

Lab04: Decision Tree and Naive Bayes

- Student ID: 21127453
- Student name: Hoang Anh Tra

How to do your homework

You will work directly on this notebook; the word `TODO` indicate the parts you need to do.

You can discuss ideas with classmates as well as finding information from the internet, book, etc...; but *this homework must be your*.

How to submit your homework

- Before submitting, save this file as `<ID>.jl`. For example, if your ID is 123456, then your file will be `123456.jl`. And export to PDF with name `123456.pdf` then submit zipped source code

Selection deleted to 123456.zip onto Moodle.

Danger

Note that you will get 0 point for the wrong submit.

Contents:

- Decision Tree
- Naive Bayes

Import library

```
1 begin
2     using Distributions, Plots, LinearAlgebra, Random, Statistics
3 end
```

```
TaskLocalRNG()
```

```
1 Random.seed!(2022)
```

Load Iris dataset

download_dataset (generic function with 2 methods)

```
1 # If you use Linux, use this function to download Iris dataset
2 function download_dataset(save_path::String="data")
3     # setup directory
4     mkpath(joinpath(dirname(@__FILE__), save_path))
5     data_dir = joinpath(dirname(@__FILE__), save_path)
6     download("https://archive.ics.uci.edu/static/public/53/iris.zip", joinpath(data_dir, "iris.zip"))
7     iris_file = joinpath(data_dir, "iris.zip")
8     cd(data_dir)
9     run(`unzip $iris_file -d $data_dir`)
10    rm(iris_file)
11    cd("..")
12 end
```

```
1 # If you have downloaded the dataset yet, please uncomment this line below and run
  this cell. Otherwise, keep it in uncomment state.
2 # download_dataset()
```

```
1 # If you don't use Linux, I have no solution for you. Please mkdir data, and goto
  https://archive.ics.uci.edu/dataset/53/iris for downloading
2 # Then, you can extract data to this dir by yourself.
3
4 # Structure:
5 # └── data
6 #     ├── bezdekIris.data
7 #     ├── Index
8 #     ├── iris.data
9 #     └── iris.names
10 # └── Lab4.jl
```

Selection deleted

iris_data_loader (generic function with 2 methods)

```
1 function iris_data_loader(data_path::String="data/iris.data")
2     # Initialize empty arrays to store data
3     sepal_length = Float64[]
4     sepal_width = Float64[]
5     petal_length = Float64[]
6     petal_width = Float64[]
7     classes = Int64[]
8
9     # open and read the data file
10    open(data_path, "r") do file
11        # read data each line
12        for line in eachline(file)
13            if line != ""
14                parts = split(line, ",")
15                push!(sepal_length, parse(Float64, parts[1]))
16                push!(sepal_width, parse(Float64, parts[2]))
17                push!(petal_length, parse(Float64, parts[3]))
18                push!(petal_width, parse(Float64, parts[4]))
19
20                if parts[5] == "Iris-setosa"
21                    push!(classes, 0)
22                elseif parts[5] == "Iris-versicolor"
23                    push!(classes, 1)
24                else
25                    push!(classes, 2)
26                end
27            end
28        end
29    end
30
31    # concat features
32    features = [sepal_length, sepal_width, petal_length, petal_width]
33    features = vcat(transpose.(features)...)
34    return features, classes
35 end
```

Selection deleted

```
1 # function change_class_to_num(y)
2 #     class = Dict("setosa"=> 0, "versicolor"=> 1, "virginica" => 2)
3 #     classnums = [class[item] for item in y]
4 #     return classnums
5 # end
```

train_test_split (generic function with 2 methods)

```
1 function train_test_split(X, y, test_ratio=0.33)
2     X = X'
3     n = size(X)[1]
4     idx = shuffle(1:n)
5     train_size = 1 - test_ratio
6     train_idx = view(idx, 1:floor(Int, train_size*n))
7     test_idx = view(idx, (floor(Int, train_size*n)+1):n)
8
9     X_train = X[train_idx,:]
10    X_test = X[test_idx,:]
11
12    y_train = y[train_idx]
13    y_test = y[test_idx]
14
15    return X_train, X_test, y_train, y_test
16 end
```

((100, 4), (50, 4), (100), (50))

```
1 begin
2     # Load features, and labels for Iris dataset
3     iris_features, iris_labels = iris_data_loader("data/iris.data")
4
5     #split dataset into training data and testing data
6     X_train, X_test, y_train, y_test = train_test_split(iris_features, iris_labels,
7     0.33)
Selection deleted
8     size(X_train), size(X_test), size(y_train), size(y_test)
9 end
```

