



Centro de Investigación
en Computación
Instituto Politécnico Nacional

INSTITUTO POLITÉCNICO NACIONAL

CENTRO DE INVESTIGACIÓN EN CÓMPUTO

Chapter 1: The learning problem

Subject: Introducción a machine learning

Alumno:

- Carpintero Mendoza Marcos Mauricio

Profesor:

Menchaca Mendez Ricardo

Índice

1. Exercises.	2
1.1. Problem setup.	2
1.1.1. Components of learning.	2
1.1.2. A simple learning model.	4
1.2. Is Learning Feasible?	5
1.2.1. Outside of the data set.	5
1.2.2. Probability to the rescue.	6
1.3. Error and noise	9
1.3.1. Noisy targets.	9
2. Problems.	10

1. Exercises.

1.1. Problem setup.

1.1.1. Components of learning.

Exercise 1.1 Express each of the following tasks in the framework of learning from data by specifying the input space \mathcal{X} , output space \mathcal{Y} , target function $f : \mathcal{X} \rightarrow \mathcal{Y}$, and the specifics of the data set that we will learn from.

Solution:

1. Medical diagnosis: A patient walks in with a medical history and some symptoms, and you want to identify the problem.

Solution:

- Let \mathcal{X} be the following finite feature vector in which each element is a boolean variable:

$$\mathcal{X} = \{\text{weight loss, tingling sensations, fatigue, red gums}\}$$

- Let \mathcal{Y} a binary variable:

$$\mathcal{Y} = \{\text{diabetes, fever}\}$$

- Suppose a target function, let $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$f = \begin{cases} \text{diabetes} & \text{if } \text{sum}(\mathcal{X}) \geq 3 \\ \text{fever} & \text{if } \text{sum}(\mathcal{X}) < 3 \end{cases}$$

- Finally, the sample training data set:

$$D = \{(\{0, 0, 0, 0\}, \text{fever}), (\{1, 0, 0, 0\}, \text{fever}), (\{0, 1, 0, 0\}, \text{fever}), \dots, (\{1, 1, 0, 1\}, \text{diabetes}), (\{1, 1, 1, 0\}, \text{diabetes}), (\{1, 1, 1, 1\}, \text{diabetes})\}$$

2. Handwritten digit recognition (for example postal zip code recognition for mail sorting)

- Let \mathcal{X} be the following finite feature vector in which each element is a boolean variable:

3. Determining if an email is spam or not.

- Let \mathcal{X} be the following finite vector in which each element is a typical spam word:

Note1: The words in \mathcal{X} were added after using stemming so that we avoid similar words.

$$\mathcal{X} = \{\text{discount, includ, cheap, free, cost, clerance, prize, low}\}$$

Note2: Keep in mind: $\forall x_i \in \mathcal{X} \exists k_i \in \mathbb{N} : k_i = \text{number of } x_i \text{ in the email.}$

- Let \mathcal{X} be the following finite vector.

$$\mathcal{Y} = \{\text{Spam, Nonspam}\}$$

- Suppose the next target function $f : \mathcal{X} \rightarrow \mathcal{Y}$, such that:

$$f = \begin{cases} \text{Spam} & \text{if } \sum_{i=0}^{\text{len}(\mathcal{X})} x_i \geq 5 \\ \text{Nonspam} & \text{otherwise} \end{cases}$$

- Finally, the data training set:

$$D = \{(1, 3, 0, 0, 0, 1, 3, 0, \text{Spam}), (\{1, 1, 1, 1, 0, 2, 0, 0\}, \text{Spam}), (\{1, 0, 0, 0, 0, 0, 1, 0\}, \text{Nonspam}), \\ (\{0, 1, 0, 1, 0, 1, 0, 0, 0\}, \text{Nonspam}), \dots, (\{1, 0, 0, 3, 0, 0, 0, 0\}, \text{Nonspam})\}$$

4. Predicting how an electrical load varies with the price, temperature and day of the week.

- Let X the feature vector which elements are the mentioned ones above.

