# SYSEN 6000: Foundations of Complex Systems

## Machine Learning

Nick Kunz [NetID: nhk37] nhk37@cornell.edu

November 18, 2022

## Data Set

Given a data set $D_{n \times m}$, where $n$ is the number of observations and $m$ is the number of feature vectors $X_1, X_2, X_3, \ldots, X_m$, with target label $Y$, where all $X_i, Y \in D_{n \times m}$, and $D_{test}$ is the test sample $D_{1 \times m}$, where $Y \notin D_{test}$, the following is:

| Employment Status | Credit Rating | Available Credit | Age | Approve Application ? |
|---|---|---|---|---|
| Unemployed | Excellent | High | Young | No |
| Unemployed | Fair | High | Young | No |
| Unemployed | Excellent | High | Middle Age | Yes |
| Unemployed | Excellent | Medium | Senior | Yes |
| Employed | Excellent | Low | Senior | Yes |
| Employed | Fair | Low | Senior | No |
| Employed | Fair | Low | Middle Age | Yes |
| Unemployed | Excellent | Medium | Young | No |
| Employed | Excellent | Low | Young | Yes |
| Employed | Fair | Medium | Young | Yes |
| Unemployed | Fair | Medium | Middle Age | Yes |
| Employed | Excellent | High | Middle Age | Yes |
| Unemployed | Fair | Medium | Senior | No |
| Unemployed | Fair | Low | Young | No |
| Unemployed | Excellent | Medium | Middle Age | Yes |
| Employed | Fair | Medium | Senior | No |
| Employed | Excellent | High | Senior | Yes |
| Employed | Fair | Medium | Senior | No |
| Employed | Fair | Medium | Middle Age | Yes |
| | | | | |
| *Unemployed* | *Excellent* | *High* | *Senior* | *???* |

where $n = 19$, $m = 4$ and:

$$
\begin{aligned}
X_1 &= \text{Employment Status} \\
X_2 &= \text{Credit Rating} \\
X_3 &= \text{Available Credit} \\
X_4 &= \text{Age} \\
Y &= \text{Approve Application?}
\end{aligned}
\tag{1}
$$

**Implementation**

The script for the give data set $D_{n \times m}$ and test sample $D_{test}$ is as follows:

```python
## features
var = [
    'Employment Status',
    'Credit Rating',
    'Available Credit',
    'Age',
    'Approve Application'
]

## observations
obs = [
    ['Unemployed', 'Excellent', 'High', 'Young', 'No'],
    ['Unemployed', 'Fair', 'High', 'Young', 'No'],
    ['Unemployed', 'Excellent', 'High', 'Middle Age', 'Yes'],
    ['Unemployed', 'Excellent', 'Medium', 'Senior', 'Yes'],
    ['Employed', 'Excellent', 'Low', 'Senior', 'Yes'],
    ['Employed', 'Fair', 'Low', 'Senior', 'No'],
    ['Employed', 'Fair', 'Low', 'Middle Age', 'Yes'],
    ['Unemployed', 'Excellent', 'Medium', 'Young', 'No'],
    ['Employed', 'Excellent', 'Low', 'Young', 'Yes'],
    ['Employed', 'Fair', 'Medium', 'Young', 'Yes'],
    ['Unemployed', 'Fair', 'Medium', 'Middle Age', 'Yes'],
    ['Employed', 'Excellent', 'High', 'Middle Age', 'Yes'],
    ['Unemployed', 'Fair', 'Medium', 'Senior', 'No'],
    ['Unemployed', 'Fair', 'Low', 'Young', 'No'],
    ['Unemployed', 'Excellent', 'Medium', 'Middle Age', 'Yes'],
    ['Employed', 'Fair', 'Medium', 'Senior', 'No'],
    ['Employed', 'Excellent', 'High', 'Senior', 'Yes'],
    ['Employed', 'Fair', 'Medium', 'Senior', 'No'],
    ['Employed', 'Fair', 'Medium', 'Middle Age', 'Yes']
]

## predictions
test = ['Unemployed', 'Excellent', 'High', 'Senior']
```