1. Decision Tree: Iterative Dichotomiser 3 (ID3)

1.1 Information Gain

Expected value of the self-information (entropy):

$$Entropy = - \sum_i^n p_i \log_2(p_i)$$

The entropy function gets the smallest value if there is a value of p_i equal to 1, reaches the maximum value if all p_i are equal. These properties of the entropy function make it is an expression of the disorder, or randomness of a system, ...

entropy

Parameters:

- `counts`: shape (n_classes): list number of samples in each class
- `n_samples`: number of data samples

Returns

- `entropy`

```
1  """
2  Parameters:
3  - `counts`: shape (n_classes): list number of samples in each class
4  - `n_samples`: number of data samples
5
6  Returns
7  - entropy
8  """
9  function entropy(counts, n_samples)
10     #TODO
11     if n_samples == 0
12         return 0.0
13     end
14
15     probs = counts ./ n_samples
16     entropy = -sum(probs .* log2.(probs))
17
18     return entropy
19 end
```

Selection deleted

entropy_of_one_division

Returns entropy of a divided group of data

Data may have multiple classes

```
1  """
2  Returns entropy of a divided group of data
3
4  Data may have multiple classes
5  """
6  function entropy_of_one_division(division)
7
8      n_samples = size(division, 1)
9      n_classes = Set(division)
10
11     counts=[]
12     # println("DIVISION: ", division)
13     # println("SAMPLES:", n_samples)
14     # println("CLASSES:", n_classes)
15
16     # count samples in each class then store it to list counts
17     #TODO:
18     for sample in n_classes
19         push!(counts, count(i->(i==sample), division))
20     end
21     entropy_value = entropy(counts, n_samples)
22
23
24     return entropy_value, n_samples
25
26 end
```

Selection deleted

get_entropy

Returns entropy of a split

y_predict is the split decision by cutoff, True/False

```
1  """
2  Returns entropy of a split
3
4  y_predict is the split decision by cutoff, True/False
5  """
6  function get_entropy(y_predict, y)
7      n = size(y,1)
8      # println("y_predict:", y_predict)
9      # println("y:", y)
10     # left hand side entropy
11     entropy_true, n_true = entropy_of_one_division(y[y_predict])
12
13     # right hand side entropy
14     entropy_false, n_false = entropy_of_one_division(y[.!y_predict])
15
16     # println("TRUE:", entropy_true, n_true)
17     # println("FALSE:", entropy_false, n_false)
18     # overall entropy
19     #TODO s=?
20     Selection deleted = n_true / n
21     p_false = n_false / n
22     s = p_true * entropy_true + p_false * entropy_false
23
24     return s
25 end
```

The information gain of classifying information set D by attribute A:

$$Gain(A) = Entropy(D) - Entropy_A(D)$$

At each node in ID3, an attribute is chosen if its information gain is highest compare to others.

All attributes of the Iris set are represented by continuous values. Therefore we need to represent them with discrete values. The simple way is to use a cutoff threshold to separate values of the data on each attribute into two part: <cutoff and > = cutoff.

To find the best cutoff for an attribute, we replace cutoff with its values then compute the entropy, best cutoff achieved when value of entropy is smallest ($\arg \min Entropy_A(D)$).

1.2 Decision tree

dtfit

Parameters:

- X: training data
- y: label of training data

Returns

- node

node: each node represented by cutoff value and column index, value and children.

- cutoff value is threshold where you divide your attribute.
- column index is your data attribute index.
- value of node is mean value of label indexes, if a node is leaf all data samples will have same label.

