)]$,
 - in regression: $\mathcal{L}(x, y) = [|y - y^*(x)| > \epsilon]$ where $y^*(x)$ — unknown target function
- Training set $X^\ell = \{x_i\}_{i=1}^{\ell} \subset X$ with known losses $\mathcal{L}(x_i, y)$
- Learning algorithm
 - $\mu: X^\ell \rightarrow f$, where $f: X \rightarrow Y$ — some function
- Average error of a function $f: X \rightarrow Y$ on a set $U \subset X$

$$\nu(f, U) = \frac{1}{|U|} \sum_{u \in U} \mathcal{L}(u, f(u))$$
- Generalization ability:
 - $\nu(\mu X^\ell, U)$ must be **sufficiently** small for **most** $U \subset X$

Now I pass to the second part of my presentation — the Learning Problem. Given a training set X^ℓ [see. . .] one must learn a function f [see. . .] that approximates the unknown target function $y^*(x)$ [see. . .] as well as possible. The approximation quality on a finite sample U is measured by the average error $\nu(f, U)$ [see. . .] also called the frequency of errors.

The most challenging problem in Learning Theory — how to guarantee that the learned function f will have a small frequency of errors [see. . .] out of the training set.

Classical bounds

- Vapnik and Chervonenkis (1974):

$$P_{\epsilon}(F) = P_{X^L} \left[\sup_{f \in F} (\nu(f, X^k) - \nu(f, X^{\ell})) > \epsilon \right] \\ \leq \Delta^F(L) \cdot 1.5 e^{-\epsilon^2 \ell}; \quad (\text{if } \ell = k)$$

where F is the entire functions set (search space);
 $\Delta^F(L)$ — *Global Shatter Coefficient* of the set F ,
the number of functions f from F distinguishable on X^L ;
 $\Delta^F(L) \leq 1.5 \frac{L^h}{h!}$, $h =$ VC dimension (growth function) of F .

◦ Vapnik and Chervonenkis (1974):

$$P(F) = P_{X^L} \left[\sup_{f \in F} |v(f, X^L) - v(f, X^h)| > \epsilon \right] \\ \leq \Delta^F(L) \cdot 1.5 e^{-2\epsilon^2}, \quad (h \ell = k)$$

where F is the entire functions set (search space).
 $\Delta^F(L)$ — Global Shatter Coefficient of the set F ,
 the number of functions f from F distinguishable on X^L ;
 $\Delta^F(L) \leq 1.5 \frac{h}{L}$, h — VC dimension (growth function) of F .

The classical generalization bound is a product of two terms: complexity term called *shatter coefficient* [see. . .] and convergence term [see. . .] that tends to zero when the sample length tends to infinity.

Let us remind that shatter coefficient of the functions set A is defined as the maximal number of functions from A pairwise indistinguishable on a set X^L .

If the set X^L is fixed then the shatter coefficient is called *local*.

If the set X^L is arbitrary then the shatter coefficient depends on sample size L only (and not on concrete objects) and is called *global* [see. . .] .

In classic VC theory only global shatter coefficient was used.

Classical bounds

- Vapnik and Chervonenkis (1974):

$$P_{\epsilon}(F) = P_{X^L} \left[\sup_{f \in F} (\nu(f, X^k) - \nu(f, X^{\ell})) > \epsilon \right] \\ \leq \Delta^F(L) \cdot 1.5 e^{-\epsilon^2 \ell}; \quad (\text{if } \ell = k)$$

where F is the entire functions set (search space);

$\Delta^F(L)$ — *Global Shatter Coefficient* of the set F ,

the number of functions f from F distinguishable on X^L ;

$\Delta^F(L) \leq 1.5 \frac{L^h}{h!}$, h = VC dimension (growth function) of F .

- The bound is extremely loose, as it doesn't take into account:
 - the distribution of objects X^{ℓ} ;
 - the target $y^*(x)$;
 - the learning algorithm μ ;
 - $1.5 e^{-\epsilon^2 \ell}$ is an asymptotic approximation.

