

# 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 result in reduction of points.
- Plots included as separate files will be considered as not provided.

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

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

### 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:  $\boldsymbol{\pi} \in [0, 1]^K, \sum \pi_j = 1$ .

- 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.
  - 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:  $\boldsymbol{\pi} \in [0, 1]^K, \sum \pi_j = 1$ .

- 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).
- 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?

Based on Lecture 5 and Recitation 5

20701750 'J' 7/1/16

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

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

Hint: Observe that  $\|w\|^2 = w^T I w$

$$V = \begin{bmatrix} w \\ b \end{bmatrix} \in \mathbb{R}^{n+1} \quad a = 0 \in \mathbb{M}_{n+1 \times n+1} \text{ maybe } 0$$

$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned} \operatorname{argmin} \|w\|^2 &= \operatorname{argmin} (w^T I w) = \\ \operatorname{argmin} (V^T \begin{bmatrix} 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} V) &= \operatorname{argmin} (\frac{1}{2} V^T Q V + 0) = \\ \operatorname{argmin} (\frac{1}{2} V^T Q V + a^T V). \end{aligned}$$

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

$$\forall i \quad y_i \quad w^T x_i + b \geq 1 \Leftrightarrow$$

$$\forall i \quad -y_i (w^T x_i - b) \leq -1 = \textcircled{*}$$

$$[A]_i = \begin{bmatrix} y_i x_i \\ y_i \end{bmatrix} \quad \forall i, \quad b = \begin{bmatrix} -1 \\ \vdots \\ -1 \end{bmatrix} \text{ מוגדר}$$

$$\textcircled{*} = Av \leq d$$

2. Consider the Soft-SVM optimization problem:

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

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:

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

$$\forall i \quad y_i \langle w, x_i \rangle \geq 1 - \xi_i \wedge \xi_i \geq 0$$

כל המסלוק צ'ה ג'יה נכון:  $y_i \langle w, x_i \rangle \geq 1$   
 אם  $\xi \geq 0$  שהמרחק יפה מונחי מרחק  
 אם  $\xi = 0$  אז המרחק נכון  
 אם  $\xi < 1$  המרחק נכון והמרחק יפה מונחי מרחק  
 אם  $\xi < 1$  המרחק נכון והמרחק יפה מונחי מרחק

$$\xi \geq 1 - y_i \langle w, x_i \rangle$$

יתקבל
עבור
כל
המניסל
החסר
מכאן

$$\xi = 1 - y_i \langle w, x_i \rangle$$

$$\forall i \quad \xi = \text{hinge}(y_i \langle w, x_i \rangle)$$

כאן

החסר

מכאן



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.

נחמד = ה"פ"ק' קוזח סמלדנאן חבוי

$$L(\theta | x, y) \stackrel{\text{joint 'lik'}}{=} \prod_{i=1}^n \mathcal{S}_{x,y}(x_i, y_i | \theta) =$$

$$= \prod_{i=1}^n \mathcal{S}_{X|Y=y_i}(x_i) \cdot \mathcal{S}_{Y|\Theta}(y_i)$$

$$= \prod_{i=1}^n \mathcal{N}(x_i | \mu_{y_i}, \sigma_{y_i}^2) \text{mult}(y_i, \pi)$$

מהנהנה  
6

לוג 2017

$$\log(L(\Theta | \mathcal{X}, \mathcal{Y})) = \log(\prod_{i=1}^n \underbrace{\mathcal{N}(x_i | \mu_{y_i}, \sigma_{y_i}^2)}_{\text{mult}(y_i, \pi)})$$

$$= \sum_{i=1}^n \log(\mathcal{N}(x_i | \mu_{y_i}, \sigma_{y_i}^2) \text{mult}(y_i, \pi))$$



DN12 10112

$$\frac{\partial l}{\partial \sigma_k^2} = -\frac{1}{2\sigma_k^2} \sum (x_i - \mu_k)^2 \Rightarrow \left( \frac{1}{\sigma_k^2} \right)^{MLE} =$$

$$\frac{1}{n_k} \sum_{x_i | y_i = k} (x_i - \mu_k^{MLE})^2$$

(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).

