

# (67577) מבוא למערכות לומדות | תרגיל 3

שם: נמרוד בר גיורא | ת"ז: 207090622

## חלק תאורטי

### Hard and Soft SVM

#### שאלה 1

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

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

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

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

**הוכחה:** נשים לב ש-hyperplane  $(w^*, b^*)$  עם  $w^* \in \mathbb{R}^d$  ו- $b^* \in \mathbb{R}$  הוא פתרון אופטימלי של הבעיה:

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

אם ורק אם הוא פתרון אופטימלי של הבעיה:

$$\underset{(w,b)}{\operatorname{argmin}} \left\| \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \right\|^2 \quad \text{s.t.} \quad \forall i \in [m] : y_i \left\langle \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \middle| x_i \right\rangle \geq 1$$

כאשר מוסיפים למטריצה  $X$  עמודת 1-דים. כלומר  $[x_i]_{d+1} = 1$  (זו הכניסה ה- $d+1$  של השורה ה- $i$  של  $X$ ) לכל  $i \in [m]$  הסיבה לכך היא שלכל  $i \in [m]$  מתקיים ש-

$$y_i (\langle w^* | x_i \rangle + b) \geq 1 \iff y_i \left( \sum_{j=1}^d w_j^* [x_i]_j + b \right) \geq 1 \iff y_i \left( \sum_{j=1}^d w_j^* [x_i]_j + b [x_i]_{d+1} \right) \geq 1 \iff y_i \left\langle \begin{bmatrix} w_1^* \\ \vdots \\ w_d^* \\ b \end{bmatrix} \middle| x_i \right\rangle \geq 1$$

ובנוסף לכל  $w \in \mathbb{R}^d$  מתקיים ש-

$$\|w^*\|^2 \leq \|w\|^2 \iff \sum_{i=1}^d (w_i^*)^2 \leq \sum_{i=1}^d (w_i)^2 \iff \sum_{i=1}^d (w_i^*)^2 + b^2 \leq \sum_{i=1}^d (w_i)^2 + b^2 \iff \left\| \begin{bmatrix} w_1^* \\ \vdots \\ w_d^* \\ b \end{bmatrix} \right\|^2 \leq \left\| \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \right\|^2$$

ולכן אם נוסיף עמודת 1-דים כנ"ל ל- $X$ , נוכל פשוט לעבוד עם הבעיה:

$$\underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \|w\|^2 \quad \text{s.t.} \quad \forall i \in [m] : y_i \langle w | x_i \rangle \geq 1$$

נגדיר  $\mathbf{a} = 0_{\mathbb{R}^{d+1}}, Q = 2I_{d+1}$

$$\mathbb{R}^m \ni \mathbf{d} = \begin{bmatrix} -1 \\ \vdots \\ -1 \end{bmatrix}, \quad \mathbb{R}^{m \times (d+1)} \ni A = -1 \cdot \begin{bmatrix} - & y_1 x_1^\top & - \\ & \vdots & \\ - & y_m x_m^\top & - \end{bmatrix}$$

ונקבל ש- $w^* \in \mathbb{R}^{d+1}$  הוא פתרון של הבעיה הנ"ל אם ורק אם לכל  $w \in \mathbb{R}^{d+1}$  מתקיים ש-

$$\begin{aligned} \|w^*\|^2 \leq \|w\|^2 &\iff w^{*\top} w^* \leq w^\top w \iff w^{*\top} I_d w^* \leq w^\top I_d w \iff \\ &\iff \frac{1}{2} w^{*\top} Q w^* + \overbrace{a w^*}^{=0} \leq \frac{1}{2} w^\top Q w + \overbrace{a w}^{=0} \iff w^* = \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} \left( \frac{1}{2} w^\top Q w + \mathbf{a}^\top w \right) \end{aligned}$$

ובנוסף:

$$\begin{aligned} \forall i \in [m] : y_i \langle w^* | x_i \rangle \geq 1 &\iff \forall i \in [m] : -y_i \langle w^* | x_i \rangle \leq -1 \iff \\ &\iff \forall i \in [m] : -y_i x_i^\top w^* \leq -1 \iff A w^* \leq \mathbf{d} \end{aligned}$$

כלומר אם ורק אם הוא פתרון אופטימלי של הבעיה:

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

■

## שאלה 2

2. Consider the Soft-SVM optimization problem:

$$\operatorname{argmin}_{\mathbf{w}, \{\xi_i\}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i \quad \text{s.t.} \quad \forall i \ y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i \wedge \xi_i \geq 0 \quad (3)$$

Denote the hinge-loss function as  $\ell^{\text{hinge}}(a) := \max\{0, 1 - a\}$ . Show that the Soft-SVM optimization problem is equivalent to the following unconstrained optimization problem:

$$\operatorname{argmin}_{\mathbf{w}, \{\xi_i\}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \ell^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \quad (4)$$

**הוכחה:** יהי  $w^* \in \mathbb{R}^d, \{\xi_i^*\}_{i=1}^m \subseteq \mathbb{R}$  - פתרון אופטימלי לבעיה (3). כלומר לכל  $w \in \mathbb{R}^d$  ולכל  $\{\xi_i\}_{i=1}^m$  מתקיים ש-

$$\frac{\lambda}{2} \|w^*\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i^* \leq \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

נשים לב שלכל  $i \in [m]$  מתקיים ש-

$$\xi_i^* = \ell^{\text{hinge}}(y_i \langle w^* | x_i \rangle) = \max\{0, 1 - y_i \langle w^* | x_i \rangle\} = \begin{cases} 0 & 1 < y_i \langle w^* | x_i \rangle \\ 1 - y_i \langle w^* | x_i \rangle & 1 \geq y_i \langle w^* | x_i \rangle \end{cases}$$

כי אם נניח בשלילה שקיים  $j$  כך ש- $\xi_j^* \neq \ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle)$  אז:

- אם  $1 < y_j \langle w^* | x_j \rangle$  אז  $\ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle) = 0$ , ולכן מכך ש- $\xi_j^* \geq 0$  נובע ש- $\ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle) < \xi_j^*$ .  
- אם  $1 \geq y_j \langle w^* | x_j \rangle$  אז  $\ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle) = 1 - y_j \langle w^* | x_j \rangle$ , ולכן מכך ש- $\xi_j^* \leq 1 - y_j \langle w^* | x_j \rangle$  נובע ש-

$$1 - y_j \langle w^* | x_j \rangle = \ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle) < \xi_j^*$$

