

# Introduction to Machine Learning (67577)

## Exercise 3 Classification

Second Semester, 2023

### Contents

<b>1</b>	<b>Submission Instructions</b>	<b>2</b>
<b>2</b>	<b>Theoretical Part</b>	<b>2</b>
2.1	Hard- & Soft-SVM .....	2
2.2	Naive Bayes Classifiers .....	2
<b>3</b>	<b>Practical Part</b>	<b>3</b>
3.1	Perceptron Classifier .....	3
3.2	Bayes Classifiers .....	4

## 1 Submission Instructions

Please make sure to follow the general submission instructions available on the course website. In addition, for the following assignment, submit a single `ex3_ID.tar` file containing:

- An `Answers.pdf` file with the answers for all theoretical and practical questions (include plotted graphs *in* the PDF file).
- The following python files (without any directories): `loss_functions.py`, `perceptron.py`, `linear_discriminant_analysis.py`, `gaussian_naive_bayes.py`, `classifiers_evaluation.py`

The `ex3_ID.tar` file must be submitted in the designated Moodle activity prior to the date specified *in the activity*.

- Late submissions will not be accepted and result in a zero mark.
- Plots included as separate files will be considered as not provided.
- Do not forget to answer the Moodle quiz of this assignment.

## 2 Theoretical Part

### 2.1 Hard- & Soft-SVM

[Based on Lecture 3 and Recitation 5](#)

1. Prove that 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  $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{I} \mathbf{w}$

2. Consider the Soft-SVM optimization problem:

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

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

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

### 2.2 Naive Bayes Classifiers

[Based on Lecture 3 and Recitation 6.](#) Let  $\mathcal{X}$  be a domain set and  $\mathcal{Y} \in [K], K \in \mathbb{N}$  the response set and let us assume there exists a joint probability distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$  with  $f_{\mathcal{D}}$  the joint probability distribution function.

Recall the Bayes Optimal Classifier which predicts the response maximizing the posterior distribution:

$$\hat{y}^{MAP} := \operatorname{argmax}_{k \in [K]} f_{Y|X=\mathbf{x}}(k) = \operatorname{argmax}_{k \in [K]} \frac{f_{X|Y=k}(\mathbf{x}) f_Y(k)}{f_X(\mathbf{x})} \quad (5)$$

*Naive Bayes* classifiers are a family of classifiers realizing the Bayes Optimal classifier where we assume that all features are *independent*. That is, for  $\mathbf{x} \sim \mathcal{P}$  then  $f_{X_i, X_j}(x_i, x_j) = f_{X_i}(x_i) f_{X_j}(x_j) \quad \forall i, j$ .

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

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

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

- (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). 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 (5) under assumptions (6). You are encouraged to use the results from (3.a).
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:  $\boldsymbol{\pi} \in [0, 1]^K, \sum \pi_j = 1$ .

- (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 (5) under assumptions (7).
- (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 (5) under assumptions (7). You are encouraged to use the results from (4.a).

### 3 Practical Part

In the following part you will implement and compare different classifiers for different data scenarios. Be sure to have pulled the latest version of the GreenGilad/IML.HUJI repository.

#### 3.1 Perceptron Classifier

Based on Lecture 3 and Recitation 5. Complete the following implementations

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

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

- To retrieve the loss at each iteration **do not** change 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 `datasets/linearly_separable.npy` dataset, what can we learn from the plot?
  2. Next run the Perceptron algorithm over the `datasets/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?

## 3.2 Bayes Classifiers

Based on Lecture 3 and Recitation 6. Complete the following implementations

- Implement the `accuracy` function in the `metrics/loss_functions.py` file as described in the function documentation.
- Implement the LDA classifier in the `learners/classifiers/linear_discriminant_analysis.py` file as described in the class documentation. Use expressions derived in class.
- Implement the `GaussianNaiveBayes` classifier in the `learners/classifiers/gaussian_naive_bayes.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 `datasets/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?

2. Repeat the procedure above (while avoiding code repetition) for `datasets/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?

