In [1]:
```python
# here, we implement knn by removing categorical variables(Type,Lifes
import numpy as np
import pandas as pd
import seaborn as sns
import random
import math
import matplotlib.pyplot as plt
from sklearn import datasets,linear_model
from sklearn.model_selection import train_test_split
import operator
```

In [2]:
```python
data1 = pd.read_csv('train_final.csv')
data2=pd.read_csv('test_final.csv')
# X_train, X_crossVal, X_test = np.split(data, [600, 800])
# X_test
# xtrain,xtest=train_test_split(data,test_size=0.25)
# train_data = data[:50]
# test_data = data[50:]
# train_data[2]
data1.head()
```

Out[2]:

|   | id | Type | LifeStyle | Vacation | eCredit | salary | property | label |
|---|----|------|-----------|----------|---------|--------|----------|-------|
| 0 | 1 | student | spend>saving | 6 | 40 | 13.62 | 3.2804 | C1 |
| 1 | 2 | student | spend>saving | 11 | 21 | 15.32 | 2.0232 | C1 |
| 2 | 3 | student | spend>saving | 7 | 64 | 16.55 | 3.1202 | C1 |
| 3 | 4 | student | spend>saving | 3 | 47 | 15.71 | 3.4022 | C1 |
| 4 | 5 | student | spend>saving | 15 | 10 | 16.96 | 2.2825 | C1 |

In [3]:
```python
data2.head()
```

Out[3]:

|   | id | Type | LifeStyle | Vacation | eCredit | salary | property | label |
|---|----|------|-----------|----------|---------|--------|----------|-------|
| 0 | 1 | student | spend<saving | 12 | 19 | 14.7900 | 3.7697 | C1 |
| 1 | 2 | student | spend>>saving | 29 | 10 | 16.1900 | 2.4839 | C1 |
| 2 | 3 | student | spend<<saving | 28 | 60 | 15.4600 | 1.1885 | C1 |
| 3 | 4 | engineer | spend>saving | 15 | 41 | 21.2600 | 1.4379 | C1 |
| 4 | 5 | librarian | spend<saving | 2 | 9 | 19.7207 | 0.6913 | C1 |

In [4]:
```python
1  data1_label=data1[['label']]
2  data1_label.head()
```
Out[4]:

| | label |
|---|---|
| 0 | C1 |
| 1 | C1 |
| 2 | C1 |
| 3 | C1 |
| 4 | C1 |

In [5]:
```python
1  data2_label=data2[['label']]
2  data2_label.head()
```
Out[5]:

| | label |
|---|---|
| 0 | C1 |
| 1 | C1 |
| 2 | C1 |
| 3 | C1 |
| 4 | C1 |

In [6]:
```python
1  data1=data1.drop(['id','Type','LifeStyle','label'],axis=1)
2  # data1.head()
3  data1.head()
```
Out[6]:

| | Vacation | eCredit | salary | property |
|---|---|---|---|---|
| 0 | 6 | 40 | 13.62 | 3.2804 |
| 1 | 11 | 21 | 15.32 | 2.0232 |
| 2 | 7 | 64 | 16.55 | 3.1202 |
| 3 | 3 | 47 | 15.71 | 3.4022 |
| 4 | 15 | 10 | 16.96 | 2.2825 |

In [7]:
```python
1  data2=data2.drop(['id','Type','LifeStyle','label'],axis=1)
2  data2.head()
```
Out[7]:

| | Vacation | eCredit | salary | property |
|---|---|---|---|---|
| 0 | 12 | 19 | 14.7900 | 3.7697 |
| 1 | 29 | 10 | 16.1900 | 2.4839 |
| 2 | 28 | 60 | 15.4600 | 1.1885 |
| 3 | 15 | 41 | 21.2600 | 1.4379 |
| 4 | 2 | 9 | 19.7207 | 0.6913 |