ובכל מקרה נקבל שהפתרון  $w^*, \{\xi'_i\}_{i=1}^m$  כאשר  $\xi'_i = \xi_i^*$  לכל  $i \neq j$  ו-  $\xi'_j = \ell^{\text{hinge}}(y_j \langle w^* | x_j \rangle)$  קטן יותר מהפתרון האופטימלי, וזו סתירה. לכן  $w^*$  הנ"ל הוא פתרון אופטימלי לבעיה (4). בכיוון השני, אם  $w^*$  הוא פתרון אופטימלי ל-(4) אז לכל  $w \in \mathbb{R}^d$  מתקיים ש-

$$\frac{\lambda}{2} \|w^*\|^2 + \frac{1}{m} \sum_{i=1}^n \ell^{\text{hinge}}(y_i \langle w^* | x_i \rangle) \leq \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^n \ell^{\text{hinge}}(y_i \langle w | x_i \rangle)$$

נסמן  $\xi_i^* = \ell^{\text{hinge}}(y_i \langle w^* | x_i \rangle)$  לכל  $i \in [m]$  ונקבל שהאילווצים של בעיה (3) מתקיימים, ובנוסף לכי שהראינו לעיל לכל  $\{\xi_i\}_{i=1}^m$  מתקיים ש-

$$\forall i \in [m] : \quad \xi_i^* = \ell^{\text{hinge}}(y_i \langle w^* | x_i \rangle) \leq \xi_i \implies \frac{1}{m} \sum_{i=1}^m \xi_i^* \leq \frac{1}{m} \sum_{i=1}^m \xi_i$$

ולכן  $w^*, \{\xi_i^*\}_{i=1}^m$  הוא פתרון אופטימלי לבעיה (3). ■

## Naive Bayes Classifiers

### שאלה 3

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

$$\begin{aligned} y &\sim \text{Multinomial}(\pi) \\ x_j | y = k &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mu_{kj}, \sigma_{kj}^2) \end{aligned} \quad (6)$$

for  $\pi$  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 (5) under assumptions (6).
- (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 (5) under assumptions (6). You are encouraged to use the results from (3.a).

(a) כאשר הדגימות בעלות פיצ'ר בודד (חד-מימדיות) ההנחות של המודל Gaussian Naive Bayes הן שלכל  $i \in [m]$ :

$$\begin{aligned} y_i &\sim \text{Mult}(\pi) \\ x_i | y_i &\sim \mathcal{N}(\mu_{y_i}, \sigma_{y_i}^2) \end{aligned}$$

כאשר  $\pi \in [0, 1]^K$  ומתקיים  $\sum_{j=1}^K \pi_j = 1$ , ובנוסף  $\mu_k \in \mathbb{R}, \sigma_k^2$  לכל  $k \in [K]$ .

לכן כדי לאמן את המודל צריך לחשב אומד של  $\{\pi_k\}_{k \in [K]}$  - ההסתברויות של המחלקות, אומד של  $\{\mu_k\}_{k \in [K]}$  - התוחלות של המחלקות, ואומד של  $\{\sigma_k^2\}_{k \in [K]}$  - השונויות של המחלקות. נחשב את ערך ה-likelihood בהינתן הדגימות:

$$\begin{aligned} \mathcal{L}(\Theta | \mathbf{x}, \mathbf{y}) &= f_{X,Y|\Theta}(\{(x_i, y_i)\}_{i=1}^m) = \prod_{i=1}^m f_{X,Y|\Theta}(x_i, y_i) = \\ &= \prod_{i=1}^m f_{X|Y=y_i}(x_i) \cdot f_{Y|\Theta}(y_i) = \prod_{i=1}^m \mathcal{N}(x_i | \mu_{y_i}, \sigma_{y_i}^2) \cdot \text{Mult}(y_i | \pi) = \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma_{y_i}^2}} \cdot \exp\left(-\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2}\right) \cdot \pi_{y_i} \end{aligned}$$

כדי למצוא  $\text{argmax}$  ל-likelihood נוכל להשתמש ב-log-likelihood:

$$\begin{aligned}
\ell(\Theta \mid \mathbf{x}, \mathbf{y}) &= \ln \left( \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma_{y_i}^2}} \cdot \exp \left( -\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right) \cdot \pi_{y_i} \right) = \\
&= \sum_{i=1}^m \ln \left( \frac{1}{\sqrt{2\pi\sigma_{y_i}^2}} \cdot \exp \left( -\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right) \cdot \pi_{y_i} \right) = \\
&= \sum_{i=1}^m \left[ \ln \left( (2\pi\sigma_{y_i}^2)^{-\frac{1}{2}} \right) - \frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} + \ln(\pi_{y_i}) \right] = \\
&= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma_{y_i}^2) - \frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right] = \\
&= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) - \frac{1}{2} \ln(\sigma_{y_i}^2) - \frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \right] - \frac{m}{2} \ln(2\pi) = \\
&= \sum_{k \in [K]} \left[ n_k \cdot \ln(\pi_k) - \frac{n_k}{2} \ln(\sigma_k^2) - \sum_{i: y_i=k} \frac{(x_i - \mu_k)^2}{2\sigma_k^2} \right] - \frac{m}{2} \ln(2\pi)
\end{aligned}$$

כאשר  $n_k = \sum_{i=1}^m \mathbb{1}_{[y_i=k]}$  לכל  $k \in [K]$ . כדי למצוא את ה- $\text{argmax}$  נגזור ביחס לפרמטרים השונים  $\pi_k, \mu_k, \sigma_k^2$  - ונשווה ל-0:  
 כדי לגזור ביחס ל- $\pi_k$  נשתמש בשיטת כופלי לגרנז' - מכיוון שיש לנו את האילוץ  $\sum_{k \in [K]} \pi_k = 1$  נגדיר:

$$g(\pi) = \sum_{k \in [K]} \pi_k - 1$$

ונגזור את  $\mathcal{L} = \ell(\Theta \mid \mathbf{x}, \mathbf{y}) - \lambda g(\pi)$  כאשר  $\lambda$  הוא משתנה חדש:

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{\partial}{\partial \pi_k} \ell(\Theta \mid \mathbf{x}, \mathbf{y}) - \lambda \frac{\partial}{\partial \pi_k} g(\pi) = \frac{n_k}{\pi_k} - \lambda$$

נשווה ל-0 ונקבל ש-

$$\frac{n_k}{\pi_k} - \lambda = 0 \iff \pi_k = \frac{n_k}{\lambda}$$

וכדי למצוא את  $\lambda$  נציב באילוץ את הערך של  $\pi_k$  שמצאנו:

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

לכן אומד להסתברויות של המחלקות הוא:

$$\forall k \in [K] : \quad \hat{\pi}_k^{\text{MLE}} = \frac{n_k}{m}$$

נגזור ביחס ל- $\mu_k$ :

$$\frac{\partial}{\partial \mu_k} \ell(\Theta \mid \mathbf{x}, \mathbf{y}) = \frac{\partial}{\partial \mu_k} \left( - \sum_{i: y_i=k} \frac{(x_i - \mu_k)^2}{2\sigma_k^2} \right) = - \sum_{i: y_i=k} \frac{\partial}{\partial \mu_k} \left( \frac{(x_i - \mu_k)^2}{2\sigma_k^2} \right) = \sum_{i: y_i=k} \frac{x_i - \mu_k}{\sigma_k^2}$$

ונשווה ל-0:

$$\sum_{i: y_i=k} \frac{x_i - \mu_k}{\sigma_k^2} = 0 \iff \sum_{i: y_i=k} (x_i - \mu_k) = 0 \iff n_k \mu_k = \sum_{i: y_i=k} x_i \iff \mu_k = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} x_i$$

לכן אומד לתוחלות של המחלקות הוא:

$$\forall k \in [K] : \quad \hat{\mu}_k^{\text{MLE}} = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} x_i$$

נגזור ביחס ל- $\sigma_k^2$ :

$$\begin{aligned} \frac{\partial}{\partial \sigma_k^2} \ell(\Theta | \mathbf{x}, \mathbf{y}) &= \frac{\partial}{\partial \sigma_k^2} \left( -\frac{n_k}{2} \ln(\sigma_k^2) \right) - \sum_{i: y_i=k} \frac{\partial}{\partial \sigma_k^2} \left( \frac{(x - \mu_k)^2}{2\sigma_k^2} \right) = \\ &= -\frac{n_k}{2\sigma_k^2} + \sum_{i: y_i=k} \frac{(x - \mu_k)^2}{2(\sigma_k^2)^2} \end{aligned}$$

ונשווה ל-0:

$$\begin{aligned} -\frac{n_k}{2\sigma_k^2} + \sum_{i: y_i=k} \frac{(x - \mu_k)^2}{2(\sigma_k^2)^2} = 0 &\iff \frac{\sum_{i: y_i=k} (x - \mu_k)^2 - \sigma_k^2 n_k}{2(\sigma_k^2)^2} = 0 \iff \\ &\iff \sum_{i: y_i=k} (x - \mu_k)^2 - \sigma_k^2 n_k = 0 \iff \\ &\iff \sigma_k^2 = \frac{1}{n_k} \sum_{i: y_i=k} (x - \mu_k)^2 \end{aligned}$$

לכן אומד לשוניויות של המחלקות הוא:

$$\forall k \in [K] : \quad \hat{\sigma}_k^{\text{MLE}} = \frac{1}{n_k} \sum_{i: y_i=k} (x - \hat{\mu}_k^{\text{MLE}})^2$$

(b) כאשר הדגימות הן בעלות  $d$  פיצ'רים ההנחות של המודל Naive Bayes Classifier הן שלכל  $i \in [m]$

$$y_i \sim \text{Mult}(\boldsymbol{\pi})$$

$$[\mathbf{x}_i]_j | y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{y_i,j}, \sigma_{y_i,j}^2)$$

כלומר הלייבל ה- $i$  מתפלג מולטינומי עם  $\boldsymbol{\pi} \in [0, 1]^K$ ,  $\sum_{k \in [K]} \pi_k = 1$ , והפיצ'ר ה- $j$  של כל דגימה  $\mathbf{x}_i \in \mathbb{R}^d$  בהינתן  $y_i$  מתפלג נורמלי עם

תוחלת  $\mu_{y_i,j} \in \mathbb{R}$  ועם שונות  $\sigma_{y_i,j}^2 \in \mathbb{R}$ .

לכן כדי לאמן את המודל צריך למצוא אומד ל- $\{\pi_k\}_{k \in [K]}$ ,  $\{\mu_{k,j}\}_{k \in [K]}^{j \in [d]}$ ,  $\{\sigma_{k,j}^2\}_{k \in [K]}^{j \in [d]}$ . פונקציית ה-likelihood היא:

$$\begin{aligned} \mathcal{L}(\Theta | \mathbf{X}, \mathbf{y}) &= f_{X,Y|\Theta}(\{(\mathbf{x}_i, y_i)\}_{i=1}^m) = \prod_{i=1}^m f_{X,Y|\Theta}(\mathbf{x}_i, y_i) = \\ &= \prod_{i=1}^m f_{X|Y=y_i}(\mathbf{x}_i) \cdot f_{Y|\Theta}(y_i) = \prod_{i=1}^m f_{X_1|Y=y_i, \dots, X_d|Y=y_i}([x_i]_1, \dots, [x_i]_d) \cdot f_{Y|\Theta}(y_i) = \\ &\stackrel{\star}{=} \prod_{i=1}^m \left[ f_{Y|\Theta}(y_i) \cdot \prod_{j=1}^d f_{X_j|Y=y_i}([x_i]_j) \right] = \\ &= \prod_{i=1}^m \left[ \text{Mult}(y_i | \boldsymbol{\pi}) \cdot \prod_{j=1}^d \mathcal{N}([x_i]_j | \mu_{y_i,j}, \sigma_{y_i,j}^2) \right] = \\ &= \prod_{i=1}^m \left[ \pi_{y_i} \cdot \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_{y_i,j}^2}} \exp\left(-\frac{([x_i]_j - \mu_{y_i,j})^2}{2\sigma_{y_i,j}^2}\right) \right] \end{aligned}$$

🌟: לפי ההנחה של המודל ש- $[\mathbf{x}_i]_j \mid y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_{y_i,j}, \sigma_{y_i,j}^2)$  כדי למצוא  $\arg\max$  ל-likelihood נשתמש ב-log-likelihood:

$$\begin{aligned} \ell(\Theta \mid \mathbf{X}, \mathbf{y}) &= \ln \left( \prod_{i=1}^m \left[ \pi_{y_i} \cdot \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_{y_i,j}^2}} \exp \left( -\frac{([\mathbf{x}_i]_j - \mu_{y_i,j})^2}{2\sigma_{y_i,j}^2} \right) \right] \right) = \\ &= \sum_{i=1}^m \ln \left( \pi_{y_i} \cdot \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_{y_i,j}^2}} \exp \left( -\frac{([\mathbf{x}_i]_j - \mu_{y_i,j})^2}{2\sigma_{y_i,j}^2} \right) \right) = \\ &= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) - \sum_{j=1}^d \left[ \frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln(\sigma_{y_i,j}^2) + \frac{([\mathbf{x}_i]_j - \mu_{y_i,j})^2}{2\sigma_{y_i,j}^2} \right] \right] = \\ &= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) - \frac{1}{2} \sum_{j=1}^d \ln(\sigma_{y_i,j}^2) - \sum_{j=1}^d \frac{([\mathbf{x}_i]_j - \mu_{y_i,j})^2}{2\sigma_{y_i,j}^2} \right] - \frac{m \cdot d}{2} \ln(2\pi) = \\ &= \sum_{k \in [K]} \left[ n_k \ln(\pi_k) - \frac{n_k}{2} \sum_{j=1}^d \ln(\sigma_{k,j}^2) - \sum_{i: y_i=k} \sum_{j=1}^d \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right] - \frac{m \cdot d}{2} \ln(2\pi) \end{aligned}$$

