



BIRZEIT UNIVERSITY

**Linux shell project**

**ENCS3130**

Student's name and id:

Amr Halahla 1200902.

Manar Shawahni 1201086.

Dr.Mohammad Jubran.

Assistent: Mahmood Abu Awwad.

- This is the main script code that contains the menu and the case statement, also it contains some logic code which do the process of the first entered option by the user, where it should be either 'r' or 'e', and if not, it will print an error message then return to the main menu again.

```

1 # Set the initial value of the verified flag to false
2 verified=false
3
4 # Start the main loop
5 while [ 1 ]; do
6     printf "\nr) read a dataset from a file\np) print the names of the features\nl) encode a feature using label encoding\nno) encode a feature using
one-hot encoding\nnm) apply MinMax scalling\ns) save the processed dataset\nne) exit\n"
7     # ask user for the main option
8     echo "Please enter your option"
9     read option
10
11     # Check if the option is 'r' or 'e'
12     if [ "$option" = "r" -o "$option" = "e" ]; then
13         # If the option is 'r' or 'e', set the verified flag to true
14         verified=true
15     fi
16
17     # If the verified flag is false, print a message and continue to the next iteration of the loop
18     # because the user should enter either r or e for the first time he run the script.
19     if [ "$verified" = false ]; then
20         echo "You must first read a dataset from a file"
21         continue
22     fi
23
24     # Use a case statement to handle each menu option
25     case "$option"
26     in
27         "r" ) # the case where the user should enter the name of the dataset source.
28             echo "Please input the name of the dataset file"
29             read dataset
30             # check if the file exists.
31             if [ ! -e "$dataset" ]
32             then
33                 echo "$dataset Not Found!"
34                 continue
35             fi
36             # if the name entered exist, then start the format checking process.
37             ./labels.sh "$dataset";;
38         "p" ) ./case3.sh;;
39         "l" ) ./case4.sh "$dataset";;
40         "o" ) ./case5.sh "$dataset";;
41         "m" ) ./case6.sh "$dataset";;
42         "s" )
43             saved=false
44             echo "Please input the name of the file to save the processed dataset"
45             read filename
46             cat newdataset.txt > "$filename"
47             saved=true
48             echo "Data Saved To The File $filename"
49             cat "$filename"
50             ;;
51         "e" )
52             if [ "$saved" ];then # check if the processed data saved or not.
53                 echo "Are you sure you want to exist"
54                 read answer
55                 if [ "$answer" = "yes" ];then
56                     exit 10
57                 else
58                     continue
59                 fi
60             else
61                 echo "The processed dataset is not saved. Are you sure you want to exist ?"
62                 read answer
63                 if [ "$answer" = "yes" ];then
64                     exit 10
65                 else
66                     continue
67                 fi
68             fi
69             ;;
70     esac
71 done

```

- This capture below shows the result if the user started the program with wrong option (neither 'r' nor 'e'):

```
amr@amr-VirtualBox: ~/Project$ ./menu.sh
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
l
You must first read a dataset from a file

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
```

- This is the first dataset we used:

```
1 id;age;gender;height;weight;active;blood;smoke;governorate;
2 1;20;male;171;55;yes;A+;no;Hebron;
3 2;20;female;158;55;yes;A-;yes;Jenin;
4 3;30;male;160;60;no;B+;no;Jerusalem;
5 4;19;female;162;55;yes;AB+;no;Ramallah;
6 5;13;female;142;37;no;O-;no;Jenin;
```

- This is the second dataset we used:

```
1 age;sex;bmi;children;smoker;region;charges;
2 19;female;27;0;yes;southwest;16884;
3 18;male;33;1;no;southeast;1725;
4 28;male;33;3;no;southeast;4449;
5 33;male;22;0;no;northwest;21984;
6 32;male;28;e;no;northwest;3866;
7 31;female;25;e;no;southeast;3756;
8 46;female;33;1;yes;southeast;8240;
```

## DATASET 1:

---

- Option 'r':

- The code:

```
"r" ) # the case where the user should enter the name of the dataset source.
echo "Please input the name of the dataset file"
read dataset
# check if the file exists.
if [ ! -e "$dataset" ]
then
    echo "$dataset Not Found!"
    continue
fi
# if the name entered exist, then start the format checking process.
./labels.sh "$dataset";;
```

- In case 'r', the user should enter the dataset file name.
- The logic of the code firstly will check if the file exists or not, if not then it will print an error file not found.
- If file exists, then the script 'labels.sh' will be executed.