$$\mathcal{X} = \{\text{Price, Temperature, Day}\}$$

- Let Y the vector with two elements:

$$\mathcal{Y} = \{\text{High, Low}\}$$

- Suppose the target function f as:

$$f = \begin{cases} \text{Low} & \text{if Day} \in \text{Weekday} \vee \text{Price} < \$15 \vee \text{Temperature} < 20C \\ \text{High} & \text{Otherwise} \end{cases}$$

- Finally, data training set:

$$D = \{(\{\text{Monday}, 20, 10\}, \text{Low}), (\{\text{Friday}, 10, 25\}, \text{Low}), \dots, \\ (\{\text{Thursday}, 20, 25\}, \text{High}), (\{\text{Wednesday}, 21, 30\}, \text{High})\}$$

5. A problem of interest to you which there is not analytic solution, but you have data from which to construct an empirical solution.

Well, the situation is: Predicting the citizens salary base on each year of high education studied (i.e. The salary increases from the first year in degree).

- Let X feature vector which elements are the years studied starting from bachelor.

$$\mathcal{X} = \{\text{Bachelor, Master, PhD, Post doctoral}\}$$

Note: If x_{i+1} exists, then x_i must be completed.

- Let Y the vector which describe the salary range.

$$\mathcal{Y} = \{\text{Super-High, High, Normal}\}$$

- Suppose the function f as:

$$f = \begin{cases} \text{Super-high} & \text{if Post-Doctoral} \geq 1 \\ \text{High} & \text{if Master} \leq 2 \\ \text{Normal} & \text{if Bachelor} \leq 3 \end{cases}$$

- Finally, data training set:

$$D = \{(\{1, 0, 0, 0\}, \text{Normal}), (\{5, 1, 0, 0\}, \text{High}), \dots, \\ (\{5, 3, 2, 0\}, \text{Super-High}), (\{5, 3, 3, 1\}, \text{Super-High})\}$$

1.1.2. A simple learning model.

Exercise 1.3 The weight update rule in (1.3) has the nice interpretation that it moves in the direction of classifying $x(t)$ correctly.

Solution:

- a) Show that $y(t)W^T(t)x(t) < 0$ [HINT: $x(t)$ is misclassified by $w(t)$]

First, by definition a correct classification:

$$\text{Sign}(W^T(t)x(t)) > 0 \rightarrow y(t) = 1$$

$$\text{Sign}(W^T(t)x(t)) \leq 0 \rightarrow y(t) = -1$$

Thus, misclassification means:

$$\text{Sign}(W^T(t)x(t)) > 0 \rightarrow y(t) = -1$$

$$\text{Sign}(W^T(t)x(t)) \leq 0 \rightarrow y(t) = 1$$

Finally, for both cases:

$$y(t)W^T(t)x(t) < 0$$

- b) Show that $y(t)W^T(t+1)x(t) > y(t)W^T(t)x(t)$ [HINT: Use (1.3)]

Since the last proof we continue:

$$\begin{aligned} y(t)W^T(t+1)x(t) &= y(t)(W(t) + y(t)x(t))^T x(t) \\ &= y(t)W(t)^T x(t) + (y(t))^2 \|x(t)\|^2 \\ &> y(t)W^T(t)x(t) \end{aligned}$$

- c) As far as classifying $x(t)$ is concerned, argue that we move from $W(t)$ to $W(t+1)$ is a move in the right direction,