Theory of Generalization Ability
 Classical Generalization Bounds
 Classical bounds

- Vapnik and Chervonenkis (1974):

$$\begin{aligned}
 P(F) = P_{X^{\ell}} \left[\sup_{f \in F} (v(f, X^{\ell}) - v(f, X^{\infty})) > \epsilon \right] \\
 \leq \Delta^F(\ell) \cdot 1.5 e^{-\epsilon^2 \ell}; \quad (\forall \ell = k)
 \end{aligned}$$

where F is the entire functions set (search space).

$\Delta^F(\ell)$ = Global Shatter Coefficient of the set F .

the number of functions f from F distinguishable on X^{ℓ} ;

$\Delta^F(\ell) \leq 1.5 \frac{h}{\ell}$, h = VC dimension (growth function) of F .

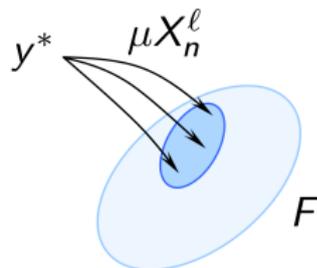
- The bound is extremely loose, as it doesn't take into account:
 - the distribution of objects X^{ℓ} ;
 - the target $y^{\ell}(\cdot)$;
 - the learning algorithm μ ;
 - $1.5 e^{-\epsilon^2 \ell}$ is an asymptotic approximation.

The main shortcoming of this bound [see. . .] is that the uniform convergence taken in VC theory as an axiom results in extremely overestimated complexity term [see. . .] . Taking supremum [see. . .] one neglect many important peculiarities [see. . .] of the given task.

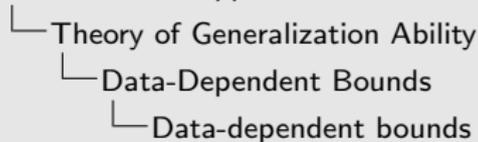
Data-dependent bounds

- *Localization effect:*

If target y^* , learning algorithm μ , and sample X^L are fixed then not all functions from F can be obtained.



- Uniform convergence bound (Vapnik & Chervonenkis, 1969)
- Theory of learnable (Valiant, 1982)
- First data-dependent bound (?)
- The most tight general VC-like bound (M. Talagrand)
- Self-bounding (Freund)
- Algorithmic luckiness ()
- Computationally tight sample complexity bounds [J. Langford]



- o Localization effect:
If target y^* , learning algorithm μ , and sample X^n are fixed then not all functions from F can be obtained.
- o Uniform convergence bound (Vapnik & Chervonenkis, 1969)
- o Theory of learnable (Valiant, 1982)
- o First data-dependent bound (?)
- o The most tight general VC-like bound (M. Talagrand)
- o Self-bounding (Freund)
- o Algorithmic luckiness ()
- o Computationally tight sample complexity bounds [J. Langford]



When a particular task is fixed, only a little part of functions [see. . .] can be obtained. This localization effect was understood long ago and several frameworks for data-dependent bound was proposed.

...

Our proposition distinguishes by the total change of Probability Axioms.

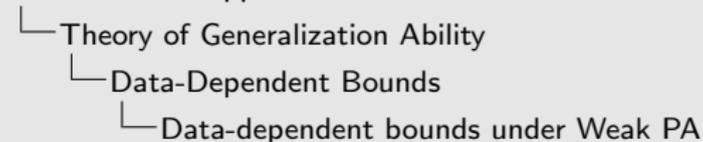
Data-dependent bounds under Weak PA

- Weak PA:

$$\begin{aligned} Q_\epsilon(\mu, X^L) &= P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \epsilon \right] \\ &\leq \Delta_L^\ell(\mu, X^L) \cdot \max_m H_{Lm}^{\ell s(\epsilon)} \\ &(\leq \Delta^F(L) \cdot 1.5 e^{-\epsilon^2 \ell}); \end{aligned}$$

where $f_n = \mu X_n^\ell$ is a result of learning;

$\Delta_L^\ell(\mu, X^L)$ — *Local Shatter Coefficient* of the set of functions obtainable by learning: $\{f_n \mid n = 1, \dots, N\}$.



