

מבוא למערכות לומדות – 67577 – תרגיל 2

אופק אבידן – 318879574

2. חלק תיאורטי:

Hard- & Soft-SVM .2.1

2.1.1

In class we saw the Hard-SVM classification model.

1. Prove that the following Hard-SVM optimization problem is a Quadratic Programming problem:

$$\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i \, y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad (1)$$

That is, find matrices Q and A and vectors \mathbf{a} and \mathbf{d} such that the above problem can be written in the following format

$$\underset{\mathbf{v} \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \mathbf{v}^\top Q \mathbf{v} + \mathbf{a}^\top \mathbf{v} \quad \text{s.t.} \quad A \mathbf{v} \leq \mathbf{d} \quad (2)$$

Hint: Observe that $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{I} \mathbf{w}$

נתבונן בבעיה:

$$\underset{w, b}{\operatorname{argmin}} \|w\|^2 \quad \text{s.t.} \quad \forall i \, y_i (\langle w, x_i \rangle + b) \geq 1$$

באופן שקול ניתן לכתוב:

$$\underset{(w, b)}{\operatorname{argmin}} w^\top I w$$

$$\begin{aligned} \text{s.t.} \quad \forall i \, y_i (\langle w, x_i \rangle + b) \geq 1 &\rightarrow y_i x_i^\top w + y_i b \geq 1 \rightarrow -y_i x_i^\top w - y_i b \leq -1 \\ &\rightarrow -y_i x_i^\top w \leq y_i b - 1 \end{aligned}$$

נגדיר:

$$Q = I, \quad A_{i' \text{th row}} = -y_i x_i^\top, \quad a = 0 \rightarrow V = \sqrt{2} W$$

ונקבל שהבעיה שקולה ל-

$$\underset{V \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} V^\top Q V + a^\top V \quad \text{s.t.} \quad A v \leq d$$

2. The **Gaussian Naive Bayes** classifier assumes a multinomial prior and independent feature-wise Gaussian likelihoods:

$$\begin{aligned} y &\sim \text{Multinomial}(\pi) \\ x_j|y=k &\overset{\text{ind.}}{\sim} \mathcal{N}(\mu_{kj}, \sigma_{kj}^2) \end{aligned} \quad (4)$$

for π a probability vector: $\pi \in [0, 1]^K, \sum \pi_j = 1$.

- (a) Suppose $x \in \mathbb{R}$ (i.e each sample has a single feature). Given a trainset $\{(x_i, y_i)\}_{i=1}^m$ fit a Gaussian Naive Bayes classifier solving (3) under assumptions (4). Fitting means finding the expressions for the maximum likelihood estimators.
- (b) Suppose $\mathbf{x} \in \mathbb{R}^d$ (i.e each sample has d feature). Given a trainset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ fit a Gaussian Naive Bayes classifier solving (3) under assumptions (4). You are encouraged to use the results from (3.a).

ראשית נמצא ביטוי ל-likelihood :

$$L(\Theta|X, y) = \prod_{i \in [m]} f_{x,y}(x_i, y_i|\Theta) \stackrel{\text{מהנחה}}{=} \prod_{i \in [m]} N(X_i|M_{y_i}, \sigma_{y_i}^2) \cdot \text{Mult}(y_i|\pi)$$

פונקציית log הינה מונוי עולה, ולכן מקסום הביטוי ימקסם גם את ה-log של הביטוי, לכן :

$$\begin{aligned} \log(L(\Theta|X, y)) &= l(\Theta|X, y) = \sum_{i \in [m]} (\log N(X_i|M_{y_i}, \sigma_{y_i}^2) + \log(\text{Mult}(y_i|\pi))) \\ &= \sum_{i \in [m]} (\log(\frac{1}{\sqrt{2\pi\sigma^2 y_i}} \exp(-\frac{(x_i - \mu y_i)^2}{2\sigma^2 y_i})) + \log(\pi y_i)) \\ &= -\frac{m}{2} \log(2\pi) + \sum_{k \in m} n_k \log(\pi k) - \frac{1}{2} n_k \log(\sigma_k^2) \\ &\quad - \frac{1}{2\sigma_k^2} \sum_{i: y_i=k} (X_i - \mu_k)^2 \end{aligned}$$

$$n_k = \sum_{i \in [m]} 1_{[y_i=k]} \quad \text{כאשר}$$

כעת נגזור את הביטוי כדי למצוא מקסימום :

