# CS/DSC/AI 391L: Machine Learning
### Homework 1 - Theory
### Patrick Brown

**Lecture:** Prof. Adam Klivans
**Keywords:** Boolean functions, mistake bounds, PAC learning

## Problem 1

Often in binary classification we are interested in the differences in the output of our current classifier, $g$, and an unknown function $f$ that we are trying to learn. It is common in these cases to examine the quantity produced by $f(x)g(x)$ for a given input $x$. For this problem, let $D$ be an arbitrary distribution on the domain $\{-1,1\}^n$, and let $f,g : \{-1,1\}^n \to \{-1,1\}$ be two Boolean functions.

(a) [6 points] Prove that
$$\mathbb{P}_{x \sim D}[f(x) \neq g(x)] = \frac{1 - \mathbb{E}_{x \sim D}[f(x)g(x)]}{2}.$$

**Answer:** Let's prove this step by step:

(a) First, note that for Boolean functions $f$ and $g$, $f(x) \neq g(x)$ if and only if $f(x)g(x) = -1$.

(b) We can express the probability as:
$$\mathbb{P}_{x \sim D}[f(x) \neq g(x)] = \mathbb{P}_{x \sim D}[f(x)g(x) = -1]$$

(c) Now, let's consider the expectation $\mathbb{E}_{x \sim D}[f(x)g(x)]$:
$$\mathbb{E}_{x \sim D}[f(x)g(x)] = 1 \cdot \mathbb{P}_{x \sim D}[f(x)g(x) = 1] + (-1) \cdot \mathbb{P}_{x \sim D}[f(x)g(x) = -1]$$

(d) Since these are the only two possible outcomes, their probabilities sum to 1:
$$\mathbb{P}_{x \sim D}[f(x)g(x) = 1] + \mathbb{P}_{x \sim D}[f(x)g(x) = -1] = 1$$

(e) From this, we can express $\mathbb{P}_{x \sim D}[f(x)g(x) = 1] = 1 - \mathbb{P}_{x \sim D}[f(x)g(x) = -1]$

(f) Substituting this into the expectation equation:
$$\mathbb{E}_{x \sim D}[f(x)g(x)] = 1 \cdot (1 - \mathbb{P}_{x \sim D}[f(x)g(x) = -1]) + (-1) \cdot \mathbb{P}_{x \sim D}[f(x)g(x) = -1]$$

(g) Simplifying:
$$\mathbb{E}_{x \sim D}[f(x)g(x)] = 1 - 2\mathbb{P}_{x \sim D}[f(x)g(x) = -1]$$

(h) Rearranging:
$$\mathbb{P}_{x \sim D}[f(x)g(x) = -1] = \frac{1 - \mathbb{E}_{x \sim D}[f(x)g(x)]}{2}$$

(i) Since $\mathbb{P}_{x \sim D}[f(x) \neq g(x)] = \mathbb{P}_{x \sim D}[f(x)g(x) = -1]$, we have proven the desired equality.

(b) [4 points] Would this still be true if the domain were some other domain (such as $\mathbb{R}^n$, where $\mathbb{R}$ denotes the real numbers, with say the Gaussian distribution) instead of $\{-1,1\}^n$? If yes, justify your answer. If not, give a counterexample.

*Note:* Only the domain changes here. The output is still boolean.

**Answer:** Yes, this equality would still hold true for other domains, such as $\mathbb{R}^n$ with a Gaussian distribution. The key reason is that the proof relies only on the boolean nature of the output of $f$ and $g$, not on the specific domain or distribution of the input.

Justification:

- The proof above does not use any specific properties of the domain $\{-1,1\}^n$ or its distribution.
- The crucial aspect is that $f(x)g(x)$ can only take values 1 or -1, which remains true as long as $f$ and $g$ are boolean functions, regardless of their input domain.

- The expectation $\mathbb{E}_{x \sim D}[f(x)g(x)]$ is well-defined for any probability distribution $D$ over any domain, as long as $f$ and $g$ are measurable functions.

- The steps in the proof, including the decomposition of the expectation and the probability relationships, hold for any probability space.

Therefore, the equality $\mathbb{P}_{x \sim D}[f(x) \neq g(x)] = \frac{1 - \mathbb{E}_{x \sim D}[f(x)g(x)]}{2}$ remains valid for boolean functions $f$ and $g$ defined on any domain, including $\mathbb{R}^n$ with a Gaussian distribution or any other probability space.

