



MONASH University

FIT5202 DATA ANALYSIS ALGORITHMS

Assignment 1

Author: Suyash Sathe

Student Id: 29279208

Email: ssat0005@student.monash.edu

Master of Data Science, Monash University

Section C. Probabilistic Machine Learning

Question 4 [Bayes Rule, 20 Marks]

Recall the simple example from Appendix A of Module 1. Suppose we have one red and one blue box. In the red box we have 2 apples and 6 oranges, whilst in the blue box we have 3 apples and 1 orange. Now suppose we randomly selected one of the boxes and picked a fruit. If the picked fruit is an apple, what is the probability that it was picked from the blue box?

Note that the chance of picking the red box is 40% and the selection chance for any of the pieces from a box is equal for all the pieces in that box.

Solution

Given: Picked fruit is an apple.

To find: Probability that the apple was picked from the blue box, i.e., $P(\text{Blue} \mid \text{Apple})$

According to Bayes' Rule,

$$P(\text{Blue} \mid \text{Apple}) = \frac{P(\text{Apple} \mid \text{Blue}) P(\text{Blue})}{P(\text{Apple})}$$

Where,

$$P(\text{Apple} \mid \text{Blue}) = \frac{3}{4}$$

$$P(\text{Blue}) = \frac{6}{10}$$

$$\begin{aligned} P(\text{Apple}) &= P(\text{Apple} \mid \text{Blue}) * P(\text{Blue}) + P(\text{Apple} \mid \text{Red}) * P(\text{Red}) \\ &= \left(\frac{3}{4} \times \frac{6}{10}\right) + \left(\frac{2}{8} \times \frac{4}{10}\right) = \frac{44}{80} = \frac{11}{20} \end{aligned}$$

$$\therefore P(\text{Blue} \mid \text{Apple}) = \frac{\frac{3}{4} \times \frac{6}{10}}{\frac{11}{20}} = \frac{3}{4} \times \frac{6}{10} \times \frac{20}{11} = \frac{9}{11}$$

$$\text{Thus, } P(\text{Blue} \mid \text{Apple}) = \frac{9}{11}$$

Question 5 [Maximum Likelihood, 20 Marks]

As opposed to a coin which has two faces, a dice has 6 faces. Suppose we are given a dataset which contains the outcomes of 10 independent tosses of a dice: $D = \{1, 4, 5, 3, 1, 2, 6, 5, 6, 6\}$. We are asked to build a model for this dice, i.e. a model which tells what the probability of each face of the dice is if we toss it. Using the maximum likelihood principle, please determine the best value for our model parameters.

Solution

Given: $D = \{1, 4, 5, 3, 1, 2, 6, 5, 6, 6\}$

Since the dice has 6 faces, let the probability associated with each face (1, 2, 3, 4, 5, 6) be $\omega_1, \omega_2, \omega_3, \omega_4, \omega_5$ and ω_6 .

Probability of ω_i where $i \in \{1, 6\}$ is $0 \leq \omega_i \leq 1$

Applying the maximum likelihood principle,

$p(D | \omega_1) = \prod_{i=1}^{10} p(x_i | \omega_1) = \omega_1^2 (1 - \omega_1)^8$, where x_i represents every event in the dataset.

$$\log(p(D | \omega_1)) = \log \omega_1^2 + \log(1 - \omega_1)^8$$

$$\frac{d \log(p(D | \omega_1))}{d \omega_1} = \frac{2}{\omega_1} - \frac{8}{(1 - \omega_1)}$$

$$0 = \frac{2}{\omega_1} - \frac{8}{(1 - \omega_1)}$$

$$\frac{2}{\omega_1} = \frac{8}{(1 - \omega_1)}$$

$$10\omega_1 = 2$$

$$\therefore \omega_1 = 0.2$$

Similarly,

$$p(D | \omega_2) = \prod_{i=1}^{10} p(x_i | \omega_2) = \omega_2 (1 - \omega_2)^9$$

$$\log(p(D | \omega_2)) = \log \omega_2 + \log(1 - \omega_2)^9$$

$$\frac{d \log(p(D | \omega_2))}{d \omega_2} = \frac{1}{\omega_2} - \frac{9}{(1 - \omega_2)}$$