$$L(\Theta | X, y) = \prod_{i=1}^m \mathcal{F}_{X, y}(\mathbf{x}_i, y_i | \Theta) =$$

$$\prod_{i=1}^m \mathcal{F}_{X | y=y_i}(\mathbf{x}_i) \mathcal{F}_{Y | \Theta}(y_i) =$$

$$= \left( \prod_{j=1}^K \pi_j \right) \prod_{i=1}^m \mathcal{N}(\mathbf{x}_{i,j} | \mu_{y_i,j}, \sigma_{y_i,j}^2) \text{mult}(y_i | \pi)$$

$$\log(L(\Theta | X, y)) = \sum_i [\log(\text{mult}(y_i | \pi)) +$$

$$\sum_j \log(\mathcal{N}(\mathbf{x}_{i,j} | \mu_{y_i,j}, \sigma_{y_i,j}^2))] =$$

$$\sum_k \pi_k \log(\pi_k) + \sum_k \sum_j \mathbb{1}_{[y_i=k]} \log(\mathcal{N}(\mathbf{x}_{i,j} | \mu_{k,j}, \sigma_{k,j}^2))$$

∴  $\mu_k = \frac{1}{n} \sum_{i: y_i = k} x_{i,j}$

$$\mu_k = \frac{n_k}{n} \quad \hat{\mu}_{k,j} = \frac{1}{n_k} \sum_{i: y_i = k} x_{i,j}$$

$$(\sigma^2)_{k,j} = \frac{1}{n_k} \sum_{i: y_i = k} (x_{i,j} - \hat{\mu}_{k,j})^2$$

∴

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

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

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).  
 (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).

$$\begin{aligned} L(\Theta | X, y) &= \prod_{i=1}^n \mathcal{F}_{x,y}(x_i, y_i | \Theta) = \\ &= \prod_{i=1}^n \mathcal{F}_{x|y=y_i}(x_i) \cdot \mathcal{F}_{y|\Theta}(y_i) \\ &= \prod_{i=1}^n \text{Poi}(x_i | \lambda_{y_i}) \cdot \text{Mult}(y_i | \pi) \end{aligned}$$

∴

הנני מציג את הפתרון:

$$\begin{aligned}
 l(\theta | X, y) &= \sum_i \log(\text{Poi}(x_i | \lambda_{y_i})) + \log(\text{mult}(y_i | \pi)) \\
 &= \sum_i \log\left(\frac{\lambda_{y_i}^{x_i} e^{-\lambda_{y_i}}}{x_i!}\right) + \log(\pi_{y_i}) \\
 &= \sum_i x_i \log(\lambda_{y_i}) - \lambda_{y_i} - \log(x_i!) + \log(\pi_{y_i}) \\
 &= \sum_k \left[ \log(\lambda_k) \sum_{i: y_i=k} x_i - n_k \lambda_k + n_k \log(\pi_k) \right] - \sum_i \log(x_i!)
 \end{aligned}$$

הנני מציג את הפתרון:

$$\frac{\partial l}{\partial \lambda_k} = \frac{1}{\lambda_k} \sum_{i: y_i=k} x_i - n_k \stackrel{!}{=} 0 \Rightarrow \lambda_k^{\text{MLE}} = \frac{1}{n_k} \sum_{i: y_i=k} x_i$$

הנני מציג את הפתרון:

$$\lambda_k^{\text{MLE}} = \frac{n_k}{n}$$

- (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).

$$L(\theta | X, y) = \prod_{i=1}^m \mathcal{S}_{X, y}(x_i, y_i | \theta)$$

$$= \prod_{i=1}^m \mathcal{S}_{X|Y=y_i}(x_i) \mathcal{S}_{Y|\theta}(y_i)$$

$$= \prod_{i=1}^m \left( \prod_{j=1}^d \text{Poi}(x_{i,j} | \lambda_{y_i,j}) \right) \text{mult}(y_i | \pi)$$

$$\text{log} \quad \text{or} \quad \text{mult} \quad \text{or} \quad \text{log}$$

$$l(\theta | X, y) = \sum_i \sum_j [\log(\text{Poi}(x_{i,j} | \lambda_{y_i,j}))] +$$

$$\sum_i \log(\text{mult}(y_i | \pi)) =$$

$$\sum_k \sum_j [\log(\lambda_{k,j}) \sum_{i: y_i=k} x_{i,j} - n_k \lambda_{k,j} + n_k \cdot \log(\pi_k)] -$$