- נגזור לפי תוחלת :

$$\frac{1}{\sigma^2 k} \sum_{i: y_i=k} (X_i - \mu_k) = 0$$

$$* \mu_k = \frac{1}{n_k} \sum_{i: y_i=k} X_i$$

- נגזרת לפי שונות :

$$-\frac{n_k}{\sigma_k^2} - \frac{1}{2\sigma_k^4} \sum_{i: y_i=k} (X_i - \mu_k)^2 = 0 \rightarrow \sigma_k^{*2} = \frac{1}{n_k} \sum_{i: y_i=k} (X_i - * \mu_k)^2$$

כעת, נגדיר לגראנז'יאן :

$$G = l(\theta|X, y) - 2(\sum_{k \in [m]} \pi_k - 1)$$

נגזור לפי π_k ונשווה ל-0:

$$0 = \left(l(\theta|X, y) - 2(\sum_{k \in [m]} \pi_k - 1) \right)' = \frac{n_k}{\pi_k} - \lambda \rightarrow \pi_k = \frac{n_k}{\lambda}$$

נוכר כי $\sum_k \pi_k = 1$ ונמצא את λ :

$$\sum_{k \in [m]} \pi_k = \sum_{k \in [m]} \frac{n_k}{\lambda} = 1 \rightarrow \lambda = m$$

לסיכום, המשעריך של כל π_k הממקסם את ה-likelihood הינו:

$$\pi_k^{MLE} = \frac{n_k}{m}$$

$$n_k = \sum_{i \in [m]} 1_{[y_i=k]} \text{ כאשר}$$

כעת נחשב את הלויקליהוד:

$$L(\theta|X, y) = \prod_{i \in [m]} f_{x,y}(x_i, y_i | \theta) \stackrel{\text{מהנחה}}{=} \prod_{j \in [m]} \prod_{i \in [m]} N(X_i | M_{y_i}, \sigma_{y_i}^2) \cdot \text{Mult}(y_i | \pi)$$

כעת, משיקולי מונוטוניות ומיקסום נפעיל לוג על כל הלויקליהוד:

$$\begin{aligned} \log(L(\theta|X, y)) &= \sum_{i \in [m]} (\log(\text{mult}(y_i, \pi_k)) + \sum_{j \in [d]} \log(N(x_{ij} | \mu_{y_{ij}}, \sigma_{y_{ij}}^2))) \\ &= \sum_k n_k \log \pi_k + \sum_k \sum_{j \in [d]} 1_{[y:=k]} \log((N(X_{ij} | \mu_y := ij, \sigma_y^2))) \end{aligned}$$

נשים לב כי נוכל להסתכל על כל פיצור j ולחשב את האומדים בעזרת החישוב שביצענו בסעיף

$$(x_j | y = k) \sim N(\mu_k, j, \sigma_k^2, j) \text{ לקבל } k \text{ עבור כל } k$$

עם אותם $\mu_{k,j}, \sigma_{k,j}^2$ שחושבו. בנוסף, מכך שהפיצרים בת"ל, לכל k אנו מקבלים מטריצה

אלכסונית שונה ממחלקה למחלקה לכן נוכל לרשום את האומדים כך:

$$\pi_k^{MLE} = \frac{n_k}{m}, \quad M_{k,j}^{MLE} = \frac{1}{n_k} \sum_{i:y_i=k} K_{i,j}, \quad \sigma_{k,j}^{2, MLE} = \frac{1}{n_k} \sum (X_{i,j} - M_{k,j}^{MLE})^2$$

3. The **Poisson Naive Bayes** classifier assumes a multinomial prior and independent feature-wise Poisson likelihoods:

$$\begin{aligned} y &\sim \text{Multinomial}(\pi) \\ x_j|y=k &\stackrel{\text{ind.}}{\sim} \text{Poi}(\lambda_{kj}) \end{aligned} \quad (5)$$

for π a probability vector: $\pi \in [0, 1]^K, \sum \pi_j = 1$.

