

began the analysis of in-sample error in Chapter 1, and we will extend this accordingly, and may not reflect the ultimate performance in a real test. We such performance has the benefit of looking at the solutions and adjusting your performance on the practice problems that you got before the final exam.

The in-sample error E_{in} , by contrast, is based on data points that have been used for training. It expresses measures training performance, similar to that have not been used for practice.

2.1 Theory of Generalization

The same distinction between training and testing happens in learning from data. In this chapter, we will develop a mathematical theory that characterizes this distinction. We will also discuss the conceptual and practical implications of the contrast between training and testing.

If the goal is for you to learn the course material, then it is merely a way to gauge how well you have learned. If the exam problems are known ahead of time, your performance on them will no longer accurately gauge how well you have learned.

Before the final exam, a professor may hand out some practice problems and solutions to the class. Although these problems are not the exact ones that will appear on the exam, studying them will help you do better. They are the training set, in your learning.

Training versus Testing

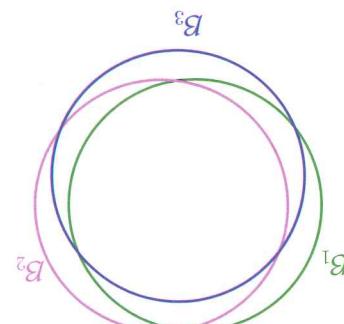
Chapter 2

effective number of hypotheses. The growth function is what will replace M . We now introduce the *growth function*, the quantity that will formalize the

2.1.1 Effective Number of Hypotheses

The mathematical theory of generalization hinges on this observation. Once we properly account for the overlaps of the different hypotheses, we will be able to replace the number of hypotheses M in (2.1) by an effective number which is finite even when M is infinite, and establish a more useful condition under which E_{out} is close to E_{in} .

Notice that $E_{\text{out}}(g)$ is very similar to $E_{\text{in}}(g)$ if ϵ is small. This is true because $E_{\text{out}}(g) - E_{\text{in}}(g) = \sum_m |E_{\text{in}}(h_m) - E_{\text{out}}(h_m)|$ is small. The union bound says that the total area covered by B_1, B_2, \dots, B_M is at most the sum of the individual areas, which is true but is a gross overestimate when the areas overlap heavily as in this example. The events “ $|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| < \epsilon$ ” and “ $|E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| < \epsilon$ ” are often strongly overlapping. If h_1 is very similar to h_2 for instance, the two events “ $|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| < \epsilon$ ” and “ $|E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| < \epsilon$ ” are likely to coincide for most data sets. In a typical learning model, many hypotheses are indeed very similar. If you take the perceptron model for instance, as you slowly vary the weight vector w , you get infinitely many hypotheses that differ from each other only infinitesimally.



$$\Pr[B_1 \text{ or } B_2 \text{ or } \dots \text{ or } B_M] \leq \Pr[B_1] + \Pr[B_2] + \dots + \Pr[B_M].$$

Then, we notice that for an example with 3 hypotheses, the right loose as illustrated in the figure to overlap, the union bound becomes particularly loose as it illustrates to include the event “ $|E_{\text{in}}(g) - E_{\text{out}}(g)| < \epsilon$ ” since g is a linear combination of the three hypotheses. The union bound says that the union bound is guaranteed to include the event “ $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| < \epsilon$ ” which is one of the hypotheses in H . We then over-estimated the probability using this way.

$$\text{(2.2)} \quad \text{“} |E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| < \epsilon \text{”},$$

⋮

$$\begin{aligned} \text{“} |E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| < \epsilon \text{” or} \\ \text{“} |E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| < \epsilon \text{” or} \end{aligned}$$

the way we got the M factor in the first place was by taking the disjunction of events:

¹Sometimes, “generalization error”, is used as another name for E_{out} , but not in this book. In order to study generalization in such models, we need to derive a counterexample to (2.1) that deals with infinite H . We would like to replace M with ∞ , the size of the hypothesis set H . If H is an infinite set, the bound goes to infinity and becomes meaningless. Unfortunately, almost all interesting learning models have simple perceptron which we discussed in Chapter 1.

on M , the size of the hypothesis set H . If H is an infinite set, the bound depends on the error bound $\sqrt{\frac{1}{2N} \ln \frac{2M}{\epsilon}}$ in (2.1), or error bar, if you will, depends the ϵ we have chosen will have a comparably higher E_{out} .

Not only do we want to know that the hypothesis with a higher E_{in} than $E_{\text{out}}(g)$. The $E_{\text{out}}(h) \leq E_{\text{in}}(h) - \epsilon$ direction of the bound assures us that with our H (no other hypothesis $h \in H$ has $E_{\text{out}}(h)$ significantly better than $E_{\text{out}}(g)$), but we also want to be sure that we did the best we could with the best training error. One with the hypothesis g that we choose (say the one with $E_{\text{in}} + \epsilon$, but we also want to do well out of sample (i.e., $E_{\text{out}} \leq E_{\text{in}} + \epsilon$), will continue to do well out of sample (say the one with $E_{\text{in}} + \epsilon$, but in a more subtle way). $E_{\text{in}} - \epsilon$ for all $h \in H$. This is important for learning, but it is also holds, that is, $E_{\text{out}} \geq E_{\text{in}} - \epsilon$ for the other side of $|E_{\text{out}} - E_{\text{in}}| \leq \epsilon$ follows.

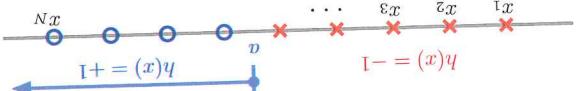
Notice that the type of inequality in (2.1) as a generalization bound because it bounds E_{out} in terms of E_{in} . To see that the Hoeffding Inequality implies this generalization bound, we rewrite (1.6) as follows: with probability ϵ it holds $E_{\text{out}} \leq E_{\text{in}} + \epsilon$, which implies $E_{\text{out}} \leq E_{\text{in}} + \epsilon$. We may now identify $\delta = 2M\epsilon - 2N\epsilon^2$, from which $\epsilon = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$, and (2.1) follows.

