

## Theoretical Part

### 2.1 Hard - & soft - SVM

1. ניתן שמתק"פ:  $\|w\|^2 = w^T I w$  כ- $I$  היא מטריצת היחידה. מתק"פ:

$$\forall i \quad y_i (\langle w, x_i \rangle + b) \geq 1 \Leftrightarrow \forall i \quad y_i (w^T x_i + b) \geq 1$$

$$\Leftrightarrow \forall i \quad -y_i (w^T x_i + b) \leq -1 \Leftrightarrow -y_i x_i^T w - y_i b \leq -1$$

אם נקבע:

$$A = \begin{bmatrix} -y_1 x_{11} & -y_1 x_{12} & \dots & -y_1 x_{1n} & -y_1 \\ -y_2 x_{21} & -y_2 x_{22} & \dots & -y_2 x_{2n} & -y_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ -y_m x_{m1} & -y_m x_{m2} & \dots & -y_m x_{mn} & -y_m \end{bmatrix} = \begin{bmatrix} -y_1 x_1^T & \dots & -y_1 \\ -y_2 x_2^T & \dots & -y_2 \\ \vdots & \ddots & \vdots \\ -y_m x_m^T & \dots & -y_m \end{bmatrix} \in M_{m \times (n+1)}$$

אם נגדיר את וקטור  $d \in \mathbb{R}^m$ :  $d = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$ . כך נקבל את המשוואה הרצויה  $Av \leq d$ .

נגדיר את  $v$  באופן הבא:  $v = \begin{bmatrix} w \\ b \end{bmatrix}$  כלומר  $v \in \mathbb{R}^{n+1}$

נגדיר את  $Q$  באופן הבא:  $Q = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}$  כאן  $Q \in M_{(n+1) \times (n+1)}$

באופן כזה כשנבצע מכפלה  $v^T Q v$  ה- $b$  יוכלו פאסר ופזר. נקבל  $\|w\|^2 = w^T I w$ .

$a = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$  :לכדן ןוללכ א ןלכ ןלכ  
פ"קןN זלכ

$$\operatorname{argmin}_{(w,b)} \|w\|^2 = \operatorname{argmin}_{(w,b)} w^T I w = \operatorname{argmin}_{(w,b)} v^T Q v$$

$$= \operatorname{argmin}_{(w,b)} \begin{bmatrix} w \\ b \end{bmatrix}^T Q \begin{bmatrix} w \\ b \end{bmatrix} = \operatorname{argmin}_{(w,b)} \frac{1}{2} \begin{bmatrix} w \\ b \end{bmatrix}^T Q \begin{bmatrix} w \\ b \end{bmatrix}$$

$$= \operatorname{argmin}_{(w,b)} \frac{1}{2} \begin{bmatrix} w \\ b \end{bmatrix}^T Q \begin{bmatrix} w \\ b \end{bmatrix} + 0^T \begin{bmatrix} w \\ b \end{bmatrix} = \operatorname{argmin}_{(w,b)} \frac{1}{2} v^T Q v + \tilde{a}^T v$$

## Naive Bayes classifiers

2.a) בהינתן  $R \ni X$  כך שלכל דגימה יש פיזור אחד וכל  $x_i$  מתפלג בהתפלגות נורמלית על  $k$  קטגוריות  $k \in N$  הוא מספר התוויות,  $\mu_k$  ו- $\sigma_k^2$  ושהתפלגות של כל תווית  $y_i$  היא  $\pi_k$ .  
ראינו כיצד שפונקציית ה-likelihood היא:

$$\begin{aligned} \mathcal{L}(\theta | X, y) &\stackrel{i.i.d}{=} \prod_{i=1}^m f_{x, y | \theta}(x_i, y_i) = \prod_{i=1}^m f_{x | y=y_i}(x_i) \cdot f_{y | \theta}(y_i) \\ &= \prod_{i=1}^m \mathcal{N}(x_i | \mu_{y_i}, \Sigma) \cdot \text{Mult}(y_i | \pi) \end{aligned}$$