## 2 Theoretical Part

## የኢትዮጵያ የስራ

209144013

## 2.1 Hard- & Soft-SVM

Based on Lecture 3 and Recitation 5

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

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

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

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

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

$V = \begin{pmatrix} \cdot & w \\ \vdots & \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$  מתקיים  $w_0 = b$  וקיים מטרית  $W \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  ו- $\lambda$  כך ש- $\lambda$  מינימיזיריה ריבועית נורמליזציה, גודל מטרית ו- $b$ .  
 מינימיזציה של מטרית  $w_0 = b - \delta$  על מנת למשוך את objective מינימום. ( $n = d + 1$ )  
 מינימיזציה של מטרית  $a = \bar{o} \in \mathbb{R}^n$  על מנת  $|Q| = 2 \cdot \begin{pmatrix} 0 & \cdots & 0 \\ \vdots & & \\ 0 & & I_{n-1} \end{pmatrix} \in \mathbb{R}^{n \times n}$ .

$$\frac{1}{2} V^T Q V + \alpha^t V = \frac{1}{2} (w_0 \dots w_n) 2 \cdot \begin{pmatrix} 0 & \dots & 0 \\ \vdots & I_{n-1} \\ 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} w_0 \\ \vdots \\ w_n \end{pmatrix} + \bar{O}^t V = (w_1 \dots w_{n-1}) I_{n-1} \begin{pmatrix} w_1 \\ \vdots \\ w_{n-1} \end{pmatrix} = \|w\|^2$$

$$y_i (\langle w, x_i \rangle + b) \geq 1 \Leftrightarrow -y_i (w_0 + \sum_{i=1}^{n-1} w_i x_i) \leq -1$$

המקרה  $A_i = -y_i(1 - x_i^T -)$  מוצג בתרשים 1.1.1. וקטור  $y$  מוגדר כ-

$$\begin{aligned} A \cdot v \leq d &\iff \forall i \in [m] \quad A_i \cdot v \leq d_i \iff -y_i(1 - x_i^t - \begin{pmatrix} w_1 \\ \vdots \\ w_{n-1} \end{pmatrix}) \leq -1 & \text{: for } i \in [m] \\ &\iff -y_i(w_0 + \sum_{i=1}^{n-1} x_i w_i) \leq -1 \end{aligned}$$

כרכ'ה

2. Consider the Soft-SVM optimization problem:

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

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 unconstraint optimization problem:

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

: פונקציית  $i \in [m]$  פ. ס.  $S = \{(\mathbf{x}_i, y_i) \mid \mathbf{x} \in \mathbb{R}^d, y_i \in \mathbb{R}\}$  פונקציית  $w \in \mathbb{R}^d$  נורמליזציה (2)

הוכיחים  $\sum_i \xi_i = \sum_i \|\mathbf{w}\|^2$  מכיון שפונקציית גזע נורמליזציה (3)-היא מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$ ; הוכיחים שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$ ; הוכיחים שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$ .

. מוכיחים כי פונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם וה רק אם  $\xi_i = 0$  עבור  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 1$  פונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \leq 1$  פונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$ ,  $\xi_i \geq 0$  מכיון שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i \Leftrightarrow \xi_i \geq 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$ ,  $\xi_i = 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$  מכיון שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \leq 1$ .

הוכיחים שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $\xi_i = \ell^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) - 1$  מכיון שפונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$  אם ורק אם  $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$ .

(3)  $\Rightarrow$  פונקציית גזע נורמליזציה מינימלית ביחס ל- $\|\mathbf{w}\|^2 + \frac{1}{m} \sum_i \xi_i$

לצורך:

**Claim 4.1** Given a training set  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  and hyperplane  $(\mathbf{w}, b)$ , the Soft-SVM optimization problem (12) is equivalent to

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} (\lambda \|\mathbf{w}\|^2 + L_S^{\text{hinge}}((\mathbf{w}, b)))$$

where  $L_S^{\text{hinge}}((\mathbf{w}, b)) := \frac{1}{m} \sum_i \ell^{\text{hinge}}(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle + b)$

*Proof.* Given a specific hyperplane  $(\mathbf{w}, b)$  consider the minimization over  $\xi_1, \dots, \xi_m$ . Since we defined the auxiliary variables to be nonnegative, the optimal assignment of  $\xi_i$  is

$$\xi_i := \begin{cases} 0 & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) > 1 \\ 1 - y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) & \text{otherwise} \end{cases}$$

Thus  $\xi_i = \ell^{\text{hinge}}(y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b))$  ■

## 2.2 Naive Bayes Classifiers

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

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

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). Fitting means finding the expressions for the maximum likelihood estimators.