- (a) Suppose $x \in \mathbb{R}$ (i.e each sample has a single feature). Given a trainset $\{(x_i, y_i)\}_{i=1}^m$ fit a Poisson Naive Bayes classifier solving (3) under assumptions (5).
 (b) Suppose $\mathbf{x} \in \mathbb{R}^d$ (i.e each sample has d feature). Given a trainset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ fit a Poisson Naive Bayes classifier solving (3) under assumptions (5). You are encouraged to use the results from (4.a).

א. ראשית נחשב את הלייקליהוד :

$$L(\Theta|X, y) = \prod_{i \in [m]} f_{x,y}(x_i, y_i|\Theta) = \prod_{i \in [m]} \text{Pois}(X_i|\lambda_i y_i) \cdot \text{Mult}(y_i|\pi)$$

כעת, משיקולי מונוטוניות ומיקסום נפעיל לוג על כל הלייקליהוד :

$$\begin{aligned} \log(L(\Theta|X, y)) &= \sum_{i \in [m]} ((\log(\text{Pois}(X_i|\lambda_i y_i)) + \sum_j \mathbb{1}_{y_i=j} \log(\text{Mult}(y_i|\pi))) \\ &= -m\lambda + \log(\lambda) \sum_i X_i - \sum \log(X_i d) \end{aligned}$$

$$\Theta^{MLE} = \lambda$$

נגזור, נשווה ל-0 ונקבל :

$$0 = (-m\lambda + \log(\lambda) \sum_i X_i - \sum \log(X_i d))' = -m + 1/\lambda * \sum_i X_i$$

לכן :

$$\pi_k^{MLE} = \frac{nk}{m}, \quad \lambda_k^{MLE} = \frac{1}{n_k} \sum_i \mathbb{1}_{[y_i=k]} X_i$$

$$n_k = \sum_i \mathbb{1}_{[y_i=k]} \text{ כאשר}$$

ב. נניח כי $X \in \mathbb{R}^d$, כלומר לכל דגימה יש d פיצירים.

נוכל להסתכל על כל פיציר j ולחשב את האומדים בעזרת החישוב שביצענו בסעיף הקודם

ועבור כל מחלקה k לקבל $(X_j, y = k) \sim \text{Pois}(\lambda_k, j)$ עד אותו λ_k, j שחישבנו.

לכן נוכל לרשום את האומדים כך :

$$\lambda_{k,j}^{MLE} = \frac{1}{n_k} \sum_i \mathbb{1}_{[y_i=k]} X_{i,j}$$

$$\pi_k^{MLE} = \frac{n_k}{m}$$

$$n_k = \sum_i \mathbb{1}_{[y_i=k]} \text{ כאשר}$$

3 Practical Part

In the following part, you will implement some classifiers you have seen in class and the recitation. You will compare the behavior of the different classifiers over different data scenarios.

3.1 Perceptron Classifier

Based on Lecture 2 and Recitation 5. Complete the following implementations

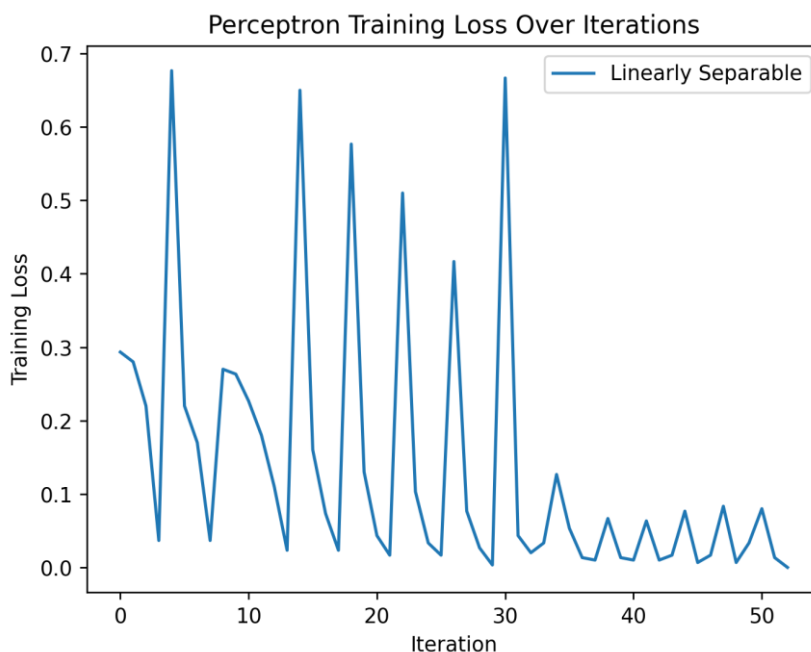
- Implement the `misclassification_error` function in the `loss_functions.py` file as described in the function documentation.
- Implement the Perceptron algorithm in the `classifiers.py` file as described in the class documentation. In the implementation use the misclassification error implemented above.