We refer to the type of inequality in (2.1) as a generalization bound because example $\delta = 0.05$, and assert with probability at least $1 - \delta$ that

$$\text{for any } \epsilon > 0. \text{ This can be rephrased as follows. Pick a tolerance level } \delta, \text{ for } \Pr[|E_{\text{in}}(g) - E_{\text{out}}(g)| < \epsilon] \geq 2M\epsilon - 2\epsilon^2 N,$$

Generalization error. We have already discussed how the value of E_{in} does not always generalize to a similar value of E_{out} . Generalization is a key issue in learning. One can define the generalization error as the discrepancy between E_{in} and E_{out} . The Hoeffding Inequality (1.6) provides a way to characterize the generalization error with a probabilistic bound, but it bounds E_{out} in terms of E_{in} . To make it easier on the not-so-mathematically inclined, we will tell you which part you can safely skip without losing the plot. The mathematical results provide fundamental insights into learning from data, and we will interpret these results in practical terms.

A word of warning: this chapter is the heaviest in this book in terms of mathematics abstracted. To make it easier on the not-so-mathematically inclined, we will tell you which part you can safely skip without losing the plot. The mathematical results provide fundamental insights into learning from data, and we will interpret these results in practical terms. Between a training set and a test set more precise. We will also make the contrast analysis to the general case in this chapter. We will also make the contrast between a training set and a test set more precise.



the right of a .

- Positive rays: H consists of all hypotheses $h: \mathbb{R} \rightarrow \{-1, +1\}$ of the form $h(x) = \text{sign}(x - a)$, i.e., the hypotheses are defined in a one-dimensional input space, and they return -1 to the left of some value a and $+1$ to the right.

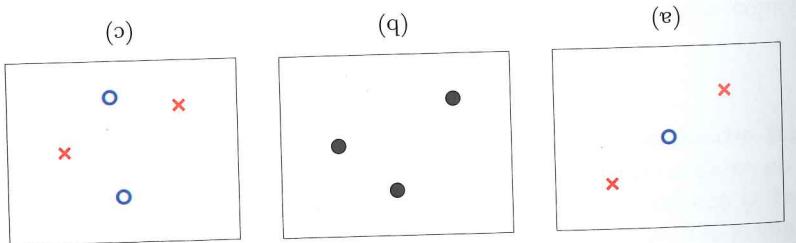
Example 2.2. Let us find a formula for $m^h(N)$ in each of the following cases.

Let us now illustrate how to compute $m^H(N)$ for some simple hypotheses sets. These examples will confirm the intuition that $m^H(N)$ grows faster when the hypothesis set \mathcal{H} becomes more complex. This is what we expect of a quantity that is meant to replace M in the generalization bound (2.1).

$$mh(4) = 14.$$

Example 2.1. If \mathcal{X} is a Euclidean plane and \mathcal{H} is a two-dimensional perceptron, what are $m_h(3)$ and $m_h(4)$? Figure 2.1(a) shows a dichotomy on 3 points that the perceptron cannot generate, while Figure 2.1(b) shows another 3 points that the perceptron can. Figure 2.1(c) is based on the maximum number of dichotomies. Because the definition of $m_h(N)$ is based on the maximum number of dichotomies, it follows that $m_h(N) = 2^N$.

Figure 2.1: Illustration of the growth function for a two-dimensional perceptron. The dichotomy of red versus blue on the 3 colinear points in part (a) cannot be generated by a perceptron, but all 8 dichotomies on the 3 points in part (b) can. By contrast, the dichotomy of red versus blue on the 4 points in part (c) cannot be generated by a perceptron. At most 14 out of the possible 16 dichotomies on any 4 points can be generated.



If H is capable of generating all possible dichotomies on x_1, \dots, x_N , then $H(x_1, \dots, x_N) = \{ -1, +1 \}_N$ and we say that H can shatter x_1, \dots, x_N . This signifies that H is as diverse as can be on this particular sample.

$$m^2 > N^{\mu}$$

In words, $m^h(N)$ is the maximum number of dichotomies that can be generated by \mathcal{H} on any N points. To compute $m^h(N)$, we consider all possible choices of N points x_1, \dots, x_N from \mathcal{X} and pick the one that gives us the most dichotomies. Like M , $m^h(N)$ is a measure of the number of hypotheses in \mathcal{H} , except that a hypothesis is now considered on N points instead of the entire \mathcal{X} . For any \mathcal{H} , since $\mathcal{H}(x_1, \dots, x_N) \subseteq \{-1, +1\}^N$ (the set of all possible dichotomies on any N points), the value of $m^h(N)$ is at most $\{-1, +1\}^N$, hence

where $| \cdot |$ denotes the cardinality (number of elements) of a set.

$$\max_{\mathcal{X} \ni \mathbf{x}_1, \dots, \mathbf{x}_N} \mathcal{H} \equiv (\Lambda^N) \mathcal{H}_{\text{sys}}$$

Definition 2.2. The growth function is defined for a hypothesis set \mathcal{H} by

One can think of the dichotomies $\mathcal{H}(x_1, \dots, x_N)$ as a set of hypotheses just like \mathcal{H} is, except that the hypotheses are seen through the eyes of N points like x_i , $i = 1, \dots, N$. The growth function is based on the number of dichotomies on x_1, \dots, x_N . A larger $\mathcal{H}(x_1, \dots, x_N)$ means \mathcal{H} is more "diverse" — generating more dichotomies on x_1, \dots, x_N . The growth function is based on the number of dichotomies.

$$(2.3) \quad \cdot \{ \mathcal{H} \ni h \mid ((N\mathbf{x})h) \cdot \dots \cdot ((I\mathbf{x})h) \} = (N\mathbf{x}, I\mathbf{x})\mathcal{H}$$

these points are defined by H on the boundaries generated by N .

The definition of the growth function is based on the number of different hypotheses that can implement, but only over a finite sample of points rather than over the entire input space \mathcal{X} . If $h \in \mathcal{H}$ is applied to a finite sample $x_1, \dots, x_N \in \mathcal{X}$, we get an N -tuple $h(x_1), \dots, h(x_N)$ of ± 1 's. Such an N -tuple is called a dichotomy since it splits x_1, \dots, x_N into two groups: those points for which h is $+1$ and those for which h is -1 . Each $h \in \mathcal{H}$ generates a dichotomy on x_1, \dots, x_N , but two different h 's may generate the same dichotomy if they happen to give the same pattern of ± 1 's on this particular sample.