ראינו שהמציאה של ה-likelihood מזהה log-likelihood:  
כל:

$$\ell(\theta | X, y) = \log \left( \prod_{i=1}^m \mathcal{N}(x_i | \mu_{y_i}, \Sigma) \cdot \text{Mult}(y_i | \pi) \right)$$

$$= \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi|\Sigma|}} \exp \left( -\frac{(x_i - \mu_{y_i})^2}{2|\Sigma|} \right) \right) + \log \pi_{y_i}$$

$$= \sum_{i=1}^m -\frac{1}{2} \log(2\pi|\Sigma|) - \frac{(x_i - \mu_{y_i})^2}{2|\Sigma|} + \log \pi_{y_i}$$

כעת, כדי למצוא את  $\argmax$  נבדוק לפי כל פרמטר בנפרד:  
למצוא MLE לפי  $\mu_k$ :

$$\frac{\partial}{\partial \mu_k} \ell(\theta | X, y) = \frac{\partial}{\partial \mu_k} \sum_{i=1}^m \left( -\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right)$$

כעת, פגשנו שאנחנו עובדים לפי  $\mu_k$ , כל הערכים מתאפסים חוץ מהאנלוגי. עכשיו  $\mu_k$  כלומר, נקרא:

$$\frac{\partial}{\partial \mu_k} \left( -\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right) = -\frac{1}{2\sigma_k^2} \cdot 2(x_i - \mu_k) \cdot -1 = \frac{x_i - \mu_k}{\sigma_k^2}$$

לכן, נקרא עבור  $y_i = k$  כל  $i$

$$\frac{\partial}{\partial \mu_k} \sum_{i=1}^m \left( -\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right) = \sum_{i=1}^m \frac{1}{\sigma_k^2} (x_i - \mu_k)$$

כלומר נקרא:

$$\mu_k^{MLE} = \frac{\sum (x_i - \mu_k) \delta(y_i = k)}{\sigma_k^2}$$

$$1 = \delta(y_i = k) \text{ אינדיקטור } y_i = k$$

כעת, נגזור לפי  $\sigma_{y_i}^2$ :

$$\frac{\partial}{\partial \sigma_k^2} \ell(\theta | X, y) = \sum_{i=1}^m \frac{\partial}{\partial \sigma_k^2} \left( -\frac{1}{2} \log(2\pi \sigma_k^2) - \frac{(x_i - \mu_{y_i})^2}{2\sigma_k^2} \right)$$

עבור לפי  $\sigma_k^2$  תאפס את כל הערכים  $y_i \neq k$  ותשאיר רק את  $y_i = k$ . נגזור ונשווה 0-ם ונקבל:

$$\frac{\partial}{\partial \sigma_k^2} \left( -\frac{1}{2} \log(2\pi \sigma_k^2) \right) = -\frac{1}{2} \cdot \frac{1}{2\pi \sigma_k^2} \cdot 2\pi = -\frac{1}{4\pi \sigma_k^2} \cdot 2\pi = -\frac{1}{2\sigma_k^2}$$

$$\frac{\partial}{\partial \sigma_k^2} \frac{(x_i - \mu_{y_i})^2}{2\sigma_k^2} = \frac{(x_i - \mu_{y_i})^2}{2\sigma_k^4}$$

$$\Rightarrow \sum_{i=1}^m \left( -\frac{1}{2\sigma_k^2} + \frac{(x_i - \mu_k)^2}{2\sigma_k^4} \right) \delta(y_i = k) = 0$$

נכפיל ב  $2\sigma_k^4$  ונקבל:

$$\sum_{i=1}^m (-\sigma_k^2 + (x_i - \mu_k)^2) \delta(y_i = k) = 0$$

$$\rightarrow \sum_{i=1}^m -\sigma_k^2 \delta(y_i = k) + \sum_{i=1}^m (x_i - \mu_k)^2 \delta(y_i = k) = 0$$

$$\rightarrow -\sigma_k^2 \cdot n_k + \sum_{i=1}^m (x_i - \mu_k)^2 \delta(y_i = k) = 0$$

$$\rightarrow \frac{\sum_{i=1}^m (x_i - \mu_k)^2 \delta(y_i = k)}{n_k} = \sigma_k^2$$

וכעת, נגזיר לפי  $\pi_k$ :

$$\frac{\partial}{\partial \pi_k} \ell(\theta | X, y) = \frac{\partial}{\partial \pi_k} \log \pi y_i$$