### 'labels.sh' script:

```
case3.sh  x  case4.sh  x  encoded.txt  x  labels.sh  x  menu.sh  x  scale.txt
1 if [ $# -lt 1 ]; then # check if the user entered a parameter that represent the dataset file name.
2   echo "Error: No dataset file provided"
3   exit 1
4 fi
5 data=$1 # extract the register value (which is the name of the file)
6 line1=$(sed -n "1p" "$data" | tr ';' '\n') # split the first line by the delimiter ';' so we can then count the number of features.
7 feature_count=$(echo "$line1" | wc -l) # count the number of features.
8 # split the second line in order to be able to check if the number of feature's values matches the number of features
9 line2=$(sed -n "2p" "$data" | tr ';' '\n')
10 values_count=$(echo "$line2" | wc -l) # count feature's values
11 #echo "$values_count"
12 if [ "$feature_count" != "$values_count" ]; then # check if they are equal, if not then print an error data formating message
13   echo "The format of the data in the dataset file is wrong"
14   exit 11
15 else
16   echo "Format is Clean"
17 fi
18 # save a copy of the main dataset file into another file to be processed.
19 cat "$data" > newdataset.txt
20 > sc.txt # file to store the name of the features that have been encoded.
21 > scale.txt # to store the features that have been encoded attatched with their codes.
22
```

- The main functionality of this script is to check the data formatting if it is correct or not.
  - The formatting is correct if the number of features in the dataset is equal to the number of values attached with it.
  - We made a copy of the dataset in order to perform all of the processes on it.
- Here the output for different tests for case 'r':

```
amr@amr-VirtualBox: ~/Project$ ./menu.sh
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
r
Please input the name of the dataset file
wrong dataset
wrong dataset Not Found!
```

The output when we entered a wrong dataset file

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
r
Please input the name of the dataset file
dataset.txt
Format is Clean
```

Here the output of the correct dataset file name, with correct formatting.

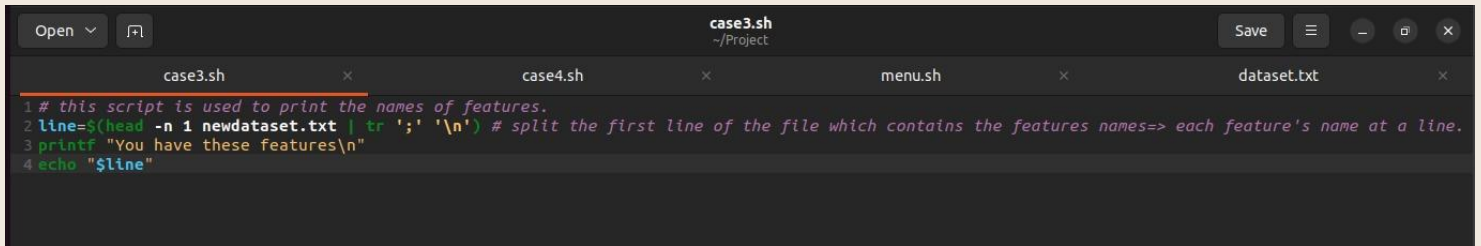
```
Please enter your option
r
Please input the name of the dataset file
dataset.txt
The format of the data in the dataset file is wrong
```

we put a wrong format in dataset file



- **Option ‘p’:**

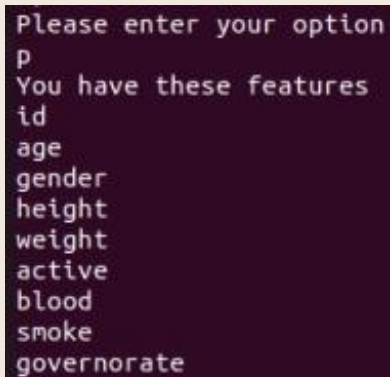
- The code:



```
1 # this script is used to print the names of features.
2 line=$(head -n 1 newdataset.txt | tr ';' '\n') # split the first line of the file which contains the features names=> each feature's name at a line.
3 printf "You have these features\n"
4 echo "$line"
```

- This script prints the name of the features that exists in the dataset file.

- **The output for this case:**



```
Please enter your option
p
You have these features
id
age
gender
height
weight
active
blood
smoke
governorate
```

- Option 'I':
  - The code:

```

1 > features.txt # file that will be used to store the name of categorical features.
2 sed -i '/^$/d' features.txt
3 dataset="$1" # get the name of the original dataset file
4 # initializing variables
5 counter=1
6 fet_index=100
7 # Count the number of lines in the dataset file
8 numoflines=$(cat newdataset.txt | wc -l)
9
10 # Ask the user to input the name of a categorical feature for label encoding
11 echo "Please input the name of the categorical feature for label encoding"
12 read feature_name
13
14 # Extract the first line of the dataset file and store it in a file called line.txt
15 line=$(sed -n "1p" newdataset.txt)
16 echo "$line" > line.txt
17
18 # Count the number of fields(features) in line.txt
19 numf=$(cat line.txt | tr ';' '\n' | wc -l)
20 # Loop through each field in line.txt
21 while [ "$numf" -gt 1 ]
22 do
23   # Extract the current field from the first line of the dataset file
24   fet=$(echo "$line" | cut -d';' -f"$counter")
25   if [ "$fet" = "$feature_name" ];then # check if the entered feature's name matches the feature at the current index.
26     fet_index=$((counter)) # if yes, store the index (at which delimiter) of the target feature.
27   fi
28   # Extract the current field from the second line of the dataset file
29   value=$(sed -n "2p" newdataset.txt | cut -d';' -f"$counter") # the value of the current feature
30
31   # Check if the value of the current field is numeric or non-numeric
32   if [[ ! "$value" =~ ^[0-9]+$ ]]; then # if its categorical, append the feature name into a file.
33     echo "$fet" >> features.txt
34   fi
35
36   numf=$((numf-1))
37   counter=$((counter+1))
38 done
39
40 if grep -qw "$feature_name" features.txt; then # verify that the entered feature is categorical or not (in order to be encoded).
41   echo "Feature $feature_name exist in the file"
42 else
43   echo "The name of categorical feature is wrong"
44   exit 3
45 fi
46
47 count_fet=$(cat features.txt | wc -l)
48 > fet_value.txt # a file in which the possible values for a specific feature are stored
49 while [ "$numoflines" -gt 1 ];do # loop through all lines in the dataset to extract the values
50   fet_value=$(sed -n "$numoflines p" newdataset.txt | cut -d';' -f"$fet_index") # the value of the feature at the current data line.
51   numoflines=$((numoflines-1))
52   echo "$fet_value" >> fet_value.txt # append the value into the file.
53 done
54
55 sort fet_value.txt | uniq > temp.txt # remove the duplicated values
56 > codes.txt # a file in which the codes for the unique possible values are stored
57 cat temp.txt > fet_value.txt
58 sed -i '/^$/d' fet_value.txt # remove blank lines from the file
59 sed -i '/^$/d' codes.txt # remove blank lines from the file
60 #printf "Unique file :\n"
61 #cat fet_value.txt
62 counter=$(cat fet_value.txt | wc -l) # the number of unique values for the current specified feature.
63 code=0
64 while [ "$counter" -gt 0 ];do # loop through all values to attach it with their codes.
65   value=$(sed -n "$counter p" fet_value.txt)
66   echo "$code" >> codes.txt # store the unique codes in a file.
67   code=$((code+1))
68   counter=$((counter-1))
69 done
70 paste -d'=' fet_value.txt codes.txt > encoded.txt # concat each value with its own code in one file (value = code).
71 #cat encoded.txt
72 sed -i '/^$/d' newdataset.txt
73 numoflines=$(cat newdataset.txt | wc -l)
74 echo "" > encodedfile.txt # a file in which the encoded data will be stored temporarily
75 sed -i '/^$/d' encodedfile.txt
76 head -n 1 newdataset.txt > encodedfile.txt # the first line will not be encoded in label-encoding.
77 echo "" > newline.txt
78 sed -i '/^$/d' newline.txt
79 counter=2
80 > scale_line.txt # This file will be used with the option "m," and it stores every feature that has been encoded with its codes.
81 echo "$feature_name" >> scale_line.txt # set feature name has been encoded.
82
83 echo "The values of the features' codes are encoded as follows :"
84 cat encoded.txt

```

```

85 echo "-----"
86 while [ "$counter" -le "$numoflines" ];do # loop through each line in the dataset.
87     newline=$(sed -n "$counter p" newdataset.txt > newline.txt) # extract current line from the dataset
88     cat newline.txt | tr ';' '\n' > tmp.txt && mv tmp.txt newline.txt # splitting
89     sed -i '/^$/d' newline.txt # remove blank lines
90     linecount=$(cat encoded.txt | wc -l) # get the number of possible encoded values.
91     while [ "$linecount" -gt 0 ];do # loop through all possible encoded values attached with their own codes.
92         var=$(sed -n "$linecount p" encoded.txt | cut -d '=' -f1) # get the name of the value
93         code=$(sed -n "$linecount p" encoded.txt | cut -d '=' -f2) # get the code of the current value.
94         base_value=$(sed -n "$fet_index p" newline.txt) # extract the pointed field(by the 'fet_index' variable) from the line of dataset.
95         sed -i "$fet_index s/\b'$var'/'$code'/" newline.txt # replace the value by its specified code.
96         linecount=$((linecount-1))
97         if [ "$base_value" = "$var" ];then # check if this is the wanted value.
98             echo "$code" >> scale_line.txt # append the code so it will be used later.
99         fi
100     done
101     cat newline.txt | tr '\n' ';' > tmp.txt && mv tmp.txt newline.txt # rearrange the encoded data line as the original format.
102     cat newline.txt >> encodedfile.txt # store each encoded line to the encoded data file.
103     printf "\n" >> encodedfile.txt
104     counter=$((counter+1))
105 done
106 cat encodedfile.txt > newdataset.txt # store the encoded data to the new dataset file.
107 printf "The dataset after label encoding for the feature $feature_name : \n"
108 cat newdataset.txt # print the final result after encode the current feature.
109 cat scale_line.txt | tr '\n' ';' > tm.txt && mv tm.txt scale_line.txt # rearrange the file in a useful form.
110 sed -i '/^$/d' scale_line.txt
111 cat scale_line.txt >> scale.txt # append the encoded feature with its codes to a file.
112 printf "\n" >> scale.txt
113 echo "$feature_name" >> sc.txt # to assign that this feature has been encoded correctly.

```

- This script performs the label encoding process.
- The label encoding can only be performed on the categorical features.
- The algorithm we used can be summarized by some mainly steps:
  - Firstly, we stored the categorical features at a file.
  - Then, we extracted the unique possible values for the wanted feature.
  - We assign a code for each value.
  - The encoding process done on each line of the dataset separately.
  - At the end, we update the encoded feature in the dataset.
- The label encoding process can be performed on different features at different times at the same dataset.



■ Some results attached for different cases:

```
Please enter your option
l
Please input the name of the categorical feature for label encoding
blood
Feature blood exist in the file
The values of the features' codes are encoded as follows :
A+=0
A-=1
AB+=2
B+=3
O-=4
-----
The dataset after label encoding for the feature blood :
id;age;gender;height;weight;active;blood;smoke;governorate;
1;20;male;171;55;yes;0;no;Hebron;
2;20;female;158;55;yes;1;yes;Jenin;
3;30;male;160;60;no;3;no;Jerusalem;
4;19;female;162;55;yes;2;no;Ramallah;
5;13;female;142;37;no;4;no;Jenin;
```

Label encoding for the  
categorical feature ' blood '

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
```

Label encoding for the  
categorical feature ' smoke '

```
l
Please input the name of the categorical feature for label encoding
smoke
Feature smoke exist in the file
The values of the features' codes are encoded as follows :
no=0
yes=1
-----
The dataset after label encoding for the feature smoke :
id;age;gender;height;weight;active;blood;smoke;governorate;
1;20;male;171;55;yes;0;0;Hebron;
2;20;female;158;55;yes;1;1;Jenin;
3;30;male;160;60;no;3;0;Jerusalem;
4;19;female;162;55;yes;2;0;Ramallah;
5;13;female;142;37;no;4;0;Jenin;
```

The dataset after label encoding  
for two features ( blood, smoke ).

- The picture below shows two different cases at same time:
  - When the feature entered is not categorical.
  - When the user request to encode the feature more than one time.

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
l
Please input the name of the categorical feature for label encoding
smoke
The name of categorical feature is wrong
```

The feature can be encoded only once, bcz since its encoded it  
will be regarded as a not categorical feature

- Option 'o':  
- The code:

```

1 > features.txt # a file in which the categorical features will be stored
2 sed -i '/^$/d' features.txt
3 dataset="s1"
4 sed -i '/^$/d' newdataset.txt
5 counter=1
6 fet_index=100
7 numoflines=$(cat newdataset.txt | wc -l)
8 # Ask the user to input the name of a categorical feature for label encoding
9 echo "Please input the name of the categorical feature for one hot encoding"
10 read feature_name
11 # Extract the first line of the dataset file and store it in a file called line.txt
12 line=$(sed -n "1p" newdataset.txt) # Extract the first line of the dataset.
13 echo "$line" > line.txt
14 numf=$(cat line.txt | tr ';' '\n' | wc -l) # count number of features in the dataset.
15 # Loop through each field in line.txt
16 while [ "$numf" -gt 1 ]
17 do
18     # Extract the current field from the first line of the dataset file
19     fet=$(echo "$line" | cut -d';' -f"$counter")
20     if [ "$fet" = "$feature_name" ];then
21         fet_index=$((counter))
22     fi
23     # Extract the current field from the second line of the dataset file
24     value=$(sed -n "2p" newdataset.txt | cut -d';' -f"$counter")
25
26     # Check if the value of the current field is numeric or non-numeric
27     if [[ ! "$value" =~ ^[0-9]+$ ]]; then
28         echo "$fet" >> features.txt # if the feature is categorical, append its name into the file.
29     fi
30
31     numf=$((numf-1))
32     counter=$((counter+1))
33 done
34
35 if grep -qw "$feature_name" features.txt; then # verify if the entered feature is one of the categorical.
36     echo "Feature $feature_name exist in the file"
37 else
38     echo "The name of categorical feature is wrong"
39     exit 3
40 fi
41
42 count_fet=$(cat features.txt | wc -l) # count number of categorical features.
43 > fet_value.txt
44 while [ "$numoflines" -gt 1 ];do
45     fet_value=$(sed -n "$numoflines p" newdataset.txt | cut -d';' -f"$fet_index")
46     numoflines=$((numoflines-1))
47     echo "$fet_value" >> fet_value.txt # store the values of the current feature into a file.
48 done
49 #-----
50 echo "$feature_name" > scale_line.txt # append the feature to the file at where the encoded features are stored.
51
52 sort fet_value.txt | uniq > temp.txt && mv temp.txt fet_value.txt # sort and remove duplicated values.
53 cat fet_value.txt > copy.txt
54 # arrange the feature values as the wanted formatting to replace the feature name in the dataset
55 cat fet_value.txt | tr '\n' ';' > tmp.txt && mv tmp.txt fet_value.txt
56 echo "" > encodedfile.txt # Each encoded line will be appended to this file.
57 sed -i '/^$/d' encodedfile.txt
58 values=$(cat fet_value.txt) # the encoding of the feature name
59 printf "\n"
60 sed -i '/^$/d' copy.txt
61 sed -i "s/\b$feature_name/;$values/" newdataset.txt # replace the the feature name with its all possible unique values arranged.
62 #-----
63 numoflines=$(cat newdataset.txt | wc -l)
64 while [ "$numoflines" -gt 1 ];do # loop through the dataset.
65     sed -n "$numoflines p" newdataset.txt > newline.txt # Extract each line from the dataset in order to be encoded.
66     cat newline.txt | tr ';' '\n' > tmp.txt && mv tmp.txt newline.txt # split the current line into multiple lines.
67     sed -i '/^$/d' newline.txt
68     linecount=$(cat copy.txt | wc -l) # count number of possible values for the current feature.
69     parameter=$(sed -n "$fet_index p" newline.txt) # the target value that must be encoding.(the value to be replaced by 1)
70     > tempencoded.txt
71     counter=1
72     while [ "$counter" -le "$linecount" ];do # loop through all possible values to determine which one is the target.
73         var=$(sed -n "$counter p" copy.txt) # extract a randomly value from the possible values.
74         if [ "$parameter" = "$var" ];then # if we reached the target value, replace it by value 1 in the code.
75             echo "1" >> tempencoded.txt # append 1 to the result of one-hot code.
76             echo "1" >> scale_line.txt # store the code beside the name of the feature so it will be used later(case'n').
77         else
78             echo "0" >> tempencoded.txt # if this is not the target, append 0 to the one-hot code.
79             echo "0" >> scale_line.txt
80         fi
81         counter=$((counter+1))
82     done
83     cat tempencoded.txt | tr '\n' ';' > tmp.txt && mv tmp.txt tempencoded.txt # rearrange the code to be in the form (0;0;1;0 ...)
84     code=$(cat tempencoded.txt) # store the final result of the code for the current value
85     #printf "code = %s\n" "$code"
86     sed -i "$numoflines s/\b$parameter/;$code/" newdataset.txt # replace the value by its own code in the dataset.
87     numoflines=$((numoflines-1))
88 done
89 cat scale_line.txt | tr '\n' ';' > tm.txt && mv tm.txt scale_line.txt
90 sed -i '/^$/d' scale_line.txt
91 cat scale_line.txt >> scale.txt
92 printf "\n" >> scale.txt
93 echo "$feature_name distinct values are :|"
94 cat copy.txt
95 cat newdataset.txt # show the dataset after encoding.
96 echo "$feature_name" >> sc.txt # append the feature name to the file where the encoded features are stored.

```

- This script performs the one-hot encoding process.
- The one-hot encoding can only be performed on the categorical features.
- The algorithm we used can be summarized by some mainly steps:
  - Firstly, we stored the categorical features at a file.
  - Then, we extracted the unique possible values for the wanted feature.
  - We replace the feature name by its unique possible values.
  - We assign '1' to the correct value, otherwise we assign '0' to the rest values.
  - The encoding process done on each line of the dataset separately.
  - At the end, we update the encoded feature in the dataset.
- The one-hot encoding process can be performed on different features at different times at the same dataset.

■ Some results attached for different cases:

```
Please enter your option
o
Please input the name of the categorical feature for one hot encoding
governorate
Feature governorate exist in the file

governorate distinct values are :
Hebron
Jenin
Jerusalem
Ramallah
id;age;gender;height;weight;active;blood;smoke;Hebron;Jenin;Jerusalem;Ramallah;
1;20;male;171;55;yes;0;0;1;0;0;0;0;
2;20;female;158;55;yes;1;1;0;1;0;0;0;
3;30;male;160;60;no;3;0;0;0;0;1;0;
4;19;female;162;55;yes;2;0;0;0;0;0;1;
5;13;female;142;37;no;4;0;0;1;0;0;0;
```

One-hot encoding for the feature governorate

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MLNMax scaling
s) save the processed dataset
e) exit
Please enter your option
o
Please input the name of the categorical feature for one hot encoding
gedner
The name of categorical feature is wrong
```

Try to encode a feature that not exist

```
Please input the name of the categorical feature for one hot encoding
gender
Feature gender exist in the file

gender distinct values are :
female
male
id;age;female;male;height;weight;active;blood;smoke;Hebron;Jenin;Jerusalem;Ramallah;
1;20;0;1;171;55;yes;0;0;1;0;0;0;0;
2;20;1;0;158;55;yes;1;1;0;1;0;0;0;
3;30;0;1;160;60;no;3;0;0;0;0;1;0;
4;19;1;0;162;55;yes;2;0;0;0;0;0;1;
5;13;1;0;142;37;no;4;0;0;1;0;0;0;
```

One-hot encoding for the feature 'gender'

The dataset after encoding some features using both encoding methods

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scaling
s) save the processed dataset
e) exit
Please enter your option
o
Please input the name of the categorical feature for one hot encoding
governorate
The name of categorical feature is wrong
```

Can not encode a feature which has been already encoded