We now use the break point k to derive a bound on the growth function $m_h(N)$ for all values of N . For example, the fact that no 4 points can be shattered by

that Example.

By inspection, find a break point k for each hypothesis set in Example 2.2 (if there is one). Verify that $m_h(k) < 2^k$ using the formulas derived in

Exercise 2.1

break point for \mathcal{H} than to compute the full growth function for that \mathcal{H} . If k is a break point, then $m_h(k) < 2^k$. Example 2.1 shows that $k = 4$ is a break point for two-dimensional perceptrons. In general, it is easier to find a

to be a break point for \mathcal{H} .

Definition 2.3. If no data set of size k can be shattered by \mathcal{H} , then k is said

point.

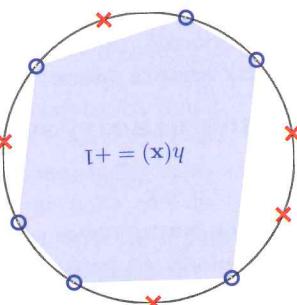
much easier than computing $m_h(N)$ itself, thanks to the notion of a break inequality in (2.1) will still hold. Getting a good bound on $m_h(N)$ will prove we can use an upper bound on $m_h(N)$ instead of the exact value, and the Fortunately, we don't have to. Since $m_h(N)$ is meant to replace M in (2.1), It is not practical to try to compute $m_h(N)$ for every hypothesis set we use.

in this case). \square

Notice that if the N points were chosen at random in the plane rather as $m_h(N)$ is concerned, since it is defined based on the maximum (2^N) as we did for the perimeter points. However, this doesn't matter as far and we wouldn't be able to shatter all the points with convex hypotheses than on the perimeter of a circle, many of the points would be "internal". To compute $m_h(N)$, we notice that given N points, the line is again

$$m_h(N) = 2^N.$$

This means that any dichotomy on these N points can be realized using a convex hypothesis, so \mathcal{H} manages to shatter these points and the growth function has the maximum possible value



2. TRAINING VERSUS TESTING
2.1. THEORY OF GENERALIZATION

than three +1 points, the convex set will be a line segment, a point, or with the dichotomy on all N points. For the diameter of a circle agrees to be convex since its vertices are on the closed interior of the polygon (which has the hypothesis made up of the N points. If you connect the +1 points with a polygon, there of \mathbb{I} 's to the N points, consider any dichotomy on these points, assigning an arbitrary pattern. Now consider N points on the N points, choose N points on the N points care- fully. Per the next figure, choose N points on the N points car-

To compute $m_h(N)$ in this case, we need to choose the N points care-

fully entirely within the set).

(a set is convex if the line segment connecting any two points in the set

$\{-1, +1\}$ that are positive inside some convex set and negative elsewhere

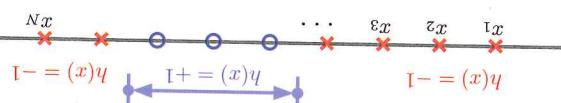
3. Convex sets: \mathcal{H} consists of all hypotheses in two dimensions $h: \mathbb{R}^2 \rightarrow$

ear $m_h(N)$ of the "simpler" positive ray case.

Notice that $m_h(N)$ grows as the square of N , faster than the lin-

$$m_h(N) = \binom{2}{N+1} + 1 = \frac{2}{1} N^2 + \frac{2}{1} N + 1.$$

regional, the resulting hypothesis is the constant -1 regardless of which region it is. Adding up these possibilities, we get $\binom{2}{N+1}$ different dichotomies. If both end values fall in the same by which two regions contain the end values of the interval, resulting split by the points into $N+1$ regions. The dichotomy we get is decided To compute $m_h(N)$, we notice that given N points, the line is again



specified by the two end values of that interval.

2. Positive intervals: \mathcal{H} consists of all hypotheses in one dimension that return +1 within some interval and -1 otherwise. Each hypothesis is

number of dichotomies.

Notice that if we picked N points where some of the points coincided not affect the value of $m_h(N)$ since it is defined based on the maximum (which is allowed), we will get less than $N+1$ dichotomies. This does

any N points, the growth function is

$$m_h(N) = N + 1.$$

To compute $m_h(N)$, we notice that given N points, the line is split by the points into $N+1$ regions. The dichotomy we get on the N points is decided by which region contains the value a . As we vary a , we will get $N+1$ different dichotomies. Since this is the most we can get for is decided by which region contains the value a . As we vary a , we will the points into $N+1$ regions. The dichotomy we get on the N points

The total number of different dichotomies on the first $N - 1$ points is given by $a + \beta$; since S_1^+ and S_2^- are identical on these $N - 1$ points, their dichotomies are redundant. Since no subset of k of these first $N - 1$ points can

$$B(N, k) = \alpha + 2\beta. \quad (2.4)$$

where x_1, \dots, x_N in the table are labels for the N points of the dichotomy. We have chosen a convenient order in which to list the dichotomies, as follows. Consider the dichotomies on x_1, \dots, x_{N-1} . Some dichotomies on these $N - 1$ points appear only once (with either +1 or -1 in the x_N column, but not on the first $N - 1$ points appear twice, once with +1 and once with -1 in both). We collect these dichotomies in the set S_1 . The remaining dichotomies appear only once (with either +1 or -1 in the x_N column, but not on the first $N - 1$ points appear twice, once with +1 and once with -1 in both). We collect these dichotomies in the set S_2 . The remaining dichotomies appear only two equal parts, S_1 and S_2 (with $+1$ and -1 in the x_N column, divided into two equal parts, S_1 and S_2 (with $+1$ and -1 in the x_N column, respectively). Let S_1 have a rows, and let S_2 have b rows each. Since the total number of rows in the table is $B(N, k)$ by construction, we have

We first assume $N \geq 2$ and $k \geq 2$ and try to develop a recursion. Consider on the one point.

one dichotomy can be allowed. A second different dichotomy must differ on at least one point and then that subset of size 1 would be shattered. $B(1, k) = 2$ for $k > 1$ since in this case there do not even exist subsets of size k ; the constraint is vacuously true and we have 2 possible dichotomies (+1 and -1)

$$B(1,k) = 2 \text{ for } k < 1.$$

$$B(N,1) = 1$$

The notation B comes from binomial, and the reason will become clear shortly. To evaluate $B(N, k)$, we start with the two boundary conditions $k = 1$ and $N = 1$.