כאשר  $n_k = \sum_{i=1}^m \mathbb{1}_{[y_i=k]}$  לכל  $k \in [K]$ . כדי למצוא את ה- $\arg\max$  נגזור ביחס לפרמטרים  $\pi_k, \mu_{k,j}, \sigma_{k,j}^2$  - ונשווה ל-0: כדי לגזור ביחס ל- $\pi_k$  נשתמש בכופלי לגרנז'. נגדיר -

$$g(\boldsymbol{\pi}) = \sum_{k \in [K]} \pi_k - 1$$

ונגזור את  $\mathcal{L} = \ell(\Theta \mid \mathbf{X}, \mathbf{y}) - \lambda g(\boldsymbol{\pi})$  כאשר  $\lambda$  הוא משתנה חדש:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \pi_k} &= \frac{\partial}{\partial \pi_k} \ell(\Theta \mid \mathbf{X}, \mathbf{y}) - \frac{\partial}{\partial \pi_k} \lambda g(\boldsymbol{\pi}) = \frac{n_k}{\pi_k} - \lambda \\ 1 &= \sum_{k \in [K]} \frac{n_k}{\lambda} \iff \lambda = m \end{aligned}$$

וע"י הצבה של הערך הזה באילוץ נקבל ש- $\pi_k = \frac{n_k}{\lambda}$

לכן אומד להסתברויות של המחלקות הוא:

$$\forall k \in [K] : \quad \hat{\pi}_k^{\text{MLE}} = \frac{n_k}{m}$$

נגזור את  $\ell(\Theta \mid \mathbf{X}, \mathbf{y})$  ביחס ל- $\mu_{k,j}$ :

$$\begin{aligned} \frac{\partial}{\partial \mu_{k,j}} \ell(\Theta \mid \mathbf{X}, \mathbf{y}) &= \frac{\partial}{\partial \mu_{k,j}} \left( - \sum_{i: y_i=k} \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right) = \\ &= - \sum_{i: y_i=k} \frac{\partial}{\partial \mu_{k,j}} \left( \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right) = \sum_{i: y_i=k} \frac{[\mathbf{x}_i]_j - \mu_{k,j}}{\sigma_{k,j}^2} \end{aligned}$$

ונשווה ל-0:

$$\sum_{i: y_i=k} \frac{[\mathbf{x}_i]_j - \mu_{k,j}}{\sigma_{k,j}^2} = 0 \iff \sum_{i: y_i=k} ([\mathbf{x}_i]_j - \mu_{k,j}) = 0 \iff \mu_{k,j} = \frac{1}{n_k} \sum_{i: y_i=k} [\mathbf{x}_i]_j$$

לכן אומד לתוחלות הוא:

$$\forall k \in [K], j \in [d] : \quad \hat{\mu}_{k,j}^{\text{MLE}} = \frac{1}{n_k} \sum_{i: y_i=k} [\mathbf{x}_i]_j$$

נגזור ביחס ל- $\sigma_{k,j}^2$ :

$$\begin{aligned}\frac{\partial}{\partial \sigma_{k,j}^2} \ell(\Theta | \mathbf{X}, \mathbf{y}) &= \frac{\partial}{\partial \sigma_{k,j}^2} \left( -\frac{n_k}{2} \sum_{j=1}^d \ln(\sigma_{k,j}^2) - \sum_{i: y_i=k} \sum_{j=1}^d \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right) = \\ &= -\frac{n_k}{2} \cdot \frac{\partial}{\partial \sigma_{k,j}^2} \ln(\sigma_{k,j}^2) - \sum_{i: y_i=k} \frac{\partial}{\partial \sigma_{k,j}^2} \left( \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right) = \\ &= -\frac{n_k}{2\sigma_{k,j}^2} + \sum_{i: y_i=k} \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2(\sigma_{k,j}^2)^2}\end{aligned}$$

ונשווה ל-0:

$$\begin{aligned}-\frac{n_k}{2\sigma_{k,j}^2} + \sum_{i: y_i=k} \frac{([\mathbf{x}_i]_j - \mu_{k,j})^2}{2(\sigma_{k,j}^2)^2} = 0 &\iff \frac{\sum_{i: y_i=k} ([\mathbf{x}_i]_j - \mu_{k,j})^2 - \sigma_{k,j}^2 n_k}{2(\sigma_{k,j}^2)^2} = 0 \iff \\ &\iff \sigma_{k,j}^2 = \frac{1}{n_k} \sum_{i: y_i=k} ([\mathbf{x}_i]_j - \mu_{k,j})^2\end{aligned}$$

לכן אומד לשוניות הוא:

$$\forall k \in [K], j \in [d]: \quad \boxed{\hat{\sigma}_{k,j}^{2\text{MLE}} = \frac{1}{n_k} \sum_{i: y_i=k} ([\mathbf{x}_i]_j - \hat{\mu}_{k,j}^{\text{MLE}})^2}$$

#### שאלה 4

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

$$\begin{aligned}y &\sim \text{Multinomial}(\boldsymbol{\pi}) \\ x_j | y = k &\stackrel{\text{ind.}}{\sim} \text{Poi}(\lambda_{kj})\end{aligned} \quad (7)$$

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

- 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 (5) under assumptions (7).
- Suppose  $\mathbf{x} \in \mathbb{R}^d$  (i.e each sample has  $d$  feature). Given a trainset  $S = \{(x_i, y_i)\}_{i=1}^m$  fit a Poisson Naive Bayes classifier solving (5) under assumptions (7). You are encouraged to use the results from (4.a).