- אוניברסיטת תל אביב קבוצת סטודנטים מילויים מודול  $\hat{\pi}_i^{\text{MLE}}, \hat{\sigma}_i^{\text{MLE}}, \hat{\mu}_i$  - תל אביב (א. 3)

$$L(\theta | x, y) = f_\theta(x, y) = f_\theta(\{x_i, y_i\}_{i=1}^m) = \prod_{i=1}^m f_{x_i, y_i | \theta}(x_i, y_i) = \prod_{i=1}^m f_{x_i | \theta}(x_i) \cdot f_{y_i | \theta}(y_i) = \\ = \prod_{i=1}^m N(x_i | \mu_{y_i}, \sigma_{y_i}^2) \cdot \text{Mult}(y_i | \pi) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma_{y_i}^2}} \cdot e^{-\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2}} \cdot \pi_{y_i}$$

•  $y_i \in \{1, \dots, K\}$

: מינימיזציה של-log-likelihood -> מינימיזציה של  $L(\theta | x, y)$  (א. 3)

$$\log(L(\theta | x, y)) = \log\left(\prod_{i=1}^m N(x_i | \mu_{y_i}, \sigma_{y_i}^2) \cdot \text{Mult}(y_i | \pi)\right) = \sum_{i=1}^m \log(N(x_i | \mu_{y_i}, \sigma_{y_i}^2)) + \log(\text{Mult}(y_i | \pi)) \\ = \sum_{i=1}^m \log(\pi_{y_i}) + \log\left(\frac{1}{\sqrt{2\pi\sigma_{y_i}^2}} \cdot e^{-\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2}}\right) = \sum_{i=1}^m \log(\pi_{y_i}) + \log\left(\frac{1}{\sqrt{2\pi\sigma_{y_i}^2}}\right) + \log\left(e^{-\frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2}}\right) = \\ \text{using } \log \exp, \log x^m = m \log x, \frac{1}{\sqrt{2\pi\sigma^2}} = (2\pi\sigma^2)^{-\frac{1}{2}}$$

$$= \sum_{i=1}^m \log(\pi_{y_i}) - \frac{1}{2} \log(2\pi\sigma_{y_i}^2) - \frac{(x_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} = -\frac{1}{2} m \log(2\pi) + \sum_{i=1}^m (n_i \log(\pi_i) - \frac{1}{2} n_i \log \sigma_i^2 - \frac{1}{2\sigma_i^2} \sum_{j \neq i} (x_j - \mu_j)^2)$$

• נספח לוג-לייליהוד:  $L(\theta | x, y) = \prod_{i=1}^m \pi_{y_i}^{n_i} \cdot \prod_{j \neq i} \frac{1}{\sqrt{2\pi\sigma_{y_j}^2}} \cdot e^{-\frac{(x_j - \mu_j)^2}{2\sigma_{y_j}^2}}$

כדי למדוד את המודל, מינימיזזציה של-log-likelihood (א. 3)

$$g(\pi) = \sum_{i=1}^K \pi_i - 1 \quad \text{מינימיזזציה של } g(\pi) \text{ מושגית ב } \pi_i = \frac{n_i}{m} \quad \text{לפניהם}$$

: מינימיזזציה של-log-likelihood (א. 3)

$$\frac{\partial L}{\partial \pi_i} = \frac{\partial}{\partial \pi_i} \left( h_i \log(\pi_i) - \lambda \pi_i \right) = \frac{h_i}{\pi_i} - \lambda = 0 \Rightarrow \pi_i = \frac{h_i}{\lambda}$$

השוויה שנותר

$$\frac{m}{\lambda} = \frac{1}{\lambda} \sum_{i=1}^k n_i = \sum_{i=1}^k \frac{n_i}{\lambda} = \sum_{i=1}^k \pi_i = \eta \Rightarrow m = \lambda$$

$$:\lambda \text{ נס. קבוצה } \pi_i = \frac{n_i}{\lambda} \text{ נס. נס.}$$

$$\hat{\pi}_i^{\text{MLE}} = \frac{n_i}{m} \quad | \sigma$$

$$\frac{\partial \log(L(\theta|x,y))}{\partial \mu_i} = \frac{\partial}{\partial \mu_i} \left( -\frac{1}{2\sigma_i^2} \sum_{j:y_j=i} (x_j - \mu_i)^2 \right) = -\frac{1}{2\sigma_i^2} \sum_{j:y_j=i} -2(x_j - \mu_i) = 0 \Rightarrow -\frac{1}{\sigma_i^2} \sum_{j:y_j=i} (x_j - \mu_i) = 0 \Rightarrow \sum_{j:y_j=i} x_j = n_i \cdot \mu_i$$