$m_h(N) \leq B(N, k)$ if k is a break point for \mathcal{H} .

The definition of $D(N, k)$ assumes a break point k , then tries to find the most dichotomies on N points without imposing any further restrictions. Since $B(N, k)$ is defined as a maximum, it will serve as an upper bound for any $m_h(N)$ that has a break point k :

Definition 2.4. $B(N, k)$ is the maximum number of dichotomies on N points such that no subset of size k of the N points can be shattered by these dichotomies.

To prove the polynomial bound, we will introduce a combinatorial quantity that counts the maximum number of dichotomies given that there is a break point, without having to assume any particular form of \mathcal{H} . This bound will therefore apply to any \mathcal{H} .

Be gentle with mistakes: If you trust our math, you can skip the following part without compromising the logical sequence. A similar green box will tell you when to regroup.

The most important fact about growth functions is that if the condition $m^h(N) = 2^N$ breaks at any point, we can bound $m^h(N)$ for all values of N by a simple polynomial based on this break point. The fact that the bound is polynomial is crucial. Absent a break point (as is the case in the convex hypothesis example), $m^h(N) = 2^N$ for all N . If $m^h(N)$ replaced M in Equation (2.1), the bound $\sqrt{\frac{1}{2M} \ln \frac{2M}{\delta}}$ on the generalization error would not go to zero regardless of how many training examples N we have. However, if $m^h(N)$ can be bounded by a polynomial — any polynomial — the generalization error will go to zero as $N \rightarrow \infty$. This means that we will generalize well given sufficient number of examples.

2.1.2 Bounding the Growth Function

$$B(N_0 + 1, k) \leq B(N_0, k) + B(N_0, k - 1).$$

Proof. The statement is true whenever $k = 1$ or $N = 1$, by inspection. The condition, we only need to worry about $k \geq 2$. By (2.7), the statement is already true when $k = 1$ (for all values of N) by the initial proof is by induction on N . Assume the statement is true for all $N \leq N_0$ and all k . We need to prove the statement is true for $N = N_0 + 1$ and all k . Since the proof is by induction on N , we have good generalization; a polynomial bound on $m^H(N)$.

The implication of Theorem 2.4 is that if H has a break point, we have for all N . The RHS is polynomial in N of degree $k - 1$.

$$m^H(N) \leq \sum_{i=0}^{k-1} \binom{i}{N}$$

Theorem 2.4. If $m^H(k) < 2^k$ for some value k , then

End case skip: Those who skipped are now rejoicing again. The next theorem states that any growth function $m^H(N)$ with a break point is bounded by a polynomial.

It turns out that $B(N, k)$ in fact equals $\sum_{i=0}^k \binom{i}{N}$ (see Problem 2.4), but we only need the inequality of Lemma 2.3 to bound the growth function. For a given break point k , the bound $\sum_{i=0}^k \binom{i}{N}$ is polynomial in N , as each term in the sum is polynomial (of degree $i \leq k - 1$). Since $B(N, k)$ is an upper bound on any $m^H(N)$ that has a break point k , we have proved

thus proved the induction step, so the statement is true for all N and k . ■

This identity can be proved by noticing that to calculate the number of ways to pick i objects from $N_0 + 1$ distinct objects, either the first object is included, in $\binom{i-1}{N_0}$ ways, or the first object is not included, in $\binom{i}{N_0}$ ways. We have to pick i objects from $N_0 + 1$ distinct objects, either the first object is included, in $\binom{i-1}{N_0}$ ways, or the first object is not included, in $\binom{i}{N_0}$ ways. We have thus proved the combinatorial identity $\binom{N_0}{N_0+1} = \binom{i}{N_0} + \binom{i-1}{N_0}$ has been used.

$$\begin{aligned} B(N_0 + 1, k) &\leq \sum_{i=0}^{k-2} \binom{i}{N_0} + \sum_{i=0}^{k-1} \binom{i}{N_0} \\ &= \sum_{i=0}^{k-1} \binom{i}{N_0 + 1} + \sum_{i=0}^{k-1} \binom{i}{N_0} \\ &= \sum_{i=0}^{k-1} \left[\binom{i}{N_0} + \binom{i-1}{N_0} \right] \\ &= \sum_{i=0}^{k-1} \binom{i+1}{N_0 + 1} = B(N_0, k). \end{aligned}$$

Applying the induction hypothesis to each term on the RHS, we get

where the first row ($N = 1$) and the first column ($k = 1$) are the boundary conditions that we already calculated. We can also use the recursion to bound $B(N, k)$ analytically.

N	1	2	3	4	5	6	...
1	1	2	2	2	2	2	...
2	1	3	4	4	4	4	...
3	1	4	7	8	8	8	...
4	1	5	11
5	1	6
6	1	7

We can use (2.7) to recursively compute a bound on $B(N, k)$, as shown in the following table.

$$B(N, k) \leq B(N - 1, k) + B(N - 1, k - 1). \quad (2.7)$$

Substituting the two inequalities (2.5) and (2.6) into (2.4), we get

$$\beta \leq B(N - 1, k - 1). \quad (2.6)$$

in this table by definition of $B(N, k)$. Therefore, points yields a subset of size k that is shattered, which we know cannot exist taking the corresponding set of dichotomies in S_2 and adding x_N to the data be shattered by the dichotomies in S_2 . If there exists such a subset, then by definition of B . Further, no subset of size $k - 1$ of the first $N - 1$ points can be shattered by the dichotomies in S_2 . If there exists such a subset, then taking the corresponding set of dichotomies in S_2 and adding x_N to the data be shattered (since no k -subset of all N points can be shattered), we deduce that be shattered (since no k -subset of all N points can be shattered), we deduce

total area of the canvas is 1.

Skech of the Proof. The data set D is the source of randomization in the original Hoefding Inequality. Consider the space of all possible data sets. Let us think of this space as a, "canvas" (Figure 2.2(a)). Each D is a point on that canvas. The probability of a point is determined by which x^n 's in \mathcal{X} happen to be in that particular D , and is calculated based on the distribution P over \mathcal{X} . Let's think of probabilities of different events as areas on that canvas, so the

in (2.12).

