

CX4240 Computing for Data Analysis - Homework 3

Name:

GTID:

Deadline: 11:59 pm EDT, March 30

- Submit your answers as one single PDF file on Gradescope.
- You will be allowed 2 total late days (48 hours) without penalty for the entire semester. Once those days are used, you will be penalized according to the following policy:
 - Homework is worth full credit before the due time.
 - It is worth 75% credit for the next 24 hours.
 - It is worth 50% credit for the second 24 hours.
 - It is worth zero credit after that.
- You are required to use Latex, or word processing software, to generate your solutions to the written questions. Handwritten solutions WILL NOT BE ACCEPTED.

1 Naive Bayes: Linear Model in Disguise [45 pts]

In this problem set, we will demonstrate that a Bernoulli Naive Bayes classifier is essentially a linear classifier. Specifically, we will show that it shares the same functional form as Logistic Regression.

Consider a binary classification problem with classes $C \in \{0, 1\}$ and a vector of n binary features $\mathbf{x} = (x_1, x_2, \dots, x_n)$ where $x_i \in \{0, 1\}$. A classifier intends to predict the class C of a given data point \mathbf{x} , according to the decision rule:

$$\text{The class of } \mathbf{x} \text{ is } \begin{cases} C = 1 & \text{if } \frac{P(C = 1|\mathbf{x})}{P(C = 0|\mathbf{x})} > 1 \\ C = 0 & \text{otherwise.} \end{cases} \quad (1)$$

We assume the class-conditional probabilities for each feature i is given:

$$\begin{aligned} P(x_i = 0|C = 0) &= \alpha_i \\ P(x_i = 0|C = 1) &= \beta_i \end{aligned} \quad (2)$$

where $\alpha_i, \beta_i \in [0, 1]$. Besides, the class prior probabilities are also given:

$$\begin{aligned} P(C = 0) &= c \\ P(C = 1) &= 1 - c \end{aligned} \quad (3)$$

Answer the following questions:

1. **[15 points]** Demonstrate that the term $\ln \frac{P(x_i|C=1)}{P(x_i|C=0)}$ can be written as a linear function of the component x_i . Write out the expression of the weight w_i and the bias b_i in terms of α_i, β_i .
2. **[10 points]** Apply Bayes' theorem and Naive Bayes assumption, demonstrate that $\ln \frac{P(C=1|\mathbf{x})}{P(C=0|\mathbf{x})}$ can be expressed as $\mathbf{w}^\top \mathbf{x} + w_0$. Write out the expression of \mathbf{w} and w_0 in terms of $\{\alpha_i\}_{i=1}^n, \{\beta_i\}_{i=1}^n$ and c . Explicitly mark where you applied the Naive Bayes assumption.
3. **[10 points]** Recall that the Logistic Regression predicts the probability of a class using the sigmoid function:

$$P_{\text{LR}}(C = 1|\mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{w}^\top \mathbf{x} + w_0))} \quad (4)$$

Using your result from the last question, show that the Naive Bayes posterior probability $P(C = 1|\mathbf{x})$ can be written in this exact same functional form.

4. **[10 pts]** Naive Bayes (Generative) and Logistic Regression (Discriminative) both result in linear decision boundaries but differ significantly in execution. Please compare these models (e.g., in terms of their learning algorithms, modeling assumptions, and etc.), and briefly discuss the benefits and weaknesses of Naive Bayes.

2 K-Means Clustering [55 pts]

In this exercise, you will implement the K-means clustering algorithm from scratch to partition data into clusters. We have prepared a Jupyter notebook to guide you through the complete K-means workflow.

Please access the notebook using [this link](#). You can also download the notebook by using File (top-left corner) → Download → Download .ipynb.

For this problem, complete all 5 sections in the notebook, execute all cells to generate the required outputs, and export the notebook into a single PDF file. For example, if you are using Google Colab, you can use File (top-left corner) → Print to generate a PDF. For submission, a separate homework entry for this problem will be created on GradeScope. Please upload the generated PDF to that entry.

Score Distribution:

1. **[0 pts] Step 1:** Data loading and visualization.
2. **[15 pts] Step 2:** Computing Euclidean distances and assigning clusters (the Assignment step). Please refer to slides of Lecture 16, Pages 13 and 20.
3. **[10 pts] Step 3:** Updating cluster centers (the Center Update step). Please refer to slides of Lecture 16, Page 13.
4. **[20 pts] Step 4:** Putting it all together — the full K-means algorithm with convergence detection. Please refer to slides of Lecture 16, Pages 13 and 27-28.
5. **[10 pts] Step 5:** Analysis — the elbow method for choosing K. Please refer to slides of Lecture 16, Page 29.

References

Solution:

Please mention any AI tools, people, post or blog etc. you used.