As we know, a correct classification (Check a)) and the definition of dot product:

$$\begin{aligned} W(t)^T x(t) &= \|W(t)\| \|x(t)\| \cos(\alpha) \\ \cos(\alpha) &= \frac{W(t)^T x(t)}{\|W(t)\| \|x(t)\|} \rightarrow \cos(\alpha) \propto W(t)^T x(t) \\ \cos(\alpha_{new}) &\propto W(t+1)^T x(t) \end{aligned}$$

$$\cos(\alpha_{new}) \propto (W(t) + x(t))^T x(t)$$

$$\propto W(t)^T x(t) + x(t)^T x(t)$$

$$\propto \cos(\alpha) + x(t)^T x(t)$$

$$\cos(\alpha_{new}) > \cos(\alpha)$$

$$\cos(\alpha_{new}) \propto (W(t) - x(t))^T x(t)$$

$$\propto W(t)^T x(t) - x(t)^T x(t)$$

$$\propto \cos(\alpha) - x(t)^T x(t)$$

$$\cos(\alpha_{new}) < \cos(\alpha)$$

Then, for each case shown above, we can identify that if we pick a misclassified point is because the $\cos \alpha$ gave a wrong value, so we must increase or decrease the value of α such makes that point a correct classification.

1.2. Is Learning Feasible?

1.2.1. Outside of the data set.

Exercise 1.7

\mathbf{x}	y	g	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8
0 0 0	○	○	○	○	○	○	○	○	○	○
0 0 1	●	●	●	●	●	●	●	●	●	●
0 1 0	●	●	●	●	●	●	●	●	●	●
0 1 1	○	○	○	○	○	○	○	○	○	○
1 0 0	●	●	●	●	●	●	●	●	●	●
1 0 1		?	○	○	○	○	●	●	●	●
1 1 0		?	○	○	●	●	○	○	●	●
1 1 1		?	○	●	○	●	○	●	○	●

For each of the following learning scenarios in the above problem, evaluate the performance of g on the three points in \mathcal{X} outside \mathcal{D} . To measure the performance, compute how many of the 8 possible target functions agree with g on all three points, on two of them, on one of them, and on none of them.

Observation: Since there are three missing vectors in \mathbf{X} , so there are eight missing functions. For each f_i we must analyze the output in it which is function of the \mathbf{X} vector and apply the following formula to get the score:

Solution:

$$Score = \sum_{i=0}^3 i(\# \text{ function agrees } i \text{ times})$$

- a) \mathcal{H} has only two hypotheses, one that always return ● and one that always returns ○. The learning algorithm picks the hypotheses that match the data set most.

$$\# \text{ times functions always returns } 1 = 1$$

$$\# \text{ times functions once returns } 1 \text{ and twice } 2 = 3$$

$$\# \text{ times functions once returns } 2 \text{ and twice } 1 = 3$$

$$\# \text{ times functions always returns } 0 = 1$$

$$Score = 0(1) + 1(3) + 2(3) + 3(1) = 9$$

- b) The same \mathcal{H} , but the learning algorithm now picks the hypotheses that matches the *least*.

$$\# \text{ times functions always returns } 1 = 1$$

$$\# \text{ times functions once returns } 1 \text{ and twice } 2 = 3$$

$$\# \text{ times functions once returns } 2 \text{ and twice } 1 = 3$$

$$\# \text{ times functions always returns } 0 = 1$$

$$Score = 0(1) + 1(3) + 2(3) + 3(1) = 9$$

- c) $\mathcal{H} = \{XOR\}$ (only one hypothesis is always picked), where XOR is defined by $XOR(x) = \bullet$ if the number of 1's in \mathcal{X} is odd and $XOR(x) = \circ$ if the number is even.

times functions agree 0 times = 1

times functions agree 1 times = 2

times functions agree 2 times = 2

times functions agree 3 times = 1

$$Score = 0(1) + 1(3) + 2(3) + 3(1) = 9$$

- d) \mathcal{H} contains all possible hypotheses (Boolean functions on three variables), and the learning algorithm picks the hypotheses that agrees with all training examples, but otherwise disagrees the most with the XOR.

times functions agree 0 times = 1

times functions agree 1 times = 2

times functions agree 2 times = 2

times functions agree 3 times = 1

$$Score = 0(1) + 1(3) + 2(3) + 3(1) = 9$$

1.2.2. Probability to the rescue.

Exercise 1.8 If $\mu = 0.9$, what is the probability that a sample of 10 marbles will have $v \leq 0.1$ [HINTS: 1. Use the binomial distribution. 2. The answer is a very small number.]