The VC generalization bound is the most important mathematical result in the theory of learning. It establishes the feasibility of learning with infinite hypothesis sets. Since the formal proof is somewhat lengthy and technical, we illustrate the main ideas in a sketch of the proof, and include the formal proof as an appendix. There are two parts to the proof; the justification that the growth function can replace the number of hypotheses in the first place, and the reason why we had to change the red items in (2.11) into the blue items as in Figure 2.1.

If you compare the blue items in (2.12) to their red counterparts in (2.11), you notice that all the items move the bound in the weaker direction. However, as long as the VC dimension is finite, the error bar still converges to zero (albeit at a slower rate), since $m_H(2N)$ is also polynomial of order $\log N$, just like $m_H(N)$. This means that, with enough data, each and every hypothesis in an infinite H with a finite VC dimension will generalize well from E_{in} to E_{out} . The key is that the effective number of hypotheses, represented by m_H , has replaced the actual number of hypotheses in the finite growth function.

$$E_{\text{in}}(a) \leq E_{\text{in}}^*(a) + \sqrt{\frac{8}{\pi} \ln \frac{4mH(2N)}{s}} \quad (2.12)$$

Theorem 2.5 (VC generalization bound). For any tolerance $\delta < 0$,

It turns out that this is not exactly the form that will hold. The quantiles in need to be technically modified to make (2.11) true. The correct bound, which is called the VC generalization bound, is given in the following theorem; it holds for any binary target function f , any hypothesis set \mathcal{H} , any learning algorithm A , and any input probability distribution P .

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{2N}{1} \ln \frac{6}{\delta}}. \quad (2.11)$$

If we treated the growth function as an effective number of hypotheses, and replaced M in the generalization bound (2.1) with $m^H(N)$, the resulting bound

2.1.4 The VC Generalization Bound

$\mathcal{X} = \{1\} \times \mathbb{F}_d^d$ is considered d -dimensional since the first coordinate $x_0 = 1$ is fixed.

Diversity is not necessarily a good thing in the context of generalization. For example, the set of all possible hypotheses is as diverse as can be, so $m^H(N) = 2^N$ for all N and $d^{vc}(H) = \infty$. In this case, no generalization at all is to be expected, as the final version of the generalization bound will show.

The VC dimension of a d -dimensional perceptron is indeed $d + 1$. This is consistent with Figure 2.1 for the case $d = 2$, which shows a VC dimension of 3. The perceptron case provides a nice intuition about the VC dimension, since $d + 1$ is also the number of parameters in this model. One can view the VC dimension as measuring the effective number of parameters. The more parameters a model has, the more diverse its hypotheses set is, which reflects in a larger value of the growth function $m^H(N)$. In the case of perceptrons, the effective parameters correspond to explicit parameters in the model, namely w_0, w_1, \dots, w_d . In other models, the effective parameters may be less obvious or implicit. The VC dimension measures these effective parameters or "degrees of freedom" that enable the model to express a hypothesis set of N points.

(b) To show that $d \leq d+1$, show that no set of $d+2$ points in \mathcal{X} can be shattered by the perceptron. [Hint: Represent each point in \mathcal{X} as a vector of length $d+1$, then use the fact that any $d+2$ vectors of length $d+1$ have to be linearly dependent. This means that some vector is a linear combination of all the other vectors. Now, if you choose the class of these other vectors carefully, then the classification of the dependent vector will be dictated. Conclude that there is some dichotomy that cannot be implemented, and therefore that for $N \geq d+2$, $m_h(N) < 2^N$.]

(a) To show that $\dim \mathcal{X} \geq d+1$, find $d+1$ points in \mathcal{X} that the perceptron can shatter. Hint: Construct a nonsingular $(d+1) \times (d+1)$ matrix whose rows represent the $d+1$ points, then use the nonsingularity to argue that the perceptron can shatter these points.

Consider the input space $\mathcal{X} = \{1\} \times \mathbb{R}^d$ (including the constant coordinate $c_0 = 1$). Show that the VC dimension of the perceptron (with $d + 1$ parameters, counting w_0) is exactly $d + 1$ by showing that it is at least $d + 1$ and at most $d + 1$, as follows.

Exercise 2.4

In this case all the powers can be shared out by π_1 . In this case, we can conclude that $d\psi^o < N$.

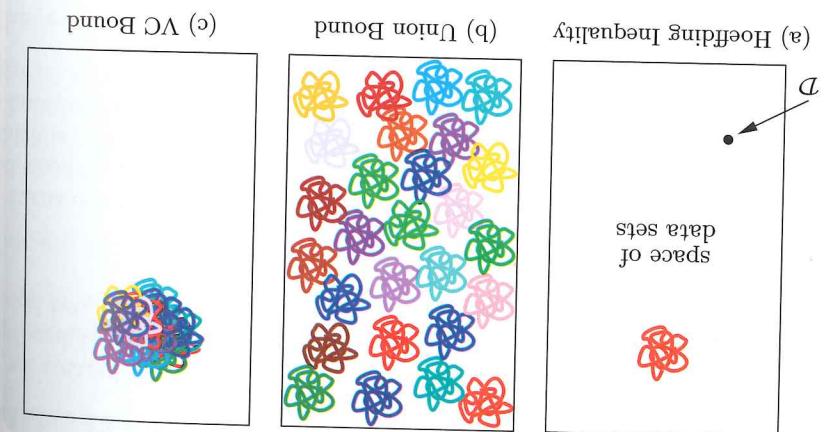
In this final iteration, we cannot conclude anything about the value of d_{VC} .

³. There is a set of N points that cannot be shattered by H . Based only on this information we cannot decide whether H is VC- d .

than enough information to conclude that $d^{y_0} \geq N$.