In the `classifiers_evaluation.py` file, implement the `run_perceptron` function as described in documentation.

- As good practice, to retrieve the loss at each iteration, we recommend not changing the previously implemented Perceptron class. Instead **specify a callback function** which receives the object and uses its `loss` function to calculate the loss over the training set. Store these values in an array to be used for plotting.

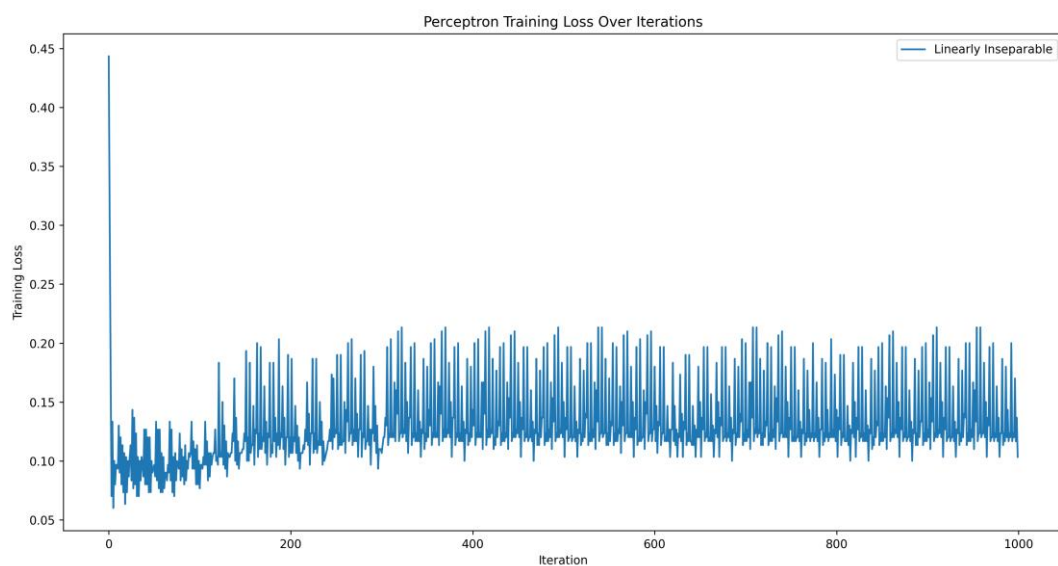
1. Fitting and plotting over the `linearly_separable.npy` dataset, what can we learn from the plot?

1. ראשית, מאחר והדאטה אכן פריד לינארית (כלומר, קו ישר (בדו מימד), מישור (בתלת מימד), או היפר מישור (במימדים גבוהים יותר) יכולים לחלק בצורה מושלמת את נקודות הנתונים של מחלקות שונות), ניתן לראות שהשגיאה האמפירית הסופית (לאחר כ-50 איטרציות) היא 0 כצפוי (אם לא היה פריד לינארית, לא היינו מגיעים ל-0). בנוסף, עד לאיטרציה 30 אפשר להבחין בקפיצות גדולות בשגיאות, זאת בשל שינוי הערכים שחוזרים מה-*misclassification error*. באופן כללי, אם הדאטה פריד לינארית, אז האלגוריתם פרספטרון אמור להתכנס. בכל איטרציה הוא "משפר" את הוקטורים w ו- b , בכיוון השיפור לפי הגודל של ה-*learning rate*. ייתכן שה-*learning rate* גדול מדי במקרים מסויימים, כלומר התזווה תהיה גדולה מדי ומכאן הקלסיפיקציה תפגע והלוס גדל משמעותית.



2. Next run the Perceptron algorithm over the `linearly_inseparable.npy` dataset and plot its loss as a function of the iterations. What is the difference between this plot and to the one in the previous question? How can we explain the difference in terms of the objective and parameter space?