(a) כאשר הדגימות הן בעלות פיצ'ר בודד ההנחות של המודל Poisson Naive Bayes הן שלכל  $i \in [m]$

$$\begin{aligned}y_i &\sim \text{Mult}(\boldsymbol{\pi}) \\ x_i | y_i &\stackrel{\text{ind.}}{\sim} \text{Poi}(\lambda_{y_i})\end{aligned}$$

כאשר  $\sum_{k \in [K]} \pi_k = 1, \pi \in [0, 1]^K$ . כדי לאמן את המודל צריך לחשב אומד של  $\{\pi_k\}_{k \in [K]}$  ושל  $\{\lambda_k\}_{k \in [K]}$ .  
נחיל בחישוב פונקציית ה-likelihood:

$$\begin{aligned}\mathcal{L}(\Theta | \mathbf{x}, \mathbf{y}) &= p_{X,Y|\Theta}(\{(x_i, y_i)\}_{i=1}^m) = \prod_{i=1}^m p_{X,Y|\Theta}(x_i, y_i) = \prod_{i=1}^m p_{X|Y=y_i}(x_i) \cdot f_{Y|\Theta}(y_i) = \\ &= \prod_{i=1}^m \text{Poi}(x_i | \lambda_{y_i}) \cdot \text{Mult}(y_i | \boldsymbol{\pi}) = \prod_{i=1}^m \frac{e^{-\lambda_{y_i}} \cdot (\lambda_{y_i})^{x_i}}{x_i!} \cdot \pi_{y_i}\end{aligned}$$

כדי למקסם את ה-likelihood אפשר למקסם את ה-log-likelihood:

$$\begin{aligned}\ell(\Theta | \mathbf{x}, \mathbf{y}) &= \ln \left( \prod_{i=1}^m \frac{e^{-\lambda_{y_i}} \cdot (\lambda_{y_i})^{x_i}}{x_i!} \cdot \pi_{y_i} \right) = \sum_{i=1}^m \left[ \ln \left( \frac{e^{-\lambda_{y_i}} \cdot (\lambda_{y_i})^{x_i}}{x_i!} \right) + \ln(\pi_{y_i}) \right] = \\ &= \sum_{i=1}^m \left[ \ln(e^{-\lambda_{y_i}}) + \ln((\lambda_{y_i})^{x_i}) - \ln(x_i!) + \ln(\pi_{y_i}) \right] = \\ &= \sum_{i=1}^m \left[ -\lambda_{y_i} + x_i \cdot \ln(\lambda_{y_i}) - \ln(x_i!) + \ln(\pi_{y_i}) \right] = \\ &= \sum_{k \in [K]} \left[ -n_k \lambda_k + \ln(\lambda_k) \cdot \sum_{i: y_i=k} x_i - \sum_{i: y_i=k} \ln(x_i!) + n_k \ln(\pi_k) \right]\end{aligned}$$

כאשר  $n_k = \sum_{i=1}^m \mathbb{1}_{[y_i=k]}$  לכל  $k \in [K]$ .

נגזור ביחס לפרמטר  $\pi_k$  בעזרת כופלי לגרנז. נגדיר  $g(\pi) = \sum_{k \in [K]} \pi_k - 1$  כגדיר את  $\ell(\Theta | \mathbf{x}, \mathbf{y}) - \gamma g(\pi)$  כאשר  $\gamma$  הוא משתנה חדש:

$$\frac{\partial}{\partial \pi_k} \ell(\Theta | \mathbf{x}, \mathbf{y}) - \gamma \frac{\partial}{\partial \pi_k} g(\pi) = \frac{n_k}{\pi_k} - \gamma$$

נשווה ל-0 ונקבל ש- $\pi_k = \frac{n_k}{\gamma}$ , וע"י הצבה של הערך הזה באילוץ נקבל ש- $\gamma = \sum_{k \in [K]} n_k = m$ . לכן אומד להסתברויות של המחלקות הוא:

$$\forall k \in [K] : \quad \hat{\pi}_k^{\text{MLE}} = \frac{n_k}{m}$$

נגזור את  $\ell(\Theta | \mathbf{x}, \mathbf{y})$  ביחס ל- $\lambda_k$ :

$$\frac{\partial}{\partial \lambda_k} \ell(\Theta | \mathbf{x}, \mathbf{y}) = -n_k + \frac{1}{\lambda_k} \cdot \sum_{i: y_i=k} x_i - 0 + 0$$

נשווה ל-0 ונקבל:

$$-n_k + \frac{1}{\lambda_k} \cdot \sum_{i: y_i=k} x_i = 0 \quad \iff \quad \lambda_k = \frac{1}{n_k} \sum_{i: y_i=k} x_i = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} \cdot x_i$$

לכן אומד מתאים לפרמטרים של ההתפלגות הוא:

$$\forall k \in [K] : \quad \hat{\lambda}_k^{\text{MLE}} = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} \cdot x_i$$

(b) כאשר הדגימות  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  הן בעלות  $d$  פיצ'רים ההנחות של המודל Poisson Naive Bayes הן שלכל  $i \in [m]$ :

$$y_i \sim \text{Mult}(\pi)$$

$$[\mathbf{x}_i]_j | y_i \stackrel{\text{ind}}{\sim} \text{Poi}(\lambda_{y_i, j})$$

כדי לאמן את המודל צריך לחשב אומדים להסתברויות של המחלקות  $\{\pi_k\}_{k \in [K]}$  ולפרמטרים של ההתפלגות פואסון  $\{\lambda_{k, j}\}_{k \in [K]}^{j \in [d]}$ . נחשב את פונקציית ה-likelihood:

$$\begin{aligned}\mathcal{L}(\Theta | \mathbf{X}, \mathbf{y}) &= p_{X, Y | \Theta}(\{(\mathbf{x}_i, y_i)\}_{i=1}^m) = \prod_{i=1}^m p_{X, Y | \Theta}(\mathbf{x}_i, y_i) = \prod_{i=1}^m p_{X | Y=y_i}(\mathbf{x}_i) \cdot p_{Y | \Theta}(y_i) = \\ &= \prod_{i=1}^m p_{X_1 | Y=y_i, \dots, X_d | Y=y_i}([\mathbf{x}_i]_1, \dots, [\mathbf{x}_i]_d) \cdot p_{Y | \Theta}(y_i) \stackrel{\star}{=} \prod_{i=1}^m \left[ p_{Y | \Theta}(y_i) \cdot \prod_{j=1}^d p_{X_j | Y=y_i}([\mathbf{x}_i]_j) \right] \\ &= \prod_{i=1}^m \left[ \text{Mult}(y_i | \pi) \cdot \prod_{j=1}^d \text{Pois}([\mathbf{x}_i]_j | \lambda_{y_i, j}) \right] = \prod_{i=1}^m \left[ \pi_{y_i} \cdot \prod_{j=1}^d \frac{e^{-\lambda_{y_i, j}} \cdot (\lambda_{y_i, j})^{[\mathbf{x}_i]_j}}{[\mathbf{x}_i]_j!} \right]\end{aligned}$$



☀: לפי ההנחה של המודל שהפיצ'רים בלתי תלויים.

