
CSE 5523: Machine Learning - Homework #1

Due: 11:59pm 02/13/2024

Homework Policy: Please submit your solutions in a single PDF file named `HW_1_name.number.pdf` (e.g., `HW_1_zhang.12807.pdf`) to Carmen. You may write your solutions on paper and scan it, or directly type your solutions and save them as a PDF file. *Submission in any other format will not be graded.* Working in groups is fine, but each member must submit their own writeup. Please write the members of your group on your solutions. **For coding problems, please append your code to your submission and report your results (values, plots, etc.) in your written solution.** You will lose points if you only include them in your code submissions.

1) Nearest neighbor & distance metrics [10 points]. Given four labeled training data instances, $([5, 2]^\top, \text{"square"})$, $([4, 0.5]^\top, \text{"cross"})$, $([6, 4]^\top, \text{"circle"})$, $([4, 5]^\top, \text{"diamond"})$, where each of them is a point in the 2D space (i.e., \mathbb{R}^2) and labeled by a “shape”, please determine the label (i.e., “shape”) of the test data instance $[3, 2]^\top$ based on the nearest neighbor rule. See [Figure 1](#) for illustration.

(a) [2 points] What is the label if L_1 distance is used?

$$\|[3, 2]^\top - [5, 2]^\top\|_1 = |3-5| + |2-2| = 2$$

Square

$$\|[3, 2]^\top - [4, 0.5]^\top\|_1 = |3-4| + |2-0.5| = 2.5$$

$$\|[3, 2]^\top - [6, 4]^\top\|_1 = |3-6| + |2-4| = 5$$

$$\|[3, 2]^\top - [4, 5]^\top\|_1 = |3-4| + |2-5| = 4$$

(b) [2 points] What is the label if L_2 distance is used?

$$\|[3, 2]^\top - [5, 2]^\top\|_2 = \sqrt{(3-5)^2 + (2-2)^2} = 2$$

CROSS

$$\|[3, 2]^\top - [4, 0.5]^\top\|_2 = \sqrt{(3-4)^2 + (2-0.5)^2} = 1.803$$

$$\|[3, 2]^\top - [6, 4]^\top\|_2 = \sqrt{(3-6)^2 + (2-4)^2} = 3.606$$

$$\|[3, 2]^\top - [4, 5]^\top\|_2 = \sqrt{(3-4)^2 + (2-5)^2} = 3.162$$

- (c) [2 points] What is the label if L_∞ distance is used? L_∞ norm of one vector is defined as the largest magnitude among each element of a vector.

$$\| [3, 2]^T - [5, 2]^T \|_\infty = (|3-5|^{\infty} + |2-2|^{\infty})^{\frac{1}{\infty}} = (|2|^{\infty})^{\frac{1}{\infty}} = 2 \quad \text{cross}$$

$$\| [3, 2]^T - [4, 0.5]^T \|_\infty = (|3-4|^{\infty} + |2-0.5|^{\infty})^{\frac{1}{\infty}} = (|1.5|^{\infty})^{\frac{1}{\infty}} = 1.5$$

$$\| [3, 2]^T - [6, 4]^T \|_\infty = (|3-6|^{\infty} + |2-4|^{\infty})^{\frac{1}{\infty}} = (|3|^{\infty})^{\frac{1}{\infty}} = 3$$

$$\| [3, 2]^T - [4, 5]^T \|_\infty = (|3-4|^{\infty} + |2-5|^{\infty})^{\frac{1}{\infty}} = (|3|^{\infty})^{\frac{1}{\infty}} = 3$$

- (d) [2 points] What is the label if cosine distance is used? Given two vectors \mathbf{x} and \mathbf{x}' , their cosine distance is defined as $1 - \frac{\mathbf{x}^\top \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}$.