Solution:

Let $R = \#$ red marbles & $G = \#$ green marbles

$$P(R) = 0.9 \rightarrow P(G) = 1 - P(R) = 0.1$$

In other hand:

$$P(v \leq 0.1) = P\left(\frac{R}{N} \leq 0.1\right) = P(R \leq 1)$$

$$R \sim B(10, 0.1)$$

Using the formula we get:

$$P(R \leq 1) = P(R = 0) + P(R = 1) = 0.1^{10} + 9 \times 10^{-9} = 10^{-9}(9 + 0.1) = 9.1 \times 10^{-9}$$

Exercise 1.9 If $\mu = 0.9$, use the Hoeffding Inequality to bound the probability that a sample of 10 marbles will have $\mu \leq 0.1$ and compare the answer to the previous exercise.

Solution:

$$v - \mu = -0.8 \rightarrow |v - \mu| \geq 0.8$$

$$\text{If } \epsilon = 0.8 \rightarrow P(|v - \mu| > 0.8) \leq 2e^{-2 \times 0.8^2 \times 10} \approx 5.52 \times 10^{-6}$$

$$P(v \leq 0.1) \leq P(|v - \mu| > 0.8) \Leftrightarrow 9.1 \times 10^{-9} \leq 5.52 \times 10^{-6}$$

Exercise 1.10 Here is an experiment that illustrates the difference between a single bin and multiple bins. Run a computer simulation for flipping 1, 000 fair coins. Flip each coin independently times. Let's focus on 3 coins as follows: c_1 is the first coin flipped; c_{rand} is a coin you choose at random; c_{min} is the coin that had the minimum frequency of heads (pick the earlier one in case of a tie). Let v_1 , v_{rand} and v_{min} be the fraction of heads you obtain for the respective three coins.

- What is μ for the three coins selected?
- Repeat this entire experiment a large number of times (e.g., 100, 000 runs of the entire experiment) to get several instances of v_1 , v_{rand} and v_{min} and plot the histograms of the distributions of v_1 , v_{rand} and v_{min} . Notice that which coins end up being c_{rand} and c_{min} may differ from one run to another.
- Using b), plot estimates for $P[|v - \mu| > \epsilon]$ as a function of ϵ , together with the Hoeffding bound $2e^{-2^2 N}$ (on the same graph).

Solution:

- All the coins are fair, so: $\mu = 0.5$

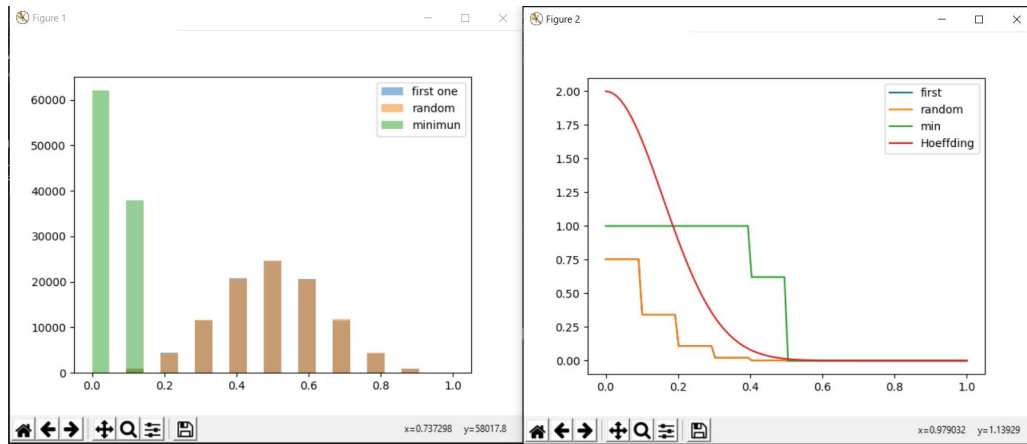


Figure 1: The histogram (left) shows the fraction of heads for the selected coins, and the plot (right) shows the probability and Hoeffding for the same ones coins.