כדי למקסם את ה-likelihood אפשר למקסם את ה-log-likelihood:

$$\begin{aligned}\ell(\Theta | \mathbf{X}, \mathbf{y}) &= \ln \left( \prod_{i=1}^m \left[ \pi_{y_i} \cdot \prod_{j=1}^d \frac{e^{-\lambda_{y_i,j}} \cdot (\lambda_{y_i,j})^{[\mathbf{x}_i]_j}}{[\mathbf{x}_i]_j!} \right] \right) = \sum_{i=1}^m \ln \left( \pi_{y_i} \cdot \prod_{j=1}^d \frac{e^{-\lambda_{y_i,j}} \cdot (\lambda_{y_i,j})^{[\mathbf{x}_i]_j}}{[\mathbf{x}_i]_j!} \right) = \\ &= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) + \ln \left( \prod_{j=1}^d \frac{e^{-\lambda_{y_i,j}} \cdot (\lambda_{y_i,j})^{[\mathbf{x}_i]_j}}{[\mathbf{x}_i]_j!} \right) \right] = \\ &= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) + \sum_{j=1}^d \left[ -\lambda_{y_i,j} + [\mathbf{x}_i]_j \cdot \ln(\lambda_{y_i,j}) - \ln([\mathbf{x}_i]_j!) \right] \right] = \\ &= \sum_{i=1}^m \left[ \ln(\pi_{y_i}) - \sum_{j=1}^d \lambda_{y_i,j} + \sum_{j=1}^d [\mathbf{x}_i]_j \cdot \ln(\lambda_{y_i,j}) - \sum_{j=1}^d \ln([\mathbf{x}_i]_j!) \right] = \\ &= \sum_{k \in [K]} \left[ n_k \ln(\pi_k) - n_k \sum_{j=1}^d \lambda_{k,j} + \sum_{j=1}^d \left[ \ln(\lambda_{k,j}) \cdot \sum_{i: y_i=k} [\mathbf{x}_i]_j \right] - \sum_{i: y_i=k} \sum_{j=1}^d \ln([\mathbf{x}_i]_j!) \right]\end{aligned}$$

כאשר  $n_k = \sum_{i=1}^m \mathbb{1}_{[y_i=k]}$  לכל  $k \in [K]$ .

נגזור ביחס לפרמטר  $\pi_k$  בעזרת כופלי לגרנז. נגדיר  $g(\pi) = \sum_{k \in [K]} \pi_k - 1$  ונגזור את  $\ell(\Theta | \mathbf{x}, \mathbf{y}) - \gamma g(\pi)$  כאשר  $\gamma$  הוא משתנה חדש:

$$\frac{\partial}{\partial \pi_k} \ell(\Theta | \mathbf{X}, \mathbf{y}) - \gamma \frac{\partial}{\partial \pi_k} g(\pi) = \frac{n_k}{\pi_k} - \gamma$$

נשווה ל-0 ונקבל ש- $\pi_k = \frac{n_k}{\gamma}$ , וע"י הצבה של הערך הזה באילוץ נקבל ש- $\gamma = \sum_{k \in [K]} n_k = m$ . לכן אומד להסתברויות של המחלקות הוא:

$$\forall k \in [K] : \boxed{\hat{\pi}_k^{\text{MLE}} = \frac{n_k}{m}}$$

נגזור את  $\ell(\Theta | \mathbf{X}, \mathbf{y})$  ביחס ל- $\lambda_{k,j}$ :

$$\frac{\partial}{\partial \lambda_{k,j}} \ell(\Theta | \mathbf{X}, \mathbf{y}) = 0 - n_k + \frac{1}{\lambda_{k,j}} \cdot \sum_{i: y_i=k} [\mathbf{x}_i]_j - 0$$

נשווה ל-0 ונקבל:

$$-n_k + \frac{1}{\lambda_{k,j}} \cdot \sum_{i: y_i=k} [\mathbf{x}_i]_j = 0 \iff -n_k \lambda_{k,j} + \sum_{i: y_i=k} [\mathbf{x}_i]_j = 0 \iff \lambda_{k,j} = \frac{1}{n_k} \cdot \sum_{i: y_i=k} [\mathbf{x}_i]_j = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} \cdot [\mathbf{x}_i]_j$$

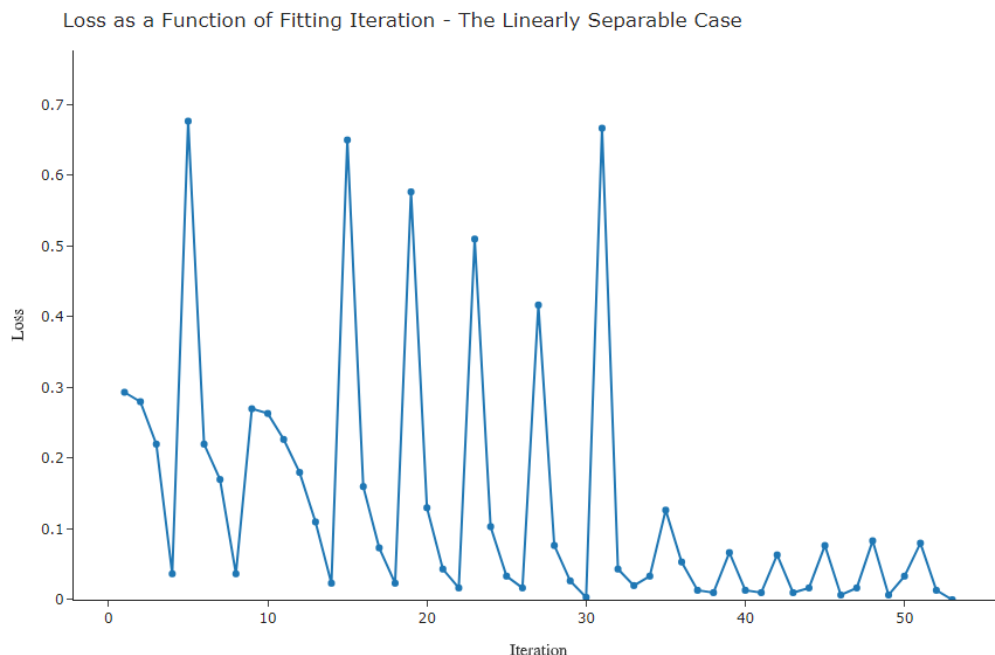
לכן אומד מתאים לפרמטרים של ההתפלגות הוא:

$$\forall k \in [K], j \in [d] : \boxed{\hat{\lambda}_{k,j}^{\text{MLE}} = \frac{1}{n_k} \sum_{i=1}^m \mathbb{1}_{[y_i=k]} \cdot [\mathbf{x}_i]_j}$$