Note that: we divide each attribute into 2 part => each node will have 2 children: left, right.

```
1  """
2  Parameters:
3  - X: training data
4  - y: label of training data
5
6  Returns
7  - node
8
9  node: each node represented by cutoff value and column index, value and children.
10 - cutoff value is threshold where you divide your attribute.
11 - column index is your data attribute index.
12 - value of node is mean value of label indexes, if a node is leaf all data samples
    will have same label.
13
14 Note that: we divide each attribute into 2 part => each node will have 2 children:
    left, right.
15 """
16 function dtfit(X, y, node=Dict(), depth=0)
17     #Stop conditions
18
19     #if all value of y are the same
20     if all(y==y[1])
21         return Dict("val"=>y[1])
22
23     else
24         # find one split given an information gain
25         col_idx, cutoff, entropy = find_best_split_of_all(X, y)
26
27         y_left = y[X[:,col_idx] .< cutoff]
28         y_right = y[X[:,col_idx] .>= cutoff]
29
30         node = Dict("index_col"=>col_idx,
31                     "cutoff"=>cutoff,
```



```

32         "val"=> mean(y),
33         "left"=> Any,
34         "right"=> Any)
35
36     left = dtfit(X[X[:,col_idx] .< cutoff, :], y_left, Dict(), depth+1)
37     right= dtfit(X[X[:,col_idx] .>= cutoff, :], y_right, Dict(), depth+1)
38
39     push!(node, "left" => left)
40     push!(node, "right" => right)
41
42     depth += 1
43 end
44 return node
45 end

```

find_best_split_of_all

Parameters:

- X: training data
- y: label of training data

Returns

- column index, cut-off value, and minimum entropy

Selection deleted

```

1  """
2  Parameters:
3  - X: training data
4  - y: label of training data
5
6  Returns
7  - column index, cut-off value, and minimum entropy
8  """
9  function find_best_split_of_all(X, y)
10     col_idx = nothing
11     min_entropy = 1
12     cutoff = nothing
13
14     for i in 1:size(X,2)
15         col_data = X[:,i]
16         entropy, cur_cutoff = find_best_split(col_data, y)
17
18         # best entropy
19         if entropy == 0
20             return i, cur_cutoff, entropy
21         elseif entropy <= min_entropy
22             min_entropy = entropy
23             col_idx = i
24             cutoff = cur_cutoff
25         end
26     end
27     return col_idx, cutoff, min_entropy
28 end

```

find_best_split

Parameters:

- col_data: data samples in column
- y: label of training data

Returns

- minimum entropy, and cut-off value

```
1  """
2  Parameters:
3  - col_data: data samples in column
4  - y: label of training data
5
6  Returns
7  - minimum entropy, and cut-off value
8  """
9  function find_best_split(col_data, y)
10     min_entropy = 10
11     cutoff = 0
12
13     Selection deleted through col_data find cutoff where entropy is minimum
14
15     for value in Set(col_data)
16         y_predict = col_data .< value
17         my_entropy = get_entropy(y_predict, y)
18
19         #TODO
20         #min entropy=?, cutoff=?
21         for entropy in my_entropy
22             if entropy <= min_entropy
23                 min_entropy = entropy
24                 cutoff = value
25             end
26         end
27
28     end
29     return min_entropy, cutoff
30 end
```

dtpredict (generic function with 1 method)

```
1  function dtpredict(tree, data)
2     pred = []
3     n_sample = size(data, 1)
4     for i in 1:n_sample
5         push!(pred, _dtpredict(tree, data[i,:]))
6     end
7     return pred
8 end
```

_dtpredict (generic function with 1 method)

```
1 function _dtpredict(tree, row)
2   cur_layer = tree
3   while haskey(cur_layer, "cutoff")
4     if row[cur_layer["index_col"]] < cur_layer["cutoff"]
5       cur_layer = cur_layer["left"]
6     else
7       cur_layer = cur_layer["right"]
8     end
9   end
10  if !haskey(cur_layer, "cutoff")
11    return get(cur_layer, "val", false)
12  end
13 end
```

1.3 Classification on Iris Dataset

tpfpntfn_cal (generic function with 2 methods)

```
1 function tpfpntfn_cal(y_test, y_pred, positive_class=1)
2   true_positives = 0
3   false_positives = 0
4   true_negatives = 0
5   false_negatives = 0
6
7   # Calculate true positives, false positives, false negatives, and true negatives
8   for (true_label, predicted_label) in zip(y_test, y_pred)
9     if true_label == positive_class && predicted_label == positive_class
10      true_positives += 1
11    elseif true_label != positive_class && predicted_label == positive_class
12      false_positives += 1
13    elseif true_label == positive_class && predicted_label != positive_class
14      false_negatives += 1
15    elseif true_label != positive_class && predicted_label != positive_class
16      true_negatives += 1
17    end
18  end
19
20  return true_positives, false_positives, true_negatives, false_negatives
21 end
```

tree =

Dict{"left" => Dict{"val" => 0}, "cutoff" => 1.0, "right" => Dict{"left" => Dict{"left"

