

CSCI 5521: Introduction to Machine Learning (Spring 2023)¹

Homework 1

Questions

- (10 points)** What is the VC dimension, $d_{\mathcal{I}}$, of intervals in \mathbb{R} ? The target function is specified by an interval $[a, b]$, and labels any example positive iff it lies inside that interval. Show how to shatter $d_{\mathcal{I}}$ points but not $d_{\mathcal{I}} + 1$ points with the function.
- (30 points)** Find the Maximum Likelihood Estimation (MLE) for the following pdf. In each case, consider a random sample of size n . Show your calculation:
 - $f(x|\theta) = \frac{1}{\theta-1} e^{-\frac{x}{\theta-1}}$, $x > 0, \theta > 1$
 - $f(x|\theta) = (\theta - 1)x^{\theta-2}$, $0 \leq x \leq 1, 1 < \theta < \infty$
 - $f(x|\theta) = \frac{1}{\theta}$, $0 \leq x \leq \theta$ (Hint: likelihood function is monotonically decreasing.)
- (30 points)** Let $P(x|C)$ denote a Bernoulli density function for a class $C \in \{C_1, C_2\}$ and $P(C)$ denote the prior,
 - Given the priors $P(C_1)$ and $P(C_2)$, and the Bernoulli densities specified by $p_1 \equiv p(x = 0|C_1)$ and $p_2 \equiv p(x = 0|C_2)$, derive the classification rules for classifying a sample x into C_1 and C_2 based on the posteriors $P(C_1|x)$ and $P(C_2|x)$. (Hint: give rules for classifying $x = 0$ and $x = 1$.)
 - Consider D -dimensional independent Bernoulli densities specified by $p_{ij} \equiv p(x_j = 0|C_i)$ for $i = 1, 2$ and $j = 1, 2, \dots, D$. Derive the classification rules for classifying a sample x into C_1 and C_2 . It is sufficient to give your rule as a function of x .
 - Follow the definition in 3(b) and assume $D = 2$, $p_{11} = 0.6$, $p_{12} = 0.1$, $p_{21} = 0.6$, and $p_{22} = 0.9$. For three different priors ($P(C_1) = 0.2, 0.6, 0.8$ and $P(C_2) = 1 - P(C_1)$), calculate the posterior probabilities $P(C_1|x)$ and $P(C_2|x)$. (Hint: Calculate the probabilities for all possible samples $(x_1, x_2) \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$).
- (30 points)** Using the provided training, validation, and test datasets, write a Python script to calculate the maximum likelihood estimation on the training set. Consider a prior function defined with respect to sigma as

$$P(C_1|\sigma) = 1 - e^{-\sigma}, \sigma > 0 \tag{1}$$

¹Instructor: Rui Kuang (kuang@umn.edu)

and $P(C_2) = 1 - P(C_1)$. Using the learned Bernoulli distributions and the given prior function, classify the samples in the validation set using your classification rules for $\sigma = 0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 2, 3, 4, 5, 6$. Finally, choose the best prior (the one that gives the lowest error rate on the validation set) and use it to classify the samples in the test set. Print to the Python console (either in terminal or PyCharm) a table of error rate of each prior on the validation set and the error rate using the best prior on the test set. (Hint: if some Bernoulli probabilities are 0, you can replace them with a small probability such as 10^{-10} to avoid the numerical problem.)

Instructions

- **Programming Questions:** All programming questions must be written in Python, no other programming languages will be accepted. Only Numpy can be used in this assignment. The code must be able to be executed from the terminal command prompt on the cselabs machines. Each function must take the inputs in the order specified and display the textual output via the Python console (either in terminal or PyCharm). For each part, you can submit additional files/functions (as needed) which will be used by the main functions specified below. Put comments in your code so that one can follow the key parts and steps. **Please follow the rules strictly. If we cannot run your code, you will receive no credit.**
 - **Question 4:**
 - * Training function in Bayes_learning.py: *Bayes_Learning*(training_data , validation_data). The function returns the outputs (p1: learned Bernoulli parameters of the first class, p2: learned Bernoulli parameters of the second class, pc1: best prior of the first class, pc2: best prior of the second class). It must also print to the terminal (sprintf) a table of error rates of all priors.
 - * Test function in Bayes_testing.py: *Bayes_Testing*(test_data, p1: the learned Bernoulli parameter of the first class, p2: the learned Bernoulli parameter of the second class, pc1: the learned prior of the first class, pc2: the learned prior of the second class). The function must print to the Python console (either in terminal or PyCharm) the error rate on the test dataset.
 - * main script main_script.py: the script loads the data and call the training and test functions to generate the results.
 - **Error rate:** Error rate is the percentage of wrongly classified data points divided by the total number of classified data points.
- **Report:** Solutions to Questions 1, 2 and 3 must be included in a report. The table of error rates on the validation set and the error rate on the test set for Question 4 must also be included in the report.
- **Things to submit:**

1. hw1_sol.pdf: A document which contains the report with solutions to all questions. Scanned answer sheets need to be clean and 100% legible.
 2. Bayes_Learning.py, Bayes_Testing.py and main_script.py: Code for Question 4.
 3. Any other files, except the data, which are necessary for your code.
- **Submit:** All material must be submitted electronically via Canvas.