כפי שהשם מרמז, הדאטה לא פריד לינארית. אנחנו רואים שהאלגוריתם פרספטרון לא מצליח להתמודד עם מידע כזה בדיוק כפי שנלמד בכיתה, והוא לא מצליח להתכנס גם אחרי 1000 איטרציות (לשם השוואה, בסעיף 1 ביצענו פחות מ-60). כלומר, לא קיים קו (או כאמור בהכללה, היפר מישור), שמחלק את הדאטה כך שהשגיאה האמפירית (לפי *misclassification error*) היא 0. מכאן שהאלגוריתם ימשיך בחיפוש אחר הפתרון עד שהוא מגיע למספר מקסימלי של איטרציות שהוגדר מראש (כאמור, במקרה זה, 1000).



3.2 Bayes Classifiers

Based on Lecture 2 and Recitation 5.

Complete the following implementations

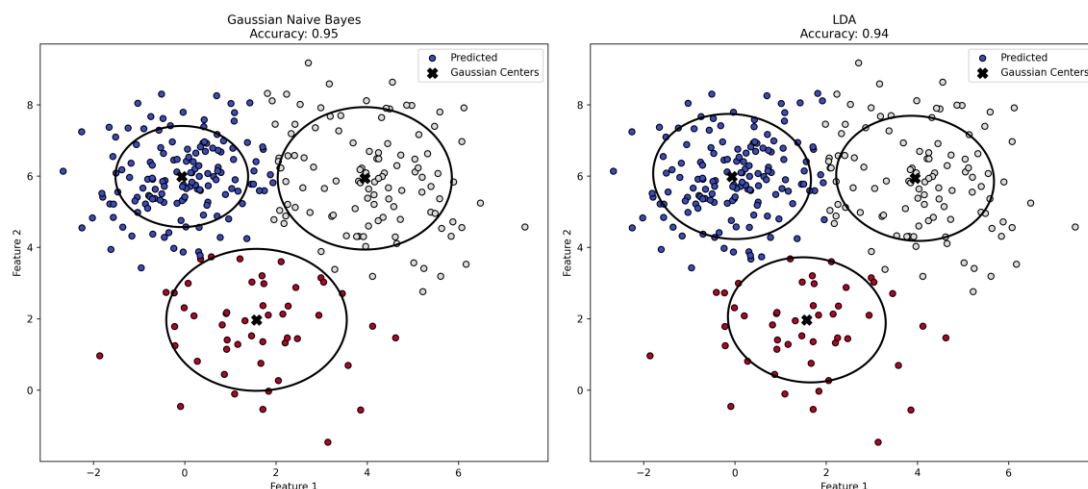
- Implement the accuracy function in the `loss_functions.py` file as described in the function documentation.
- Implement the LDA classifier in the `classifiers.py` file as described in the class documentation. Use expressions derived in class.
- Implement the GaussianNaiveBayes classifier in the `classifiers.py` file as described in the class documentation. Use expressions derived in question 3b of the theoretical part.

Then, implement and answer the following questions:

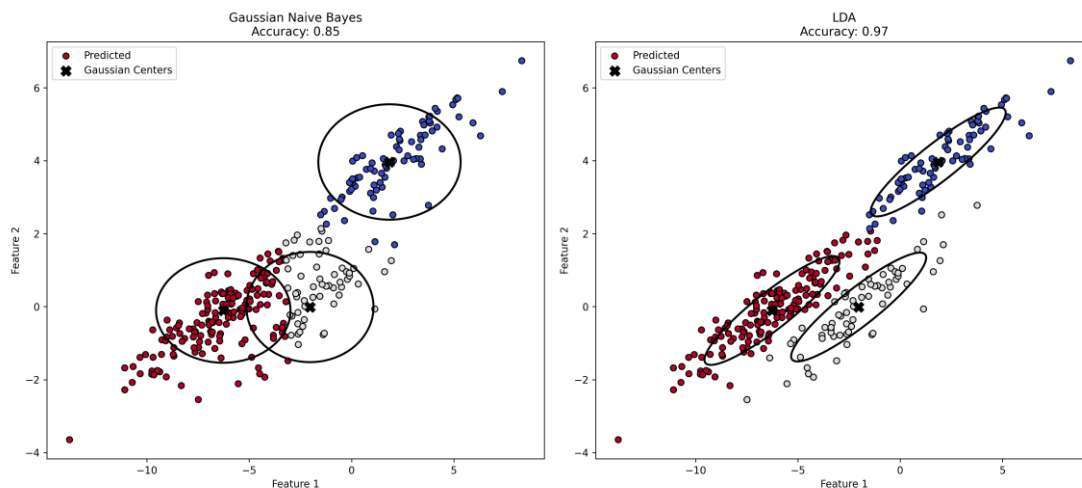
1. In the `compare_gaussian_classifiers` function, `classifiers_evaluation.py` file, load the `gaussians1.npy` dataset. Fit both the Gaussian Naive Bayes and LDA algorithms previously implemented. Plot the following:
 - A single figure with two subplots:
 - (a) 2D scatter-plot of samples, with marker color indicating Gaussian Naive Bayes *predicted* class and marker shape indicating *true* class.
 - (b) 2D scatter-plot of samples, with marker color indicating LDA *predicted* class and marker shape indicating *true* class.
 - (c) Provide classifier name and accuracy (over train) in sub-plot title
 - For both subplots add:
 - (a) Markers (colored black and shaped as 'X') indicating the center of fitted Gaussians.
 - (b) An ellipsis (colored black) centered in Gaussian centers and shape dictated by fitted covariance matrix.
 - Specify dataset name in figure title.

Explain what can be learned from the plots above regarding the distribution used to sample the data?

1. ניתן לראות שהגרפים שהתקבלו משני המודלים דומים מאוד, עם יתרון קטן לטובת ה-GNB. זאת מאחר וההתפלגות שממנה דגמו את אוסף הנתונים (הדאטאסט שקיבלנו) היא בעלת שונות משותפת אפס (כלומר, הדגימות בלתי תלויות). הדבר נובע מההנחה של GNB שהנתונים בת'ל ולכן התוצאה של שני המודלים תהיה זהה. כמו כן, אליפסה עגולה (שבה הצירים דומים באורכם) מעידה על כך שהשונות במשתנים השונים דומה. אם האליפסה מאורכת בכיוון מסוים, זה מצביע על כך שהשונות במשתנה אחד גבוהה יותר מאשר במשתנה אחר.



2. Repeat the procedure above (while avoiding code repetition) for `gaussians2.npy`. What is the difference between the two scenarios? What can be learned regarding the distribution used to sample the data? Which of the two classifiers better matches this dataset and why?



ניתן לראות שהדיוק ממודל LDA גבוה מהדיוק במודל GNB בניגוד לגרף הקודם. התופעה נובעת ככל הנראה מהיותם של הפיצורים תלויים זה בזה בניגוד להנחה של GNB שהם בת'ל. בנוסף, ניתן להבחין שהאליפסה בגרף ה LDA מתוחה בכיוון הדגימות מה שמסביר את הדיוק הגבוה יותר (כאמור, אם האליפסה מאורכת בכיוון מסוים, זה מצביע על כך שהשונות במשתנה אחד גבוהה יותר מאשר במשתנה אחר). גם כיוון האליפסה מראה את השונות בין המשתנים. אליפסה מוטה מעידה על כך שיש קשר (קורלציה) בין המשתנים. כיוון האליפסה יצביע על הכיוון שבו המשתנים מתפזרים יחד. אם האליפסה מוטה בכיוון מסוים, זה מעיד על כך ששני המשתנים משתנים באותו כיוון. אם האליפסה ישרה, זה מעיד על כך שאין קורלציה בין המשתנים. נשים לב שהאליפסות אינן מקבילות לצירים ב- LDA , כלומר האליפסות מוטות, בניגוד לסעיף הקודם, וזו עשויה להיות עדות נוספת לכך שהפיצורים תלויים.

שימוש ב-Chat GPT

השתמשתי בצ'אט ג'יפיטי ברוב החלק התכנותי. לצערי אין לי ניסיון ב-`numpy` ולכן נעזרתי בו לא מעט, בעיקר בהדפסת הגרפים. אני כן מבין לאחר מכן את מה שהוא ביצע, ואני משתדל לפרט המון בדו"ח.