```
Please enter your option
p
You have these features
id
age
female
male
height
weight
active
blood
smoke
Hebron
Jenin
Jerusalem
Ramallah
```

Finally, the name of features in the dataset after one-hot encoding.



- Option 'm':
  - The code:

```

1 > features.txt
2 > numeric_features.txt
3 sed -i '/^$/d' features.txt
4 dataset='s1'
5 sed -i '/^$/d' "$dataset"
6 counter=1
7 fet_index=100
8 numoflines=$(cat "$dataset" | wc -l)
9 echo "Please input the name of the feature to be scaled"
10 read feature_name
11
12 # The code below is used to check whether the feature entered is categorical in its original case before determining
13 # whether it has been encoded or not.
14 line=$(sed -n "1p" "$dataset") # extract first line from the dataset
15 echo "$line" > line.txt
16 numf=$(cat line.txt | tr ';' '\n' | wc -l) # count number of features in order to able looping through it.
17 # Loop through each field in line.txt
18 while [ "$numf" -gt 1 ]
19 do
20   fet=$(echo "$line" | cut -d';' -f"$counter") # extract feature name.
21   if [ "$fet" = "$feature_name" ];then # verify if its the wanted feature.
22     fet_index=$((counter))
23     fi
24     value=$(sed -n "2p" "$dataset" | cut -d';' -f"$counter")
25
26     # Check if the value of the current field is numeric or non-numeric
27     if [[ ! "$value" =~ ^[0-9]+$ ]]; then
28       echo "$fet" >> features.txt
29     else
30       echo "$fet" >> numeric_features.txt
31     fi
32     numf=$((numf-1))
33     counter=$((counter+1))
34 done
35 #-----
36 # loop to determine if the feature entered isn't categorical, also to save its index.
37 if ! grep -qw "$feature_name" features.txt;then
38   line1=$(head -n 1 newdataset.txt) # Obtain the first line of the dataset in order to determine which numerical feature is the target.
39   echo "$line1" | tr ';' '\n' > line1.txt
40   line1=$(cat line1.txt | wc -l)
41   counter=1
42   while [ "$counter" -le "$line1" ];do
43     name=$(sed -n "$counter p" line1.txt)
44     if [ "$name" = "$feature_name" ];then
45       numeric_index=$((counter)) # save the index of the numeric feature index.
46     fi
47     counter=$((counter+1))
48   done
49 fi
50
51 #-----
52 # If the feature is categorical, we must determine whether it has been encoded or not.
53 if grep -qw "$feature_name" features.txt; then
54   echo "Feature $feature_name exist in the file"
55   if grep -qw "$feature_name" sc.txt; then # grep through the file containing the names of encoded features.
56     echo "The Feature $feature_name Has Been Encoded"
57     # Extract only the codes (without the feature name) that are attached after the feature name.
58     line=$(grep "$feature_name" scale.txt | cut -d';' -f2-)
59     echo "$line" | tr ';' '\n' > scaledcode.txt # Save the codes in a file, one for each line.
60     sed -i '/^$/d' scaledcode.txt # remove blanked lines.
61     numofcodes=$(cat scaledcode.txt | wc -l)
62     #cat scaledcode.txt
63     #echo "$numofcodes"
64     counter=1
65     # Sort the codes to make it easier to find the minimum and maximum values.
66     # Extract the min and max numbers
67     min=$(sort -n scaledcode.txt | head -n 1)
68     max=$(sort -n scaledcode.txt | tail -n 1)
69     > vector.txt
70     printf "[" >> vector.txt
71     printf "max value = %d \nmin value = %d\n" $max $min # print the min and max values.
72     while [ "$counter" -le "$numofcodes" ];do # loop through all codes to be encoded
73       xl=$(sed -n "$counter p" scaledcode.txt) # extract each code separately
74       upper=$((xl-min)) #calculate the numerator.
75       lower=$((max-min)) # calculat denominator.
76       xiscale=$(echo "scale=1;$upper/$lower" | bc) # calculate the float result of dividing upper by lower.
77       printf "%.2f" $xiscale >> vector.txt # insert the scaled code into the vector
78       if [ "$counter" -ne "$numofcodes" ]; then
79         printf ", " >> vector.txt
80       fi
81       printf "x$counter scaled into %.2f \n" $xiscale
82       counter=$((counter+1))
83     done
84     printf "]" >> vector.txt
85     #print the scaled vector
86     printf "x-sacled = "
87     cat vector.txt
88   else

```