Weak PA:

$$\begin{aligned}
 Q(\mu, X^N) &= P_{\mu}[\sigma(\ell_n, X^N) - \sigma(\ell_n, X^N) > \epsilon] \\
 &\leq \Delta_{\mu}^{\sigma}(\mu, X^N) \cdot \max_{\mu} H(\ell_n, \mu) \\
 &(\leq \Delta^{\sigma}(\epsilon) \cdot 1.5 \epsilon^{-2})
 \end{aligned}$$

where $\ell_n = \mu X_n^N$ is a result of learning.
 $\Delta_{\mu}^{\sigma}(\mu, X^N)$ — Local Shatter Coefficient of the set of functions obtainable by learning: $\{\ell_n \mid n = 1, \dots, N\}$.

The bound can be obtained under Weak PA [see. . .] . We will see later that this bound is still very loose. The only advantage of this bound is that this functional [see. . .] can be measured effectively to understand the causes of bound looseness.

Note that the older data-independent VC bound can be derived from this one.

The next idea is to eliminate maximum on m [see. . .] in convergence term.

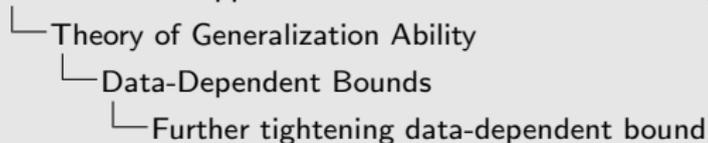
Further tightening data-dependent bound

- **Idea:** the scalar complexity contains too small information about the learning process.
- Splitting the local shatter coeff into *Shatter Profile* $\{D_m\}_{m=0}^L$:

$$\Delta_L^\ell(\mu, X^L) = \sum_{m=1}^L D_m(\mu, X^L)$$

$D_m(\mu, X^L)$ — local shatter coefficient of the set of functions having m errors on X^L : $\{f_n \mid \nu(f_n, X^L) = \frac{m}{L}, n = 1, \dots, N\}$.

Combinatorial Approach to Generalization Bounds Tightening



Further tightening data-dependent bound

- Idea: the scalar complexity contains too small information about the learning process.
- Splitting the local shatter coeff into Shatter Profile $\{D_m\}_{m=0}^L$:

$$\Delta_L(\mu, X^L) = \sum_{m=0}^L D_m(\mu, X^L)$$

$D_m(\mu, X^L)$ — local shatter coefficient of the set of functions having m errors on X^L : $\{f_n \mid v(f_n, X^L) = m, n = 1, \dots, N\}$.

This idea has yet another interpretation. One scalar characteristic of complexity contains too small information about learning process. It's not sufficient to know, how many different functions can be obtained as a result of learning. Also it is worth to take into account how many functions of a given quality can be obtained.

For this reason we split the local shatter coefficient into $L + 1$ components [see. . .] . Each component can be considered as a local shatter coefficient of the set of functions that make exactly m errors on the full sample X^L . So we obtain the non-scalar characteristic of complexity that we call *Shatter Profile* [see. . .] .

Further tightening data-dependent bound

- **Idea:** the scalar complexity contains too small information about the learning process.
- Splitting the local shatter coeff into *Shatter Profile* $\{D_m\}_{m=0}^L$:

$$\Delta_L^\ell(\mu, X^L) = \sum_{m=1}^L D_m(\mu, X^L)$$

$D_m(\mu, X^L)$ — local shatter coefficient of the set of functions having m errors on X^L : $\{f_n \mid \nu(f_n, X^L) = \frac{m}{L}, n = 1, \dots, N\}$.

- Weak PA:

$$Q_\epsilon(\mu, X^L) \leq \sum_{m=1}^L D_m(\mu, X^L) \cdot H\left(\frac{\ell}{L} \frac{s(\epsilon)}{m}\right);$$

Combinatorial Approach to Generalization Bounds Tightening

- Theory of Generalization Ability

- Data-Dependent Bounds

- Further tightening data-dependent bound

Further tightening data-dependent bound

- **Idea:** the scalar complexity contains too small information about the learning process.
- Splitting the local shatter coeff into Shatter Profile $\{D_m\}_{m=0}^L$:

$$\Delta_L(\mu, X^L) = \sum_{m=1}^L D_m(\mu, X^L)$$

$D_m(\mu, X^L)$ — local shatter coefficient of the set of functions having m errors on X^L : $\{f_n \mid v(f_n, X^L) = m, n = 1, \dots, N\}$.