מכיון שבניכוי ההתחלה קאול  $\sum_{i=1}^k \pi_i = 1$  נשתמש בלגראנז' ונקבל:

$$\mathcal{L} = \ell(\theta | X, y) + \lambda \left( \sum_{i=1}^k \pi_i - 1 \right)$$

פיתרנו כאן

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{n_k}{\pi_k} - \lambda = 0 \Leftrightarrow \pi_k = \frac{n_k}{\lambda}$$

$$m = \lambda \quad \text{כיון ש-} \sum_{i=1}^k \pi_i = 1 \quad \text{נקבל ש-}$$

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

כאשר  $m$  הוא מספר הנתונים.

b.2 עבור  $d$  פיצורים יתקיים :

$$\mathcal{L}(\theta | x, y) = \prod_{i=1}^m f_{x|y=y}(x_i) \cdot f_{y|\theta}(y_i)$$

$$= \prod_{j=1}^d \prod_{i=1}^m N(x_{ij} | \mu_{y_{ij}}, \sigma_{y_{ij}}^2) \cdot \text{Mult}(y_i | \pi)$$

נשים לב שלכל פיצור הנצטרך שנקרא  $\mu_k, \pi_k, \sigma_k^2$  תאוסם את כל האופציות עם פיצורים אחרים. כלומר, הפיצורים הם קתגוריה ולכן לכל  $[d]$  נקרא :

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

$$\sigma_{kj}^{2, MLE} = \frac{1}{n_k} \sum (x_{ij} - \mu_{kj})^2$$

$$\mu_{kj}^{MLE} = \frac{1}{n_k} \sum (x_{ij} - \mu_{kj})$$

: likelihood-ის სიზუსტის ზოგადი (a.3)

$$L(\theta | x, y) = \prod_{i=1}^m f_{x,y|\theta}(x_i, y_i) = \prod_{i=1}^m f_{x|y=y}(x_i) \cdot f_{y|\theta}(y_i)$$

$$= \prod_{i=1}^m \text{Pois}(x_i | \lambda_{y_i}) \cdot \text{Mult}(y_i)$$

: log-likelihood სიზუსტის ზოგადი

$$\ell(\theta | x, y) = \log \left( \prod_{i=1}^m \text{Pois}(x_i | \lambda_{y_i}) \cdot \text{Mult}(y_i | \pi) \right)$$

$$= \sum_{i=1}^m \log(\text{Pois}(x_i | \lambda_{y_i}) + \log(\text{Mult}(y_i | \pi)))$$

$$= \sum_{i=1}^m \log \left( \frac{e^{-\lambda_{y_i}} \cdot \lambda_{y_i}^{x_i}}{x_i!} \right) + \log(\pi_{y_i})$$

$$= \sum_{i=1}^m \log(e^{-\lambda_{y_i}} \cdot \lambda_{y_i}^{x_i}) - \log(x_i!) + \log(\pi_{y_i})$$

$$= \sum_{i=1}^m -\lambda_{y_i} + \log(\lambda_{y_i}^{x_i}) - \log(x_i!) + \log(\pi_{y_i})$$

$$= \sum_{i=1}^m -\lambda_{y_i} + x_i \log(\lambda) - \log(x_i!) + \log(\pi_{y_i})$$

$$= \sum_{y_i=k} \left[ \log(\lambda_k) \cdot \sum_{i | y_i=k} x_i - n_k \cdot \lambda_k + n_k \cdot \log(\pi_k) \right] - \sum_{i=1}^m \log(x_i!)$$

כעת, נבחר את  $\lambda$  כאלו שהסתברותם  $0 < \delta$ .  
 צביעה על  $\lambda_k$ :

$$\frac{\partial}{\partial \lambda_k} \ell(\theta | x, y) = \frac{\partial}{\partial \lambda_k} \sum_{y_i=k} [\log(\lambda_k) \cdot \sum_{i|y_i=k} x_i - n_k \cdot \lambda_k]$$

$$= \frac{1}{\lambda_k} \cdot \sum_{i|y_i=k} x_i - n_k$$

$$\frac{1}{\lambda_k} \cdot \sum_{i|y_i=k} x_i - n_k = 0 \Leftrightarrow \frac{1}{\lambda_k} \cdot \sum_{i|y_i=k} x_i = n_k$$

$$\Leftrightarrow \lambda_k = \frac{1}{n_k} \cdot \sum_{i|y_i=k} x_i$$

אם  $\pi_k$  הוא ההסתברות של  $y_i = k$  (א.2) - נ

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

(b.3)  $\text{pen}$  ו-  $\text{likelihood}$

$$\ell(\theta | x, y) = \prod_{i=1}^m f_{x|y=y_i}(x_i) \cdot f_{y|\theta}(y_i)$$

$$= \prod_{i=1}^m \prod_{j=1}^d \text{pois}(x_{ij} | \lambda_{y_i,j}) \cdot \text{Mult}(y_i | \pi)$$



