# 1  The accuracy of learning decision boundaries

This problem exercises your basic probability (e.g. from 70) in the context of understanding why lots of training data helps to improve the accuracy of learning things.

For each $\theta \in (1/3, 2/3)$, define $f_\theta : [0, 1] \to \{0, 1\}$, such that

$$f_\theta(x) = \begin{cases} 1 \text{ if } x > \theta \\ 0 \text{ otherwise.} \end{cases}$$
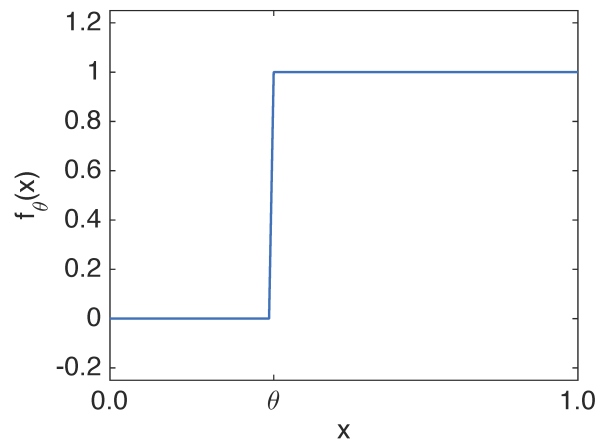
The function is plotted in Figure 1.



Figure 1: Plot of function $f_\theta(x)$ against $x$.

We draw samples $X_1, X_2, \ldots, X_n$ uniformly at random and i.i.d. from the interval $[0, 1]$. Our goal is to learn an estimate for $\theta$ from $n$ random samples $(X_1, f_\theta(X_1)), (X_2, f_\theta(X_2)), \ldots, (X_n, f_\theta(X_n))$.

Let $T_{min} = \max(\{\frac{1}{3}\} \cup \{X_i | f_\theta(X_i) = 0\})$. We know that the true $\theta$ must be larger than $T_{min}$.

Let $T_{max} = \min(\{\frac{2}{3}\} \cup \{X_i | f_\theta(X_i) = 1\})$. We know that the true $\theta$ must be smaller than $T_{max}$.

The gap between $T_{min}$ and $T_{max}$ represents the uncertainty we will have about the true $\theta$ given the training data that we have received.

(a) **What is the probability that $T_{max} - \theta > \epsilon$ as a function of $\epsilon$? And what is the probability**

**that** $\theta - T_{min} > \epsilon$ **as a function of** $\epsilon$?

(b) Suppose that you would like the estimator $\hat{\theta} = (T_{max} + T_{min})/2$ for $\theta$ that is $\epsilon$-close (defined as $|\hat{\theta} - \theta| < \epsilon$, where $\hat{\theta}$ is the estimation and $\theta$ is the true value) with probability at least $1 - \delta$. Both $\epsilon$ and $\delta$ are some small positive numbers. **Please bound or estimate how big of an** $n$ **do you need?** You do not need to find the optimal lowest sample complexity $n$, an approximation using results of question (a) is fine.

(c) Let us say that instead of getting random samples $(X_i, f(X_i))$, we were allowed to choose where to sample the function, but you had to choose all the places you were going to sample in advance. **Propose a method to estimate** $\theta$. **How many samples suffice to achieve an estimate that is** $\epsilon$**-close as above?** (**Hint:** You need not use a randomized strategy.)

(d) Suppose that you could pick where to sample the function adaptively — choosing where to sample the function in response to what the answers were previously. **Propose a method to**

**estimate $\theta$. How many samples suffice to achieve an estimate that is $\epsilon$-close as above?**

(e) In the three sampling approaches above: random, deterministic, and adaptive, **compare the scaling of $n$ with $\epsilon$ (and $\delta$ as well for the random case).**

(f) **Why do you think we asked this series of questions? What are the implications of those results in a machine learning application?**

# 2 The Classical Bias-Variance Tradeoff

Consider a random variable $X$, which has unknown mean $\mu$ and unknown variance $\sigma^2$. Given $n$ iid realizations of training samples $X_1 = x_1, X_2 = x_2, \ldots, X_n = x_n$ from the random variable, we wish to estimate the mean of $X$. We will call our estimate of $\mu$ the random variable $\hat{X}$, which has mean $\hat{\mu}$. There are a few ways we can estimate $\mu$ given the realizations of the $n$ samples:

1. Average the $n$ samples: $\frac{x_1+x_2+\ldots+x_n}{n}$.

2. Average the $n$ samples and one sample of 0: $\frac{x_1+x_2+\ldots+x_n}{n+1}$.

3. Average the $n$ samples and $n_0$ samples of 0: $\frac{x_1+x_2+\ldots+x_n}{n+n_0}$.

4. Ignore the samples: just return 0.

In the parts of this question, we will measure the *bias* and *variance* of each of our estimators. The *bias* is defined as

$$\mathbb{E}[\hat{X} - \mu]$$

and the *variance* is defined as

$$\text{Var}[\hat{X}].$$

(a) What is the bias of each of the four estimators above?

(b) What is the variance of each of the four estimators above?

(c) Suppose we have constructed an estimator $\hat{X}$ from some samples of $X$. We now want to know how well $\hat{X}$ estimates a new independent sample of $X$. Denote this new sample by $X'$. Derive a general expression for $\mathbb{E}[(\hat{X} - X')^2]$ in terms of $\sigma^2$ and the bias and variance of the estimator $\hat{X}$. Similarly, derive an expression for $\mathbb{E}[(\hat{X} - \mu)^2]$. Compare the two expressions and comment on the differences between them.

(d) It is a common mistake to assume that an unbiased estimator is always "best." Let's explore this a bit further. Compute $E[(\hat{X} - \mu)^2]$ for each of the estimators above.

(e) Demonstrate that the four estimators are each just special cases of the third estimator, but with different instantiations of the hyperparameter $n_0$.

(f) What happens to bias as $n_0$ increases? What happens to variance as $n_0$ increases?