$$\| [3, 2]^T - [5, 2]^T \|_{\cos} = 1 - \frac{[3 \ 2] [5 \ 2]^T}{\sqrt{3^2+2^2} \sqrt{5^2+2^2}} = 0.0215$$

$$\| [3, 2]^T - [4, 0.5]^T \|_{\cos} = 1 - \frac{[3 \ 2] [4 \ 0.5]^T}{\sqrt{3^2+2^2} \sqrt{4^2+0.5^2}} = 0.106$$

$$\| [3, 2]^T - [6, 4]^T \|_{\cos} = 1 - \frac{[3 \ 2] [6 \ 4]^T}{\sqrt{3^2+2^2} \sqrt{6^2+4^2}} = 0$$

$$\| [3, 2]^T - [4, 5]^T \|_{\cos} = 1 - \frac{[3 \ 2] [4 \ 5]^T}{\sqrt{3^2+2^2} \sqrt{4^2+5^2}} = 0.047$$

- (e) [2 points] What is the label if the following Mahalanobis distance is used? Given two vectors \mathbf{x} and \mathbf{x}' , their Mahalanobis distance is defined as $(\mathbf{x} - \mathbf{x}')^\top \mathbf{A} (\mathbf{x} - \mathbf{x}')$. Here, please consider \mathbf{A} as $\begin{bmatrix} 9 & -3 \\ -3 & 1 \end{bmatrix}$.

$$\| [3, 2]^T - [5, 2]^T \|_{\text{Mahalanobis}} = [3-5 \ 2-2] \begin{bmatrix} 9 & -3 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 3-5 \\ 2-2 \end{bmatrix} = 81 \quad \text{diamond}$$

$$\| [3, 2]^T - [4, 0.5]^T \|_{\text{Mahalanobis}} = [3-4 \ 2-0.5] \begin{bmatrix} 9 & -3 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 3-4 \\ 2-0.5 \end{bmatrix} = 20.25$$

$$\| [3, 2]^T - [6, 4]^T \|_{\text{Mahalanobis}} = [3-6 \ 2-4] \begin{bmatrix} 9 & -3 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 3-6 \\ 2-4 \end{bmatrix} = 49$$

$$\| [3, 2]^T - [4, 5]^T \|_{\text{Mahalanobis}} = [3-4 \ 2-5] \begin{bmatrix} 9 & -3 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 3-4 \\ 2-5 \end{bmatrix} = 0$$

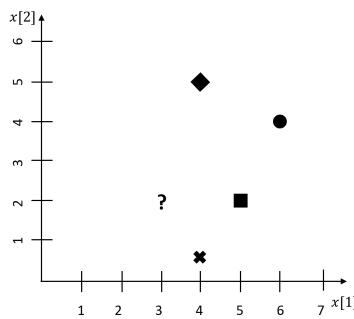


Figure 1: A training set with four labeled data instances and a test data instance (i.e., the question mark).

2) **Maximum Likelihood Estimation [10 points].** Consider a random variable \mathbf{X} (possibly a vector) whose distribution may be written as $f(\mathbf{x}; \theta)$, where θ is called the parameter. *Maximum likelihood estimation* is one of the most important parameter estimation techniques. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be iid (independent and identically distributed) random variables distributed according to $f(\mathbf{x}; \theta)$. By independence, the joint distribution of the observations is the product

$$\prod_{i=1}^n f(\mathbf{X}_i; \theta) \quad (1)$$

Viewed as a function of θ , this quantity is called the likelihood of θ . It is often more convenient to work with the *log-likelihood*,

$$\sum_{i=1}^n \log f(\mathbf{X}_i; \theta) \quad (2)$$

A maximum likelihood estimate (MLE) of θ is any parameter

$$P(\mathbf{x}; \lambda) = \frac{e^{-\lambda} \lambda^{\mathbf{x}}}{\mathbf{x}!} \quad \hat{\theta} \in \arg \max_{\theta} \sum_{i=1}^n \log f(\mathbf{X}_i; \theta) \quad (3)$$

where “ $\arg \max$ ” denotes the set of all values achieving the maximum. If there is a unique maximizer, it is called the maximum likelihood estimate. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be iid Poisson random variables with intensity parameter λ . Determine the maximum likelihood estimator of λ .

