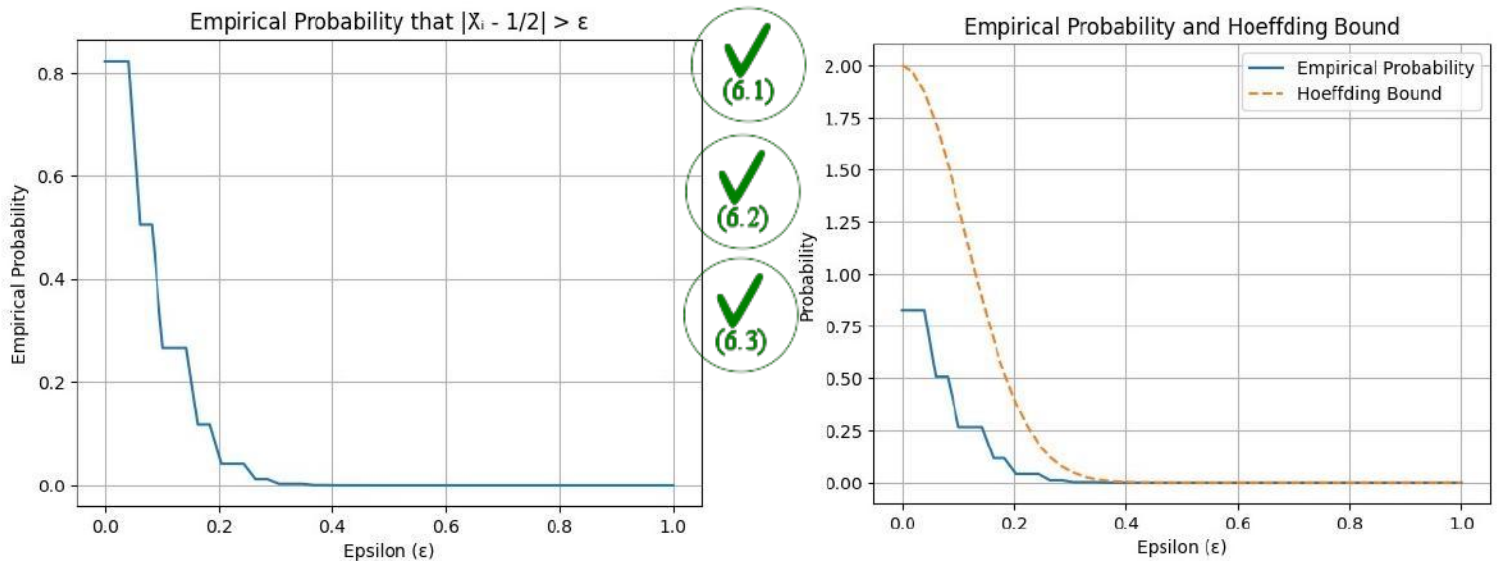


Programming Assignment

1. Visualizing the Hoeffding bound (10 pts).

- Use **numpy** to generate an $N \times n$ matrix of samples from $Bernoulli(1/2)$. Calculate for each row i the empirical mean, $\bar{X}_i = \frac{1}{n} \sum_{j=1}^n X_{i,j}$, where $N = 200000$ and $n = 20$.
- Take 50 values of $\epsilon \in [0, 1]$ (**numpy.linspace(0,1, 50)**), and calculate the empirical probability that $|\bar{X}_i - 1/2| > \epsilon$. Plot the empirical probability as a function of ϵ .
- Add to your plot the Hoeffding bound of that probability, as a function of ϵ .

Submit your plots (no need to submit code for this question).



- Nearest Neighbor (20 pts).** In this question, we will study the performance of the Nearest Neighbor (NN) algorithm on the MNIST dataset. The MNIST dataset consists of images of handwritten digits, along with their labels. Each image has 28×28 pixels, where each pixel is in gray-scale, and can get an integer value from 0 to 255. Each label is a digit between 0 and 9. The dataset has 70,000 images. Although each image is square, we treat it as a vector of size 784.

- Run the algorithm using the first $n = 1000$ training images, on each of the test images, using $k = 10$. What is the accuracy of the prediction (i.e. the percentage of correct classifications)? What would you expect from a completely random predictor?

• ככל שכמות הדגימות שאנחנו משתמשים בהם היא גדלה, הדיוק של הקבוצה 86%.

• המסלול הנבחר הוא יחיד מספר בין 0-9 באופן שרירותי ולכן לא צפוי להיות 10% דיוק.

- Plot the prediction accuracy as a function of k , for $k = 1, \dots, 100$ and $n = 1000$. Discuss the results. What is the best k ?