In [8]:
```
1  data1_norm = (data1-data1.min())/(data1.max()-data1.min())
2  data1_norm.head()
```
Out[8]:

|   | Vacation | eCredit | salary | property |
|---|----------|---------|--------|----------|
| 0 | 0.079365 | 0.107558 | 0.219960 | 0.183167 |
| 1 | 0.158730 | 0.052326 | 0.293102 | 0.112797 |
| 2 | 0.095238 | 0.177326 | 0.346023 | 0.174200 |
| 3 | 0.031746 | 0.127907 | 0.309882 | 0.189984 |
| 4 | 0.222222 | 0.020349 | 0.363663 | 0.127311 |

In [9]:
```
1  data2_norm = (data2-data2.min())/(data2.max()-data2.min())
2  data2_norm.head()
```
Out[9]:

|   | Vacation | eCredit | salary | property |
|---|----------|---------|--------|----------|
| 0 | 0.20 | 0.058824 | 0.104637 | 0.398926 |
| 1 | 0.54 | 0.021008 | 0.175059 | 0.243041 |
| 2 | 0.52 | 0.231092 | 0.138339 | 0.085992 |
| 3 | 0.26 | 0.151261 | 0.430086 | 0.116229 |
| 4 | 0.00 | 0.016807 | 0.352657 | 0.025714 |

In [10]:
```
1  data1_norm['label']=data1_label['label']
2  data1_norm.head()
3
```
Out[10]:

|   | Vacation | eCredit | salary | property | label |
|---|----------|---------|--------|----------|-------|
| 0 | 0.079365 | 0.107558 | 0.219960 | 0.183167 | C1 |
| 1 | 0.158730 | 0.052326 | 0.293102 | 0.112797 | C1 |
| 2 | 0.095238 | 0.177326 | 0.346023 | 0.174200 | C1 |
| 3 | 0.031746 | 0.127907 | 0.309882 | 0.189984 | C1 |
| 4 | 0.222222 | 0.020349 | 0.363663 | 0.127311 | C1 |

In [11]:
```
1  data2_norm['label']=data2_label['label']
2  data2_norm.head()
3
```

Out[11]:

| | Vacation | eCredit | salary | property | label |
|---|---|---|---|---|---|
| 0 | 0.20 | 0.058824 | 0.104637 | 0.398926 | C1 |
| 1 | 0.54 | 0.021008 | 0.175059 | 0.243041 | C1 |
| 2 | 0.52 | 0.231092 | 0.138339 | 0.085992 | C1 |
| 3 | 0.26 | 0.151261 | 0.430086 | 0.116229 | C1 |
| 4 | 0.00 | 0.016807 | 0.352657 | 0.025714 | C1 |

