

Amin Fadaee

Machine Learning HW3

Python 3.5, Libraries: Numpy, Pandas, Random, Matplotlib, Math, Scipy Documentation: HTML, Bootstrap 4, jqMath, Graphs: Xmind

Problem 1

- By inspecting the plot, it can be seen that based on the distributions that we would derive for each of the classes, the probability of $p(B|x)$ is greater than $p(A|x)$.
- The instance in the second plot is clearly an outlier whether it belongs to A or B, however it would have higher probability to belong to B ($q(B|x) = p(A|x)$). This is due to the fact that the variance for the class A is really small comparing to B and the probability of the specified deviation considering this little amount of variance is highly unlikely.
- The decision boundary for this data is going to be a circle which includes the 4 points belonging to B, therefore all instances in the circle is going to be classified as B and all the ones out of the circle as A, so the label is B.

Problem 2

Assuming we have n classes and each instance has n features if we were to classify a new instance based on a generative model (let's say a gaussian bayesian classifier) we would need to compute the covariance matrix of $n \times n$ and the mean vector μ of length n for each of the n classes:

- If all the features are correlated, then C would have no non-zero elements and therefore finding C and μ would require finding the values of $n^2 \times n$ parameters for each class. Also the priors should be obtained as well (be priors), so in total we would need to compute the values of $n(n^2 + n) + n$.
- If the features are independent then C would become a diagonal matrix and therefore we need to find n parameters for each class and n priors which add up to $2n + n$ parameters.
- Gaussian bayes classifiers are not constrained by the values of each attribute, so if instances have real values and features are independent, then we still need to derive $2n + n$ parameters.

Problem 3

Here is the joint probability distribution of the network:

$$P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)$$

Now if we like to obtain $P(B|D=1)$ we can do this using the conditional probability:

$$P(B|D=1) = \frac{P(A,B,C,D=1)}{P(D=1)} = \frac{P(A)P(B|A)P(C|A,B)P(D=1|A,B,C)}{P(D=1|A,B,C)}$$

However, dealing with the posterior would be much easier:

$$P(B|D=1) = \frac{P(D=1|B)P(B)}{P(D=1)}$$

Now we are going to compute the probability of each of the components in the above formula. First $P(B)$:

$$\begin{aligned} P(B) &= \sum_A P(A,B) = \sum_A P(A)P(B|A) = P(A=1)P(B|A=1) + P(A=2)P(B|A=2) \\ P(B=1) &= P(A=1)P(B=1|A=1) + P(A=2)P(B=1|A=2) = 0.2 \times 0.37 + 0.8 \times 0.21 \\ P(B) &= 0.347 \\ P(B=2) &= 0.153 \end{aligned}$$

Now let's do $P(D=1|B)$:

$$\begin{aligned} P(D=1|B) &= \sum_{A,C} P(A,B,C,D=1) = \sum_{A,C} P(A)P(B|A)P(C|A,B)P(D=1|A,B,C) \\ &= P(A=1)P(B=1|A=1)P(C=1|A=1,B=1)P(D=1|A=1,B=1,C=1) + P(A=1)P(B=1|A=1)P(C=2|A=1,B=1)P(D=1|A=1,B=1,C=2) \\ &\quad + P(A=2)P(B=1|A=2)P(C=1|A=2,B=1)P(D=1|A=2,B=1,C=1) + P(A=2)P(B=1|A=2)P(C=2|A=2,B=1)P(D=1|A=2,B=1,C=2) \\ P(D=1|B=1) &= 0.2 \times 0.37 \times 0.3 \times 0.5 + 0.2 \times 0.37 \times 0.7 \times 0.15 + 0.8 \times 0.21 \times 0.25 \times 0.5 + 0.8 \times 0.21 \times 0.75 \times 0.15 \\ P(D=1|B=1) &= 0.2 \times 0.6075 + 0.2 \times 0.03975 + 0.8 \times 0.2625 + 0.8 \times 0.23625 = 0.65775 \\ P(D=1|B=2) &= 0.8 \times 0.21 \times 0.3 \times 0.5 + 0.8 \times 0.21 \times 0.7 \times 0.15 + 0.2 \times 0.37 \times 0.25 \times 0.5 + 0.2 \times 0.37 \times 0.75 \times 0.15 \\ P(D=1|B=2) &= 0.42825 + 0.20325 + 0.04675 + 0.085125 = 0.763375 \\ P(D=1) &= P(D=1|B=1)P(B=1) + P(D=1|B=2)P(B=2) \\ P(D=1) &= 0.65775 \times 0.347 + 0.763375 \times 0.153 = 0.322 + 0.1168 = 0.4388 \end{aligned}$$