: log-likelihood -η σ³ρηθ ηκ ρθη

$$Q(\theta|x,y) = \log \left( \prod_{i=1}^m \prod_{j=1}^d \text{Pois}(x_{ij} | \lambda_{y_{ij}}) \cdot \text{Mult}(y_i | \pi) \right)$$

$$= \sum_{i=1}^m \sum_{j=1}^d \log(\text{Pois}(x_{ij} | \lambda_{y_{ij}})) + \log(\text{Mult}(y_i | \pi))$$

$$= \sum_{i=1}^m \sum_{j=1}^d \log(\text{Pois}(x_{ij} | \lambda_{y_{ij}})) + \sum_{i=1}^m \log(\text{Mult}(y_i | \pi))$$

$$= \sum_{k=1}^K \sum_{j=1}^d \log \left( \frac{e^{-\lambda_{kj}} \cdot \lambda_{kj}^{x_{ij}}}{x_{ij}!} \right) + n_k \cdot \log(\pi_k)$$

$$= \sum_{k=1}^K \sum_{j=1}^d \left[ n_k \cdot (-\lambda_{kj} + \log(\lambda_{kj})) \sum_{i|y_i=k} x_{ij} + n_k \cdot \log(\pi_k) \right] - \sum_{i=1}^m \log(x_{ij}!)$$

-e ρ³ρηη (b.2)δ ηηηρ

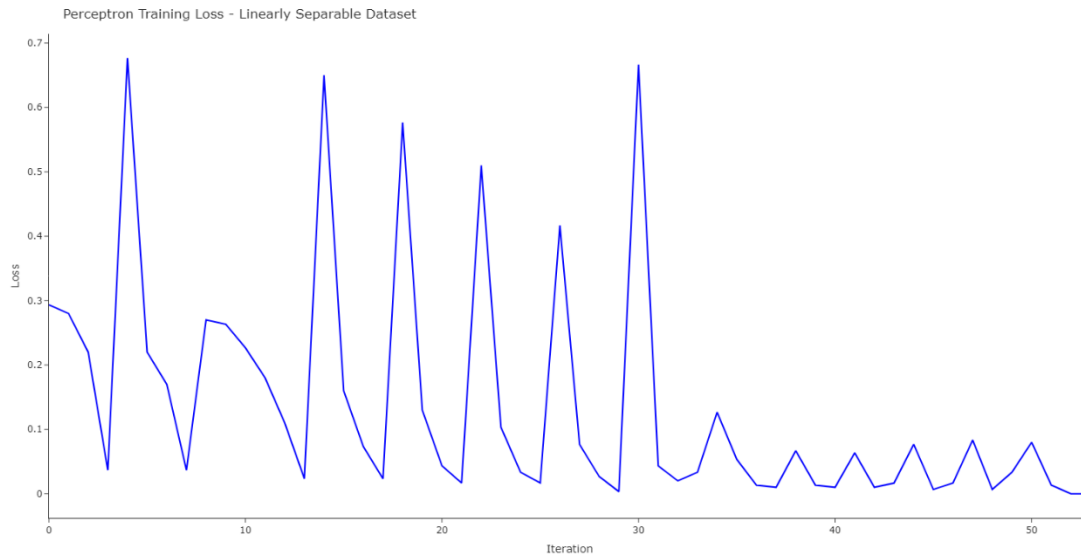
$$\lambda_k^{MLE} = \frac{1}{n_k} \cdot \sum_{i|y_i=k} x_i$$