Figure 2.2: Illustration of the proof of the VC bound, where the canvas represents the space of all data sets, with areas corresponding to probability of given hypotheses $h \in \mathcal{H}$, the event " $|E_m(h) - E_{out}(h)| < \epsilon$ " consists of all points D for which the statement is true. For a particular h , let us paint all these "bad" points using one color. What the basic Hoeffding Inequality tells us is that the colored area on the canvas will be small (Figure 2.2(a)). For a given hypothesis $h \in \mathcal{H}$, the event " $|E_m(h) - E_{out}(h)| < \epsilon$ " depends on D and \mathcal{D} . Now, if we take another $h \in \mathcal{H}$, the event " $|E_m(h) - E_{out}(h)| < \epsilon$ " may contain different points, since the event depends on h . Let us paint these points with a different color. The area covered by all the points we colored will be at most the sum of the two individual areas, which is the case only if the two with a different color overlap with previous colors. If we keep throwing in new colored areas for each $h \in \mathcal{H}$, and never consider two different colors, then the colored area will soon be mostly covered in color and not taking the overlaps of the union bound in the Hoeffding Inequality (Figure 2.2(b)). Even if each h contributed very little, the sheer number of hypotheses will eventually make the colored area cover the whole canvas. This was the problem with binary target functions (because of 100 times (because of 100 different h 's), then the total colored area is now $1/100$ of what it would have been if the colored points had not overlapped at all. This is the essence of the VC bound as illustrated in (Figure 2.2(c)). The argument goes as follows.

The bulk of the VC proof deals with how to account for the overlaps. Here is the idea. If you were told that the hypotheses in \mathcal{H} are such that each point on the canvas that is colored will be colored 100 times (because of 100 different h 's), then the total colored area is now $1/100$ of what it would have been if the colored points had not overlapped at all. This is the essence of the VC bound as illustrated in (Figure 2.2(c)). The argument goes as follows. The bulk of the VC proof deals with how to account for the overlaps. Here is the idea. If you were told that the hypotheses in \mathcal{H} are such that each point on the canvas that is colored will be colored 100 times (because of 100 different h 's), then the total colored area is now $1/100$ of what it would have been if the colored points had not overlapped at all. This is the essence of the VC bound as illustrated in (Figure 2.2(c)). The argument goes as follows. The bulk of the VC proof deals with how to account for the overlaps. Here is the idea. If you were told that the hypotheses in \mathcal{H} are such that each point on the canvas that is colored will be colored 100 times (because of 100 different h 's), then the total colored area is now $1/100$ of what it would have been if the colored points had not overlapped at all. This is the essence of the VC bound as illustrated in (Figure 2.2(c)). The argument goes as follows.



Many hypotheses share the same dichotomy on a given D , since there are finitely many dichotomies even with an infinite number of hypotheses. Any finite growth function enables us to do this to account for this kind of hypotheses false for all the hypotheses that look the same on that particular D . What statement based on D alone will be simultaneously true or simultaneously redundant many dichotomies even with an infinite number of hypotheses. Any redundancy in a precise way, so we can get a factor similar to the 100, in the above example.

When \mathcal{H} is infinite, the redundancy factor will also be infinite since the hypotheses will be divided among a finite number of dichotomies. Therefore, the reduction in the total colored area when we take the redundancy into consideration will be dramatic. If it happens that the number of dichotomies is only a polynomial, the reduction will be so bring the total probability down to a very small value. This is the essence of the proof of Theorem 2.5.

The reason $m_h(2N)$ appears in the VC bound instead of $m_h(N)$ is that the proof uses a sample of $2N$ points instead of N points. Why do we need $2N$ points? The event " $|E_m(h) - E_{out}(h)| < \epsilon$ " depends not only on D , but also on the entire \mathcal{X} because $E_{out}(h)$ is based on \mathcal{X} . This breaks the main premise of the proof, since we have to justify why the two-sample condition " $|E_m(h) - E_{out}(h)| < \epsilon$ " holds.

Of course we have to replace the original condition " $|E_m(h) - E_{out}(h)| < \epsilon$ " by a factor of 4, and also end up doing so, we end up having to shrink the ϵ 's by a factor of 4, and also end up getting so, we end up having to shrink the ϵ 's by a factor of 4, and also end up you get the formula in (2.12). \square

The VC generalization bound (2.12) is a universal result in the sense that since the same bound has to cover a lot of different cases. Indeed, the bound that the bound it provides may not be particularly tight in any given case, target functions as well. Given the generality of the result, one would suspect distributions, and binary target functions. It can be extended to other types of it applies to all hypothesis sets, learning algorithms, input spaces, probability distributions to all hypothesis sets, learning algorithms, input spaces, probability distributions, and binary target functions. It can be extended to other types of

4The term ‘complexity’ comes from a similar metaphor in computational complexity.

overestimate; a more practical constant of proportionality is closer to 10. □
practise. The constant of proportionality it suggests is 10,000, which is a gross is approximately proportional to the VC dimension, as has been observed in You can see that the inequality suggests that the number of examples needed in similar calculation will find that $N \approx 40,000$. For $d_{\text{vc}} = 5$, we get $N \approx 50,000$. process, rapidly converging to an estimate of $N \approx 30,000$. If $d_{\text{vc}} = 4$, a We then try the new value $N = 21,193$ in the RHS and continue this iterative

$$N \geq \frac{0.1^2}{8} \ln \left(\frac{4(2 \times 1000)^3 + 4}{0.1} \right) \approx 21,193.$$

Trying an initial guess of $N = 1,000$ in the RHS, we get

$$N \geq \frac{0.1^2}{8} \ln \left(\frac{4(2N)^3 + 4}{0.1} \right).$$

$\epsilon = 0.1$ and $\delta = 0.1$). How big a data set do we need? Using (2.13), we need would like the generalization error to be at most 0.1 with confidence 90% (so Example 2.6. Suppose that we have a learning model with $d_{\text{vc}} = 3$ and simple iterative methods.

which is again implicit in N . We can obtain a numerical value for N using

$$N \geq \frac{\epsilon^2}{8} \ln \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right), \quad (2.13)$$

upper bound in (2.10) which is based on the the VC dimension, we get a both sides of the inequality. If we replace $m_h(2N)$ in (2.12) by its polynomial This gives an implicit bound for the sample complexity N , since N appears on suffices to obtain generalization error at most ϵ (with probability at least $1 - \delta$).

$$N \geq \frac{\epsilon^2}{8} \ln \left(\frac{4m_h(2N)}{\delta} \right)$$

$\sqrt{\frac{8}{N} \ln \frac{4m_h(2N)}{\delta}}$, and so it suffices to make $\sqrt{\frac{8}{N} \ln \frac{4m_h(2N)}{\delta}} \leq \epsilon$. It follows that be at most ϵ . From Equation (2.12), the generalization error is bounded by learning model. Fix $\delta < 0$, and suppose we want the generalization error to We can use the VC bound to estimate the sample complexity for a given indicates how much data is needed to get good generalization.