```

89     echo "this feature $feature_name is categorical feature and must be encoded first"
90     exit 3
91 fi
92 # case of numeric features.
93 elif grep -qw "$feature_name" numeric_features.txt; then
94     > numeric_codes.txt # file at where the number feature's values are stored.
95     echo "The Feature is Numeric"
96     #echo "$numeric_index"
97     counter=2
98     numoflines=$(cat newdataset.txt | wc -l) # loop through the dataset to extract the values.
99     while [ "$counter" -le "$numoflines" ];do
100         value=$(sed -n "$counter p" newdataset.txt | cut -d';' -f$numeric_index)
101         echo "$value" >> numeric_codes.txt # store the current value
102         counter=$((counter+1))
103     done
104     counter=1
105     # Extract the min and max numbers
106     min=$(sort -n numeric_codes.txt | head -n 1)
107     max=$(sort -n numeric_codes.txt | tail -n 1)
108     printf "max value = %d \nmin value = %d\n" $max $min
109     > vector.txt
110     printf "[" >> vector.txt
111     numoflines=$(cat numeric_codes.txt | wc -l)
112     while [ "$counter" -le "$numoflines" ];do # loop through all codes.
113         xi=$(sed -n "$counter p" numeric_codes.txt) # extract the current code
114         upper=$((xi-min)) # calculate the numerator
115         lower=$((max-min)) # calculate the denominator
116         xiscale=$(echo "scale=1; $upper/$lower" | bc) # calculate the scale of the code.
117         printf "%.2f" $xiscale >> vector.txt # insert the scaled value to the vector.
118         if [ "$counter" -ne "$numoflines" ]; then
119             printf ", " >> vector.txt
120         fi
121         printf "x$counter scaled into %.2f\n" $xiscale
122         counter=$((counter+1))
123     done
124     printf "]" >> vector.txt
125     printf "x-sacled = "
126     cat vector.txt # print the scaled vector.
127 else
128     echo "the feature does not exist in the dataset."
129 fi

```

- First step, we checked if the feature is categorical or numeric.
- If the feature is numeric, we start the process of minmax scaling.
- If the feature is categorical, we should first check if it has been encoded.
  - o If yes, then calculate the scaling vector of the feature.
  - o If not, print that the feature must be encoded first.
- Print the scaling vector.

- Some results attached for different cases:

```
blood
Feature blood exist in the file
The Feature blood Has Been Encoded
max value = 4
min value = 0
x1 scaled into 0.00
x2 scaled into 0.20
x3 scaled into 0.70
x4 scaled into 0.50
x5 scaled into 1.00
x-scaled = [0.00, 0.20, 0.70, 0.50, 1.00].
```

Apply MinMax scaling on the feature blood, and print the scaled vector.

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scaling
s) save the processed dataset
e) exit
```

```
Please enter your option
m
Please input the name of the feature to be scaled
smoke
Feature smoke exist in the file
The Feature smoke Has Been Encoded
max value = 1
min value = 0
x1 scaled into 0.00
x2 scaled into 1.00
x3 scaled into 0.00
x4 scaled into 0.00
x5 scaled into 0.00
x-scaled = [0.00, 1.00, 0.00, 0.00, 0.00].
```

MinMax scaling for the feature 'smoke', which has been encoded before.

```
Please enter your option
m
Please input the name of the feature to be scaled
id
The Feature is Numeric
max value = 5
min value = 1
x1 scaled into 0.00
x2 scaled into 0.20
x3 scaled into 0.50
x4 scaled into 0.70
x5 scaled into 1.00
x-scaled = [0.00, 0.20, 0.50, 0.70, 1.00].
```

MinMax scaling on a numeric feature 'id'.

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scaling
s) save the processed dataset
e) exit
```

Try to get a MinMax scaling of a categorical feature that has not been encoded before.

```
Please enter your option
m
Please input the name of the feature to be scaled
gender
Feature gender exist in the file
this feature gender is categorical feature and must be encoded first
```

- If the feature that we entered is not in the dataset:

```
Please enter your option
m
Please input the name of the feature to be scaled
wrongfeature
the feature does not exist in the dataset.

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
```

- Option 's':

- The code:

```
42     "s" )
43         saved=false
44         echo "Please input the name of the file to save the processed dataset"
45         read filename
46         cat newdataset.txt > "$filename"
47         saved=true
48         echo "Data Saved To The File $filename"
49         cat "$filename"
50     ;;
```

- The output:

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scaling
s) save the processed dataset
e) exit
Please enter your option
s
Please input the name of the file to save the processed dataset
saved.txt
Data Saved To The File saved.txt
1d;age;gender;height;weight;active;A+;A-;AB+;B+;0-;smoke;governorate;
1;20;male;171;55;yes;1;0;0;0;0;0;0;
2;20;female;158;55;yes;0;1;0;0;0;0;1;1;
3;30;male;160;60;no;0;0;0;1;0;0;2;
4;19;female;162;55;yes;0;0;1;0;0;0;3;
5;13;female;142;37;no;0;0;0;0;1;0;1;
```

- Option 'e':
  - The code:

```

51 "e" )
52     if [ "$saved" ];then
53         echo "Are you sure you want to exist"
54         read answer
55         if [ "$answer" = "yes" ];then
56             exit 10
57         else
58             continue
59         fi
60     else
61         echo "The processed dataset is not saved. Are you sure you want to exist ?"
62         read answer
63         if [ "$answer" = "yes" ];then
64             exit 10
65         else
66             continue
67         fi
68     fi
69 ;;
70 esac
71 done