- Weak PA:

$$Q(\mu, X^L) \leq \sum_{m=1}^L D_m(\mu, X^L) \cdot H\left(\frac{L}{m}\right).$$

With shatter profile we obtain a more tight bound [see...]. The previous one could be obtained from this one taken the maximum of convergence term by m [see...].

The Effective Local Shatter Profile and Coefficient

- **Ideally accurate but “unfair” bound** answers a question: what would be the shatter profile D_m to bound be exact?
- **Theorem**

$$\begin{aligned} Q_{\epsilon, m}(\mu, X^L) &= P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \epsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right] \\ &\leq D_m(\mu, X^L) \cdot H_{Lm}^{\ell s(\epsilon)}; \end{aligned}$$

Let us change “ \leq ” by “ $=$ ” and express D_m :

Combinatorial Approach to Generalization Bounds Tightening

- └ Theory of Generalization Ability

- └ Measuring Effective Local Shatter

- └ The Effective Local Shatter Profile and Coefficient

- Ideally accurate but "unfair" bound answers a question: what would be the shatter profile D_m to bound be exact?

- Theorem

$$Q_{\mu, \sigma}(m, X^m) = \mathbb{P}_\sigma \left[\left| \bar{r}(f_m, X_m^*) - \bar{r}(f_m, X_m) \right| > \sigma \right] \leq \left[\bar{r}(f_m, X^m) - \frac{\sigma}{2} \right] \leq D_m(\mu, X^m) \cdot H_{\sigma}^{\left(\frac{4\sigma^2}{m} \right)}$$

Let us change " \leq " by " $=$ " and express D_m :

To compare proposed bounds empirically we must have an ideally accurate bound as a reference point. Then we introduce a subsidiary functional Q_m [see. . .] which helps to estimate all components of the Shatter Profile separately [see. . .]. We call this bound *unfair* because usually one estimate the complexity term to get an upper bound of the quality functional. Whereas here we do a contrary thing: we estimate the quality functional via Cross-Validation [see. . .] in order to answer a question: what would be the shatter profile D_m [see. . .] to bound be exact?

The Effective Local Shatter Profile and Coefficient

- **Ideally accurate but “unfair” bound** answers a question: what would be the shatter profile D_m to bound be exact?

- **Theorem**

$$Q_{\epsilon, m}(\mu, X^L) = P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \epsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right] \\ \leq D_m(\mu, X^L) \cdot H_{Lm}^{\ell s(\epsilon)};$$

Let us change “ \leq ” by “ $=$ ” and express D_m :

- *Effective Local Shatter Profile* \hat{D}_m , $m = 0, \dots, L$:

$$\hat{D}_m = \frac{\frac{1}{|N'|} \sum_{n \in N'} \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \epsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right]}{H_{Lm}^{\ell s(\epsilon)}}.$$

- *Effective Local Shatter Coefficient*: $\hat{\Delta}_L^\ell = \hat{D}_0 + \dots + \hat{D}_L$.

Combinatorial Approach to Generalization Bounds Tightening

- Theory of Generalization Ability

- Measuring Effective Local Shatter

- The Effective Local Shatter Profile and Coefficient

The Effective Local Shatter Profile and Coefficient

- Ideally accurate but "unfair" bound answers a question: what would be the shatter profile D_m to bound be exact?

- Theorem

$$Q_{\epsilon, m}(p, X^L) = P_m \left[|v(f_m, X_m^L) - v(f_m, X_m^L) > \epsilon \right] |v(f_m, X^L) - \frac{p}{2}| \\ \leq D_m(p, X^L) \cdot H_{\epsilon}^{(L, m)}$$

Let us change " \leq " by " $=$ " and express D_m :

- Effective Local Shatter Profile \tilde{D}_m , $m = 0, \dots, L$:

$$\tilde{D}_m = \frac{1}{H_{\epsilon}^{(L, m)}} \sum_{\substack{X^L \\ |v(f_m, X_m^L) - v(f_m, X_m^L) > \epsilon}} |v(f_m, X^L) - \frac{p}{2}|$$

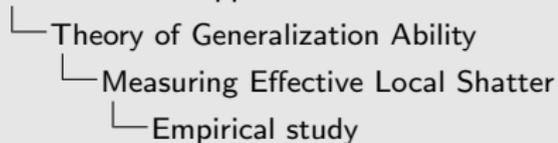
- Effective Local Shatter Coefficient: $\Delta_L^{\epsilon} = D_0 + \dots + D_L$.