↓

$$\frac{1}{n_i} \sum_{j:y_j=i} x_j = \mu_i$$

$$\hat{\mu}_i^{\text{MLE}} = \frac{1}{n_i} \sum_{j:y_j=i} x_j \quad | \sigma$$

$$\frac{\partial \log(L(\theta|x,y))}{\partial \sigma_i^2} = \frac{\partial}{\partial \sigma_i^2} \left( -\frac{1}{2\sigma_i^2} \sum_{j:y_j=i} (x_j - \mu_i)^2 - \frac{n_i}{2} \log \sigma_i^2 \right) = \frac{1}{2(\sigma_i^2)^2} \sum_{j:y_j=i} (x_j - \mu_i)^2 - \frac{n_i}{2\sigma_i^2} = 0 \Rightarrow$$

$\left(\frac{1}{x}\right)' = -\frac{1}{x^2}$   
 $\log'(x) = \frac{1}{x}$

$$\Rightarrow \frac{1}{\sigma_i^2} \sum_{j:y_j=i} (x_j - \mu_i)^2 = h_i \quad \Rightarrow \sigma_i^2 = \frac{1}{h_i} \sum_{j:y_j=i} (x_j - \mu_i)^2$$

$$\hat{\sigma}_i^{\text{MLE}} = \frac{1}{h_i} \sum_{j:y_j=i} (x_j - \mu_i)^2 \quad | \sigma$$

- (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).

$$- y \sim \text{Multinomial}(\pi)$$

(b) כ"ל רצוי מינימיזציה של פונקציית האפסון-טולסן (epsilon-constraint function) על מנת למצוא את הערך המרבי של  $\hat{\mu}_{ij}^{\text{MLE}}$ .

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

$$= \prod_{i=1}^m f(x_1=x_1, \dots, x_d=x_d | y_i=y_i) \cdot \text{Mult}(y_i, \pi) = \prod_{\substack{i=1 \\ j \in \mathbb{N}^d \setminus \mathbb{B}^d}}^m \left( \frac{d}{\prod_{j=1}^d} f(x_j=x_j | y_i=y_i) \right) \text{Mult}(y_i, \pi) =$$

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

: פונקציית log-likelihood - היא פונקציה מוגדרת כפונקציית האנרגיה המינימלית של מודל.

$$\log(L(\theta|x,y)) = \log\left(\prod_{i=1}^m \prod_{j=1}^d N(x_{ij} | \mu_{y,j}, \sigma_{y,j}^2) \cdot \text{Mult}(y_i, \pi)\right)$$

$$= \sum_{i=1}^m \left( \sum_{j=1}^d \log(N(x_{ij} | \mu_{y_{ij}}, \sigma_{y_{ij}}^2)) + \log(\text{Mult}(y_i, \pi)) \right) = \sum_{i=1}^m h_i \log(\pi_i) + \sum_{m:y_m=i} \sum_{j=1}^d \log(N(x_{mj} | \mu_{y_{ij}}, \sigma_{y_{ij}}^2))$$

.  $\forall i \in [k] \quad h_i := \sum_{j=1}^m \mathbf{1}_{y_m=j}$

$$\hat{\pi}_i^{\text{MLE}} = \frac{n_i}{m}$$

הנחיות רלוונטיות נקבעו כמפורט לעיל. מטרת הבדיקה היא למדוד את השוני בין הנקודות המבוקשות ונקודות המבוקשיות.

$N(x_{2;j} | \mu_{y_{2;j}}, \sigma_{y_{2;j}})$  表示针对特征  $x_{2;j}$  的 log-likelihood 分布，其中  $j \in [d]$  表示类别， $i \in [k]$  表示样本。

$$\widehat{\sigma}_{ij}^{MLF} = \frac{1}{h_i} \sum_{m:y_m = i} (x_{mj} - \mu_{ij})^2$$

$$\hat{\mu}_{ij}^{MLE} = \frac{1}{n_i} \sum_{m:y_m=i} x_{mj}$$

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

Poisson likelihoods:

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

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 Poisson Naive Bayes classifier solving (5) under assumptions (7).