In [12]:
```
1  xtrain=data1_norm.copy()
2  xtest=data2_norm.copy()
```

```
In [13]:   1  k_list=[]
           2  acc_list=[]
           3  for k in range(1,26,2):
           4      k_list.append(k)
           5      predict=[]
           6      def euc_distance(testrow,trainrow,length):
           7          distance=0
           8  #        for i in range(2):
           9  #            if(testrow[i]==trainrow[i]):
          10  #                distance+=1
          11  #        for i in range(2):
          12  #            if(testrow[i]==trainrow[i]):
          13  #                distance+=1
          14          for i in range(length):
          15              distance+=pow((testrow[i]-trainrow[i]),2)
          16          return math.sqrt(distance)
          17      def getNeighbours(traindata,testRow,k):
          18          distance_with_train=[]
          19          length=len(testRow)-1
          20          for x in range(len(traindata)):
          21              dist=euc_distance(testRow,traindata[x],length)
          22              distance_with_train.append((traindata[x],dist))
          23          distance_with_train.sort(key=operator.itemgetter(1))
          24          neighbors = []
          25          for x in range(k):
          26              neighbors.append(distance_with_train[x][0])
          27          return neighbors
          28      def getResponse(neighbors):
          29          votes = {}
          30          for x in range(len(neighbors)):
          31              response = neighbors[x][-1]
          32              if response in votes:
          33                  votes[response] += 1
          34              else:
          35                  votes[response] = 1
          36          sortedVotes = sorted(votes.items(), key=operator.itemgetter(1
          37          return sortedVotes[0][0]
          38      def getAccuracy(xtest, predict):
          39          correct = 0
          40          for x in range(len(xtest)):
          41              if xtest[x][-1] == predict[x]:
          42                  correct += 1
          43          return (correct/float(len(xtest))) * 100.0
          44      for i in range(len(xtest)):
          45          neighbour=getNeighbours(xtrain.values,xtest.values[i],k)
          46  #        print(neighbour)
          47          result = getResponse(neighbour)
          48          predict.append(result)
          49  #        print('> predicted=' + repr(result) + ', actual=' + repr(x
          50      accuracy = getAccuracy(xtest.values, predict)
          51      acc_list.append(accuracy)
          52      print('Accuracy: ' + repr(accuracy) + '%','with k=',k)

Accuracy: 28.57142857142857% with k= 1
Accuracy: 28.57142857142857% with k= 3
Accuracy: 33.33333333333333% with k= 5
```
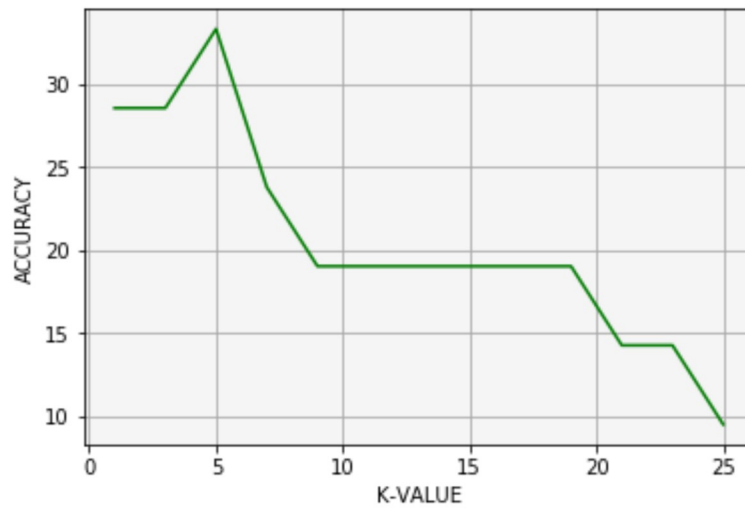
```
In [14]:    1  plt.plot(k_list,acc_list,color='green')
            2  plt.xlabel('K-VALUE')
            3  plt.ylabel('ACCURACY')
            4  plt.grid(True)
            5  plt.show()
```



```
In [16]:    1  bigdata = pd.concat([xtrain, xtest], ignore_index=True)
            2  bigdata.head()
```

Out[16]:

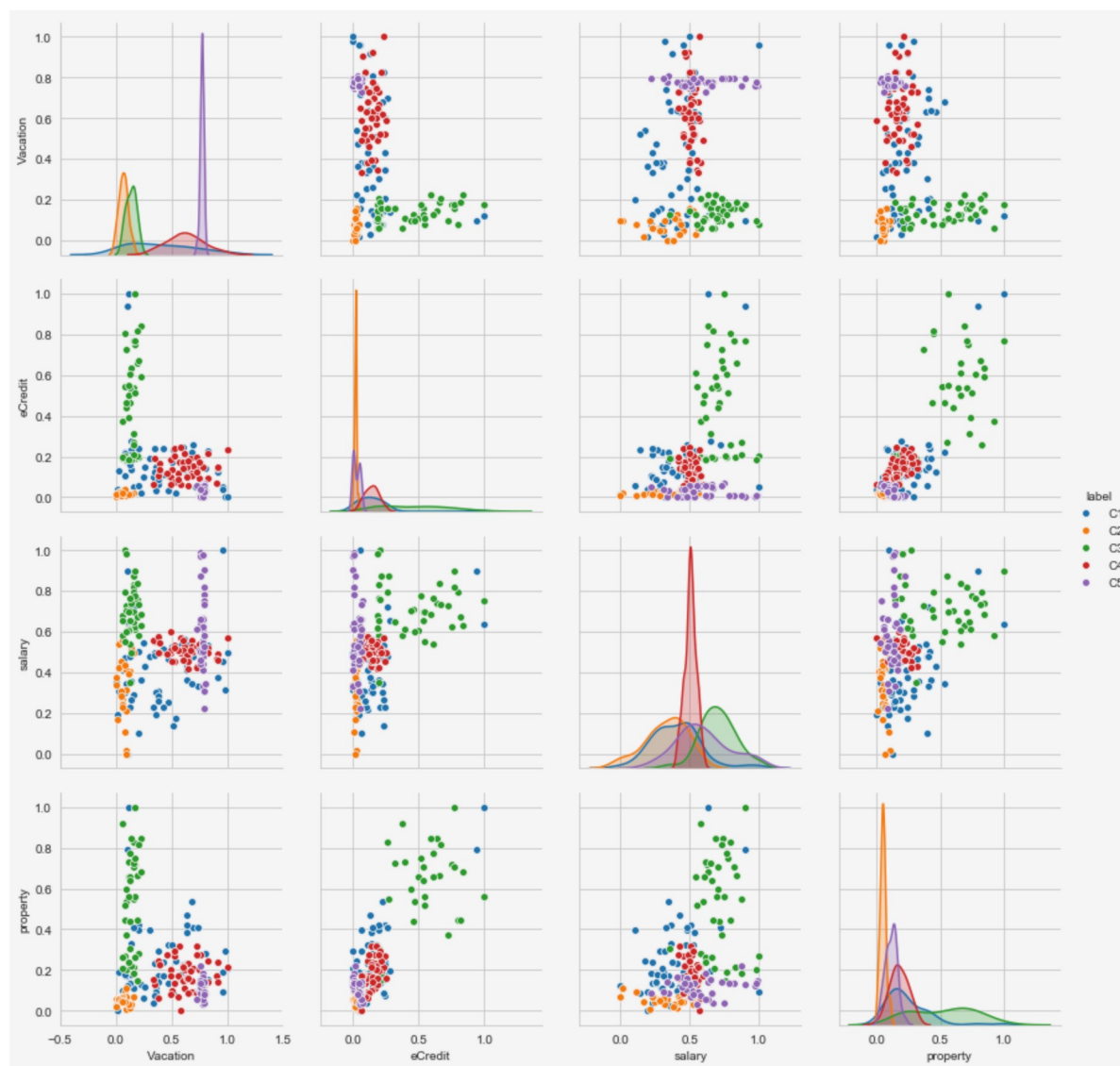|   | Vacation | eCredit | salary | property | label |
|---|----------|---------|--------|----------|-------|
| 0 | 0.079365 | 0.107558 | 0.219960 | 0.183167 | C1 |
| 1 | 0.158730 | 0.052326 | 0.293102 | 0.112797 | C1 |
| 2 | 0.095238 | 0.177326 | 0.346023 | 0.174200 | C1 |
| 3 | 0.031746 | 0.127907 | 0.309882 | 0.189984 | C1 |
| 4 | 0.222222 | 0.020349 | 0.363663 | 0.127311 | C1 |

In [17]:
```python
import seaborn as sns
plt.close();
sns.set_style("whitegrid");
sns.pairplot(bigdata, hue="label", size=3);
plt.show()
```

```
C:\Users\prem\Anaconda3\lib\site-packages\seaborn\axisgrid.py:2065: Us
erWarning: The `size` parameter has been renamed to `height`; pleaes u
pdate your code.
  warnings.warn(msg, UserWarning)
C:\Users\prem\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: F
utureWarning: Using a non-tuple sequence for multidimensional indexing
is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the fut
ure this will be interpreted as an array index, `arr[np.array(seq)]`,
which will result either in an error or a different result.
  return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



In [ ]:
```
1
```

In [ ]:
```
1
```