And finally $P(D=1)$:

$$P(D=1) = \sum_B P(B,D=1) = \sum_B P(B)P(D=1|B) = P(B=1)P(D=1|B=1) + P(B=2)P(D=1|B=2) = 0.347 \times 0.65775 + 0.153 \times 0.763375 = 0.228 + 0.1168 = 0.3448$$

And now for the ultimate result:

$$\begin{aligned} P(B|D=1) &= \frac{P(D=1|B)P(B)}{P(D=1)} = \frac{0.65775 \times 0.347}{0.3448} = 0.65775 \\ P(B=1|D=1) &= \frac{0.65775 \times 0.347}{0.3448} = 0.65775 \\ P(B=2|D=1) &= \frac{0.763375 \times 0.153}{0.3448} = 0.34225 \end{aligned}$$

$$P(B=1|D=1) = 0.65775$$

$$P(B=2|D=1) = 0.34225$$

Problem 4

In finding the probability that x belongs to class w_i ($p(w_i|x)$) there are two general practices. The generative and discriminative methods. In discriminative methods the values of $p(w_i|x)$ are obtained directly from the data (logistic regression is an example of these models), whereas in the generative models they are derived with an additional step of generating the conditional probabilities ($p(x|w_i)$). Bayesian classifiers are examples of generative models.

Problem 5

Logistic Regression directly estimates the parameters of $P(Y|X)$ (discriminative), whereas Naive Bayes directly estimates parameters for $P(X)$ and $P(Y|X)$ (generative). Asymptotically (as the number of training examples grows toward infinity) the GNB and Logistic Regression converge toward identical classifiers, however, GNB and Logistic Regression converge toward their asymptotic accuracies at different rates and Logistic Regression parameter estimates converge more slowly.

Logistic Regression

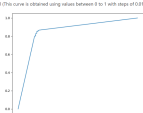
The following functions has been implement to build the logistic regression classifier for this problem:

- `predict_proba(X, w)`: This function takes in as argument the features matrix X and the weights w and computes the probability $P(Y = 1|x)$, w for x being a row in X it uses the sigmoid function $\frac{1}{1 + e^{-w \cdot x}}$.
- `features_derivative(X, labels, w, regularization)`: This function computes the derivative of log likelihood and takes as input X , The feature matrix (rows), list of labels of X rows, w , the list of coefficients and λ , penalty, the parameter λ for L2 regularization. This function implements the following formula:
$$\frac{\partial}{\partial w_j} \sum_{i=1}^n \log(P(Y_i = 1|x_i)) = \sum_{i=1}^n (Y_i - P(Y_i = 1|x_i)) x_{ij}$$
where 1 is the indicator function (1 if $Y_i = 1$, 0 if $Y_i = 0$) and the $-2\lambda w_j$ will be omitted for w_0 and $\lambda = 0$ for no regularization.
- `compute_log_likelihood(X, labels, w)`: This function computes the log likelihood which is the quality measure of logistic regression and can be used to check the proper functioning of our algorithm and printing its output throughout the algorithm runtime in order to check whether our gradient ascent algorithm is raising the log likelihood in each iteration.
- `logistic_regression(X, labels, w, step_size, max_iter, regularization)`: This function implements the logistic regression model training via gradient ascent by updating w using the following formula:
$$w_j^{t+1} = w_j^t + \text{step_size} \times \left(\frac{\partial \log L}{\partial w_j} - 2\lambda w_j \right)$$
- `false_positive_rate(labels, predictions)`: Computes the false positive rate of the prediction
- `true_positive_rate(labels, predictions)`: Computes the true positive rate of the prediction