$$\begin{aligned} \lambda \in \arg \max_{\lambda} \sum_{i=1}^n \log f(X_i; \lambda) &= \arg \max_{\lambda} \sum_{i=1}^n \log \frac{e^{-\lambda} \lambda^{X_i}}{X_i!} \\ &= \arg \max_{\lambda} \sum_{i=1}^n (\log e^{-\lambda} + \log \lambda^{X_i} - \log X_i!) \\ &= \arg \max_{\lambda} \sum_{i=1}^n (-\lambda + X_i \log \lambda - \log X_i!) \\ &= \arg \max_{\lambda} \sum_{i=1}^n \frac{\partial}{\partial \lambda} (-\lambda + X_i \log \lambda - \log X_i!) \\ &= \arg \max_{\lambda} \sum_{i=1}^n \left(-1 + \frac{X_i}{\lambda} \right) \end{aligned}$$

$$\sum_{i=1}^n \left(\frac{X_i}{\lambda} - 1 \right) = 0$$

$$\sum_{i=1}^n \frac{X_i}{\lambda} = \sum_{i=1}^n 1 \Rightarrow \sum_{i=1}^n \frac{X_i}{\lambda} = n$$

$$\lambda = \frac{1}{n} \sum_{i=1}^n X_i$$

3) Linear regression with input dependent noise [20 points].

From the probabilistic perspective of linear regression, we have a model in which the output is linearly dependent on the input with respect to the parameters.

$$y = \mathbf{w}^\top \mathbf{x} + \epsilon \quad (4)$$

We will see in the class if the noise is from a normal distribution (i.e., $\epsilon \sim \mathcal{N}(0, \sigma^2)$), maximizing likelihood estimation leads to the sum of squares loss (or residual sum of square loss) formulation.

Here, we will consider a simple case where y and x are scalars (i.e., $x, y \in \mathbb{R}$) but the noise is dependent on the input such that $\epsilon \sim \mathcal{N}(0, \sigma^2 x^2)$. That is, the standard deviation scales linearly with the input. We observe N independent training examples denoted as $D_{\text{tr}} = \{(x_i, y_i)\}_{i=1}^N$. The model is

$$y = wx + \epsilon.$$

Let's assume x is IID sampled from a known distribution with a probability density $p(x)$.

- (a) [5 points] Write an expression for the likelihood of data $p(D_{\text{tr}}; w)$ in terms of $p(x)$, σ , x_i , y_i , w , N , and π .

$$Y = W^\top X$$

$$P(Y|X) \sim \mathcal{N}(W^\top X, \sigma^2)$$

$$Y = W^\top X + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 x^2)$$

$$P(Y|X) \sim \mathcal{N}(W^\top X, \sigma^2 x^2)$$

$$P(Y|X; w) = \frac{1}{\sqrt{2\pi\sigma^2 x^2}} e^{-\frac{(Y-W^\top X)^2}{2\sigma^2 x^2}}$$

$$\begin{aligned} P(D_{\text{tr}}; w) &= \prod_{i=1}^N P(x_i, y_i; w) \\ &= \prod_{i=1}^N P(x_i) P(y_i|x_i; w) \\ &= \sum_{i=1}^N \log P(x_i) P(y_i|x_i; w) \\ &= \sum_{i=1}^N P(x_i) \log \frac{1}{\sqrt{2\pi\sigma^2 x_i^2}} e^{-\frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2}} \end{aligned}$$

- (b) [5 points] Find the maximum likelihood estimate w^{ML} of weights w .