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

## Practical part – IML

### 3.1 Perceptron Classifier

#### 1. Fitting and plotting over the linearly\_separable.npy dataset, what can we learn from the plot?



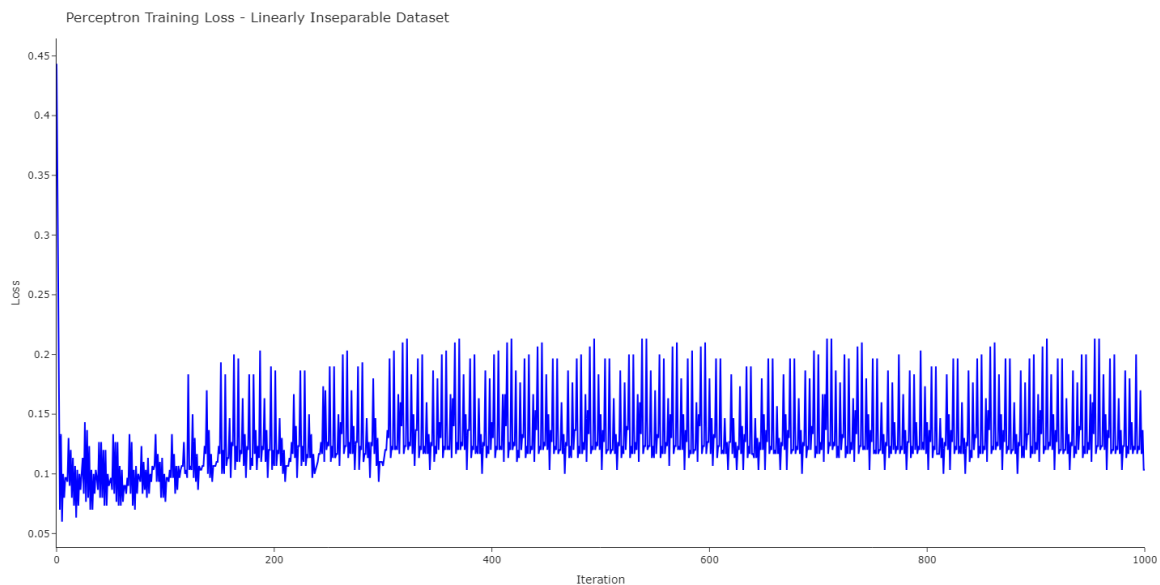
הגרף ממחיש את התהליך האיטרטיבי של אלגוריתם של Perceptron על מערך נתונים הניתן להפרדה ליניארית. זה מראה את יכולתו של האלגוריתם למצוא hyperplane מפריד בהצלחה, זאת ניתן לראות על ידי התייצבותו של הגרף בנקודה אפס לאחר 50+ איטרציות. התנהגות זו מדגימה את יעילותו ותכונת ההתכנסות של ה-Perceptron על מערכי נתונים הניתנים להפרדה ליניארית.

נשים לב שבתחילה, Perceptron מתחיל עם מספר גבוה של סיווגים שגויים, אך ככל מספר האיטרציות עולה, מתעדכן וקטור המשקולים ומספר הסיווגים השגויים הולך ופוחת משמעותית עד התייצבותו באפס.

התנודות החדות בגרף נובעות ככל הנראה מאופי עדכוני המשקל של ה-Perceptron, כאשר הוא מבצע תיקונים עבור נקודות שסווגו בצורה שגויה. אך החל מאיטרציה 30 הוא מצליח לסווג את הנקודות בצורה נכונה עד שלבסוף הוא מצליח להביא את ה MLE לאפס.

בנוסף, מכיוון שהמערך נתונים הוא מופרד ליניארית, מובטח שהמסווג ימצא hyperplane שיחלק אותו באופן ליניארי ומוצלח ואכן הגרף מאשר זאת על ידי התייצבותו החל מאיטרציה 50+ סביב נקודת האפס.

2. 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?



מהגרף ניתן לראות שעבור מערך נתונים שאינו ניתן להפרדה האלגוריתם אינו מתכנס ומתייצב סביב נקודת האפס כלל מה שממחיש את חוסר ההצלחה שלו בלמצוא hyperplane מפריד בהצלחה.