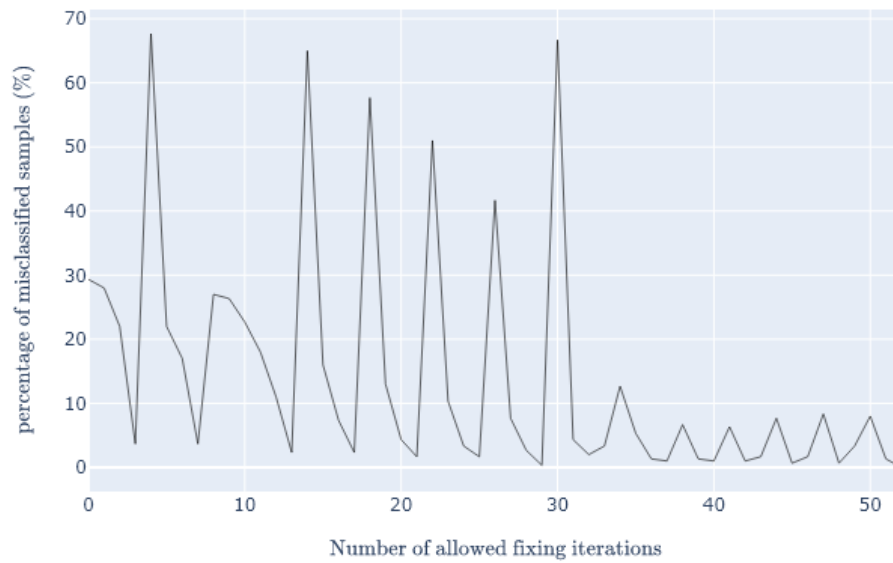
$$\sum_{i,j} \log(x_{i,j}!) \Rightarrow$$

$$\frac{1}{\lambda_{k,j}} = \frac{1}{n_k} \sum_{i: y_i=k} x_{i,j}, \quad \frac{1}{\pi_k} = \frac{n_k}{n}$$