We call the unfair estimate of shatter profile the *Effective Local Shatter Profile* [see...]. Then we define the *Effective Local Shatter Coefficient* as the sum of all profile components [see...].

Empirical study

- The goal of experiment — to understand:
 - how tight these bounds are?
 - what effect is more important to bound tightening?
- Bounds to be compared:
 - The worst: classical VC bound
 - Local shatter coefficient bound
 - Local shatter profile bound
 - The best: effective (unfair) local shatter profile bound
- Testing area — the Rule Induction algorithm, because:
 - Global SC is well known
 - Local SC can be easily estimated during rule search
- Testing area — 7 tasks from UCI repository

Combinatorial Approach to Generalization Bounds Tightening



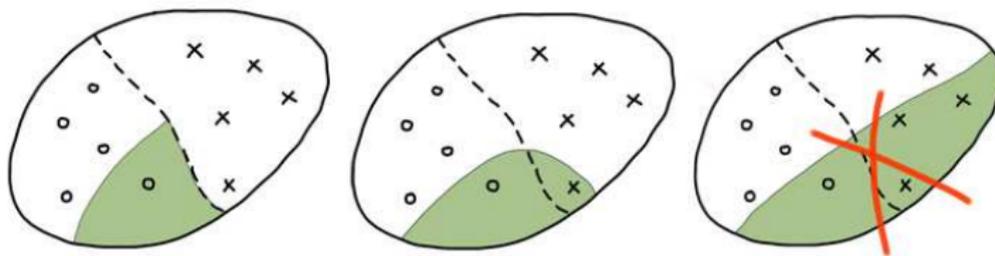
Empirical study

- The goal of experiment — to understand:
 - how tight these bounds are?
 - what effect is more important to bound tightening?
- Bounds to be compared:
 - The worst: classical VC bound
 - Local shatter coefficient bound
 - Local shatter profile bound
 - The best: effective (sofar) local shatter profile bound
- Testing area — the Rule Induction algorithm, because:
 - Global SC is well known
 - Local SC can be easily estimated during rule search
- Testing area — 7 tasks from UCI repository

The goal of our empirical study was to check the looseness of our bounds and to understand, what effects must be taken into account additionally? We compare all three bounds [see. . .] with an ideal [see. . .] . We choose a Rule induction System as a Testing Area because global and local shatter coefficients can be easily obtained for this kind of learning algorithms.

The rule definition

- Features $f_1(x), \dots, f_n(x)$.
- Rule ϕ is a conjunction: $\phi(x, \theta) = \bigwedge_{f \in \Omega} [f(x) \leq \theta_f]$,
 where $\Omega \subset \{f_1, \dots, f_n\}$, θ_f — threshold for feature $f(x)$.
- Rule $\phi(x, \theta)$ is well interpretable while $|\Omega| \lesssim 5$.
- $\phi(x) = 1 \iff$ Rule $\phi(x)$ covers object $x \in \mathbb{X}$.
- Rule $\phi_y(x)$ of class y covers many objects from y and none or a few objects from $\mathbb{Y} \setminus y$:



Experiments with Rule Induction System

The Rule Induction Classifier

The rule definition

The rule definition

- Features $f_1(x), \dots, f_n(x)$.
- Rule ϕ is a conjunction: $\phi(x, \theta) = \bigwedge_{f \in \Omega} [f(x) \leq \theta_f]$,
where $\Omega \subset \{f_1, \dots, f_n\}$, θ_f — threshold for feature $f(x)$.
- Rule $\phi(x, \theta)$ is well interpretable while $|\Omega| \lesssim 5$.
- $\phi(x) = 1 \iff$ Rule $\phi(x)$ covers object $x \in X$.
- Rule $\phi_y(x)$ of class y covers many objects from y and none or a few objects from $X \setminus y$.



A few words about what is rule induction. A rule is a well interpretable predicate, usually a conjunction, that covers many objects of one class and none [see. . .] or a few [see. . .] objects of other classes. For example this [see. . .] is not a rule because it covers both classes significantly.