Code:

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  def simulate(ncoin, ntoss):
5      outcome = np.random.choice((0., 1.), size=(ncoin, ntoss))
6      heads = np.sum(outcome, axis=1)
7      heads /= 10
8      return heads[0], np.random.choice(heads), np.min(heads)
9
10 def plott1():
11     bins = np.linspace(0, 1, 22)
12     plt.hist(experiments[:, 0], bins, alpha=0.5, label="first one")
13     plt.hist(experiments[:, 1], bins, alpha=0.5, label="random")
14     plt.hist(experiments[:, 2], bins, alpha=0.5, label="minimun")
15     plt.legend()
16     plt.show()
17
18 def P(outcomes, mu, epsilon):
19     less = np.sum(outcomes < mu - epsilon)
20     greater = np.sum(outcomes > mu + epsilon)
21     return (less + greater) / len(outcomes)
22
23 def hoeffding(epsilon, N):
24     return 2 * np.exp(-2 * epsilon * epsilon * N)
25
26 def plott2():
27     epsilons = np.linspace(0, 1, 100)
28     P_first = [P(experiments[:, 0], mu, epsilon) for epsilon in epsilons]
29     P_random = [P(experiments[:, 1], mu, epsilon) for epsilon in epsilons]
30     P_min = [P(experiments[:, 2], mu, epsilon) for epsilon in epsilons]
31     P_hoeffding = [hoeffding(epsilon, ntoss) for epsilon in epsilons]
32     plt.plot(epsilons, P_first, label="first")
33     plt.plot(epsilons, P_random, label="random")
34     plt.plot(epsilons, P_min, label="min")
35     plt.plot(epsilons, P_hoeffding, label="Hoeffding")
36     plt.legend()
37     plt.show()
38
39 if __name__ == '__main__':
40     ncoin = 1000 # Number of fair coins for each experiment
41     ntoss = 10 # Number of toss for each coin
42     mu = 0.5 # Probability
43     nexperiment = 100000 # Number of experiments
44     experiments = np.empty((nexperiment, 3))
45     for i in range(nexperiment):
46         first, rand, minimun = simulate(ncoin, ntoss)
47         experiments[i, 0], experiments[i, 1], experiments[i, 2] = first, rand, minimun
48     plott1()
49     plott2()

```

1.3. Error and noise

1.3.1. Noisy targets.

Exercise 1.13 Consider the bin model for a hypothesis h that makes an error with probability μ in approximating a deterministic target function (both h and f are binary functions). If we use the same h to approximate a noisy version of f given by:

$$P(y \mid \mathbf{x}) = \begin{cases} \lambda & y = f(\mathbf{x}), \\ 1 - \lambda & y \neq f(\mathbf{x}). \end{cases}$$

- a) What is the probability of error that h makes in approximating y ?
- b) At what value of λ will the performance of h be independent of μ ? [Hint: The noisy target will look completely random.]

Solution:

- a) First, we just know that $P[h(\mathbf{x}) \neq f(\mathbf{x})] = \mu$ and we want to know the $P[h(\mathbf{x}) \neq y]$, so there are two cases:

1. $P[h(\mathbf{x}) \neq y] = P[h(\mathbf{x}) = f(\mathbf{x})]P[f(\mathbf{x}) \neq y]$
2. $P[h(\mathbf{x}) \neq y] = P[h(\mathbf{x}) \neq f(\mathbf{x})]P[f(\mathbf{x}) = y]$

If we sum the two parts we obtain the answer:

$$\begin{aligned} &= (1 - \mu)(1 - \lambda) + (\mu)(\lambda) \\ &= 1 + 2\mu\lambda - \mu - \lambda \end{aligned}$$