(4)

הנחתה מוגדרת ב- $y$  כ  $K$  ו- $\hat{\pi}_i, \hat{\lambda}_i$  הם ה-MLE (א)

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

$$\prod_{i=1}^m \text{Mult}(y_i | \pi) \cdot \text{Pois}(x_i | \lambda_{y_i}) = \prod_{i=1}^m \pi_{y_i} \cdot \frac{e^{-\lambda_{y_i}} \lambda_{y_i}^{x_i}}{(x_i)!}$$

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

$$x | y = y_i \sim \text{Poi}(\lambda_{y_i}) \text{ נסובב כפונקציית כפיפה } x \in \mathbb{N} - 0, 1, 2, \dots \text{ נסובב כפונקציית}$$

$$x \sim \text{Poi}(\lambda) \Rightarrow P(x=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

ה-MLN log-likelihood - ה-MLN מוגדרת כ- $\log L(\theta | x, y)$  ו- $L(\theta | x, y)$  מוגדרת כ- $\prod_{i=1}^m \pi_{y_i} \cdot \frac{e^{-\lambda_{y_i}} \lambda_{y_i}^{x_i}}{(x_i)!}$

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

כעת נוכחים  $y_1, \dots, y_m$  ו- $\lambda_1, \dots, \lambda_m$

לעתה  $\forall j \in K \quad h_j := \sum_{i:y_i=j} 1$

$$\sum_{j=1}^K h_j \log(\pi_j) - h_j \lambda_j + \log(\lambda_j) \leq x_i - \sum_{i:y_i=j} \log(x_i)!$$

$$L = \log(L(\theta | x, y)) = \log(\prod_{i=1}^m \pi_{y_i} \cdot \frac{e^{-\lambda_{y_i}} \lambda_{y_i}^{x_i}}{(x_i)!}) = \sum_{i=1}^m \log(\pi_{y_i}) - \lambda_{y_i} + x_i \log(\lambda_{y_i}) - \log((x_i)!)$$

$$\frac{\partial L}{\partial \pi_j} = \frac{\partial}{\partial \pi_j} \left( h_j \log(\pi_j) - \lambda_j \pi_j \right) = \frac{h_j}{\pi_j} - \lambda_j = 0 \Rightarrow \pi_j = \frac{h_j}{\lambda_j} \Rightarrow \frac{m}{\lambda} = \frac{1}{\lambda} \sum_{j=1}^K h_j = \sum_{j=1}^K \pi_j = 1 \Rightarrow \lambda = m$$

. ב- $\lambda = m$  מוגדרת כ- $\hat{\pi}_j = \frac{h_j}{m}$

$$\frac{\partial L}{\partial \lambda_j} = \frac{\partial}{\partial \lambda_j} \left( -h_j \lambda_j + \sum_{i:y_i=j} x_i \log(\lambda_j) \right) = -h_j + \sum_{i:y_i=j} \frac{x_i}{\lambda_j} = 0 \Rightarrow h_j = \sum_{i:y_i=j} \frac{x_i}{\lambda_j} \Rightarrow \lambda_j = \frac{1}{h_j} \sum_{i:y_i=j} x_i$$

$$\hat{\lambda}_j = \frac{1}{h_j} \sum_{i:y_i=j} x_i$$

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

$$\forall i \in [k], \forall j \in [d] \quad \hat{\pi}_i^{MLE}, \hat{\lambda}_{ij}^{MLE} \quad \text{-} \text{איך יכול נזיר (b)}$$

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

$$= \prod_{i=1}^m \pi_{y_i} \cdot \prod_{j=1}^d f(x_j = x_{ij} | y_i = y_i) = \prod_{i=1}^m \pi_{y_i} \cdot \prod_{j=1}^d \frac{e^{-\lambda_{y_i j}} \lambda_{y_i j}^{x_{ij}}}{(x_{ij})!}$$

$y_i \sim \text{Mult}(\pi)$   
 $x_{ij} \sim \text{Pois}(\lambda_{y_i j})$

: מינימיזציה של log-likelihood  $\rightarrow$  מינימיזציה של  $L(\theta | \mathbf{x}, \mathbf{y})$   $\rightarrow$  מינימיזציה של  $\log(L(\theta | \mathbf{x}, \mathbf{y}))$

$$\log(L(\theta | \mathbf{x}, \mathbf{y})) = \log\left(\prod_{i=1}^m \pi_{y_i} \cdot \prod_{j=1}^d \frac{e^{-\lambda_{y_i j}} \lambda_{y_i j}^{x_{ij}}}{(x_{ij})!}\right) = \sum_{i=1}^m \log(\pi_{y_i}) + \log\left(\prod_{j=1}^d \frac{e^{-\lambda_{y_i j}} \lambda_{y_i j}^{x_{ij}}}{(x_{ij})!}\right) =$$

$$= \sum_{i=1}^m \log(\pi_{y_i}) + \sum_{j=1}^d (-\lambda_{y_i j} + x_{ij} \log(\lambda_{y_i j}) - \log(x_{ij}!)) = \sum_{i=1}^m \log(\pi_i) + \sum_{j=1}^d -n_i \lambda_{ij} + \log(\lambda_{ij}) \sum_{t:y_t=j} x_{ij} - \sum_{i=1}^m \log(x_{ij}!)$$