```

## ▪ The output:

```

Please enter your option
5
Please input the name of the file to save the processed dataset
saved.txt
Data Saved To The File saved.txt
id;age;gender;height;weight;active;A+;A-;AB+;B+;0-;smoke;governorate;
1;20;male;171;55;yes;1;0;0;0;0;0;0;
2;20;female;158;55;yes;0;1;0;0;0;1;1;
3;30;male;160;60;no;0;0;0;1;0;0;2;
4;19;female;162;55;yes;0;0;1;0;0;0;3;
5;13;female;142;37;no;0;0;0;0;1;0;1;

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
e
Are you sure you want to exist
no

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
e
Are you sure you want to exist
yes
anr@anr-VirtualBox:~/Project$

```

The case when the processed dataset is saved to file

```

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
e
The processed dataset is not saved. Are you sure you want to exist ?
no

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
e
The processed dataset is not saved. Are you sure you want to exist ?
yes
anr@anr-VirtualBox:~/Project$

```

The case when the processed dataset is not saved

## DATASET 2:

---

### ❖ The output of option ‘r’ & ‘p’:

```
r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
r
Please input the name of the dataset file
tempdataset.txt
Format is Clean

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
p
You have these features
age
sex
bmi
children
smoker
region
charges
```

### ❖ The output of option ‘l’:

```
Please enter your option
l
Please input the name of the categorical feature for label encoding
sex
Feature sex exist in the file
The values of the features' codes are encoded as follows :
female=0
male=1
-----
The dataset after label encoding for the feature sex :
age;sex;bmi;children;smoker;region;charges;
19;0;27;0;yes;southwest;16884;
18;1;33;1;no;southeast;1725;
28;1;33;3;no;southeast;4449;
33;1;22;0;no;northwest;21984;
32;1;28;e;no;northwest;3866;
31;0;25;e;no;southeast;3756;
46;0;33;1;yes;southeast;8240;
```

### ❖ The outputs of option ‘o’:

```
Please enter your option
o
Please input the name of the categorical feature for one hot encoding
region
Feature region exist in the file
region distinct values are :
northwest
southeast
southwest
age;sex;bmi;children;smoker;northwest;southeast;southwest;charges;
19;0;27;0;yes;0;0;1;16884;
18;1;33;1;no;0;1;0;1725;
28;1;33;3;no;0;1;0;4449;
33;1;22;0;no;1;0;0;21984;
32;1;28;e;no;1;0;0;3866;
31;0;25;e;no;0;1;0;3756;
46;0;33;1;yes;0;1;0;8240;
```

```
Please enter your option
p
You have these features
age
sex
bmi
children
smoker
northwest
southeast
southwest
charges
```



## ❖ The outputs of option ‘m’:

```
Please enter your option
m
Please input the name of the feature to be scaled
sex
Feature sex exist in the file
The Feature sex Has Been Encoded
max value = 1
min value = 0
x1 scaled into 0.00
x2 scaled into 1.00
x3 scaled into 1.00
x4 scaled into 1.00
x5 scaled into 1.00
x6 scaled into 0.00
x7 scaled into 0.00
x-scaled = [0.00, 1.00, 1.00, 1.00, 1.00, 0.00, 0.00].
```

```
Please enter your option
m
Please input the name of the feature to be scaled
region
Feature region exist in the file
The Feature region Has Been Encoded
max value = 1
min value = 0
x1 scaled into 0.00
x2 scaled into 1.00
x3 scaled into 0.00
x4 scaled into 0.00
x5 scaled into 1.00
x6 scaled into 0.00
x7 scaled into 1.00
x8 scaled into 0.00
x9 scaled into 0.00
x10 scaled into 1.00
x11 scaled into 0.00
x12 scaled into 0.00
x13 scaled into 0.00
x14 scaled into 1.00
x15 scaled into 0.00
x16 scaled into 0.00
x17 scaled into 1.00
x18 scaled into 0.00
x19 scaled into 0.00
x20 scaled into 0.00
x21 scaled into 1.00
x-scaled = [0.00, 1.00, 0.00, 0.00, 1.00, 0.00, 1.00, 0.00, 0.00, 1.00, 0.00, 0.00, 0.00, 1.00, 0.00, 0.00, 1.00, 0.00, 0.00, 0.00, 1.00].
```

```
Please enter your option
m
Please input the name of the feature to be scaled
sex
Feature sex exist in the file
this feature sex is categorical feature and must be encoded first

r) read a dataset from a file
p) print the names of the features
l) encode a feature using label encoding
o) encode a feature using one-hot encoding
m) apply MinMax scalling
s) save the processed dataset
e) exit
Please enter your option
m
Please input the name of the feature to be scaled
age
The Feature is Numeric
max value = 46
min value = 18
x1 scaled into 0.00
x2 scaled into 0.00
x3 scaled into 0.30
x4 scaled into 0.50
x5 scaled into 0.50
x6 scaled into 0.40
x7 scaled into 1.00
x-scaled = [0.00, 0.00, 0.30, 0.50, 0.50, 0.40, 1.00].
```

**THANK  
YOU**