## Perceptron Classifier

### שאלה 1

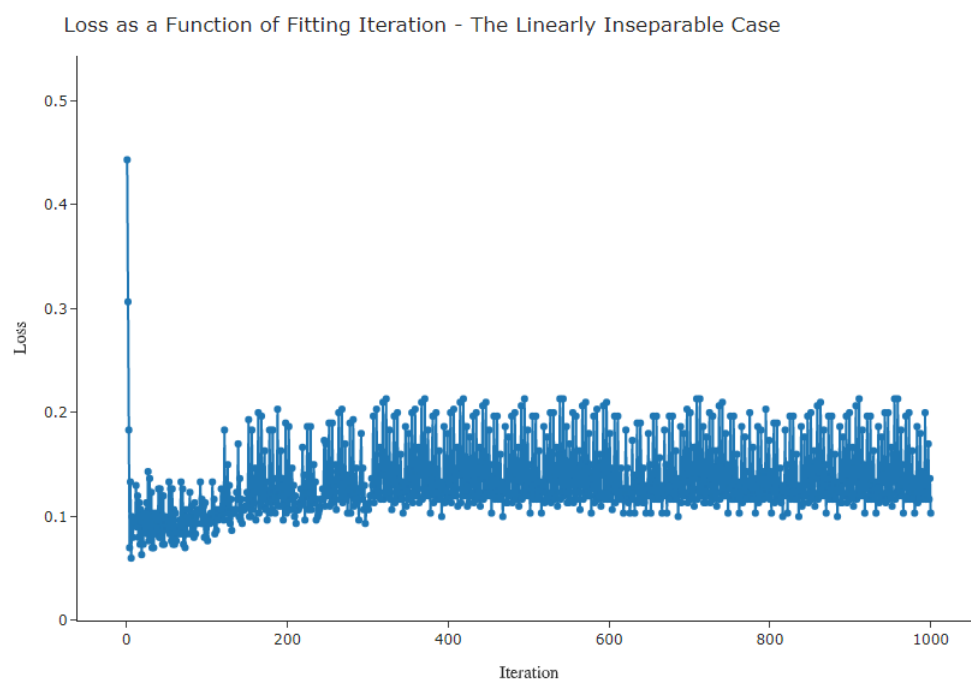
גרף המציג את ההתקדמות של האלגוריתם על הדאטא שניתן להפרדה ליניארית - מציג את ערך ה-Loss על ה-training data כפונקציה של מספר האיטרציות של האלגוריתם:



ניתן לראות שלאחר 50 איטרציות של האלגוריתם הוא מתכנס - ערך ה-Loss שווה ל-0 ממש. כלומר האלגוריתם מצליח להפריד ליניארית את הדאטא.

### שאלה 2

גרף המציג את ההתקדמות של האלגוריתם על הדאטא שלא ניתן להפרדה ליניארית:



ההבדל בין שני הגרפים ניכר - רואים בגרף השני שהאלגוריתם **לא מתכנס** על הדאטא שלא ניתן להפרדה ליניארית - גם אחרי 1000 איטרציות של האלגוריתם ערך ה-Loss נשאר גדול מ-0, כלומר האלגוריתם לא מצליח להפריד ליניארית את הדאטא. במונחים פורמליים - רק עבור הדאטא שניתן להפרדה ליניארית האלגוריתם מצא וקטור  $w \in \mathbb{R}^{d+1}$  כך שלכל  $i \in [m]$  מתקיים ש-

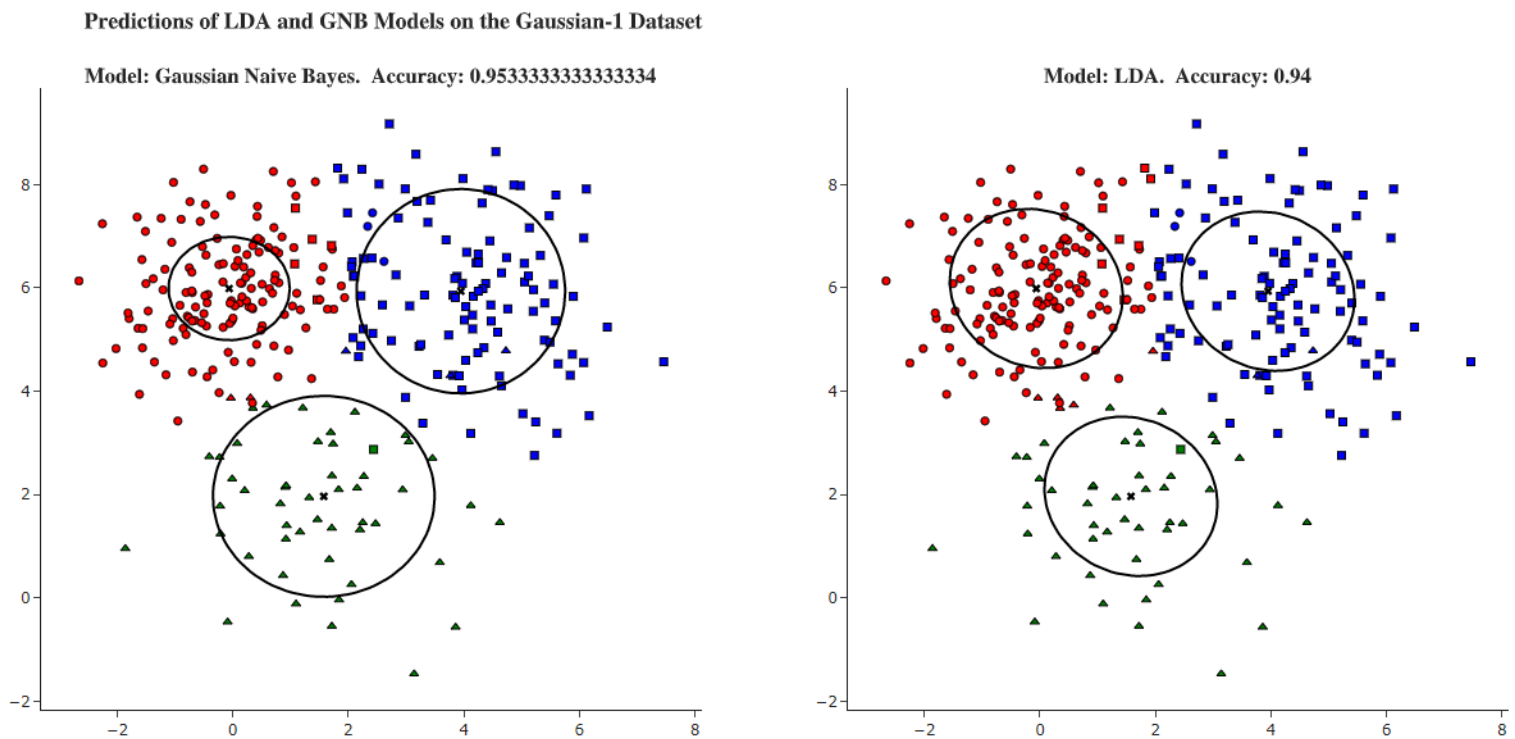
$$y_i \cdot \langle w | x_i \rangle > 0$$