$$= \alpha \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2 x_i^2}} + \log e^{-\frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2}} \right]$$

$$= \alpha \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2 x_i^2}} - \frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2} \right]$$

$$w = \arg \max_w \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2 x_i^2}} - \frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2} \right]$$

$$w = \arg \min_w \sum_{i=1}^N \frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2}$$

$$\frac{\partial}{\partial w} \left[\frac{(y_i - w^\top x_i)^2}{2\sigma^2 x_i^2} \right] = 0$$

$$= \frac{-2(y_i - w^\top x_i)}{2\sigma^2 x_i^2} x_i = 0$$

$$= \sum_{i=1}^N \frac{-y_i}{\sigma^2 x_i} + \sum_{i=1}^N \frac{w^\top x_i}{\sigma^2 x_i}$$

$$w^\top \sum_{i=1}^N \frac{1}{\sigma^2 x_i} = \sum_{i=1}^N \frac{y_i}{\sigma^2 x_i}$$

$$w^\top = \sum_{i=1}^N \frac{y_i}{\sigma^2 x_i} \frac{\sigma^2}{N} = \frac{\sum_{i=1}^N \frac{y_i}{x_i}}{N}$$

- (c) [5 points] Assuming that the prior distribution of w is $\mathcal{N}(0, \sigma^2)$, express the posterior distribution of w , i.e., $p(w|D_{\text{tr}})$ in terms of $p(x)$, σ , α , x_i , y_i , w , N , and π .

$$P(D_{\text{tr}}|w) = P(x_i) \frac{1}{\sqrt{2\pi\sigma^2} X_i^2} e^{-\frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2}}$$

$$P(w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(w-\mu)^2}{2\sigma^2}}$$

$$P(w|D_{\text{tr}}) = \frac{P(D_{\text{tr}}|w) P(w)}{P(D_{\text{tr}})}$$

$$= \frac{\prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi\sigma^2} X_i^2} e^{-\frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2}} \right) \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{w^2}{2\sigma^2}} \right) P(w)}{P(D_{\text{tr}})}$$

- (d) [5 points] Find the MAP estimate w^{MAP} of weights w .

$$\begin{aligned} & \log P(w|D_{\text{tr}}) \\ &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2} X_i^2} e^{-\frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2}} \right) \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{w^2}{2\sigma^2}} \right) \\ &= \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2} X_i^2} - \frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2} + \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{w^2}{2\sigma^2} \right] \end{aligned}$$

$$\underset{w}{\operatorname{argmax}} \log P(w|D_{\text{tr}}) = \underset{w}{\operatorname{argmin}} \left[\frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2} + \frac{w^2}{2\sigma^2} \right]$$

$$\frac{\partial}{\partial w} \left[\frac{(y_i - w^T X_i)^2}{2\sigma^2 X_i^2} + \frac{w^2}{2\sigma^2} \right] = 0$$

$$= \frac{2(y_i - w^T X_i)}{2\sigma^2 X_i^2} X_i + \frac{2w}{2\sigma^2}$$

$$= \sum_{i=1}^N \frac{-y_i}{\sigma^2 X_i} + \sum_{i=1}^N \frac{w^T X_i}{\sigma^2 X_i} + \frac{w}{\sigma^2}$$

$$W \sum_{i=1}^N \frac{1}{\sigma^2} + W \frac{1}{\sigma^2} = \sum_{i=1}^N \frac{y_i}{X_i \sigma^2}$$

$$W \left(\frac{N}{\sigma^2} + \frac{1}{\sigma^2} \right) = \sum_{i=1}^N \frac{y_i}{X_i \sigma^2}$$

$$W = \frac{\sum_{i=1}^N \frac{y_i}{X_i}}{1+N}$$

4) **Weighted Linear Regression [10 points].** Consider a linear regression problem in which we want to weigh different training examples differently. Specifically, suppose we want to minimize

$$E_D(w) = \frac{1}{2} \sum_{i=1}^N r_i (w^\top x_i - y_i)^2. \quad (5)$$

In class, we worked out what happens for the case where all the weights (r_i 's) are one. In this problem, we will generalize some of those ideas to the weighted setting. In other words, we will allow the weights r_i to be different for each of the training examples.