Python 3: Data Set & Test Sample

# Decision Tree Classifier

The following is an exhibit of a Decision Tree classifier utilizing Entropy and Information Gain for the given data set $D_{n \times m}$ and $D_{test}$, where Entropy is computed as:

$$H(X) = - \sum_{j=1}^{c} p_j \log_2(p_j) \tag{2}$$

and $p_j$ is the probability of observing class $C$, which can also be expressed as:

$$H(X) = \sum_{j=1}^{c} \log_2 \left( \frac{1}{p_j} \right) p_j \tag{3}$$

Information Gain is computed as:

$$IG(Y, X) = H(Y) - H(Y|X) \tag{4}$$

Building the Decision Tree classifier utilizing Information Gain occurs in these steps.

---

**Algorithm 1** Decision Tree Classifier

---

**Input:** $D_{n \times m}$, $D_{test}$
**Output:** class $C$ of $D_{test}$

1: **for** each $i$ **do**
2:     compute $IG(Y, X_i)$
3: **end for**
4: create node $N_i$ with highest $IG$
5: **if** all $Y_j \in N_i$ are the same class $C$, **then**
6:     **return** leaf node labeled class $C$ to tree
7: **else** step 1
8: **end if**
9: **if** tree is **done**, **then**
10:     traverse $N_i$ in tree with $X_{j \times m} \in D_{test}$ until $C$
11:     **return** leaf node as class $C$
12: **end if**

---

**Implementation**

The scripts for the Decision Tree classifier are as follows:

```python
## libraries
import math
import pprint

## feat vect of data
def feat_vect(x_y, y_i):

    n = len(x_y)

    return [x_y[i][y_i] for i in range(0, n)]


## feat subset of data
def feat_subs(x_y, y_i):

    n = len(x_y)

    return [x_y[i][0:y_i] for i in range(0, n)]


## count unique vals of feat
def counter(y_j):

    ## find unique vals
    uni_val = list(
        set([j for j in y_j])
    )

    ## count unique vals
    uni_val_cts = dict.fromkeys(
        uni_val, 0
    )

    n = len(y_j)
    i_j = [None] * n

    for j in range(0, n):
        i_j[j] = y_j[j]
```

```python
39
40      for i in i_j:
41          uni_val_cts[i] += 1
42
43      return uni_val_cts
44

45
46  ## entropy of feat
47  def entropy(y_j):
48
49      ## count unique vals
50      cts = counter(
51          y_j = y_j
52      )
53
54      ## compute probs
55      n = len(y_j)
56
57      prob = [(j / n) for j in cts.values()]
58
59      ## compute entropy
60      return sum(
61          [-p * math.log(p, 2) for p in prob]
62      )
63

64
65  ## info gain of data
66  def info_gain(x_y, x_i, y_i):
67
68      ## trgt feat
69      y_j = feat_vect(
70          x_y = x_y,
71          y_i = y_i
72      )
73
74      ## attr feat
75      x_a = feat_vect(
76          x_y = x_y,
77          y_i = x_i
78      )
79
80      ## count unique vals
81      cts = counter(
82          y_j = x_a
83      )
84
85      ## comp entropy
86      entp = entropy(
87          y_j = y_j
88      )
89
90      ## cond entropy of feat
91      entp_cond = 0
92
93      for i, uni in enumerate(cts):
94
95          ## subset data by unique vals in trgt feat
96          x_y_uni = [j for j in x_y if j[x_i] == uni]
97
98          ## attr feat of subset
99          y_v = feat_vect(
100             x_y = x_y_uni,
101             y_i = y_i
102         )
```

```python
103
104          ## cond entropy of attr feat
105          entp_cond += (cts[uni] / len(x_a)) * entropy(
106              y_j = y_v
107          )
108
109      ## info gain
110      return entp - entp_cond
111
112
113  ## decision tree on data
114  def grow_tree(x_y, y_i, tree = None):
115
116      ## feat info gain
117      m = len(x_y[0]) - 1
118      info = [None] * m
119
120      for i in range(0, m):
121          info[i] = info_gain(
122              x_y = x_y,
123              x_i = i,
124              y_i = y_i
125          )
126
127      ## best feat split
128      x_i_star = info.index(max(info))
129
130      x_i_feat = feat_vect(
131          x_y = x_y,
132          y_i = x_i_star
133      )
134
135      x_i_splt = counter(
136          y_j = x_i_feat
137      )
138
139      if tree is None:
140          tree = dict()
141          tree[x_i_star] = dict()
142
143      else:
144          pass
145
146      ## cont feat split
147      n = len(x_y)
148
149      for i in x_i_splt:
150          x_y_splt = list()
151
152          for j in range(0, n):
153              if x_y[j][x_i_star] == i:
154                  x_y_splt.append(x_y[j])
155
156          y_j = [x_y_splt[k][y_i] for k in range(0, len(x_y_splt))]
157          y_j_freq = counter(
158              y_j = y_j
159          )
160
161          ## leaf node
162          if len(y_j_freq) == 1:
163              tree[x_i_star][i] = y_j_freq
164
165          ## recursion
166          else:
```

```
167              tree[x_i_star][i] = grow_tree(
168                  x_y = x_y_splt,
169                  y_i = y_i
170              )
171
172      return tree
173
174
175 ## predict on decision tree dict
176 def pred_tree(tree, test, verb = False):
177
178      k = list(tree.keys())[0]
179      v = list(tree.values())[0]
180
181      if isinstance(v, dict) == True:
182          if verb == True:
183              print(test[k])
184
185          return pred_tree(
186              tree = v[test[k]],
187              test = test,
188              verb = verb
189          )
190      else:
191          return k
192
193 ## decision tree
194 tree = grow_tree(
195      x_y = obs,
196      y_i = 4
197 )
198
199 ## prediction
200 pred_tree(
201      tree = tree,
202      test = test
203 )
```