Classifier is a combination of rules

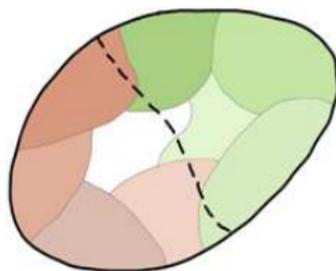
- Decision List of rules:

If $\phi_{y_1}^1(x) \rightarrow f(x) := y_1$;

If $\phi_{y_2}^2(x) \rightarrow f(x) := y_2$;

...

If $\phi_{y_T}^T(x) \rightarrow f(x) := y_T$;

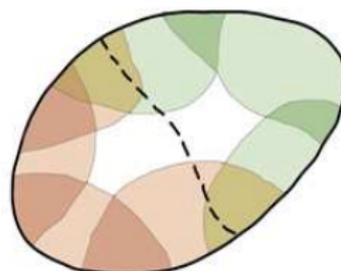


- Weighted Voting of rules:

$$f(x) = \arg \max_{y \in \mathbb{Y}} \sum_{t=1}^{T_y} w_y^t \phi_y^t(x),$$

where

$\phi_y^t(x)$ — t -th rule of class y ;
 w_y^t — its weight.



Combinatorial Approach to Generalization Bounds Tightening

- Experiments with Rule Induction System

- The Rule Induction Classifier

- Classifier is a combination of rules

Classifier is a combination of rules

Decision List of rules:

If $\phi_1^y(x) \rightarrow f(x) := y_1;$
 If $\phi_2^y(x) \rightarrow f(x) := y_2;$
 ...
 If $\phi_r^y(x) \rightarrow f(x) := y_r;$

Weighted Voting of rules:

$$f(x) = \arg \max_{y \in Y} \sum_{s=1}^{T_0} w_s^y \phi_s^y(x),$$

where
 $\phi_s^y(x)$ — s -th rule of class y ;
 w_s^y — its weight.



A rule-based classifier can be considered as an ensemble of rules based on the principle of seniority voting [see. . .] or majority voting [see. . .] .

Results

Task	L	Glob.SC	Loc.SC	Loc.S.Profile	Eff.Loc.SC
crx	690	$2.8 \cdot 10^8$	$3.5 \cdot 10^4$	$1.1 \cdot 10^4$	21 ± 11
german	1000	$5.2 \cdot 10^8$	$3.1 \cdot 10^4$	$1.8 \cdot 10^4$	47 ± 38
hepatitis	155	$5.5 \cdot 10^6$	$1.8 \cdot 10^4$	$8.4 \cdot 10^3$	58 ± 46
horse-colic	300	$1.9 \cdot 10^6$	$1.3 \cdot 10^4$	$6.3 \cdot 10^3$	5 ± 3
hypothyroid	3163	$5.3 \cdot 10^8$	$2.2 \cdot 10^4$	$9.2 \cdot 10^3$	43 ± 28
liver	345	$1.5 \cdot 10^7$	$2.9 \cdot 10^4$	$9.5 \cdot 10^3$	12 ± 8
promoters	106	$4.4 \cdot 10^9$	$5.3 \cdot 10^4$	$2.4 \cdot 10^4$	13 ± 4

Interpretation

In all tasks effective local SC $\ll L$.

Then the "effective local VC dimension" never exceeds 1 !

Combinatorial Approach to Generalization Bounds Tightening

Experiments with Rule Induction System

Experimental results: shatter coefficients

Results

Results

Task	L	Glob. SC	Loc. SC	Loc. SC Profile	Eff. Loc. SC
crx	690	$2.8 \cdot 10^6$	$3.5 \cdot 10^4$	$1.1 \cdot 10^4$	21 ± 11
german	1020	$5.2 \cdot 10^6$	$3.1 \cdot 10^4$	$1.8 \cdot 10^4$	47 ± 30
hepatitis	155	$5.5 \cdot 10^6$	$1.8 \cdot 10^4$	$8.4 \cdot 10^3$	50 ± 46
horse-colic	300	$1.9 \cdot 10^6$	$1.3 \cdot 10^4$	$6.3 \cdot 10^3$	5 ± 3
hypothyroid	3163	$5.3 \cdot 10^6$	$2.2 \cdot 10^4$	$0.2 \cdot 10^4$	43 ± 20
liver	340	$1.5 \cdot 10^6$	$2.9 \cdot 10^4$	$9.5 \cdot 10^3$	12 ± 8
promoters	106	$4.4 \cdot 10^6$	$5.3 \cdot 10^4$	$2.4 \cdot 10^4$	13 ± 4

Interpretation

In all tasks effective local SC $< L$.