Suppose we have a training set $\{(x_i, y_i); i = 1, \dots, N\}$ of N independent examples, but in which the y_i 's were observed with different variances. Specifically, suppose that

$$p(y_i|x_i; w) = \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left(-\frac{(y_i - w^\top x_i)^2}{2(\sigma_i)^2}\right). \quad (6)$$

In other words, y_i has mean $w^\top x_i$ and variance $(\sigma_i)^2$ (where the σ_i 's are fixed, known constants). Show that finding the maximum likelihood estimate of w reduces to solving a weighted linear regression problem. State clearly what the r_i 's are in terms of the σ_i 's.

$$\begin{aligned} P(D|w) &= \prod_{i=1}^N P(y_i|x_i; w) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(y_i - w^\top x_i)^2}{2(\sigma_i)^2}} \\ \log P(D|w) &= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(y_i - w^\top x_i)^2}{2(\sigma_i)^2}} \right] \\ \arg \max_w \log P(D|w) &= \arg \min_w \sum_{i=1}^N \left[\frac{(y_i - w^\top x_i)^2}{2(\sigma_i)^2} \right] \end{aligned}$$

$$\text{Since we want to } \arg \min_w \frac{1}{2} \sum_{i=1}^N r_i (w^\top x_i - y_i)^2$$

$$\text{so } r_i = \frac{1}{(\sigma_i)^2}$$

5) Naive Bayes as linear classifiers [15 points].

Given a training set $D_{\text{tr}} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{n=1}^N$ for binary classification, let us fit a Naive Bayes model to it; i.e., $p(y|\mathbf{x}) \propto p(y) \prod_{d=1}^D p(x[d]|y)$, where $p(x[d]|y)$ is a Gaussian distribution $\mathcal{N}(\mu_{d,y}, \sigma_d^2)$. That is, we assume that for each dimension d , the two one-dimensional Gaussian distributions (one for each class) share the same variance. For $p(y)$, let $p(+1) = \lambda$. Please show the followings.

(a) [6 points] $p(y|\mathbf{x})$ can be re-written as $\frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}$ ($y \in \{+1, -1\}$)

$$\begin{aligned} P(y=1|\mathbf{x}) &= \frac{P(x|y=1) P(y=1)}{P(x)} \\ &= \frac{P(x|y=1) P(y=1)}{P(x, y=1) + P(x, y=-1)} \\ &= \frac{P(x|y=1) P(y=1)}{P(x|y=1) P(y=1) + P(x|y=-1) P(y=-1)} \\ &= \frac{1}{1 + \frac{P(x|y=-1) P(y=-1)}{P(x|y=1) P(y=1)}} \end{aligned}$$

$$= \frac{1}{1 + e^{\log \left[\frac{P(x|y=-1) P(y=-1)}{P(x|y=1) P(y=1)} \right]}}$$

$$\begin{aligned} \text{Let } Z &= \log \frac{P(x|y=1) P(y=1)}{P(x|y=-1) P(y=-1)} = y(\mathbf{w}^\top \mathbf{x} + b) \\ &= \frac{1}{1 + e^{-Z}} \end{aligned}$$

$$P(y=1|\mathbf{x}) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}$$

(b) [6 points] What are the corresponding \mathbf{w} and b ?