Python 3: Decision Tree Classifier

```
"Yes"
```

```
1 ## view tree
2 pprint.pprint(tree)
```

Python 3: Decision Tree

```
{3: {'Middle Age': {'Yes': 6},
     'Senior': {1: {'Excellent': {'Yes': 3}, 'Fair': {'No': 4}}},
     'Young': {0: {'Employed': {'Yes': 2}, 'Unemployed': {'No': 4}}}}}
```

# Naive Bayes Classifier

The following is an exhibit of a Naive Bayes classifier assuming independence for the given data set $D_{n \times m}$ and $D_{test}$, where Bayes theorem is computed as:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)} \tag{5}$$

which assumes independence, such that:

$$P(X_1, X_2, X_3, \ldots, X_m|Y) = \prod_{i=1}^{m} P(X_i|Y) \tag{6}$$

where $X_i$ and $X_j$ are conditionally independent on $Y$, such that:

$$\forall\, i \neq j \tag{7}$$

which implies that a classifier can be expressed as:

$$P(C_k|X) = \frac{P(C_k)P(X|C_k)}{P(X)} = \frac{P(C_k) \prod_{i=1}^{m} P(X_i|C_k)}{P(X)} \tag{8}$$

where $k$ is a class $C$, and $C$ can be estimated by:

$$\hat{C} = \underset{k \in K}{\operatorname{argmax}}\, P(C_k) \prod_{i=1}^{m} P(X_i|Y) \tag{9}$$

Building the Naive Bayes classifier assuming independence occurs in these steps.

---

**Algorithm 2** Naive Bayes Classifier

---

**Input:** $D_{n \times m}$, $D_{test}$
**Output:** class $C$ of $D_{test}$

1: **for** each $k$ **do**
2:     compute $P(C_k)$
3:     **for** each $i$ **do**
4:         compute $P(X_i|C_k)$
5:     **end for**
6: **end for**
7: **return** $\operatorname{argmax}_{k \in K} P(C_k) \prod_{i=1}^{m} P(X_i|C_k)$