$$0 = \frac{1}{\omega_2} - \frac{9}{(1-\omega_2)}$$

$$\frac{1}{\omega_2} = \frac{9}{(1-\omega_2)}$$

$$10\omega_2 = 1$$

$$\therefore \omega_2 = 0.1$$

$$p(D | \omega_3) = \prod_{i=1}^{10} p(x_i | \omega_3) = \omega_3(1 - \omega_3)^9$$

$$\log(p(D | \omega_3)) = \log \omega_3 + \log(1 - \omega_3)^9$$

$$\frac{d\log(p(D | \omega_3))}{d\omega_3} = \frac{1}{\omega_3} - \frac{9}{(1-\omega_3)}$$

$$0 = \frac{1}{\omega_3} - \frac{9}{(1-\omega_3)}$$

$$\frac{1}{\omega_3} = \frac{9}{(1-\omega_3)}$$

$$10\omega_3 = 1$$

$$\therefore \omega_3 = 0.1$$

$$p(D | \omega_4) = \prod_{i=1}^{10} p(x_i | \omega_4) = \omega_4(1 - \omega_4)^9$$

$$\log(p(D | \omega_4)) = \log \omega_4 + \log(1 - \omega_4)^9$$

$$\frac{d\log(p(D | \omega_4))}{d\omega_4} = \frac{1}{\omega_4} - \frac{9}{(1-\omega_4)}$$

$$0 = \frac{1}{\omega_4} - \frac{9}{(1-\omega_4)}$$

$$\frac{1}{\omega_4} = \frac{9}{(1-\omega_4)}$$

$$10\omega_4 = 1$$

$$\therefore \omega_4 = 0.1$$

$$p(D | \omega_5) = \prod_{i=1}^{10} p(x_i | \omega_5) = \omega_5^2(1 - \omega_5)^8$$

$$\log(p(D | \omega_5)) = \log \omega_5^2 + \log(1 - \omega_5)^8$$

$$\frac{d\log(p(D | \omega_5))}{d\omega_5} = \frac{2}{\omega_5} - \frac{8}{(1-\omega_5)}$$

$$0 = \frac{2}{\omega_5} - \frac{8}{(1-\omega_5)}$$

$$\frac{2}{\omega_5} = \frac{8}{(1-\omega_5)}$$

$$10\omega_5 = 2$$

$$\therefore \omega_5 = 0.2$$

$$p(D | \omega_6) = \prod_{i=1}^{10} p(x_i | \omega_6) = \omega_6^3 (1 - \omega_6)^7$$

$$\log(p(D | \omega_6)) = \log \omega_6^3 + \log(1 - \omega_6)^7$$

$$\frac{d\log(p(D | \omega_6))}{d\omega_6} = \frac{3}{\omega_6} - \frac{7}{(1-\omega_6)}$$

$$0 = \frac{3}{\omega_6} - \frac{7}{(1-\omega_6)}$$

$$\frac{3}{\omega_6} = \frac{7}{(1-\omega_6)}$$

$$10\omega_6 = 3$$

$$\therefore \omega_6 = 0.3$$

Thus, the best values of our model parameters are $\omega_1 = 0.2$, $\omega_2 = 0.1$, $\omega_3 = 0.1$, $\omega_4 = 0.1$, $\omega_5 = 0.2$ and $\omega_6 = 0.3$.

Section D. Ridge Regression

Question 6 [Ridge Regression, 25 Marks]

Given the gradient descent algorithms for linear regression (discussed in Chapter 2 of Module 2), derive weight update steps of stochastic gradient descent (SGD) as well as batch gradient descent (BGD) for linear regression with L2 regularisation norm. Show your work with enough explanation in your PDF report; you should provide the steps of SGD and BGD, separately.

Solution

Derivation of weight update steps for SGD and BGD

For Ridge Regression,

Objective function = Error + Regularization

Objective function = $E(w) + \lambda \cdot \Omega(w)$

Where,

$E(w) \rightarrow$ error term

$\lambda \rightarrow$ regularization parameter

$\Omega(w) \rightarrow$ penalty (regularization function)

For Linear Regression,

$$E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w \cdot \phi(x_n))^2$$

For L2 Regularization, regularization function,

$$\Omega(w) = \frac{\lambda}{2} \sum_{j=0}^{M-1} w_j^2$$

Therefore,

$$\text{Objective function} = \frac{1}{2} \sum_{n=1}^N (t_n - w \cdot \phi(x_n))^2 + \frac{\lambda}{2} \sum_{j=0}^{M-1} w_j^2$$