Then the "effective local VC dimension" never exceeds 1!

We measured complexity terms for seven real tasks from UCI repository [see. . .] and four types of bounds.

The result is choking.

None of the bounds can be called tight!

Moreover, the effective local shatter coefficient [see. . .] is always significantly less than the sample size [see. . .]. This means that the effective local VC dimension never exceeds 1. So the VC dimension bears no relation to generalization bounds.

Directions of further investigation

- So, what causes of bound looseness did we miss?
I have two ideas:

Directions of further investigation

- So, what causes of bound looseness did we miss?

I have two ideas:

- Functions from $\{f_n \mid n = 1, \dots, N\}$ have different chance to be obtained as a result of learning.
- Two similar functions from $\{f_n \mid n = 1, \dots, N\}$ really are not 2 but ≈ 1 function.

What can we do today

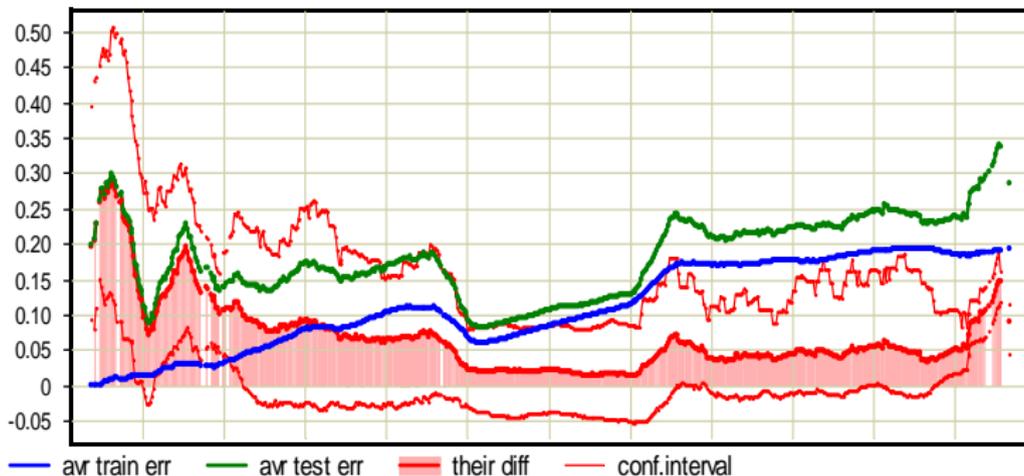
- Define the *overfitting* of a rule $\phi(x)$ as

$$\delta(\phi, X^\ell, X^k) = \nu(\phi, X^k) - \nu(\phi, X^\ell)$$

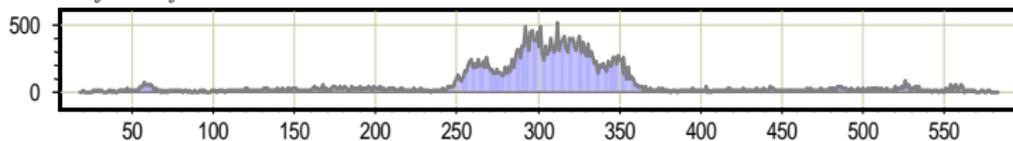
- Estimate (nonparametric) regression:
how δ depends on rules properties measured on training set X_n^ℓ during inductive search:
 - number of errors made on X_n^ℓ ;
 - number of objects covered on X_n^ℓ ;
 - number of terms in a rule;
 - entropy of a rule on X_n^ℓ ;
 - breadth and width of the inductive search;
- We use these technique for:
 - Controlling rules overfitting
 - Risk Assessment — estimation of PD (probability of default) in Credit Scoring application: $PD(x) = \nu(\phi, X^\ell) + \hat{\delta}(\phi)$