---

### Implementation
The scripts for the Naive Bayes classifier are as follows:

```
## class prob
def class_prob(y_j):

    ## count unique vals
    y_j_cts = counter(
```

```python
         y_j = y_j
     )

     y_i_prob = y_j_cts.copy()

     n = len(y_j)

     for i in y_j_cts.keys():
         y_i_prob[i] = round(
             number = (y_j_cts[i] / n),
             ndigits = 3
         )

     return y_i_prob


## conditional prob
def cond_prob(x_y, y_i):

     ## trgt feat
     y_j = feat_vect(
         x_y = x_y,
         y_i = y_i
     )

     ## count unique vals
     cts = counter(
         y_j = y_j
     )

     ## cond prob
     n = len(x_y)
     m = len(x_y[0]) - 1

     y_i_splt = list()
     x_y_cond = cts.copy()

     for i in cts.keys():
         x_y_splt = list()
         x_y_cond[i] = dict()

         for j in range(0, n):
             if x_y[j][y_i] == i:
                 x_y_splt.append(x_y[j])

         n_splt = len(x_y_splt)

         for j in range(0, m):
             x_splt_subs = feat_vect(
                 x_y = x_y_splt,
                 y_i = j
             )

             cts_splt_subs = counter(
                 y_j = x_splt_subs
             )

             x_splt_prob = {
                 k: round(v / n_splt, 3) for k, v in cts_splt_subs.items()
             }

             x_y_cond[i].update(x_splt_prob)

     return x_y_cond
```

```python
70
71
72  ## posterior prob
73  def post_prob(p_clss, p_cond, test):
74
75      post_prob = p_clss.copy()
76
77      for i in p_clss.keys():
78          test_prob = 1
79
80          for j in test:
81              test_prob = test_prob * p_cond[i][j]
82
83          post_prob[i] = round(
84              number = p_clss[i] * test_prob,
85              ndigits = 3
86          )
87
88      return post_prob
89
90
91  ## predict on post prob dict
92  def pred_post(p_post):
93
94      ## arg max func
95      return max(
96          p_post,
97          key = p_post.get
98      )
99
100 ## class prob
101 y_j = feat_vect(
102     x_y = obs,
103     y_i = 4
104 )
105
106 prob_clss = class_prob(
107     y_j = y_j
108 )
109
110 ## conditional prob
111 prob_cond = cond_prob(
112     x_y = obs ,
113     y_i = 4
114 )
115
116 ## posterior prob
117 prob_post = post_prob(
118     p_clss = prob_clss,
119     p_cond = prob_cond,
120     test = test
121 )
122
123 ## prediction
124 pred_post(
125     p_post = prob_post
126 )
```

Python 3: Naive Bayes Classifier

```
"Yes"
```

# k-Nearest Neighbors Classifier

The following is an exhibit of a k-Nearest Neighbors classifier utilizing Jaccard Similarity for the given data set $D_{n \times m}$ and $D_{test}$, where the Jaccard Similarity distance metric is computed as:

$$J(A|B) = \frac{|A \cap B|}{|A \cup B|} \tag{10}$$

Building the k-Nearest Neighbors classifier with Jaccard Similarity occurs in these steps.
$k$ should be set to an odd number to avoid cases where the number of $C$ in $Y_j \in D_{k \times m}$ are equal.

---
**Algorithm 3** k-NN Classifier

---
**Input:** $D_{n \times m}$, $D_{test}$
**Output:** class $C$ of $D_{test}$
 1: **set** $k$
 2: **for** each $j$ **do**
 3:     compute distance $I_j = J(X_{j \times m}|D_{test})$
 4: **end for**
 5: **sort** distances $I_j$ in ascending order
 6: **truncate** $I_j$ by $k$ to $I_k$
 7: **subset** $D_{n \times m}$ to $D_{k \times m}$ indexed by $I_k$
 8: **return** majority of $Y_j \in D_{k \times m}$ as $C$

---

### Implementation
The scripts for the k-NN classifier are as follows:

```python
## jaccard similiarity
def jaccard(a, b):

    ## intersect
    inter = len(list(set(a).intersection(b)))

    ## union
    n_a, n_b = len(a), len(b)
    union = (n_a + n_b) - inter

    ## proportion
    return round(
        number = inter / union,
        ndigits = 2
    )

## dist matrix on data
def dist_matrix(x_y, y_i, test = False):

    x_i = feat_subs(
        x_y = x_y,
        y_i = y_i
    )

    n = len(x_i)
    dist_mtrx = [None] * n

    if test is False:
        for i in range(0, n):

            dist = [None] * n