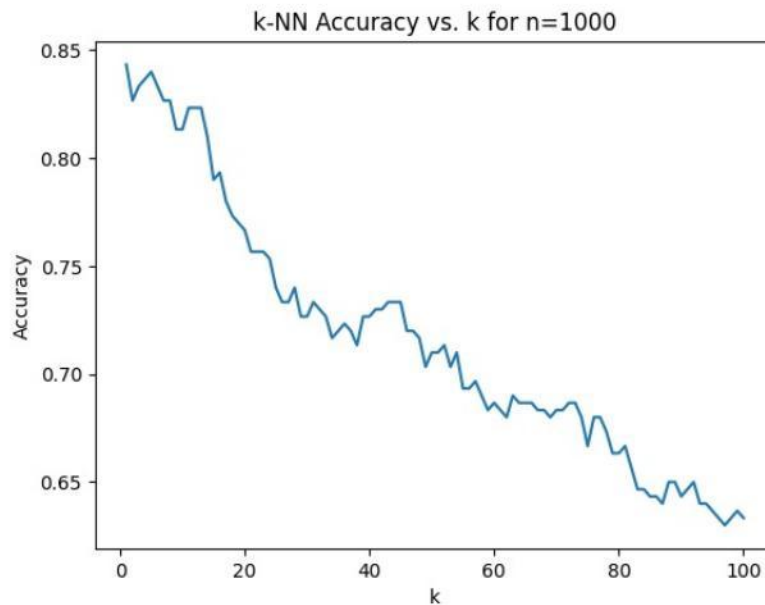
לפי התוצאות נראה כי היא הכי טובה היא $k=5$ או $k=1$. בעצם היא שניהם הדיוק הסביר ביותר.

נראה כי ערכים קטנים של k (1-5) הם הדיוק הטוב ביותר. זה כי ככל שהיא גדלה, המסלול נבחר יותר שרירותי.

חוקים אידיאליים קרובים למוקטור (פאזיזט שטוקר). עם זאת, כחומר $k=1$ הוא יכולה לפעמים להיות $overfitting$. לכן 1 נראה

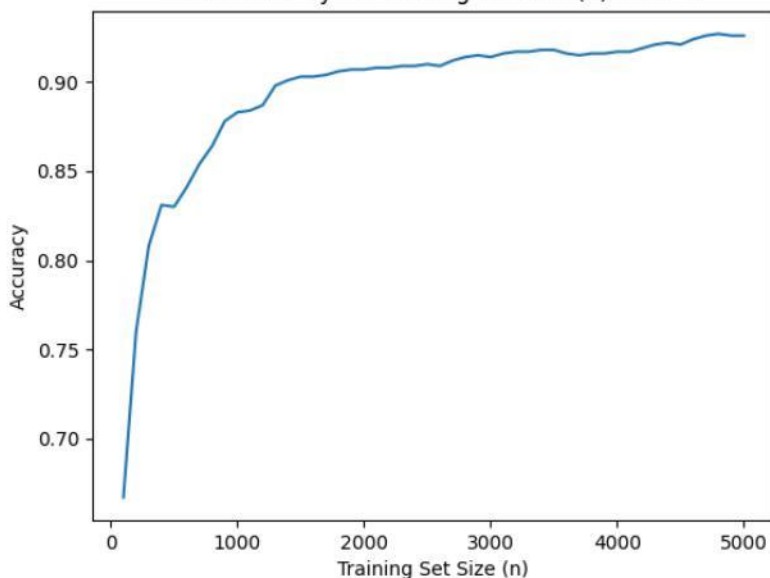
אומנם מדויק יותר אבל הפסד בין 1 ל-5 מעט ואולי היה אולי משמעותי אם כמות ה- $test Data$ הייתה קטנה יותר. הנה שני דוגמאות נוספות:

✓
(7.3)



(d) Using $k = 1$, run the algorithm on an increasing number of training images. Plot the prediction accuracy as a function of $n = 100, 200, \dots, 5000$. Discuss the results.

k-NN Accuracy vs. Training Set Size (n) for k=1



✓
(7.4)

ב' שניתן לראות באיור ככל שגודל המידע

ה- Training Data גדול יותר, עולה דיוק הדיוק $k=1$

אבל בשל כמות המידע (באזור 1200) יש האטה בקצב

הדיוק שגדל לאט בעצמו ה-90% מסתוות:

1. תסיבות ה- Training Data - העליה המהירה בהתחלה

מגיעה במהרה תשובה למחזוריות באזור גדול לצורך דיוק

האופטימלי.

2. עקומת המידה: בהתחלה המידע לומד במהירות אבל

לבסוף חווה ירידה בדיוק עקב חוויה של מידע חקלא.

3. Overfitting: עבור n קטנים: ה- $n=100$ הוא מידע ספחית $k=1$ מדיק.

4. מודל אופטימלי של קבוצת Training Data: בעוד שניתן כי יותר נתונים עצמים, תוצאה זאת מראה Trade off , יש כמות אופטימלית של נתונים.

שבה אחר מנסה דיוק מדי. והעבר לנק' האופטימלית זו ייתכן שגודל עיבוד ואיסוף הנתונים אף מועדקת למידה האופטימלית.

✗
(8)