(כשמוסיפים ל- $x$  את ה-intercept כדי שהעל-מישור יוכל לא לעבור בראשית הצירים). עבור הדאטא שאינו ניתן להפרדה ליניארית לא קיים וקטור  $w$  כזה, ולכן האלגוריתם לא התכנס.

## Bayes Classifiers

### שאלה 1

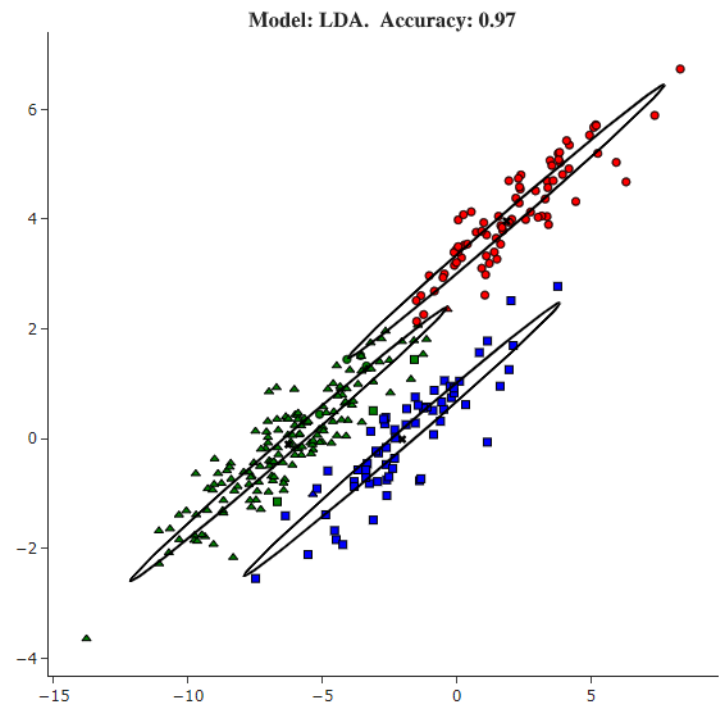
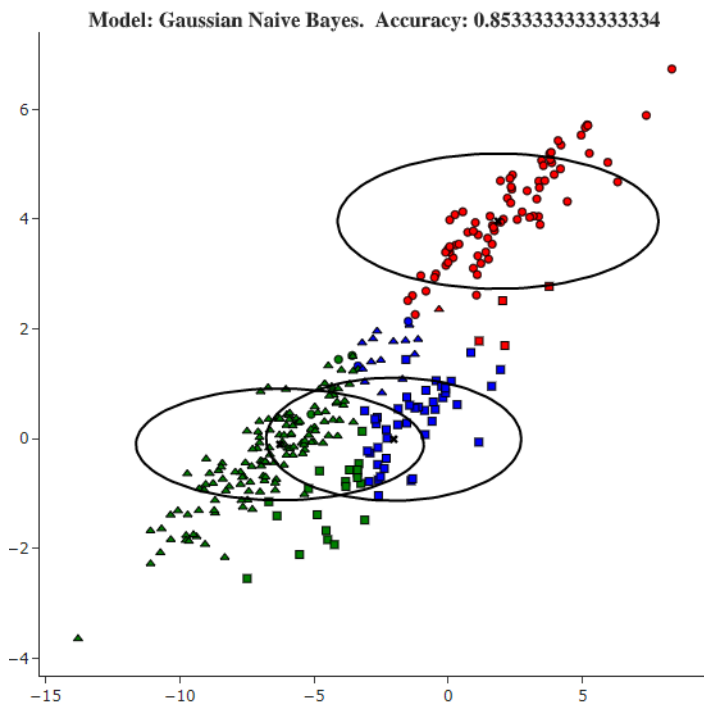
הגרף הנדרש:



לפי הגרפים הללו אפשר להסיק שההתפלגות שממנה נדגם הדאטא דומה יותר להתפלגות שהמודל Gaussian Naive Bayes מניח, כי הוא הצליח להפריד את הדאטא קצת יותר טוב מה-LDA. כלומר - הפיצ'רים בדאטא הם בלתי תלויים, ו-3 המחלקות נדגמו עם התפלגות נורמלית סביב 3 הנקודות (שמסומנות ב- $\times$ ) שניתן לראות בגרף.

להס

Predictions of LDA and GNB Models on the Gaussian-2 Dataset



בתרחיש הזה, דווקא המודל LDA הפריד את הדאטא יותר טוב מאשר ה-GNB.

באמת ניתן לראות שבדאטא הזה דווקא יש תלות בין הפיצ'רים. דרך אחת לראות זאת היא ע"י האליפסות שבגרפים המציגות את מטריצת השונות המשותפת של הפיצ'רים - ה-LDA לא מניח אי-תלות בין הפיצ'רים והוא באמת בונה מודל שבו השונות המשותפת שלהם שונה מ-0. לעומת זאת ההנחה של ה-GNB שהשונות המשותפת היא 0 גורמת לכך שהוא חוזה את הפיזור של הנקודות באופן פחות מדויק. בנוסף ניתן לראות גם שכאשר הערך בציר ה- $x$  גדל - כלומר הערך של הפיצ'ר הראשון גדל - גם הערך בציר ה- $y$  גדל - כלומר הערך של הפיצ'ר השני גדל.

מעבר לכך מכיוון שה-LDA זיהה ברמת דיוק גבוהה את הנקודות ניתן להסיק שההתפלגות שממנה נדגם הדאטא היא מהתפלגות נורמלית עם 3 מוקדים שקרובים לנקודות (שמסומנות ב- $\times$ ) שניתן לראות בגרף.