

Machine learning 2 - Final Assignment

Frans de Boer

frans.deboer

5661439

June 28, 2022

The git repository for this assignment (including code and the report) can be found at github.com. I will make this repository public after the Brightspace deadline has passed.

1 PAC Learning

1.a Extend the proof that was given in the slides for the PAC-learnability of hyper-rectangles: show that axis-aligned hyper-rectangles in n -dimensional feature spaces ($n > 2$) are PAC learnable.

We can extend the proof by taking the number of sides/faces of the hyper-rectangle into account.

- A hyper-rectangle in 2-dimensional space is simply a rectangle and has 4 sides.
- A hyper-rectangle in 3-dimensional space is a rectangular box and has 6 sides.

For now I will assume the number of sides on a n -dimensional hyper-rectangle to be $F(n)$, and I will get back later to the calculation of $F(n)$.

To start we can follow the steps from the lecture. We have the ground truth n -dimensional hyper-rectangle R , and the current tightest fit rectangle R' . The error $(R - R')$ can be split into $F(n)$ strips T' . On each of these strips we can now grow a new strip T such that the probability mass of T is $\frac{\epsilon}{F(n)}$.

If T covers T' for all strips then the probability of an error is

$$P[\text{error}] \leq \sum_{i=0}^{F(n)-1} P[T_i] = F(n) \frac{\epsilon}{F(n)} = \epsilon \quad (1)$$

We can now estimate the probability that T does not cover T'

$$\begin{aligned}
P[\text{random } x \text{ hits } T] &= \frac{\epsilon}{F(n)} \\
P[\text{random } x \text{ misses } T] &= 1 - \frac{\epsilon}{F(n)} \\
P[m \text{ random } x\text{'s misses } T] &= \left(1 - \frac{\epsilon}{F(n)}\right)^m
\end{aligned}$$

Since we have $F(n)$ strips:

$$\begin{aligned}
P[m \text{ random } x\text{'s miss any } Ts] &\leq F(n) \left(1 - \frac{\epsilon}{F(n)}\right)^m \\
P[R' \text{ has larger error than } \epsilon] &\leq F(n) \left(1 - \frac{\epsilon}{F(n)}\right)^m < \delta
\end{aligned}$$

Bounding the chance that our R' has an error larger than ϵ by δ :

$$F(n) \left(1 - \frac{\epsilon}{F(n)}\right)^m < \delta$$

Using $e^{-x} \geq (1 - x)$

$$F(n) e^{-m\epsilon/F(n)} \geq F(n) \left(1 - \frac{\epsilon}{F(n)}\right)^m$$

So instead we can use:

$$\begin{aligned}
F(n) e^{-m\epsilon/F(n)} &< \delta \\
-m\epsilon/F(n) &< \log(\delta/F(n)) \\
m\epsilon/F(n) &> \log(F(n)/\delta) \\
m &> (F(n)/\epsilon) \log(F(n)/\delta)
\end{aligned}$$

So any n -dimensional hyper-rectangle is learnable.

1.b Assume we have a 2-dimensional feature space \mathbb{R}^2 , and consider the set of concepts that are L1-balls: $c = \{(x, y) : |x| + |y| \leq r\}$ (basically, all L1-balls centered around the origin). Use a learner that fits the tightest ball. Show that this class is PAC-learnable from training data of size m .

Same thing? what. literally just question a but now rotated 45 degrees

1.c Now we extend the previous class by allowing the varying center: $c = \{(x, y) : |x - x_0| + |y - y_0| \leq r\}$. Is this class PAC-learnable? Proof if it is, or otherwise disproof it.

maybe, maybe not. can't you still just do the tightest fit L1 ball? consistent learner so PAC learnable?

2 VC Dimension

Let us assume we are dealing with a two-class classification problem in d -dimensions. Let us start off with refreshing our memories. Consider the class of linear hypotheses (i.e., all possible linear classifiers) in d -dimensional space.

2.a Given N data points, how many possible labelings does this data set have?

Each data point can have 2 labels, so 2^n

2.b What is the dimensionality d we need, at the least, for our class of linear hypotheses to be able to find perfect solutions to all possible labelings of a data set of size N ? Stated differently, what is the smallest d that allows us to shatter N points?

$d - 1$

2.c Consider $d = 1$ and a data set of size N . At maximum, how many different labelings of such a data set can decision stumps solve perfectly, i.e., with zero training error?

This would be all cases where there is no overlap, i.e. all points with class 1 (or 2) are either all on the left or right side. See Figure 1 for an example with $N=3$. Note that the first 3 and last 3 examples are the same, just with inverted labelings.