Your answers should be based on $\mu_{d,y}$, σ_d , and λ . For $\mathbf{w} \in \mathbb{R}^D$, you may simply write the expression of $w[d]$ for a specific d .

$$\begin{aligned} y(\mathbf{w}^\top \mathbf{x} + b) &= \log \frac{P(x|y=1) P(y=1)}{P(x|y=-1) P(y=-1)} = \log \frac{P(y=1) \prod_{d=1}^D P(x[d]|y)}{P(y=-1) \prod_{d=1}^D P(x[d]|y)} = \log \frac{\lambda}{1-\lambda} + \sum_{d=1}^D \log \frac{\frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-\frac{(x[d]-\mu_{d,y=1})^2}{2\sigma_d^2}}}{\frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-\frac{(x[d]-\mu_{d,y=-1})^2}{2\sigma_d^2}}} \\ &= \log \frac{\lambda}{1-\lambda} + \sum_{d=1}^D \left[\log e^{-\frac{(x[d]-\mu_{d,y=1})^2}{2\sigma_d^2}} - \log e^{-\frac{(x[d]-\mu_{d,y=-1})^2}{2\sigma_d^2}} \right] = \log \frac{\lambda}{1-\lambda} + \sum_{d=1}^D \left(-\frac{(x[d]-\mu_{d,y=1})^2}{2\sigma_d^2} + \frac{(x[d]-\mu_{d,y=-1})^2}{2\sigma_d^2} \right) \\ &= \log \frac{\lambda}{1-\lambda} + \sum_{d=1}^D \frac{x_d^2 + \mu_{d,y=1}^2 - 2x_d\mu_{d,y=1} - x_d^2 - \mu_{d,y=-1}^2 + 2x_d\mu_{d,y=-1}}{2\sigma_d^2} \\ W[d] &= \sum_{d=1}^D \frac{2x_d\mu_{d,y=1} - 2x_d\mu_{d,y=-1}}{2\sigma_d^2} = x[d] \frac{\mu_{d,y=1} - \mu_{d,y=-1}}{\sigma^2} \\ b &= \log \frac{\lambda}{1-\lambda} + \sum_{d=1}^D \frac{\mu_{d,y=1}^2 - \mu_{d,y=-1}^2}{2\sigma^2} \end{aligned}$$

6) Naive Bayes Classifier [35 points].

Download the files **spambase.train** and **spambase.test**.

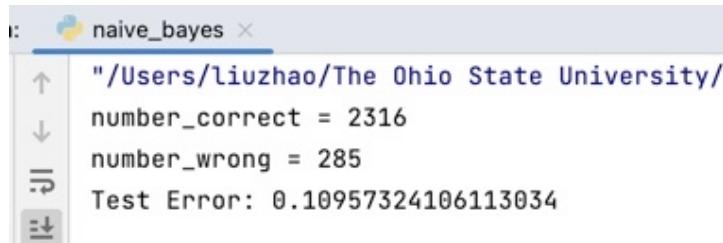
The file spambase.train contains 2000 training data and spambase.test has 2601 test data. Both datasets have 58 columns: the first 57 columns are input features, corresponding to different properties of an email, and the last column is an output label indicating spam (1) or non-spam (0). Please fit the Naive Bayes model using the training data.

As a pre-processing step, quantize each variable to one of two values, say 1 and 2, so that values below the median map to 1, and others map to 2. Please aggregate the training and test test to obtain the median value of each variable, and then use the median to quantize both data sets.

- (a) Report the test error (misclassification percentage) of Naive Bayes classifier. As a sanity check, what would be the test error if you always predicted the same class, namely, the majority class from the training data?