After training the model on the provided data with no regularization ($\lambda = 0$) the training and test accuracy was as the following:

- Training Accuracy: 100%
- Test Accuracy: 84%

And here the ROC curve associated with this model (This curve is obtained using values between 0 to 1 with steps of 0.01 in size):



Now let's train the model with regularization. The training has been done with values of 10^{-1} , 10^{-2} , 10^{-3} , 10^{-4} , 10^{-5} , 10^{-6} , 10^{-7} , 10^{-8} , 10^{-9} , 10^{-10} for λ and each of the values has been assessed using 10-fold cross validation. Here are the results for each λ :

λ	1.00E-05	0.0001	0.001	0.01	0.1	1	10	100
Accuracy	0.60875	0.60875	0.60875	0.60675	0.60475	0.60275	0.60125	0.60075

Based on the table above the best λ is 0.01, so we train the whole training set with this parameter. Here is the performance and the ROC curve of this model:

- Training Accuracy: 100%
- Test Accuracy: 91%



As can be seen the performance dropped tremendously with regularization.

General Naive Bayes

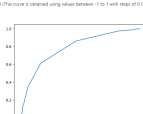
The following functions has been implement to build the Naive Bayes classifier for this problem:

- `stats(data, labels)`: Finds μ and diagonal covariance matrix and the prior probability for each class (assuming 2, one of which being 0) of data
- `conditional(mu, cov, x)`: Finds the conditional probability $p(x|y)$ assuming a gaussian underlying distribution and using the multivariate normal distribution formula. For not dealing with the hassles of high precision data, the probability is returned in the logarithmic form.
- `gib_classification(x, cov, priors, mu, covs, priors, x, threshold)`: Classifies x into 1 or zero using the stats provided via gaussian naive bayes procedure:
$$\begin{cases} g(A_1 - A_2) > 0 & \text{Class} = 1 \\ \text{else} & \text{Class} = 0 \\ A_i = & \log(p(x|w_i) + \log(prior)) \end{cases}$$
- `false_positive_rate(labels, predictions)`
- `true_positive_rate(labels, predictions)`

After using the training data to obtain the stats, here is the accuracy for each set of data after they have been classified:

- Training Accuracy: 0.80475
- Test Accuracy: 0.7225

And here the ROC curve associated with this model (This curve is obtained using values between -1 to 1 with steps of 0.1 in size as threshold, after normalizing the posteriors by their sum):



Bayesian Network

Using the `pgmpy` python library the following function has been implemented:

- `bayes_network(train, test, columns, labels, edges)`

This function creates a bayesian network using the edges provided as input, then test the model using the test data and returns the accuracy. Using 10-Fold cross validation and the Cars data set, different models has been tested. First let's take a look at the data:

The cars data set has 6 attributes in total:

- Buying Price (bought) with values `high`, `high`, `med`, `low`
- Maintenance Cost (mcost) with values `high`, `high`, `med`, `low`
- Number of Doors (doors) with values 2, 3, 4, `four`
- Passenger Capacity (persons) with values 2, 4, `more`
- Safety Level (safety) with values `low`, `med`, `high`
- Size of Luggage Boot (lugg_boot) with values `small`, `med`, `big`

Now let's see the different network structures in action:

- Naive structure with all of the features:



- Naive structure with features `bought`, `persons`, `safety` (On these 3 by intuition tend to have the most influence):



- Naive structure with features `mcost`, `doors`, `safety`, `bought`:



- Complex structure with features `bought`, `mcost`, `doors`, `persons`:



As the number of doors influences the number of persons able to use the car and the buying price influence the maintenance cost, we relate the two features.

- Complex structure with features `bought`, `mcost`, `doors`, `persons`, `safety`:



- Complex structure with features `bought`, `mcost`, `doors`, `persons`, `safety`, `bought`:



Here is the table of results:

	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
Network						
Accuracy	0.2312	0.8125	0.6302	0.6105	0.6267	0.2222

As can be seen the second model which we pick by intuition is the best among the naive structures and the accuracy even surpasses both models with all features (model 1 and model 6) and the best accuracy overall belongs to the 5th model where the logical assumption in the relations between 4 variables has been made added by safety.