1 tree = dtfit(X_train, y_train)

```

1 begin
2   pred = dtpredict(tree, X_test)
3
4   acc = 0
5   precision = 0
6   recall = 0
7   f1 = 0
8
9   for i ∈ [0, 1, 2]
10     # Calculate true positives, false positives, false negatives, and true
        negatives
11     true_positives, false_positives, true_negatives, false_negatives =
        tpfpntnfn_cal(y_test, pred, i)
12
13     # Calculate precision, recall, and F1-score
14     acc += (true_positives + true_negatives) / (true_positives +
        false_positives + true_negatives + false_negatives)
15     precision += true_positives / (true_positives + false_positives)
16     recall += true_positives / (true_positives + false_negatives)
17   end
18
19   acc = acc / 3
20   precision = precision / 3
21   recall = recall / 3
22   f1 = 2 * precision * recall / (precision + recall)
23   print(" acc: $acc\n precision: $precision\n recall: $recall\n f1_score: $f1\n")
24 end

```

```

acc: 0.9733333333333333
precision: 0.9649122807017544
recall: 0.9393939393939394
f1_score: 0.9519821328866556

```



2. Bayes Theorem

Bayes formulation $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

If \mathcal{B} is our data \mathcal{D} , \mathcal{A} and w are parameters we need to estimate:

$$\underbrace{P(w|\mathcal{D})}_{\text{Posterior}} = \frac{1}{\underbrace{P(\mathcal{D})}_{\text{Normalization}}} \underbrace{P(\mathcal{D}|w)P(w)}_{\text{Likelihood Prior}}$$

Naive Bayes

To make it simple, it is often assumed that the components of the \mathcal{D} random variable (or the features of the \mathcal{D} data) are independent with each other, if w is known. It mean:

$$P(\mathcal{D}|w) = \prod_{i=1}^d P(x_i|w)$$

- d: number of features

2.1. Probability Density Function

update (generic function with 1 method)

```
1 #update histogram for new data
2 function update(_hist, _mean, _std, data)
3     """
4     P(hypo/data)=P(data/hypo)*P(hypo)*(1/P(data))
5     """
6     hist = copy(_hist)
7     #P(hypo/data)=P(data/hypo)*P(hypo)*(1/P(data))
8
9     #Likelihood * Prior
10    #TODO
11    for hypo in keys(hist)
12        hist[hypo] = likelihood(_mean, _std, data, hypo) * hist[hypo]
13    end
14
15    #Normalization
16
17    #TODO: s=P(data)
18    #s=?
19    s = sum(values(hist))
20
21    for hypo in keys(hist)
22        hist[hypo] = hist[hypo]/s
23    end
24    return hist
25 end
```

maxHypo (generic function with 1 method)

```
1 function maxHypo(hist)
2     #find the hypothesis with maximum probability from dictionary hist
3     #TODO
4     max_prob = maximum(values(hist))
5     max_hypo = 0
6     for (key, val) in hist
7         if val == max_prob
8             max_hypo = key
9             break
10        end
11    end
12
13
14    return max_hypo
15 end
```

2.2 Classification on Iris Dataset

Gaussian Naive Bayes

- Naive Bayes can be extended to use on continuous data, most commonly by using a normal distribution (Gaussian distribution).
- This extension called Gaussian Naive Bayes. Other functions can be used to estimate data distribution, but Gauss (or the normal distribution) is the easiest to work with since we only need to estimate the mean and standard deviation from the training data.

Define Gauss function

$$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Gauss (generic function with 1 method)

```
1 function Gauss(std, mean, x)
2     #Compute the Gaussian probability distribution function for x
3     #TODO
4     return (1 ./ (std * sqrt( 2 * pi ))) * exp( -((x .- mean)^2) / (2 * std^2))
5 end
```

likelihood (generic function with 5 methods)

```
1 function likelihood(_mean=nothing, _std=nothing ,data=nothing, hypo=nothing)
2     """
3     Returns: res=P(data/hypo)
4     -----
5     Naive bayes:
6         Atributes are assumed to be conditionally independent given the class value.
7     """
8
9     std=_std[hypo]
10    mean=_mean[hypo]
11    res=1
12    #TODO
13    #res=res*P(x1/hypo)*P(x2/hypo)...
14    for i in 1:length(data)
15        res *= pdf(Normal(mean[i], std[i]), data[i])
16    end
17
18    return res
19 end
```