- b) When we say independent of μ it means that whichever its value could be does not matter in the original prediction: $h(x) = y$. So, in what we are really interested is in the value for λ , in the previous result, such that μ disappears:

$$2\mu\lambda - \mu = 0$$

$$2\mu\lambda = \mu$$

$$\lambda = \frac{1}{2}$$

2. Problems.

- 1.1) We have 2 opaque bags, each containing 2 balls. One bag has 2 black balls and the other has a black and a white ball. You pick a bag at random and then pick one of the balls in that bag at random. When you look at the ball it is black. You now pick the second ball from that same bag. What is the probability that this ball is also black?

Solution:

Let be, \mathbf{B}_k : k-Ball is black, so we want to know:

$$\begin{aligned}\mathbb{P}(B_2|B_1) &= \frac{\mathbb{P}(B_1 B_2)}{\mathbb{P}(B_1)} \\ &= \frac{\sum_k^n \mathbb{P}(B_1 B_2 | Bag_k) \mathbb{P}(Bag_k)}{\sum_k^n \mathbb{P}(B_1 | Bag_k) \mathbb{P}(Bag_k)} \\ &= \frac{1(\frac{1}{2}) + \frac{1}{2}(0)}{1(\frac{1}{2}) + \frac{1}{2}(\frac{1}{2})} \\ &= \frac{2}{3}\end{aligned}$$

- 1.2) Consider the perceptron in two dimensions: $h(x) = \text{sgn}(\mathbf{w}^T \mathbf{x})$ where $\mathbf{w} = [w_0, w_1, w_2]^T$ and $\mathbf{x} = [1, x_1, x_2]^T$. Technically, \mathbf{x} has three coordinates, but we call this perceptron two-dimensional because the first coordinate is fixed at 1.
- 1) Show that the regions on the plane where $h(x) = +1$ and $h(x) = -1$ are separated by a line. If we express this line by the equation $x_2 = ax_1 + b$, what are the slope a and intercept b in terms of w_0, w_1, w_2 ?
 - 2) Draw a picture for the cases $w = [1, 2, 3]$ and $w = -[1, 2, 3]^T$.

Solution:

- 1) We have the following perceptron: $h(x) = \text{sgn}(\mathbf{w}^T \mathbf{x})$, so:

$$\begin{aligned}h(x) = +1 &\rightarrow w^T x > 0 \\ h(x) = -1 &\rightarrow w^T x < 0 \\ \Rightarrow w^T x = 0 &\text{ is the separator line}\end{aligned}$$

The equation is:

$$\begin{aligned}w_0 + w_1 x_1 + w_2 x_2 &= 0 \\ w_0 + w_1 x_1 &= -w_2 x_2 \\ \frac{w_0 + w_1 x_1}{-w_2} &= \frac{-w_2 x_2}{-w_2} \\ -\frac{w_0}{w_2} - \frac{w_1}{w_2} x_1 &= x_2\end{aligned}$$

The constants are next to the x_i , so:

$$b = -\frac{w_0}{w_2}; a = -\frac{w_1}{w_2}$$

2) We obtain the following plots for each vector:

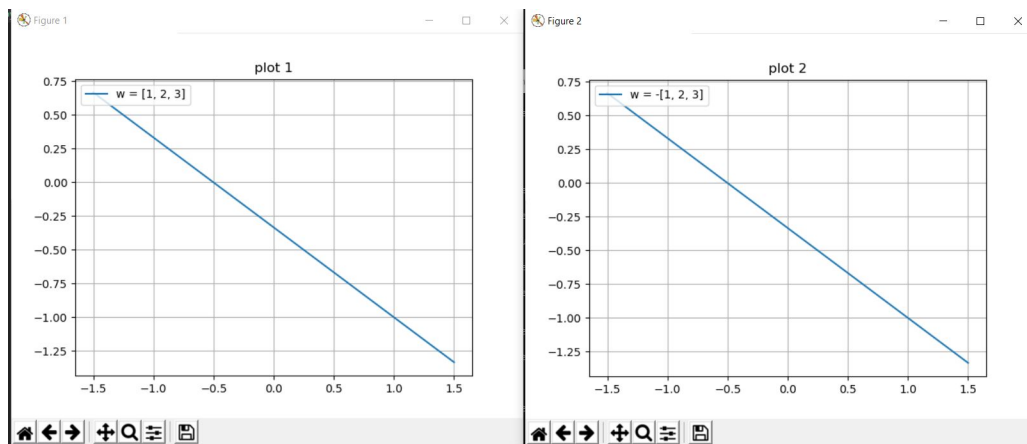


Figure 2: The separator line is identical in both plots, but the area in which $h(x) = +1$ puts each x_i in the left is on top and in the right is on bottom.

Code:

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 if __name__ == '__main__':
5     x = np.linspace(-1.5, 1.5, 80)
6     w1 = [1, 2, 3]
7     w2 = [-1, -2, -3]
8     y1 = -(w1[0]/w1[2]) - (w1[1]*x)/w1[2]
9     plot1 = plt.figure(1)
10    plt.plot(x, y1, label='w = [1, 2, 3]')
11    plt.title('plot 1')
12    plt.legend(loc='upper left')
13    plt.grid()
14    plot1.show()
15
16    y2 = -(w2[0]/w2[2]) - (w2[1]*x)/w2[2]
17    plot2 = plt.figure(2)
18    plt.plot(x, y2, label='w = -[1, 2, 3]')
19    plt.title('plot 2')
20    plt.legend(loc='upper left')
21    plt.grid()
22    plot2.show()
23
24    input()