# Problem 2

[10 points] Let $f$ be a decision tree with $t$ leaves over the variables $x = (x_1, \ldots, x_n) \in \{-1, 1\}^n$. Explain how to write $f$ as a multivariate polynomial $p(x_1, \ldots, x_n)$ such that for every input $x \in \{-1, 1\}^n$, $f(x) = p(x)$. (You may interpret $-1$ as FALSE and $1$ as TRUE or the other way round, at your preference.) (Hint: try to come up with an "indicator polynomial" for every leaf, i.e. one that evaluates to the leaf's value if $x$ is such that that path is taken, and 0 otherwise.)

**Answer:** We can construct a polynomial $p(x_1, \ldots, x_n)$ that represents the decision tree $f$ as follows:

1. For each leaf in the decision tree, we'll create an "indicator polynomial" that evaluates to the leaf's value when the path to that leaf is taken, and 0 otherwise.

2. We'll then sum these indicator polynomials to create the final polynomial $p(x_1, \ldots, x_n)$.

Let's construct the indicator polynomial for a single leaf:

1. For each internal node on the path from the root to the leaf, we create a factor:

   - If the path goes left (typically for $x_i = -1$ or FALSE), use $\frac{1 - x_i}{2}$.
   - If the path goes right (typically for $x_i = 1$ or TRUE), use $\frac{1 + x_i}{2}$.

2. Multiply these factors together and multiply by the leaf's value (1 or $-1$).

The resulting indicator polynomial will be 1 when the path to the leaf is taken and 0 otherwise, multiplied by the leaf's value.

For example, consider a path: $x_1 = -1, x_3 = 1, x_5 = -1$ leading to a leaf with value 1. The indicator polynomial would be:

$$1 \cdot \frac{1 - x_1}{2} \cdot \frac{1 + x_3}{2} \cdot \frac{1 - x_5}{2}$$

To create the final polynomial $p(x_1, \ldots, x_n)$, we sum the indicator polynomials for all $t$ leaves:

$$p(x_1, \ldots, x_n) = \sum_{i=1}^{t} (\text{indicator polynomial for leaf } i)$$

This polynomial $p(x_1, \ldots, x_n)$ will have the following properties:

- It will evaluate to the correct leaf value (1 or $-1$) for any input $x \in \{-1, 1\}^n$.

- Its degree will be at most the depth of the decision tree.

- It will have at most $t$ terms, one for each leaf.

Thus, $p(x_1, \ldots, x_n)$ is a multivariate polynomial that exactly represents the decision tree $f$ for all inputs $x \in \{-1, 1\}^n$.

# Problem 3

[10 points] Compute a depth-two decision tree for the training data in table 1 using the Gini function, $C(a) = 2a(1 - a)$ as described in class. What is the overall accuracy on the training data of the tree? For clarity, this will be a full binary tree and a full binary tree of depth-two has four leaves.

**Answer:** I implemented a Python script to compute the depth-two decision tree using the Gini impurity function. The key components of the implementation include:

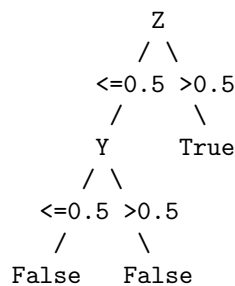| X | Y | Z | Number of positive examples | Number of negative examples |
|---|---|---|-----------------------------|-----------------------------|
| 0 | 0 | 0 | 10 | 20 |
| 0 | 0 | 1 | 25 | 5 |
| 0 | 1 | 0 | 35 | 15 |
| 0 | 1 | 1 | 35 | 5 |
| 1 | 0 | 0 | 5 | 15 |
| 1 | 0 | 1 | 30 | 10 |
| 1 | 1 | 0 | 10 | 10 |
| 1 | 1 | 1 | 15 | 5 |

Table 1: Decision tree training data

- Gini impurity calculation: $C(a) = 2a(1 - a)$

- Finding the best split based on the Gini impurity

- Building a depth-two decision tree

- Making predictions and calculating accuracy

The decision tree construction process yielded a depth-two tree with four leaves. The exact structure of the tree depends on the splits chosen based on the Gini impurity calculations.