Example 1

error rate



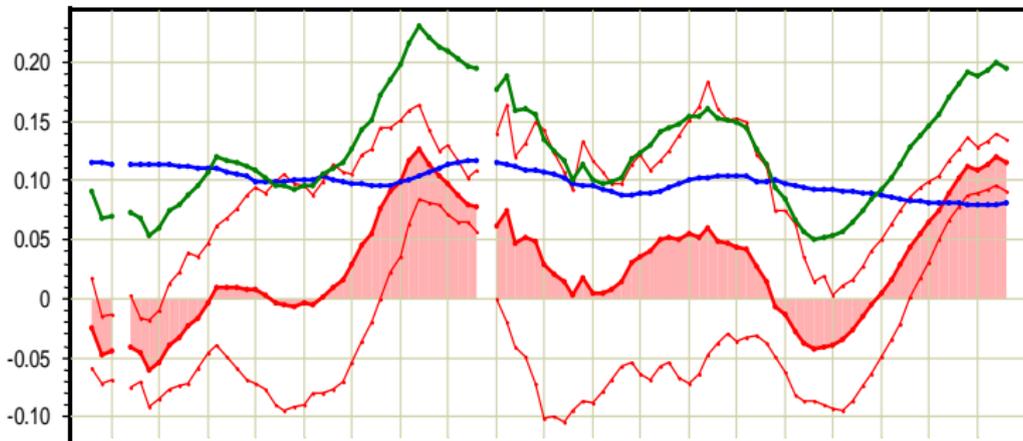
number of rules found



objects covered on training set

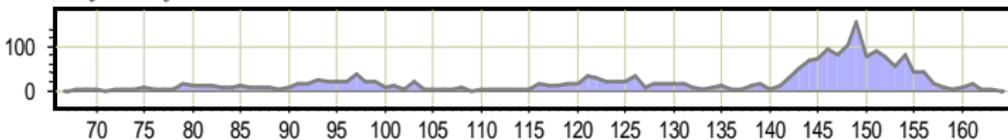
Example 2

error rate



— av train err — av test err — their diff — conf.interval

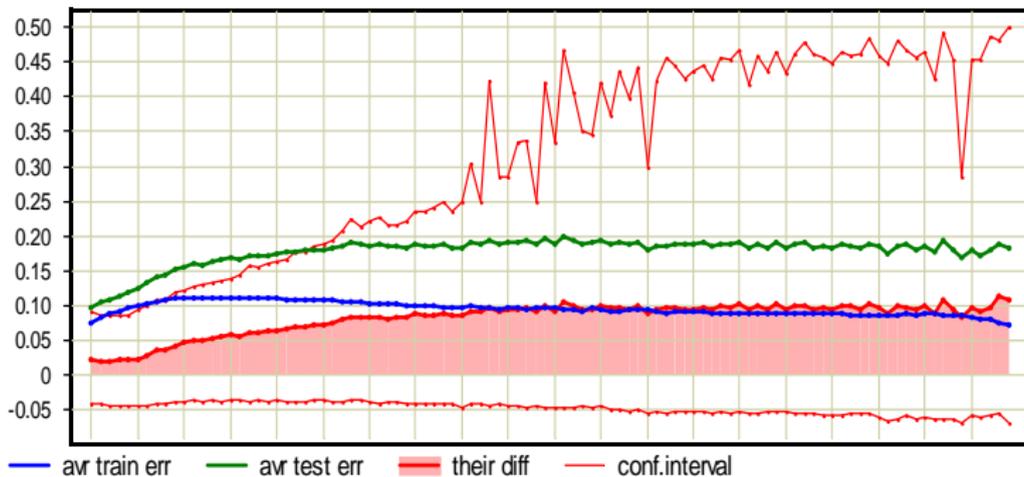
number of rules found



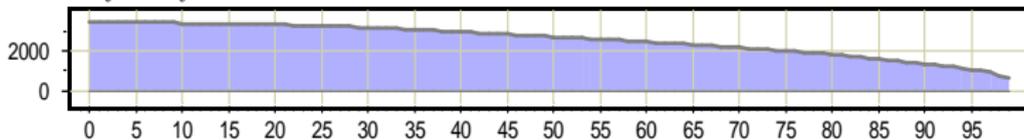
entropy of rules on training set

Example 3

error rate



number of rules found



the rating of the rule