נשים לב לשונות בין שני הגרפים הבאה לידי ביטוי נקודות הבאות:

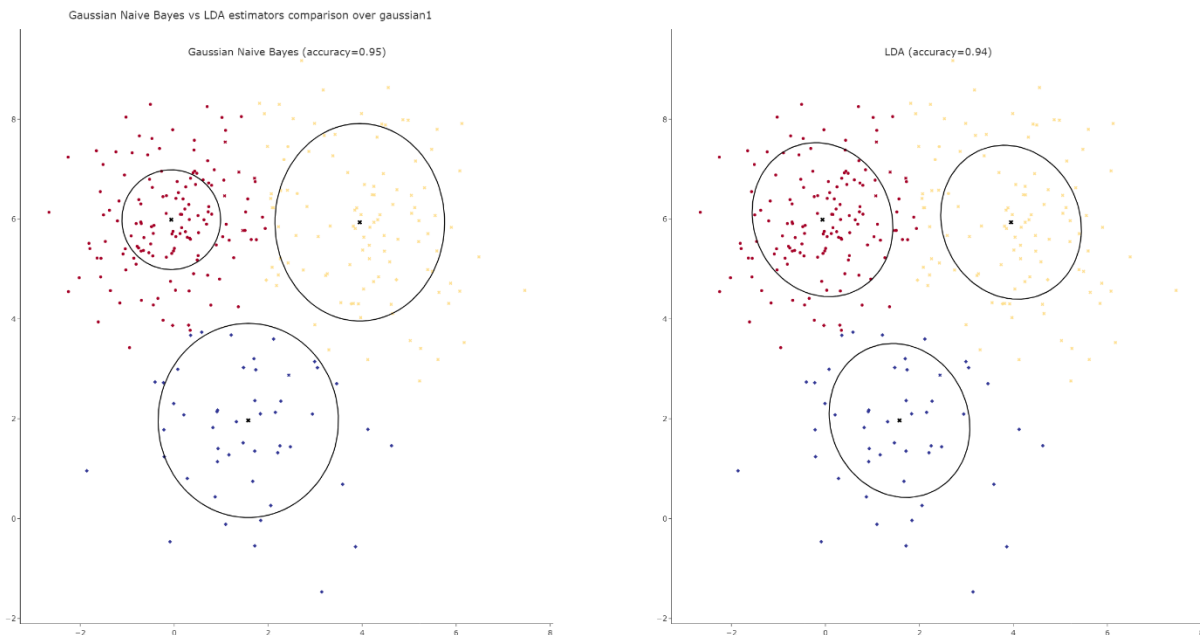
התנהגות הפונקציה: במקרה של מערך נתונים הניתן להפרדה, הפונקציה הולכת ומתייצבת על אפס בעוד שבגרף של מערך נתונים בלתי ניתן להפרדה הפונקציה עולה ויורדת ולא מתכנסת בשום שלב עד לעצירת האלגוריתם אחרי max\_iter איטרציות.

ערכי loss נמוכים: במקרה של מערך נתונים הניתן להפרדה, ערך הloss הולך וקטן במהלך האיטרציות ואילו בגרף מערך נתונים בלתי ניתן להפרדה, ערכי הloss נשארים גבוהים ומתנודדים בערכי טווח מסוים, מה שמשקף סיווגים שגויים מתמשכים.

ההבדל בנתונים מדגיש את ההשפעה של הפרדת הנתונים על יכולתו של ה-Perceptron לייעל את המטרה שלו. עבור נתונים הניתנים להפרדה ליניארית, ה-Perceptron מוצא פתרון מושלם במרחב הפרמטרים שממזער את פונקציית המטרה לאפס. עבור נתונים בלתי ניתנים להפרדה באופן ליניארי, ה-Perceptron נאבק למצוא פתרון יציב, מה שמוביל לעלייה וירידה מתמשכת כאשר הוא מנסה להתייצב ולמזער ללא הצלחה סיווגים שונים.