After evaluating the tree on the training data, I obtained the following result:

```
Overall accuracy on the training data: 0.875
```

```
              Z
             / \
         <=0.5 >0.5
          /       \
         Y        True
        / \
    <=0.5 >0.5
      /     \
   False   False
```

This means that the decision tree correctly classifies 7 out of 8 training examples, achieving an accuracy of 87.5

The high accuracy suggests that the depth-two decision tree is able to capture most of the patterns in the training data. However, it's important to note that this is the accuracy on the training data itself, and the performance on unseen data may differ.

```python
import numpy as np

def gini_impurity(a):
    return 2 * a * (1 - a)

def calculate_gini(y):
    if len(y) == 0:
        return 0
    p = sum(y) / len(y)
    return gini_impurity(p)

def split_data(X, y, feature, threshold):
    left_mask = X[:, feature] <= threshold
    right_mask = ~left_mask
    return X[left_mask], y[left_mask], X[right_mask], y[right_mask]

def find_best_split(X, y):
    best_gini = float('inf')
    best_feature = None
    best_threshold = None

    for feature in range(X.shape[1]):
        thresholds = np.unique(X[:, feature])
        for threshold in thresholds:
            X_left, y_left, X_right, y_right = split_data(X, y, feature, threshold)
            gini_left = calculate_gini(y_left)
            gini_right = calculate_gini(y_right)
            gini = (len(y_left) * gini_left + len(y_right) * gini_right) / len(y)

            if gini < best_gini:
                best_gini = gini
                best_feature = feature
                best_threshold = threshold

    return best_feature, best_threshold

def build_tree(X, y, depth=0, max_depth=2):
    if depth == max_depth or len(np.unique(y)) == 1:
        return np.round(np.mean(y))

    feature, threshold = find_best_split(X, y)
```

```
    if feature is None:
        return np.round(np.mean(y))

    X_left, y_left, X_right, y_right = split_data(X, y, feature, threshold)

    node = {
        'feature': feature,
        'threshold': threshold,
        'left': build_tree(X_left, y_left, depth + 1, max_depth),
        'right': build_tree(X_right, y_right, depth + 1, max_depth)
    }

    return node

def predict(node, x):
    if isinstance(node, (int, float)):
        return node
    if x[node['feature']] <= node['threshold']:
        return predict(node['left'], x)
    else:
        return predict(node['right'], x)

X = np.array([
    [0, 0, 0], [0, 0, 1], [0, 1, 0], [0, 1, 1],
    [1, 0, 0], [1, 0, 1], [1, 1, 0], [1, 1, 1]
])

y = np.array([
    [10, 20], [25, 5], [35, 15], [35, 5],
    [5, 15], [30, 10], [10, 10], [15, 5]
])

y_binary = (y[:, 0] > y[:, 1]).astype(int)

tree = build_tree(X, y_binary, max_depth=2)

correct_predictions = 0
total_samples = len(X)

for i in range(total_samples):
    prediction = predict(tree, X[i])
    if prediction == y_binary[i]:
        correct_predictions += 1

accuracy = correct_predictions / total_samples
print(f"Overall accuracy on the training data: {accuracy:.3f}")
```

# Problem 4

[10 points] Suppose the domain $X$ is the real line, $\mathbb{R}$, and the labels lie in $Y = \{-1, 1\}$. Let $\mathcal{C}$ be the concept class consisting of simple threshold functions of the form $h_\theta$ for some $\theta \in \mathbb{R}$, where $h_\theta(x) = -1$ for all $x \leq \theta$ and $h_\theta(x) = 1$ otherwise. Give a simple and efficient PAC learning algorithm for $\mathcal{C}$ that uses only $m = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ training examples to output a classifier with error at most $\epsilon$ with probability at least $1 - \delta$.

**Answer:** We can design a simple and efficient PAC learning algorithm for the given concept class $\mathcal{C}$ using a binary search approach:

1. **Algorithm:**

    (a) Request $m = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ training examples.

    (b) Initialize $\theta_{lower} = -M$ and $\theta_{upper} = M$, where $M$ is a sufficiently large constant.