error tolerance ϵ is violated. How fast N grows as ϵ and δ become smaller generalization error, and the confidence parameter d determines how often the by two parameters, ϵ and δ . The error tolerance ϵ determines the allowed to achieve a certain generalization performance. The performance is specified The sample complexity denotes how many training examples N are needed

2.2.1 Sample Complexity

the bound is used in practice. It not absolutely. With this understanding, let us look at the different ways

generalization is a popular rule of thumb. Thus, the VC bound can be used as a guideline for generalization, relatively useful in practice, and some rules of thumb have emerged in terms of the VC dimension. For instance, requiring that N be at least $10 \times d_{\text{vc}}$ to get decent generalization. Because of this observation, the VC analysis proves those with higher d_{vc} . These models with lower d_{vc} tend to generalize better than applications, learning models with d_{vc} from practical experience, not a mathematical statement. This is an observation from real-life data sets to compare the generalization performance of these models. It tends to be equally loose for different learning models, and hence is useful sets, the only kind we use in practice. Second, although the bound is loose, analysis is what establishes the feasibility of learning for infinite hypothesis space rather to go through the analysis itself? Two reasons. Why did we bother to use the VC line of analysis leads to a very loose bound. The reality is that the VC line of analysis have resulted in only diminishing returns. Technical attempts in the literature have led to a very loose bound. Some effort could be put into tightening the VC bound, but many highly

will contribute further slack to the VC bound.
3. Bounding $m_h(N)$ by a simple polynomial of order d_{vc} , as given in (2.10),

far fewer dichotomies than 2^N , while $m_h(N) = 2^N$.
the case of convex sets in two dimensions, which we examined in Example 2.2, if you pick N points at random in the plane, they will likely have its expected value instead of the upper bound $m_h(N)$. For instance, in bound if we considered specific x_1, \dots, x_N and used $|\mathcal{H}(x_1, \dots, x_N)|$ or ability distribution P over \mathcal{X} . However, we would get a more tuned estimate. This does allow the bound to be independent of the problem size of which N points are in the data set, gives us a worst-case 2. Using $m_h(N)$ to quantify the number of dichotomies on N points, re-

in some slack.
1. The basic Hoeffding Inequality used in the proof already has a slack. The inequality gives the same bound whether E_{out} is close to 0.5 or close to zero. However, the variance of E_{in} is quite different in these two cases. Therefore, having one bound capture both cases will result

Why is the VC bound so loose? The slack in the bound can be attributed to a number of technical factors. Among them,

Exercise 2.5
Suppose we have a simple learning model whose growth function is $m_h(N) = N + 1$, hence $d_{\text{vc}} = 1$. Use the VC bound (2.12) to estimate the probability that E_{out} will be within 0.1 of E_{in} given 100 training examples. [Hint: The estimate will be ridiculous.]

One way to think of $\mathcal{O}(N, \mathcal{H}, \delta)$ is that it is a penalty for model complexity. The penalties us by worsening the bound on E^{out} when we use a more complex H (larger $d_{\mathcal{H}}$). If someone manages to fit a simpler model with the same training data, then it is a penalty for model complexity. The result is that the bound on E^{out} is just a sample estimate like E^{in} . How do we know our estimate of E^{out} , we are in fact asserting that E^{test} generalizes very well to E^{out} . After all, E^{test} is just a sample estimate like E^{in} .

Let us call the error we get on the test set E^{test} . When we report E^{test} as all that approach, we are in fact asserting that E^{test} generalizes very well to E^{out} . We would like to now take a closer look at this approach.

result is taken as an estimate of E^{out} . We would like to now take a closer look training process. The final hypothesis g is evaluated on the test set, and the is to estimate E^{out} by using a test set, a data set that was not involved in the An alternative approach that we alluded to in the beginning of this chapter system is expected to perform.

you need a more accurate estimate so that your customer knows how well the an accurate forecast of E^{out} . If you are developing a system for a customer, a guideline for the training process, it is next to useless if the goal is to get out-of-sample error E^{out} based on E^{in} . While the estimate can be useful as As we have seen, the generalization bound gives us a loose estimate of the

2.2.3 The Test Set

combinatation of the two terms, as illustrated informally in Figure 2.3. Although $\mathcal{O}(N, \mathcal{H}, \delta)$ goes up when \mathcal{H} has a higher VC dimension, E^{in} is and hurt $\mathcal{O}(N, \mathcal{H}, \delta)$. The optimal model is a compromise that minimizes a to fit the data. Therefore, we have a tradeoff: more complex models help E^{in} likely to go down with a higher VC dimension as we have more choices within \mathcal{H} we have more training examples, as we would expect.

gets worse if we insist on higher confidence (lower δ), and it gets better when

error, they will get a more favorable estimate for E^{out} . The penalty $\mathcal{O}(N, \mathcal{H}, \delta)$

ment ($\delta = 0.1$). We could ask what error bar can we offer with this confidence-

Example 2.7. Suppose that $N = 100$ and we have a 90% confidence require-

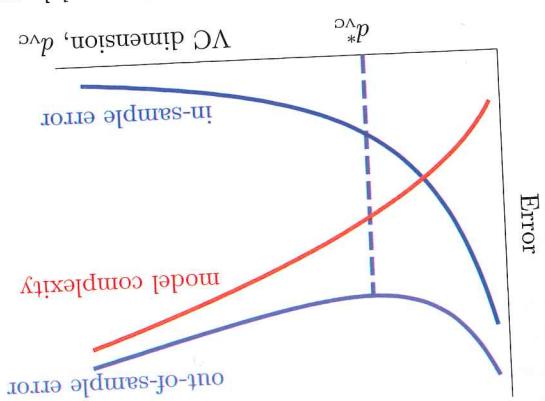


Figure 2.3: When we use a more complex learning model, one that has higher VC dimension $d_{\mathcal{H}}$, we are likely to fit the training data better relative error, thus attains a minimum at some intermediate $d_{\mathcal{H}}$. A combination of the two, which estimates the out-of-sample complexity. A combination of the two, but we pay a higher penalty for model fitting in a lower-in-sample error, but we pay a higher penalty for model complexity. In most practical situations, however, we are given a fixed data set D , so N is also fixed. In this case, the relevant question is what performance can we expect given this particular N . The bound in (2.12) answers this question:

Sample complexity fixes the performance parameters ϵ (generalization error) and δ (confidence parameter) and estimates how many examples N are needed. In most practical situations, however, we are given a fixed data set D , so N is also fixed. In this case, the relevant question is what performance can we expect given this particular N . The bound in (2.12) answers this question:

2.2.2 Penalty for Model Complexity

If we use the polynomial bound based on $d_{\mathcal{H}}$ instead of $m_{\mathcal{H}}(2N)$, we get

$$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question: in most practical situations, however, we are given a fixed data set D , so N is also fixed. In this case, the relevant question is what performance can we expect given this particular N . The bound in (2.12) answers this question:

Let us look more closely at the two parts that make up the bound on E^{out} in (2.12). The first part is E^{in} , and the second part is a term that increases as the VC dimension of \mathcal{H} increases.