```

- 1.6) Consider a sample of 10 marbles drawn independently from a bin that holds red and green marbles. The probability of a red marble is μ . For $\mu = 0.05$, $\mu = 0.5$, and $\mu = 0.8$, compute the probability of getting no red marbles ($v = 0$) in the following cases.
- We draw only one such sample. Compute the probability that $v = 0$.
 - We draw 1,000 independent samples. Compute the probability that (at least) one of the samples has $v = 0$.
 - Repeat (b) for 1,000,000 independent samples.

Solution:

- a) These experiments have a Binomial distribution, so:

$$X = \# \text{ red marbles in sample. } X \sim B(10, \mu)$$

- For $\mu = 0.05$

$$P[X = 0] = 0.5987$$

- For $\mu = 0.5$

$$P[X = 0] = 9.8 \times 10^{-4}$$

- For $\mu = 0.8$

$$P[X = 0] = 1.02 \times 10^{-7}$$

- b) These are new experiments which a Binomial distribution based on the previous experiment:

$$Y = \# \text{ samples with 0 red marbles } Y \sim B(1000, \lambda)$$

- For $\lambda = 0.5987$

$$P[Y \geq 1] = 1$$

- For $\lambda = 9.8 \times 10^{-4}$

$$P[Y \geq 1] = 0.62487$$

- For $\lambda = 1.02 \times 10^{-7}$

$$P[Y \geq 1] = 1 \times 10^{-4}$$

- c) Here is the same experiment, therefore is the almost the same distribution. We just have to increase the number of samples. For samples $n = 1,000,000$.

$$Y = \# \text{ samples with 0 red marbles} \quad Y \sim B(1000000, \lambda)$$

- For $\lambda = 0.5987$

$$P[Y \geq 1] = 1$$

- For $\lambda = 9.8 \times 10^{-4}$

$$P[Y \geq 1] = 1$$

- For $\lambda = 1.02 \times 10^{-7}$

$$P[Y \geq 1] = 0.09697$$

- 1.10) Assume that $X = \{x_0, x_1, \dots, x_N, x_{N+1}, \dots, x_{N+M}\}$ and $Y = \{-1, +1\}$ with an unknown target function $f : X \rightarrow Y$. The training data set $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$. Define the off-training-set error of a hypothesis h with respect to f by

$$E_{\text{off}}(h, f) = \frac{1}{M} \sum_{m=1}^M M[h(x_{N+m}) \neq f(x_{N+m})]$$

- a) Say $f(x) = +1$ for all x and

$$h(x) = \begin{cases} +1 & \text{for } x = x_k \text{ and } k \text{ is odd and } 1 \leq k \leq M + N \\ -1 & \text{otherwise} \end{cases}$$

What is $E_{\text{off}}(h, f)$?