custom_std (generic function with 1 method)

```
1 # Custom standard deviation function
2 function custom_std(data, mean_val)
3     return sqrt(sum((data .- mean_val) .^ 2) / length(data))
4 end
```

custom_mean (generic function with 1 method)

```
1 # Custom mean function
2 function custom_mean(data)
3     return sum(data) / length(data)
4 end
```

gfit (generic function with 4 methods)

```
1 function gfit(X, y, _std=nothing, _mean=nothing, _hist=nothing)
2     """Parameters:
3     X: training data
4     y: labels of training data
5     """
6     n=size(X,1)
7     #number of iris species
8     #TODO
9     #n_species=???
10    n_species = length(Set(y))
11    # println(n)
12    # println(Set(y))
13
14    hist=Dict()
15    mean=Dict()
16    std=Dict()
17
18    #separate dataset into rows by class
19    for hypo in Set(y)
20        #rows have hypo label
21        #TODO rows=
22        rows = X[y .== hypo, :]
23        # println("hypo=", hypo, rows)
24        #histogram for each hypo
25        #TODO probability=?
26        probability = size(rows, 1) / n
27        hist[hypo]=probability
28        # println(hist)
29
30        #Each hypothesis represented by its mean and standard derivation
31        # """mean and standard derivation should be calculated for each column (or
each attribute)"""
32        #TODO mean[hypo]=?, std[hypo]=?
33
34        mean[hypo] = [custom_mean(rows[:, col]) for col in 1:size(rows, 2)]
35        std[hypo] = [custom_std(rows[:, col], mean[hypo][col]) for col in
361:size(rows, 2)]
37    end
38    _mean=mean
39    _std=std
40    _hist=hist
41    return _hist, _mean, _std
42 end
```


`_gpredict` (generic function with 2 methods)

```
1 function _gpredict(_hist, _mean, _std, data, plot=true)
2     """
3     Predict label for only 1 data sample
4     -----
5     Parameters:
6     data: data sample
7     plot: True: draw histogram after update new record
8     -----
9     return: label of data
10    """
11    hist = update(_hist, _mean, _std, data)
12    # println(hist)
13    if (plot == true)
14        plt = bar(collect(keys(hist)), collect(values(hist)))
15    end
16    return maxHypo(hist)
17 end
```

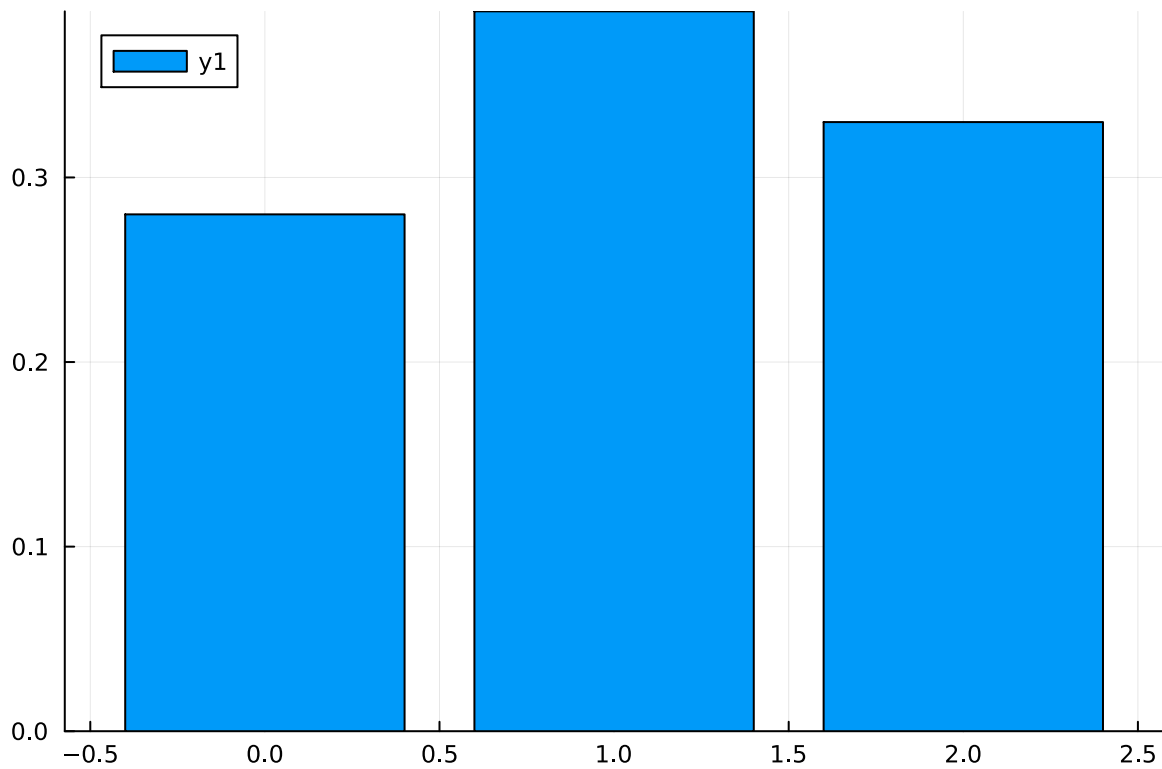
`plot_pdf` (generic function with 1 method)

```
1 function plot_pdf(_hist)
2     bar(collect(keys(_hist)), collect(values(_hist)))
3 end
```