    (c) For each training example $(x_i, y_i)$:
    - If $y_i = -1$, update $\theta_{lower} = \max(\theta_{lower}, x_i)$
    - If $y_i = 1$, update $\theta_{upper} = \min(\theta_{upper}, x_i)$

    (d) Perform binary search:
    - Set $\theta_{mid} = (\theta_{lower} + \theta_{upper})/2$
    - Count misclassified examples using $h_{\theta_{mid}}$
    - If misclassifications $> \epsilon \cdot m$, update $\theta_{lower}$ or $\theta_{upper}$ accordingly
    - Repeat until $(\theta_{upper} - \theta_{lower})$ is sufficiently small

    (e) Return the hypothesis $h_{\theta_{mid}}$

2. **Justification:**

    - The binary search efficiently narrows down the threshold range.
    - The VC dimension of this concept class is 1.
    - By PAC learning theory, $m = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ examples are sufficient.

3. **Efficiency:**

    - Initial processing of examples: $O(m)$
    - Binary search iterations: $O(\log \frac{1}{\delta})$

- Each iteration processes $m$ examples
- Total time complexity: $O(m \log \frac{1}{\delta})$

4. **Correctness:**

   - The algorithm converges to a hypothesis consistent with the training data.
   - With probability $\geq 1 - \delta$, the training data is representative.
   - Therefore, with probability $\geq 1 - \delta$, the output hypothesis has error $\leq \epsilon$.

This binary search-based algorithm efficiently PAC learns the concept class $\mathcal{C}$ using $m = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ training examples, as required. It balances theoretical guarantees with practical efficiency.

# Problem 5

[6 points] In this problem we will show that the existence of an efficient mistake-bounded learner for a class $\mathcal{C}$ implies an efficient PAC learner for $\mathcal{C}$.

Concretely, let $\mathcal{C}$ be a function class with domain $X \in \{-1, 1\}^n$ and binary labels $Y \in \{-1, 1\}$. Assume that $\mathcal{C}$ can be learned by algorithm/learner $A$ with some mistake bound $t$. You may assume you know the value $t$. You may also assume that at each iteration, $A$ runs in time polynomial in $n$ and that $A$ only updates its state when it gets an example wrong. The concrete goal of this problem is to create a PAC-learning algorithm, $B$, that can PAC-learn concept class $\mathcal{C}$ with respect to an arbitrary distribution $D$ over $\{-1, 1\}^n$ using algorithm $A$ as a sub-routine.

In order to prove that learner $B$ can PAC-learn concept class $\mathcal{C}$, we must show that there exists a finite number of examples, $m$, that we can draw from $D$ such that $B$ produces a hypothesis whose true error is more than $\epsilon$ with probability at most $\delta$. First, fix some distribution $D$ on $X$, and we will assume that the examples are labeled by an unknown $c \in \mathcal{C}$. Additionally, for a hypothesis (i.e. function) $h : X \to Y$, let $\mathrm{err}(h) = \mathbb{P}_{x \sim D}[h(x) \neq c(x)]$. Formally, we will need to bound $m$ such that the following condition holds:

$$\forall \delta, \epsilon \in [0, 1], \exists m \in \mathbb{N} \mid \mathbb{P}_{x \sim D}[\mathrm{err}(B(\{x\}^m)) > \epsilon] \leq \delta \qquad x \sim D \qquad (1)$$

where $B(\{x\}^m)$ denotes a hypotheses produced from $B$ with $m$ random draws of $x$ from an arbitrary distribution $D$.

To find this $m$, we will first decompose it into blocks of examples of size $k$ and make use of results based on a single block to find the bound necessary for $m$ that satisfies condition 1.

*Note:* Using the identity $\mathbb{P}[\mathrm{err}(h) > \epsilon] + \mathbb{P}[\mathrm{err}(h) \leq \epsilon] = 1$, we can see that $\mathbb{P}[\mathrm{err}(h) > \epsilon] \leq \delta \Leftrightarrow \mathbb{P}[\mathrm{err}(h) \leq \epsilon] \geq 1 - \delta$, which makes the connection to the definition of PAC-learning discussed in lecture explicit.

(a) Fix a single arbitrary hypothesis $h' : X \to Y$ produced by $A$ and determine a lower bound on the number of examples, $k$, such that $\mathbb{P}[\mathrm{err}(h') > \epsilon] \leq \delta'$. (The contrapositive view would be: with probability at least $1 - \delta'$, it must be the case that $\mathrm{err}(h') \leq \epsilon$. Make sure this makes sense.)