Now, gradient of the objective function,

$$\begin{aligned}\nabla(\text{Objective function}) &= \nabla \left(\frac{1}{2} \sum_{n=1}^N (t_n - w \cdot \phi(x_n))^2 + \frac{\lambda}{2} \sum_{j=0}^{M-1} w_j^2 \right) \\ &= \nabla \left(\frac{1}{2} \sum_{n=1}^N (t_n - w \cdot \phi(x_n))^2 \right) + \nabla \left(\frac{\lambda}{2} \sum_{j=0}^{M-1} w_j^2 \right)\end{aligned}$$

$$\therefore \nabla(\text{Objective function}) = \left[- \sum_{n=1}^N (t_n - w \cdot \phi(x_n)) \right] \cdot \phi(x_n) + \lambda \cdot w$$

where, $\sum_{j=0}^{M-1} w_j = w$

For SGD and BGD,

$$\text{new parameter} = \text{old parameter} - \eta \nabla(\text{Objective function})$$

where, η is the learning rate.

Therefore,

$$w^{(T)} = w^{(T-1)} - \eta^{(T)} \nabla(\text{Objective function})$$

$$\therefore w^{(T)} = w^{(T-1)} - \eta^{(T)} \left(\left[- \sum_{n=1}^N (t_n - w^{(T-1)} \cdot \phi(x_n)) \right] \cdot \phi(x_n) + \lambda \cdot w^{(T-1)} \right)$$

Therefore, $w^{(T)}$ is the updated weight for stochastic gradient descent (SGD) as well as batch gradient descent (BGD) for linear regression with L2 regularisation norm.

Stochastic Gradient Descent (SGD)

In SGD algorithm, the weight parameters are calculated and updated for each row/value (X_i) of the training data. The weight parameter is updated as per the above derivation of $w^{(T)}$. This updated weight parameter is used to update the next weight parameter. This process is continued till the termination condition is reached. Following is the algorithm for the above steps:

1. Initialize the weight parameters to some random initial values.

$$w^{(1)} = [w_0, w_1, w_2, \dots, w_m]$$

2. Plug the parameter values into the weight update function for every point in the dataset to get the gradient ($\nabla(\text{Objective function})$).
3. Calculate the step size: $\eta \cdot \nabla(\text{Objective function})$

$$step\ size = \eta^{(1)} \left(\left[- \sum_{n=1}^N (t_n - w^{(0)} \cdot \phi(x_n)) \right] \cdot \phi(x_n) + \lambda \cdot w^{(0)} \right)$$

4. Calculate the new parameter:

$$new\ parameter = old\ parameter - step\ size$$

$$new\ parameter = old\ parameter - \eta \cdot \nabla(Objective\ function)$$

5. Repeat steps 3 to 5 until the step size becomes very small ,i.e., the algorithm converges, or maximum number of iterations are performed.

In this algorithm, the dataset is resampled to get a batch of the dataset and for every data point in the batch, the weight parameters are updated until it converges, or threshold of iterations is reached.

Batch Gradient Descent (BGD)

In the BGD algorithm, the weight parameters are calculated and updated using entire dataset as a batch and not every data point in the dataset. The calculation of weight parameters is similar to SGD, but instead of every point in the dataset, the gradient is calculated using the mean of all data points (X_i) in the batch. Following is the algorithm for the above steps:

1. Initialize the weight parameters to some random initial values like SGD.

$$w^{(1)} = [w_0, w_1, w_2, \dots, w_m]$$
2. Plug the parameter values into the weight update function and use the mean of predicted values and true values from the dataset to get the gradient ($\nabla(Objective\ function)$).

(Note that mean of true and predicted values for the dataset is taken and not the individual data points).

3. Calculate the step size: $\eta \cdot \nabla(Objective\ function)$

$$step\ size = \eta \left(\left[- \sum_{n=1}^N (mean(t_n) - w^{(0)} \cdot mean(\phi(x_n))) \right] \cdot mean(\phi(x_n)) + \lambda \cdot w^{(0)} \right)$$

4. Calculate the new parameter:

$$new\ parameter = old\ parameter - step\ size$$

5. Repeat steps 3 to 5 until the step size becomes very small ,i.e., the algorithm converges, or maximum number of iterations are performed.

Thus, the implementation and function of SGD and BGD is quite similar. The only difference is that SGD calculates weight parameters for each point in the dataset whereas BGD considers the mean of the batch of the dataset.