`gpredict` (generic function with 1 method)

```
1 function gpredict(_hist, _mean, _std, data)
2     """Parameters:
3     Data: test data
4     -----
5     return labels of test data
6     """
7     pred=[]
8     n_sample = size(data, 1)
9     for i in 1:n_sample
10        push!(pred, _gpredict(_hist, _mean, _std, data[i,:]))
11    end
12    return pred
13 end
```

Show histogram of training data



```

1 begin
2   # print(X_train, y_train)
3   _hist, _mean, _std = gfit(X_train, y_train)
4   plt = plot_pdf(_hist)
5 end

```

Test with 1 data record

```

1 begin
2   #label of test_y[10]
3   print("Label of X_test[10]: ", y_test[20])
4
5   #update model and show histogram with X_test[10]:
6   print("\nOur histogram after update X_test[10]: ", _gpredict(_hist, _mean,
7     _std, X_test[20,:], true))
8 end

```

```

Label of X_test[10]: 1
Our histogram after update X_test[10]: 1

```



Evaluate your Gaussian Naive Bayes model

```

1 begin
2   _pred = gpredict(_hist, _mean, _std, X_test)
3
4   _acc = 0
5   _p = 0
6   _r = 0
7   _f1 = 0
8
9   #TODO: Self-define and calculate accuracy, precision, recall, and f1-score
10
11   for i ∈ [0, 1, 2]
12     # Calculate true positives, false positives, false negatives, and true
negatives
13     tp, fp, tn, fn = tpfptnfn_cal(y_test, _pred, i)
14
15     # Calculate precision, recall, and F1-score
16     _acc += (tp + tn) / (tp + fp + tn + fn)
17     _p += tp / (tp + fp)
18     _r += tp / (tp + fn)
19   end
20
21   _acc = _acc / 3
22   _p = _p / 3
23   _r = _r / 3
24   _f1 = 2 * _p * _r / (_p + _r)
25
26   print(" acc: $_acc\n precision: $_p\n recall: $_r\n f1_score: $_f1\n")
27 end

```

```

acc: 0.9866666666666667
precision: 0.9814814814814815
recall: 0.9696969696969697
f1_score: 0.9755536381938454

```



TODO: F1, Recall and Precision report

The Gaussian Naive Bayes model demonstrates strong performance across multiple evaluation metrics:

- **Accuracy (acc):** The model achieves an accuracy of approximately 98.67%, signifying the overall correctness of its predictions. This metric is calculated as the ratio of correctly predicted instances to the total number of instances.
- **Precision:** With a precision value of approximately 98.15%, the model excels in making accurate positive predictions. Precision measures the accuracy of the model's positive predictions, reflecting its ability to avoid false positives.
- **Recall:** The model achieves a recall, or sensitivity, of approximately 96.97%. This metric indicates the model's capability to capture a high percentage of actual positive instances. A higher recall value suggests that the model is effective at identifying positive cases.
- **F1-score:** The F1-score, computed as approximately 97.56%, represents the harmonic mean of precision and recall. This balanced metric considers both false positives and false negatives, providing a comprehensive assessment of the model's performance.

In summary, the Gaussian Naive Bayes model exhibits robust performance, demonstrating high accuracy, precision, recall, and F1-score across various aspects of its predictions.