1570 p. 17

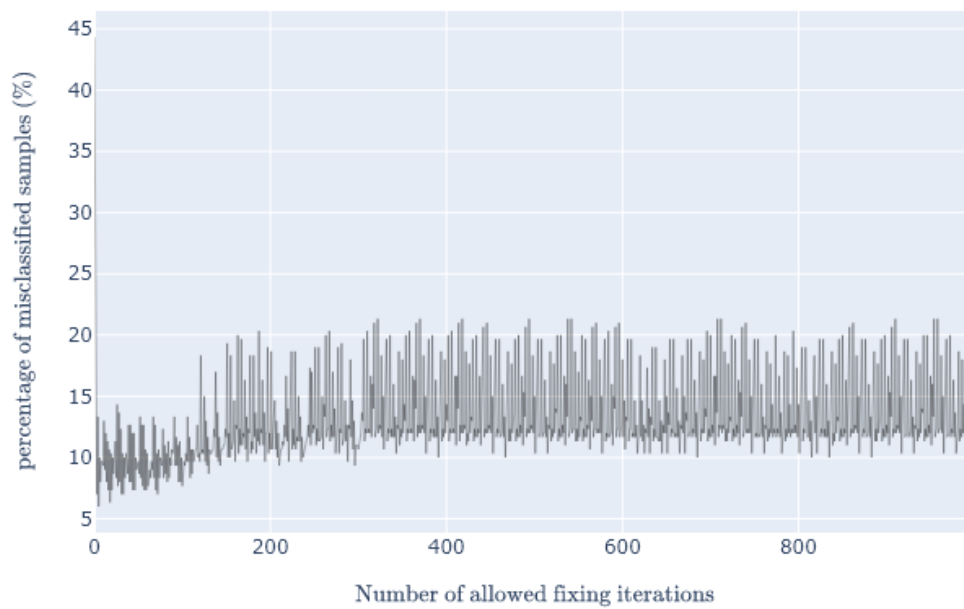
2.1

Perceptron model error percentage over - linearly separable dataset



2.2

Perceptron model error percentage over - linearly inseparable dataset

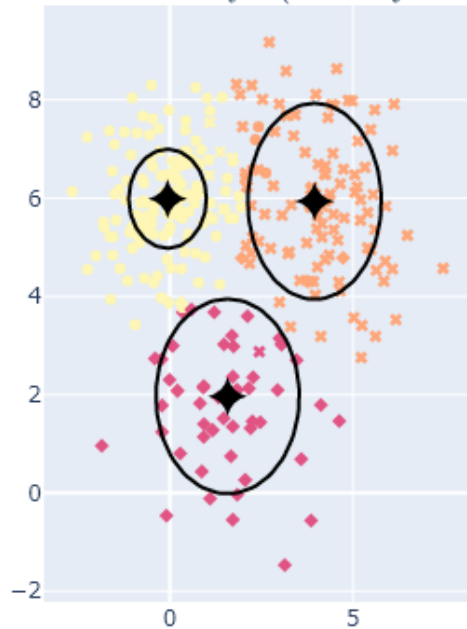




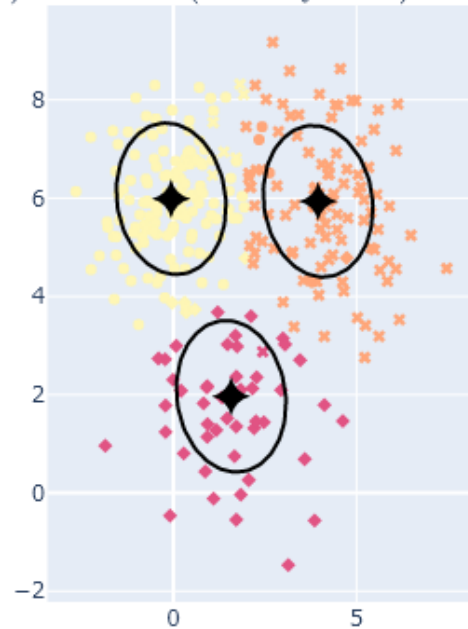
מביט בנקודות שבהן אנו רואים  
 את ההבדלים בין המצבים השונים  
 ונראה שהם לא תמיד כאלו שאנו חושבים  
 שהם. למשל, ההבדלים שהם חסרי  
 כלל, אלא הם תלויים בהם. זהו  
 דבר שיש לו.

Naive Bayse vs. LDA estimators comparison over gaussian1

Gaussian Naive Bayes (accuracy=95.33%)



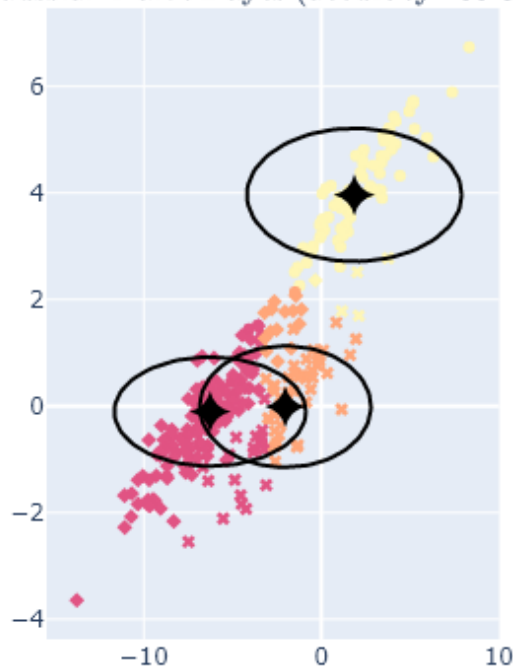
LDA (accuracy=94%)



לערוך טאג, און דאס איז דער  
 וואס ער האט געטון מיט דעם  
 אונזערע דאטעס, און דאס  
 שטענדיג אונזערע דאטעס, LDA  
 דערנאך האט ער געטון  
 און דאס איז דער  
 אונזערע דאטעס, און דאס

Naive Bayse vs. LDA estimators comparison over gaussian2

Gaussian Naive Bayes (accuracy=85.33%)



LDA (accuracy=97%)