: מינימיזציה של  $\sum_{i=1}^m \log(\pi_i) + \sum_{j=1}^d (-n_i \lambda_{ij} + \log(\lambda_{ij}) \sum_{t:y_t=j} x_{ij} - \sum_{i=1}^m \log(x_{ij}!))$

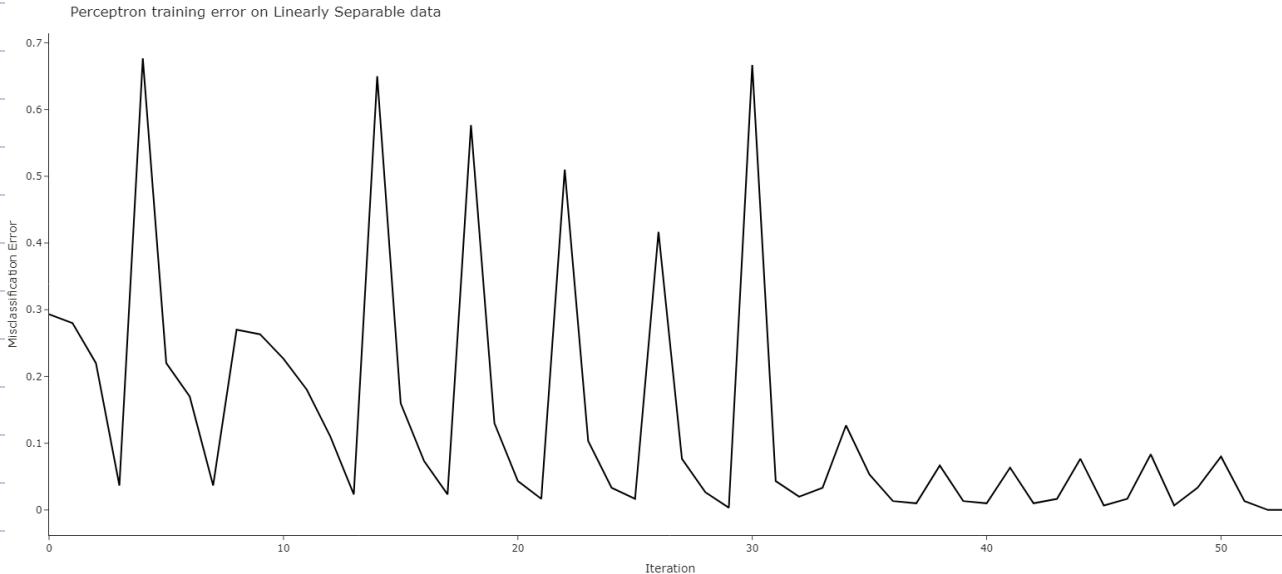
, מינימיזציה של  $\sum_{i=1}^m \log(\pi_i)$  מושג בהמינימום של הפונקציית האנרגיה, כלומר  $\hat{\pi}_i^{MLE} = \frac{n_i}{m}$ .

$$\hat{\lambda}_{ij}^{MLE} = \frac{1}{n_i} \sum_{m:y_m=j} x_{mj}$$

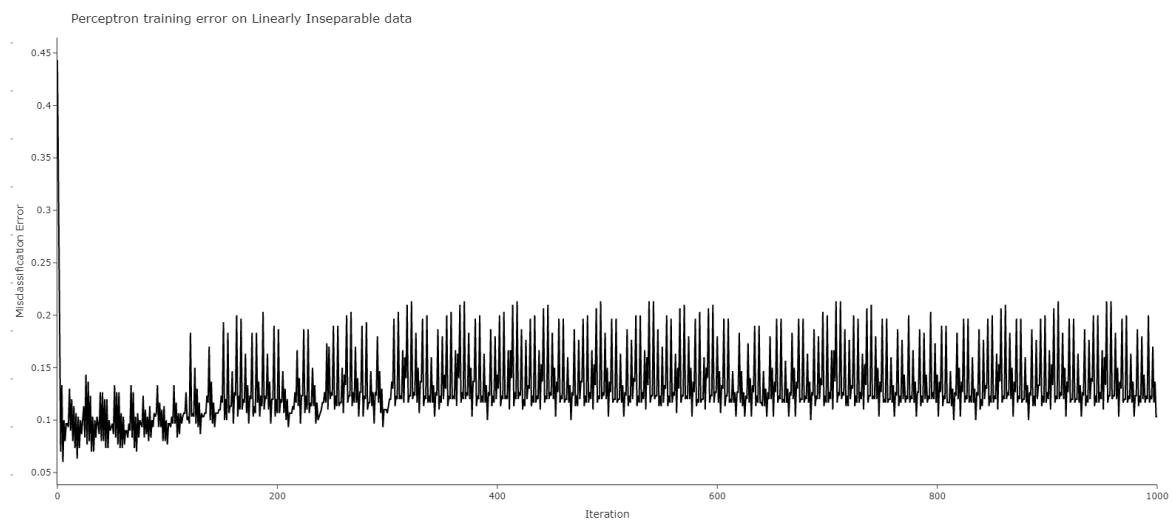
. ( $j$  הינה מינימום של  $\sum_{i=1}^m \log(\lambda_{ij}) \sum_{t:y_t=j} x_{ij} - \sum_{i=1}^m \log(x_{ij}!)$ )