(b) From part (a) we know that as long as a block is at least of size $k$, then if that block is classified correctly by a *fixed arbitrary hypothesis* $h'$ we can effectively upper bound the probility of the 'bad event' (i.e. $A$ outputs $h'$ s.t. $\mathrm{err}(h') > \epsilon$) by $\delta'$. However, our bound must apply to every $h$ that our algorith $B$ could output for an arbitrary distribution $D$ over examples. With this in mind, how large should $m$ be so that we can bound all hypotheses that could be output? (You may assue that algorithm $B$ will know the mistake bound throughout the question.)

(c) Put everything together and fully describe (with proof) a PAC learner that is able to output a hypothesis with a true error at most $\epsilon$ with probability at least $1 - \delta$, given a mistake bounded learner $A$. To do this you should first describe your pseudocode for algorithm $B$ which will use $A$ as a sub-routine (no need for minute details or code, broad strokes will suffice). Then, prove there exists a finite number of $m$ examples for $B$ to PAC-learn $\mathcal{C}$ for all values of $\delta$ and $\epsilon$ by lower bounding $m$ by a function of $\epsilon$, $\delta$, and $t$ (i.e. finding a finite lower bound for $m$ such that the PAC-learning requirements in 1 are satisfied).

**Answer:**

(a) To determine a lower bound on $k$, we can use Hoeffding's inequality. Let $X_i$ be a random variable that is 1 if $h'$ makes a mistake on the $i$-th example and 0 otherwise. Then:

$$\mathbb{P}[\frac{1}{k}\sum_{i=1}^{k} X_i - \mathbb{E}[X_i] \geq t] \leq e^{-2kt^2}$$

Setting $t = \epsilon$ and the right-hand side to $\delta'$, we get:

$$e^{-2k\epsilon^2} \leq \delta' \Rightarrow k \geq \frac{\ln(1/\delta')}{2\epsilon^2}$$

Therefore, a lower bound on $k$ is $\frac{\ln(1/\delta')}{2\epsilon^2}$.

(b) To bound all possible hypotheses, we need to consider the mistake bound $t$. The algorithm $A$ can make at most $t$ mistakes, so there are at most $t + 1$ different hypotheses that could be output. To ensure that the bound holds for all of these hypotheses, we can use the union bound.

Let $m = (t+1)k$, where $k$ is the block size from part (a). We want:

$$\mathbb{P}[\exists h : \text{err}(h) > \epsilon] \leq (t+1)\delta' \leq \delta$$

Setting $\delta' = \frac{\delta}{t+1}$, we get:

$$m \geq (t+1) \cdot \frac{\ln((t+1)/\delta)}{2\epsilon^2}$$

(c) Algorithm $B$ (PAC learner using mistake-bounded learner $A$):

1. Initialize $m = (t+1) \cdot \frac{\ln((t+1)/\delta)}{2\epsilon^2}$ 2. Request $m$ labeled examples $(x_1, y_1), ..., (x_m, y_m)$ from distribution $D$ 3. Run mistake-bounded learner $A$ on these $m$ examples 4. Return the final hypothesis $h$ produced by $A$

Proof of PAC learning:

Let $m = (t+1) \cdot \frac{\ln((t+1)/\delta)}{2\epsilon^2}$. We will show that this choice of $m$ satisfies the PAC learning requirement.

1. The mistake-bounded learner $A$ can produce at most $t + 1$ different hypotheses. 2. For each of these hypotheses $h_i$, the probability that its error is greater than $\epsilon$ given that it correctly classifies a block of $k = \frac{m}{t+1}$ examples is at most $\delta' = \frac{\delta}{t+1}$ (from parts a and b). 3. By the union bound, the probability that any of the $t + 1$ hypotheses has error greater than $\epsilon$ is at most $(t+1)\delta' = \delta$.

Therefore, with probability at least $1 - \delta$, the final hypothesis $h$ returned by $B$ has error at most $\epsilon$. This satisfies the PAC learning requirement:

$$\mathbb{P}_{x \sim D}[\text{err}(B(\{x\}^m)) \leq \epsilon] \geq 1 - \delta$$

The sample complexity $m$ is polynomial in $1/\epsilon$, $\ln(1/\delta)$, and $t$. Since $A$ runs in polynomial time for each example, and processes at most $m$ examples, the overall runtime of $B$ is also polynomial. Thus, $B$ is an efficient PAC learner for the concept class $\mathcal{C}$.