Baseline Test error : 38.56 %

Naive Bayes test error : 10.957 %



```
naive_bayes
"/Users/liuzhao/The Ohio State University/
number_correct = 2316
number_wrong = 285
Test Error: 0.10957324106113034
```

```
1 import numpy as np
2 import pandas as pd
3
4 SPAM = 1
5 NON_SPAM = 0
6
7
8 if __name__ == '__main__':
9     train = pd.read_csv("spambase.train", header=None)
10    test = pd.read_csv("spambase.test", header=None)
11
12    train = pd.concat([train, test], ignore_index=True, sort=False)
13
14    number_correct = 0
15    number_wrong = 0
16
17    # number of rows in Test dataset
18    test_rows = test.shape[0]
19    # number of columns in Test dataset
20    test_columns = test.shape[1]
21
22    # number of rows in Training dataset
23    train_rows = train.shape[0]
24    # number of columns in Training dataset
25    train_columns = train.shape[1]
26
27    # Training dataset with Y = 1
28    train_y_1 = train[train[57].isin([1])]
29    number_train_y_1 = train_y_1.shape[0]
30
31    # Training dataset with Y = 0
32    train_y_0 = train[train[57].isin([0])]
33    number_train_y_0 = train_y_0.shape[0]
34
35    # Probability(Y = 1)
36    probability_y_1 = number_train_y_1 / train_rows
37    # Probability(Y = 0)
38    probability_y_0 = number_train_y_0 / train_rows
```

```
39
40     medians = []
41     for column in range(train_columns - 1):
42         medians.append(train[column].median())
43         # print(f"column = {column}, median = {train[column].median()}")
44
45     """
46     median_condition_probability = np.array([[1, 2
47 , 3]])
48
49     for column in range(train_columns - 1):
50         median = medians[column]
51
52         number_match_1 = train_y_1[train_y_1[column
53 ] <= median].shape[0]
54         theta_1 = number_match_1 / number_train_y_1
55
56         number_match_0 = train_y_0[train_y_0[column
57 ] <= median].shape[0]
58         theta_0 = number_match_0 / number_train_y_0
59
60         row_median_proba = np.array([[median,
61             theta_1, theta_0]])
62
63         if column == 0:
64             median_condition_probability =
65             row_median_proba
66         else:
67
68             median_condition_probability = np.
69             concatenate((median_condition_probability,
70             row_median_proba))
71
72         """
73
74         #print(median_condition_probability)
75         #number_match_1 = train_y_1[train_y_1[11] < 0.
76         #1450000000000002].shape[0]
77         #print(f"number_match_1 = {number_match_1}")
78         #print(f"number_train_y_1 = {number_train_y_1
79         })")
80         #print(f"number_train_y_1 = {number_train_y_0
```

```
69  }")
70
71     for row in range(test_rows):
72         predict_value = 0
73         real_value = test.iloc[row, test_columns
74 - 1]
75
76         # For spam (Y=1) which indicates label = 1
77         probability_spam = 1
78         # For Non-spam (Y=0) which indicates label
79 = 0
80         probability_non_spam = 1
81
82         for column in range(test_columns - 1):
83             new_value = test.iloc[row, column]
84
85             median = medians[column]
86
87             number_match_1 = train_y_1[train_y_1[
88             column] <= median].shape[0]
89             theta_1 = number_match_1 /
90             number_train_y_1
91
92             number_match_0 = train_y_0[train_y_0[
93             column] <= median].shape[0]
94             theta_0 = number_match_0 /
95             number_train_y_0
96
97             if new_value > median:
98                 theta_1 = 1 - theta_1
99                 theta_0 = 1 - theta_0
100
101             probability_spam = probability_spam *
102             theta_1
103             probability_non_spam =
104             probability_non_spam * theta_0
105
106             probability_spam = probability_spam *
107             probability_y_1
108             probability_non_spam =
109             probability_non_spam * probability_y_0
```

```
100
101      #print(f"row = {row}, probability_spam = {probability_spam}, probability_non_spam = {probability_non_spam}")
102
103      if probability_spam >=
104          probability_non_spam:
105              predict_value = 1
106          else:
107              predict_value = 0
108
109      #print(f"row = {row}, real_value = {real_value}, predict_value = {predict_value}")
110
111      if predict_value == real_value:
112          number_correct = number_correct + 1
113      else:
114          number_wrong = number_wrong + 1
115
116      print(f"number_correct = {number_correct}")
117      print(f"number_wrong = {number_wrong}")
118
119      print(f"Test Error: {number_wrong / test_rows}")
120
```