Let us look more closely at the two parts that make up the bound on E^{out} in (2.12).

□ $E^{\text{out}}(g) \leq E^{\text{in}}(g) + \mathcal{O}(N, \mathcal{H}, \delta)$

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + 0.301$

with confidence $\geq 90\%$. This is a pretty poor bound on E^{out} . Even if $E^{\text{in}} = 0$,

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{100} \ln \left(\frac{4(201)}{0.1} \right)} \approx E^{\text{in}}(g) + 0.848$

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + 0.848$

if \mathcal{H} has $d_{\mathcal{H}} = 1$. Using (2.14), we have

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4(2N)^{d_{\mathcal{H}}} + 1}{\delta} \right)}$

another valid bound on the out-of-sample error,

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N)}{\delta} \right)}$

with probability at least $1 - \delta$, given this particular N . The bound in (2.12) answers this question:

$E^{\text{out}}(g) \leq E^{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \left(\frac{4m_{\mathcal{H}}(2N$

These definitions make E_m a sample estimate of E_{out} , just as it was in the case of binary functions. In fact, the error measure used for binary functions can also be expressed as a squared error.

$$\cdot \zeta (({}^u x)f - ({}^u x)u) \sum_{N=1}^{l=u} \frac{N}{l} = (u) {}^{ul} H$$

data set,

$$\cdot \left[\zeta((\mathbf{x})f - (\mathbf{x})y) \right] \mathbb{E} = (y)^{\text{uno}} E$$

An error measure that is commonly used in this case is the squared error $e(h(x), f(x)) = (h(x) - f(x))^2$. We can define in-sample and out-of-sample versions of this error measure. The out-of-sample error is based on the expected value of the error measure over the entire input space \mathcal{X} ,

In order to deal with real-valued functions, we need to adapt the definitions of E_{in} and E_{out} that have so far been based on binary functions. We defined E_{in} and E_{out} in terms of binary error; either $h(x) = f(x)$ or else $h(x) \neq f(x)$. If $f(x)$ and $h(x)$ are from each other, rather than just whether their values are and h are real-valued, a more appropriate error measure would gauge how far apart they are.

Although the VC analysis was based on binary target functions, it can be extended to real-valued functions, as well as to other types of functions. The proof in those cases are quite technical, and they do not add to the insight that the VC analysis of binary functions provides. Therefore, we will introduce an alternative approach that covers real-valued functions and provides new insights into generalization. The approach is based on bias-variance analysis,

2.2.4 Other Target Types

In some of the learning literature, E_{test} is used as synonymous with E_{out} . When we report experimental results in this book, we will often treat E_{test} based on a large test set as if it was E_{out} because of the closeness of the two quantities.

of the data points provided by the customer as a test set, we end up using fewer examples for training. Since the training set is used to select one of the hypotheses in \mathcal{H} , training examples are essential to finding a good hypothesis. If we take a big chunk of the data for testing and end up with too few examples, we may not get a good hypothesis from the training part even if we can relatively evaluate it in the testing part. We may end up reporting to the customer, with high confidence mind you, that the g we are delivering is terrible. There is thus a tradeoff to setting aside test examples. We will address this tradeoff in more detail and learn some clever tricks to get around it.

There is a price to be paid for having a test set. The test set does not affect the outcome of our learning process, which only uses the training set. The test set just tells us how well we did. Therefore, if we set aside some

Another aspect that distinguishes the test set from the training set is that the test set is not biased. Both sets are finite samples that are bound to have some variance due to sample size, but the test set doesn't have an optimistic or pessimistic bias in its estimate of F_{out} . The training set has an optimistic bias, since it was used to choose a hypothesis that looks good on it. The VC generalization bound implicitly takes that bias into consideration, and that's why it gives a huge error bar. The test set just has straight finite-sample variance, but no bias. When you report the value of E_{test} to your customer and they try your system on new data, they are as likely to be pleasantly surprised as unpleasantly surprised, though quite likely not to be surprised at all.

(d) Is there any reason why you shouldn't reserve even more examples for testing?

(a) Using a 5% error tolerance ($\delta = 0.05$), which estimate has the higher error bar?

A data set has 600 examples. To properly test the performance of the final hypothesis, you set aside a randomly selected subset of 200 examples which are never used in the training phase; these form a test set. You use a learning model with 1,000 hypotheses and select the final hypothesis based on the 400 training examples. We wish to estimate $E_{out}(g)$. We have access to two estimates: $E_{in}(g)$, the in-sample error on the 400 training examples; and, $E_{test}(g)$, the test error on the 200 test examples that were set aside.

Exercise 2.6

Therefore, the generalization bound that applies to E_{test} is the simple Hoeffding Inequality with one hypothesis. This is a much tighter bound than the VC bound. For example, if you have 1,000 data points in the test set, E_{te} will be within 5% of E_{out} with probability $\geq 98\%$. The bigger the test set you use, the more accurate E_{test} will be as an estimate of E_{out} .

The alternative number of hypotheses that matters in the generalization behavior of E_{test} is 1. There is only one hypothesis as far as the test set concerned, and that's the final hypothesis g that the training phase produce. This hypothesis would not change if we used a different test set as it would be used a different training set. Therefore, the simple Hoefding Inequality valid in the case of a test set, it would not be considered a test set any more as set in any shape or form, it wouldn't be affected by the test set any more at all.

that we have developed the theory of generalization in concrete mathematical terms.