### 3.1 Perceptron Classifier

1. Fitting and plotting over the datasets/linearly\_separable.npy dataset, what can we learn from the plot?
  2. Next run the Perceptron algorithm over the datasets/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?



ב-1972 (בכינור נ-60 קיבוצי צ'רנוביץ), פירוטו תיכון כ' ב-1973 ו-1974 נקבעו נספח ז' ונספח י' (ב-1974).

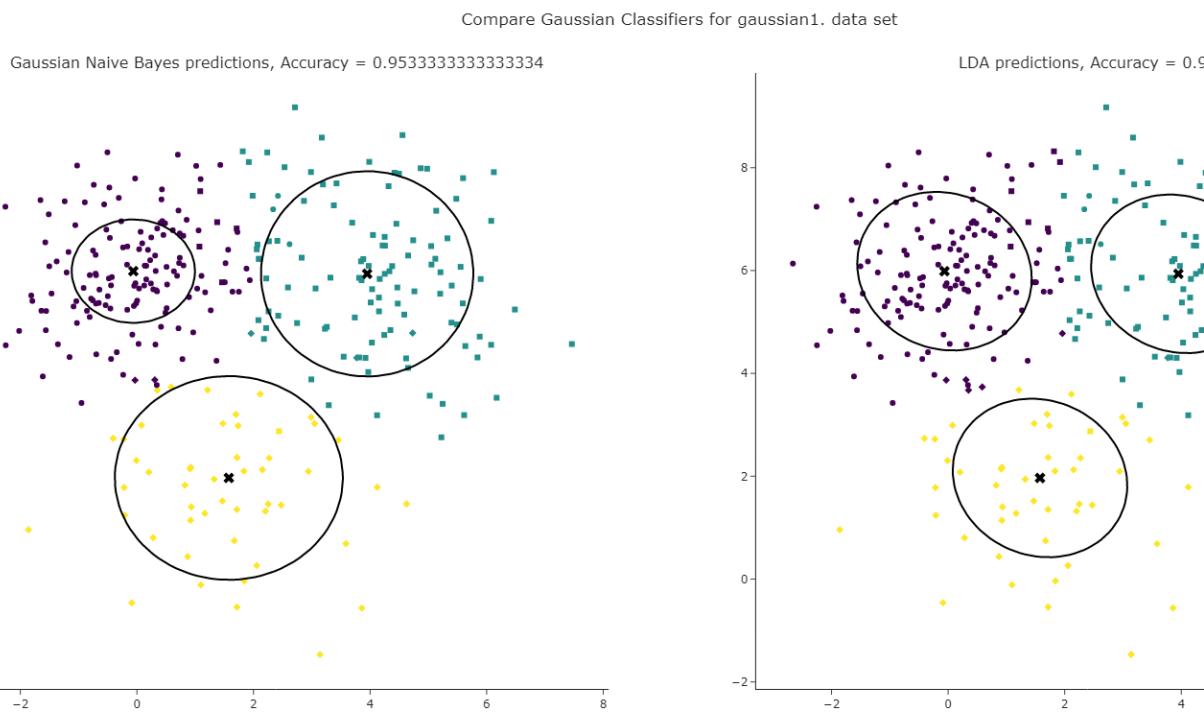


$$\begin{aligned} & \text{minimize} && 0 \\ & \text{subject to} && y_i \cdot \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 0 \quad i = 1, \dots, m \end{aligned}$$

soft SVM (soft margin SVM) objective function is:

## 3.2 Bayes Classifiers

6

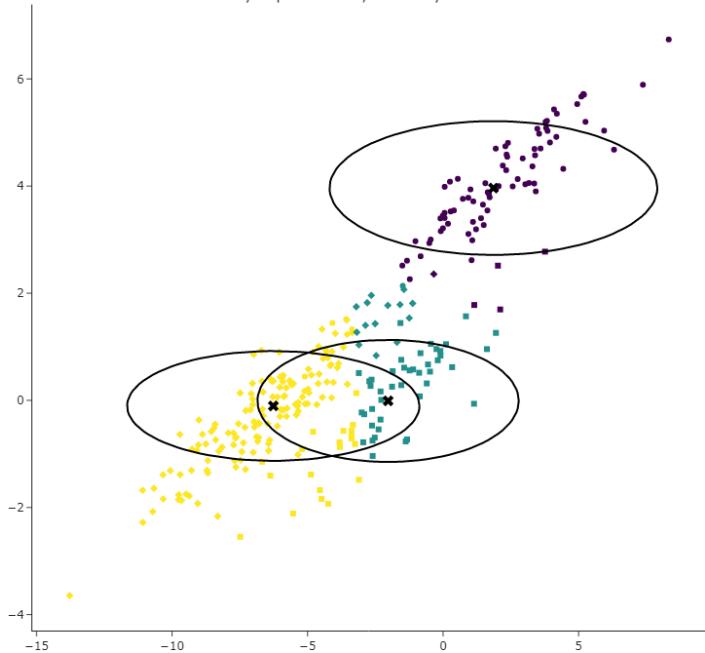


- Explain what can be learned from the plots above regarding the distribution used to sample the data?

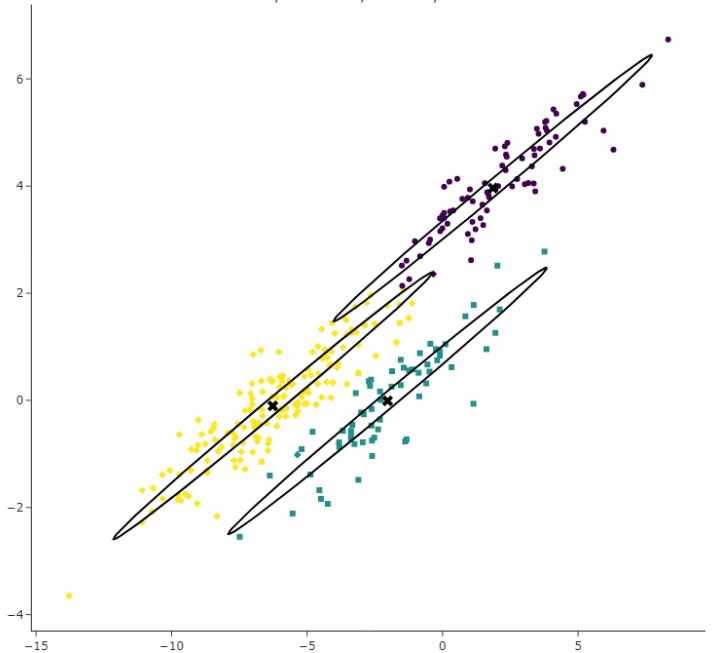
העדרת תוצאות מדויקות מ-0.95 ל-0.94 מראה שמשתמש בcovariance matrix (cov(x,y)=cov(y,x)=0) מוביל ל결לים לא-טראנספורמציוניים (transformations). מטרת ה-LDA היא לפרק דimensionality ו- $\text{cov}(x,y)=\text{cov}(y,x)=0$  מונע מכך. מטרת ה-GNB היא לפרק דimensionality ו- $\text{cov}(x,y)=\text{cov}(y,x)=0$  מונע מכך. מטרת ה-LDA היא לפרק Dimensionality ו- $\text{cov}(x,y)=\text{cov}(y,x)=0$  מונע מכך. מטרת ה-GNB היא לפרק Dimensionality ו- $\text{cov}(x,y)=\text{cov}(y,x)=0$  מונע מכך.

## Compare Gaussian Classifiers for gaussian2. data set

Gaussian Naive Bayes predictions, Accuracy = 0.8533333333333334



LDA predictions, Accuracy = 0.97



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?

- מונע מ- $\text{GNB}$  לטעות ב- $\text{LDA}$  על ידי שימוש ב- $\text{Cov}$  (בנוסף ל- $\text{mean}$  ו- $\text{std}$ )