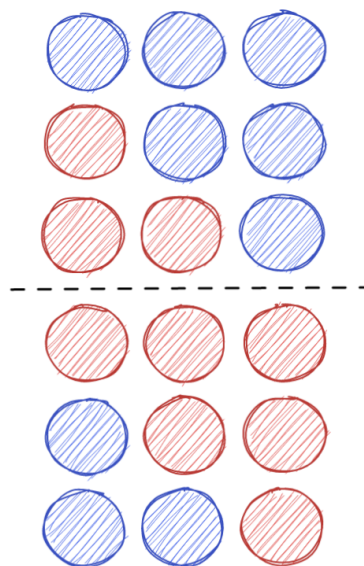


Figure 1: Example of perfectly solvable labels with $N = 3$

With a dataset of size N , there are N ways in which we can have class 1 be on the left side and have a decision stump be able to solve it perfectly. There can be $1, 2, \dots, N$ points with class 1 in a row. Because each of these possible label assignments has a complement (just invert the classes), we end up with a total of $2N$ possible labelings that a decision stump can solve perfectly.

2.d With the answer from 3.c, give an upper-bound on the maximum number of different labelings a data set of size N in d dimensions can get by means of decision stumps.

We can once again look at Figure 1, but now instead of true labels we can view the labels (red and blue) as labels predicted by the decision stump. The decision stump either assigns all points to the left of it as class 1 or as class 2. So in $d = 1$ we know the maximum number of different labelings a data set of size N can get is $2N$.

Now to generalize this case to any dimension d . A possible problem with higher dimension spaces is that not each stump can differentiate all N points. See Figure 2 as an example. Since we are looking for the upper bound from now on I will make the assumption that all points in the data set can be differentiated by the decision stump, i.e. for $d = 2$

$$x_i \neq x_j \wedge y_i \neq y_j \forall i, j$$

See Figure 3 for an example.

Figure 2: Example of dataset that can't be differentiated by a stump in $d = 2$ (The current stump can't differentiate 1 and 2)

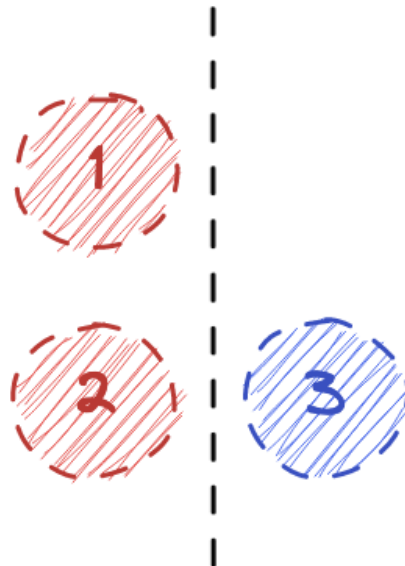
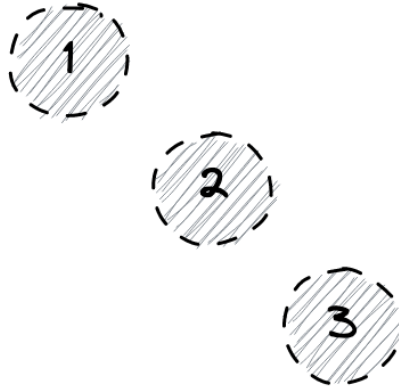


Figure 3: Example of dataset that can be differentiated by a stump in $d = 2$



A stump can only choose 1 variable (dimension) to make its decision. This means that per variable a stump can create $2N$ possible labelings (Section 2.c). In d dimensions we can create d possible stumps, so in total there are $2Nd$ possible labelings.

So the upper-bound on the maximum number of different labels a data set of size N in d dimensions can get by means of decision stumps is $2Nd$.

2.e Using the bound from Exercise 3.d, determine the smallest upper-bound on the VC-dimension for decision stumps for dimensionalities $d \in \{1, 2, 3, 4, 5, 6, 7, 8, 1024, 2^{100}\}$

VC-Dimension: The largest number of nodes N that a classifier can shatter in all possible 2^N ways for any possible data set. **Smallest upper-bound:** First size N in dimension d where decision stumps cannot shatter any data set

As was shown in Section 2.d the largest number of different labels a set of size N in d dimensions can get by means of decision stumps is $2Nd$. This means that for a decision stump to be able to shatter a data set of size N in dimension d the number of possible label assignments 2^N needs to be less than or equal to $2Nd$, so

$$2^N \leq 2Nd$$

For $d = 1$ the VC dimension is 2. Any more than this and the decision stump won't be able to shatter the dataset by simply alternating the ground truth class assignments.

Assume this is an N where $2^N = 2Nd$ for $d = 2$