- b) We say that a target function f can 'generate' D in a noiseless setting if $y_n = f(x_n)$ for all $(x_n, y_n) \in D$. For a fixed D of size N , how many possible $f : X \rightarrow Y$ can generate D in a noiseless setting?
- c) For a given hypothesis h and an integer k between 0 and M , how many of those f in (b) satisfy $E_{\text{off}}(h, f) = \frac{k}{M}$?
- d) For a given hypothesis h , if all those f that generate D in a noiseless setting are equally likely in probability, what is the expected off training-set error $E_f[E_{\text{off}}(h, f)]$?
- e) A deterministic algorithm A is defined as a procedure that takes D as an input, and outputs a hypothesis $h = A(D)$. Argue that for any two deterministic algorithms A_1 and A_2 .

$$E_f[E_{\text{off}}(A_1(D, f))] = E_f[E_{\text{off}}(A_2(D, f))]$$

,

Solution:

- a) If we rewrite the conditions, we can sum up and observe that the only case when $f(x_{N+i}) \neq h(x_{N+i}) \quad \forall 1 \leq i \leq M$ is when x_{N+i} is odd. So:

$$E_{\text{off}}(h, f) = \frac{1}{M} [\text{Total number of odd between 1 and } M]$$

There are two cases:

$$\begin{aligned} M \text{ is even} & \rightarrow E_{\text{off}}(h, f) = \frac{1}{M} \frac{M}{2} = \frac{1}{2} \\ M \text{ is odd} & \rightarrow E_{\text{off}}(h, f) = \frac{1}{M} \left(\frac{M-1}{2} + 1 \right) = \frac{1}{2} - \frac{1}{2M} + \frac{1}{M} = \frac{1}{2} + \frac{1}{2M} \end{aligned}$$

- b) If $|D| = N$ there are 2^n possible target functions that can generate D in a noiseless setting.
 c) What we really want to know is:

$$\frac{1}{M} \sum_{m=1}^M M[h(x_{N+m}) \neq f(x_{N+m})] = k$$

And this can sum up to the number of ways we have k errors in M trials. And this is a combination in the entire set:

$$C_k^M = \frac{M!}{k!(M-k)!}$$

d)

$$\begin{aligned} E_f[E_{\text{off}}(h, f)] &= \frac{1}{M} E_f[\text{Number of } h(x_{N+m}) \neq f(x_{N+m})] \\ &= \frac{1}{M} \times M \times \frac{1}{2^M} = \frac{1}{2^M} \end{aligned}$$

- 1.11) The matrix which tabulates the cost of various errors for the CIA and Supermarket applications in Example 1.1 is called a risk or loss matrix. For the two risk matrices in Example 1.1, explicitly write down the in sample error E_{in} that one should minimize to obtain g . This in-sample error should weight the different types of errors based on the risk matrix.

Solution:

a) Supermarket:

$$\begin{aligned} E_{\text{in}}(h) &= \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n)) \\ &= \frac{1}{N} \sum_{y_n=1} e(h(x_n, 1)) + \sum_{y_n=-1} e(h(x_n, -1)) \\ &= \frac{1}{N} \sum_{y_n=1} 10[[h(x_n \neq 1)]] + \sum_{y_n=-1} [[h(x_n \neq -1)]] \end{aligned}$$

b) CIA:

$$\begin{aligned} E_{\text{in}}(h) &= \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n)) \\ &= \frac{1}{N} \sum_{y_n=1} e(h(x_n, 1)) + \sum_{y_n=-1} e(h(x_n, -1)) \\ &= \frac{1}{N} \sum_{y_n=1} [[h(x_n \neq 1)]] + \sum_{y_n=-1} 1000[[h(x_n \neq -1)]] \end{aligned}$$