# Lab04: Decision Tree and Naive Bayes

- Student ID:
- Student name:

**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 and pdf into `123456.zip` onto Moodle.

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

**Contents:**

- Decision Tree
- Naive Bayes

## Import library

☁ Code not executed in *Safe preview*

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

☁ Code not executed in *Safe preview*

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

## Load Iris dataset

```julia
# If you use Linux, use this function to download Iris dataset
function download_dataset(save_path::String="data")
    # setup directory
    mkpath(joinpath(dirname(@__FILE__), save_path))
    data_dir = joinpath(dirname(@__FILE__), save_path)
    download("https://archive.ics.uci.edu/static/public
    /53/iris.zip",joinpath(data_dir, "iris.zip"))
    iris_file = joinpath(data_dir, "iris.zip")
    cd(data_dir)
    run(`unzip $iris_file -d $data_dir`)
    rm(iris_file)
    cd("..")
end
```

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

```julia
# If you don't use Linux, I have no solution for you. Please makedir data, and goto
https://archive.ics.uci.edu/dataset/53/iris for dowloading
# Then, you can extract data to this dir by yourself.

# Structure:
# ├── data
# │   ├── bezdekIris.data
# │   ├── Index
# │   ├── iris.data
# │   └── iris.names
# ├── Lab4.jl
```

```julia
1  function iris_dataloader(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
```

```julia
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
```

```julia
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
```

Safe preview ⓘ

```
1 begin
2     # Load features, and labels for Iris dataset
3     iris_features, iris_labels = iris_dataloader("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,
      0.33)
7
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, ...

```
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
12 end
```

Safe preview ⓘ

```
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
13     # count samples in each class then store it to list counts
14     #TODO:
15
16 end
```

```
1  """
2  Returns entropy of a split
3
4  y_predict is the split decision by cutoff, True/Fasle
5  """
6  function get_entropy(y_predict, y)
7      n = size(y,1)
8      # left hand side entropy
9      entropy_true, n_true = entropy_of_one_division(y[y_predict])
10
11     # right hand side entropy
12     entropy_false, n_false = entropy_of_one_division(y[.~y_predict])
13
14     # overall entropy
15      #TODO s=?
16
17     return s
18 end
```

Safe preview ⓘ

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

$$Gain(A) = Entrophy(D) - Entrophy_A(D)$$

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 $(\arg\min Entrophy_A(D))$.

```
1  md"""
2  The information gain of classifying information set D by attribute A:
3
4  $$Gain(A)=Entrophy(D)-Entrophy_{A}(D)$$
5
6  At each node in ID3, an attribute is chosen if its information gain is highest
   compare to others.
7
8  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`.
9
10 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  $\left
   (\arg \min Entrophy_ {A} (D) \right)$.
11 """
```

## 1.2 Decision tree

```julia
"""
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 thresold 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.
"""
function dtfit(X, y, node=Dict(), depth=0)
    #Stop conditions

    #if all value of y are the same
    if all(y.==y[1])
        return Dict("val"=>y[1])

    else
        # find one split given an information gain
        col_idx, cutoff, entropy = find_best_split_of_all(X, y)

        y_left = y[X[:,col_idx] .< cutoff]
        y_right = y[X[:,col_idx] .>= cutoff]

        node = Dict("index_col"=>col_idx,
                    "cutoff"=>cutoff,
                    "val"=> mean(y),
                    "left"=> Any,
                    "right"=> Any)

        left = dtfit(X[X[:,col_idx] .< cutoff, :], y_left, Dict(), depth+1)
        right= dtfit(X[X[:,col_idx] .>= cutoff, :], y_right, Dict(), depth+1)

        push!(node, "left" => left)
        push!(node, "right" => right)

        depth += 1
    end
    return node
end
```

Safe preview ⓘ

```
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
```

```
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     #Loop 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
22
23     end
24     return min_entropy, cutoff
25 end
```

Safe preview ⓘ

```
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
```

```
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

```
1  function tpfptnfn_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
```

```
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 =
       tpfptnfn_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
```

# 2. Bayes Theorem

Bayes formulation $P\left(A|B\right) = \dfrac{P\left(B|A\right)P\left(A\right)}{P\left(B\right)}$

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

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

## Naive Bayes

To make it simple, it is often assumed that the components of the $D$ 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

```
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
12     #Normalization
13
14     #TODO: s=P(data)
15     #s=?
16     s = 0
17
18
19     for hypo in keys(hist)
20         hist[hypo] = hist[hypo]/s
21     end
22     return hist
23  end
```

```
1  function maxHypo(hist)
2      #find the hypothesis with maximum probability from hist
3      #TODO
4
5      return max_hypo
6  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

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

Safe preview ⓘ

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

```
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
15     return res
16 end
```

```
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
12     hist=Dict()
13     mean=Dict()
14     std=Dict()
15
16     #separate  dataset into rows by class
17     for hypo in Set(y)
18         #rows have hypo label
19         #TODO rows=
20
21         #histogram for each hypo
22         #TODO probability=?
23
24         hist[hypo]=probability
25
26         #Each hypothesis represented by its mean and standard derivation
27         """mean and standard derivation should be calculated for each column (or each
   attribute)"""
28         #TODO mean[hypo]=?, std[hypo]=?
29
30     end
31     _mean=mean
32     _std=std
33     _hist=hist
34     return _hist, _mean, _std
```

```
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     if (plot == true)
13         plt = bar(collect(keys(hist)), collect(values(hist)))
14     end
15     return maxHypo(hist)
16 end
```

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

```
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      _hist, _mean, _std = gfit(X_train, y_train)
3      plt = plot_pdf(_hist)
4  end
```

## Test wih 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, _std,
       X_test[20,:], true))
7
8  end
```

## 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     print(" acc: $_acc\n precision: $_p\n recall: $_r\n f1_score: $_f1\n")
12 end
```

**TODO**: F1, Recall and Precision report