## 3.2 Bayes Classifiers

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

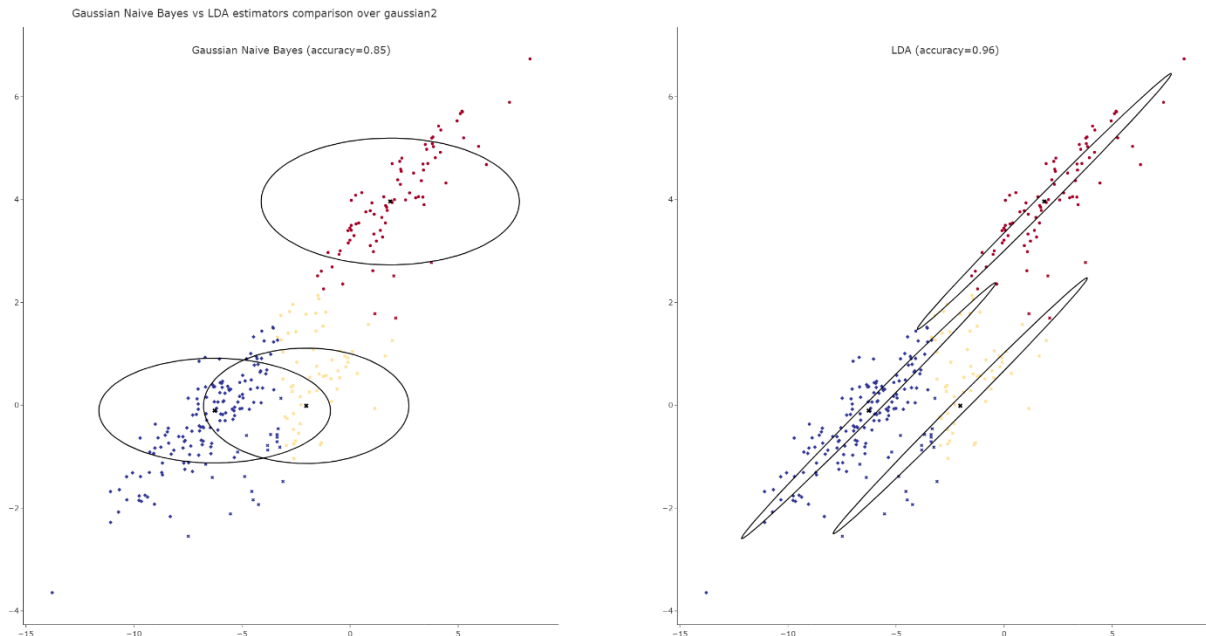


מהגרף ניתן להסיק כי ההתפלגות המשותפת שממנה נדגמו הדגימות היא נמוכה. נסתכל על תוצאות הגרף של ה-LDA classifier, מודל זה משתמש בחישוביו בהנחה שיש תלות בין הפיצ'רים. כמו כן, עבור כל מחלקה הוא משתמש באותה מטריצת שונות משותפת בחישוב (לכן נקבל שהאליפסות הן זהות עבור כל מחלקה) בעוד שה GNB classifier אינו משתמש בהנחה שיש תלות בין הפיצ'רים ועבור כל מחלקה הוא משתמש בשונות שלה עצמה (לכן נקבל שכל אליפסה היא בגודל שונה).

לכן, מבחינת ביצועים שני ה-classifiers משיגים תוצאות יחסית דומות וקרובות כיוון שהשונות המשותפת היא נמוכה מאוד ולכן ב-LDA אף שהוא משתמש בהנחה שיש תלות בין הפיצ'רים, זה לא משפיע על החישובים.

ניתן להסיק גם שההתפלגות המשותפת היא נמוכה על ידי פיזור הנקודות. נשים לב שעבור נקודה כלשהי, ערך ה- $x$  שלה לא משפיע על ערך ה- $y$ .

2. 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?



מהגרף ניתן להסיק שההתפלגות שממנה נדגמו הדגימות היא גבוהה משמעותית. ניתן לראות זאת מהגרף של LDA classifier שמניח שישנה תלות בין הפיצ'רים ואכן רואים שישנה התפלגות משותפת על ידי צורת האליפסה העקומה שנוצרה. בנוסף, ניתן לראות זאת על ידי פיזור הנקודות, שעבור נקודה כלשהי, ערך ה  $x$  משפיע על ערך ה  $y$  (ככל שערך האיקס גבוה, ערך ה  $y$  גבוה גם).

לעומת זאת, GNB classifier לא מצליח לתפוס את ההתפלגות, שכן הוא לא משתמש בהנחה שיש תלות בין הפיצ'רים, לכן, עבור הסט נתונים הזה, הביצועים שלו יהיו נמוכים לעומת ה LDA שמשיג תוצאה ממש גבוהה.

לכן נסיק שעבור סט הנתונים הזה, ההתפלגות המשותפת היא גבוהה ולכן המודל לחיזוי שיהיה מתאים יותר הוא LDA classifier.