

Title: Uniform convergence in probability

URL: https://en.wikipedia.org/wiki/Uniform_convergence_in_probability

PageID: 22999791

Categories: Category:Combinatorics, Category:Machine learning, Category:Theorems in probability theory

Source: Wikipedia (CC BY-SA 4.0).

Uniform convergence in probability is a form of convergence in probability in statistical asymptotic theory and probability theory . It means that, under certain conditions, the empirical frequencies of all events in a certain event-family converge to their theoretical probabilities . Uniform convergence in probability has applications to statistics as well as machine learning as part of statistical learning theory .

The law of large numbers says that, for each single event A , its empirical frequency in a sequence of independent trials converges (with high probability) to its theoretical probability. In many application however, the need arises to judge simultaneously the probabilities of events of an entire class S from one and the same sample. Moreover it, is required that the relative frequency of the events converge to the probability uniformly over the entire class of events S [1] The Uniform Convergence Theorem gives a sufficient condition for this convergence to hold. Roughly, if the event-family is sufficiently simple (its VC dimension is sufficiently small) then uniform convergence holds.

Definitions

For a class of predicates H defined on a set X and a set of samples $x = (x_1, x_2, \dots, x_m)$, where $x_i \in X$, the empirical frequency of $h \in H$ on x is

The theoretical probability of $h \in H$ is defined as $Q_P(h) = P\{y \in X : h(y) = 1\}$.

The Uniform Convergence Theorem states, roughly, that if H is "simple" and we draw samples independently (with replacement) from X according to any distribution P , then with high probability , the empirical frequency will be close to its expected value , which is the theoretical probability. [2]

Here "simple" means that the Vapnik–Chervonenkis dimension of the class H is small relative to the size of the sample. In other words, a sufficiently simple collection of functions behaves roughly the same on a small random sample as it does on the distribution as a whole.

The Uniform Convergence Theorem was first proved by Vapnik and Chervonenkis [1] using the concept of growth function .

Uniform convergence theorem

The statement of the uniform convergence theorem is as follows: [3]

If H is a set of $\{0, 1\}$ -valued functions defined on a set X and P is a probability distribution on X then for $\epsilon > 0$ and m a positive integer, we have:

And for any natural number m , the shattering number $\Pi_H(m)$ is defined as:

From the point of Learning Theory one can consider H to be the Concept/Hypothesis class defined over the instance set X . Before getting into the details of the proof of the theorem we will state Sauer's Lemma which we will need in our proof.

Sauer–Shelah lemma

The Sauer–Shelah lemma [4] relates the shattering number $\Pi_h(m)$ to the VC Dimension.

Lemma: $\Pi_H(m) \leq (em/d)^d$, where d is the VC Dimension of the concept class H .

Corollary: $\Pi_H(m) \leq md^m$.

Proof of uniform convergence theorem

[1] and [3] are the sources of the proof below. Before we get into the details of the proof of the Uniform Convergence Theorem we will present a high level overview of the proof.

Symmetrization: We transform the problem of analyzing $|Q_P(h) - Q^x(h)| \geq \epsilon$ into the problem of analyzing $|Q^r(h) - Q^s(h)| \geq \epsilon/2$, where r and s are i.i.d samples of size m drawn according to the distribution P . One can view r as the original randomly drawn sample of length m , while s may be thought as the testing sample which is used to estimate $Q_P(h)$.

Permutation: Since r and s are picked identically and independently, so swapping elements between them will not change the probability distribution on r and s . So, we will try to bound the probability of $|Q^r(h) - Q^s(h)| \geq \epsilon/2$ for some $h \in H$ by considering the effect of a specific collection of permutations of the joint sample $x = r || s$. Specifically, we consider permutations $\sigma(x)$ which swap x_i and x_{m+i} in some subset of $1, 2, \dots, m$. The symbol $r || s$ means the concatenation of r and s . [citation needed]

Reduction to a finite class: We can now restrict the function class H to a fixed joint sample and hence, if H has finite VC Dimension, it reduces to the problem to one involving a finite function class.

We present the technical details of the proof.

Symmetrization

Lemma: Let $V = \{x \in X^m : |Q_P(h) - Q^x(h)| \geq \epsilon \text{ for some } h \in H\}$ and

Then for $m \geq \frac{2}{\epsilon^2}$, $P_m(V) \leq 2 P_{2m}(R)$.

Proof:

By the triangle inequality, if $|Q_P(h) - Q^r(h)| \geq \epsilon$ and $|Q_P(h) - Q^s(h)| \leq \epsilon/2$ then $|Q^r(h) - Q^s(h)| \geq \epsilon/2$.

Therefore,

since r and s are independent.

Now for $r \in V$ fix an $h \in H$ such that $|Q_P(h) - Q^r(h)| \geq \epsilon$. For this h , we shall show that

Thus for any $r \in V$, $A \geq P_m(V)^2$ and hence $P_{2m}(R) \geq P_m(V)^2$. And hence we perform the first step of our high level idea.

Notice, $m \cdot \widehat{Q}_s(h)$ is a binomial random variable with expectation $m \cdot Q_P(h)$ and variance $m \cdot Q_P(h)(1 - Q_P(h))$. By Chebyshev's inequality we get

for the mentioned bound on m . Here we use the fact that $x(1-x) \leq 1/4$ for x .

Permutations

Let Γ_m be the set of all permutations of $\{1, 2, 3, \dots, 2m\}$ that swaps i and $m+i$ $\forall i$ in some subset of $\{1, 2, 3, \dots, 2m\}$.

Lemma: Let R be any subset of X^{2m} and P any probability distribution on X . Then,

where the expectation is over x chosen according to P^{2m} , and the probability is over σ chosen uniformly from Γ_m .

Proof:

For any $\sigma \in \Gamma_m$,

(since coordinate permutations preserve the product distribution P^{2m} .)

The maximum is guaranteed to exist since there is only a finite set of values that probability under a random permutation can take.

Reduction to a finite class

Lemma: Basing on the previous lemma,

Proof:

Let us define $x = (x_1, x_2, \dots, x_{2m})$ and $t = |H| \cdot x$ which is at most $|H| \cdot (2m)$. This means there are functions $h_1, h_2, \dots, h_t \in H$ such that for any $h \in H$, $\exists i$ between 1 and t with $h(x_k) = h_i(x_k)$ for $1 \leq k \leq 2m$.

We see that $\sigma(x) \in R$ iff for some $h \in H$ satisfies, $|\{1 \leq i \leq m : h(x_{\sigma(i)}) = 1\}| - |\{m+1 \leq i \leq 2m : h(x_{\sigma(i)}) = 1\}| \geq \epsilon/2$.

Hence if we define $w_{ij} = 1$ if $h_j(x_i) = 1$ and $w_{ij} = 0$ otherwise.

For $1 \leq i \leq m$ and $1 \leq j \leq t$, we have that $\sigma(x) \in R$ iff for some j in $1, \dots, t$ satisfies $|\sum_i w_{\sigma(i)}(j) - \sum_i w_{(m+i)}(j)| \geq \epsilon/2$. By union bound we get

Since, the distribution over the permutations σ is uniform for each i , so $w_{\sigma(i)}(j) - w_{(m+i)}(j)$ equals $\pm |w_{ij} - w_{m+i,j}|$, with equal probability.

Thus,

where the probability on the right is over β_i and both the possibilities are equally likely. By Hoeffding's inequality, this is at most $2e^{-m\epsilon^2/8}$.

Finally, combining all the three parts of the proof we get the Uniform Convergence Theorem.

References