```

```python
33              for j in range(0, n):
34                  dist[j] = jaccard(
35                      a = x_i[i],
36                      b = x_i[j]
37                  )
38
39              dist_mtrx[i] = dist
40
41      else:
42          for i in range(0, n):
43              dist_mtrx[i] = jaccard(
44                  a = test,
45                  b = x_i[i]
46              )
47
48      return dist_mtrx
49
50  ## knn on data
51  def knn(x_y, y_i, k, test = False):
52
53      dist_mtrx = dist_matrix(
54          x_y = x_y,
55          y_i = y_i,
56          test = test
57      )
58
59      ## find near neigh index
60      n = len(dist_mtrx)
61      near_negh = dict()
62
63      if test is False:
64          for i in range(0, n):
65
66              ## neigh sorted by index
67              idx = sorted(
68                  range(len(dist_mtrx[i])),
69                  key = lambda k: dist_mtrx[i][k]
70              )
71
72              ## subset near by k
73              near_negh[i] = list(
74                  reversed(idx[-k - 1:][:-1])
75              )
76      else:
77          ## neigh sorted by index
78          idx = sorted(
79              range(len(dist_mtrx)),
80              key = lambda k: dist_mtrx[k]
81          )
82
83          ## subset near by k
84          near_negh = list(
85              reversed(idx[-k - 1:][:-1])
86          )
87
88      return near_negh
89
90  ## predict on knn test list
91  def knn_pred(x_y, y_i, k, test):
92
93      idx_knn = knn(
94          x_y = x_y,
95          y_i = y_i,
96          k = k,
```

```
 97            test = test
 98        )
 99
100        n = len(idx_knn)
101        x_y_knn = [None] * n
102
103        for i in range(0, n):
104            x_y_knn[i] = obs[idx_knn[i]]
105
106        y_j_knn = feat_vect(
107            x_y = x_y_knn,
108            y_i = y_i
109        )
110
111        y_j_hat = counter(
112            y_j = y_j_knn
113        )
114
115        return max(
116            y_j_hat,
117            key = y_j_hat.get
118        )
```
Python 3: k-NN Classifier

$k = 1$

```
1  ## prediction
2  knn_pred(
3      x_y = obs,
4      y_i = 4,
5      k = 1,
6      test = test
7  )
```
Python 3: 1-NN Classifier

"Yes"

$k = 3$

```
1  ## prediction
2  knn_pred(
3      x_y = obs,
4      y_i = 4,
5      k = 3,
6      test = test
7  )
```
Python 3: 3-NN Classifier

"No"

$k = 5$

```
1  ## prediction
2  knn_pred(
3      x_y = obs,
4      y_i = 4,
5      k = 5,
6      test = test
7  )
```
Python 3: 5-NN Classifier

"Yes"

# Dependency Support Classifiers

The following verifies the previous details of the given classifiers utilizing common dependencies for machine learning to include: *pandas* and *scikit-learn*. The outputs of the scripts follow the figures.

## Data Set

```python
## libraries
import pandas as pd
from sklearn.preprocessing import OrdinalEncoder
from sklearn.preprocessing import MinMaxScaler

## data
data = pd.DataFrame(
    data = obs,
    columns = var
)

data_pred = pd.DataFrame(
    data = [test],
    columns = var[0:4]
)

## pre-process
enc = OrdinalEncoder()
mms = MinMaxScaler()

for i in data.columns:
    data[i] = enc.fit_transform(data[[i]])

data = mms.fit_transform(data)
data = pd.DataFrame(
    data = data,
    columns = var
)
data_pred = enc.fit_transform(data_pred)
```

Python 3: Pandas Dataframe

## Decision Tree Classifier

```python
## libraries
from sklearn import tree

## train
clf = tree.DecisionTreeClassifier(
    criterion = 'entropy'
)

clf = clf.fit(
    X = data.drop(['Approve Application'], axis = 1).values,
    y = data['Approve Application']
)

## predict
if int(clf.predict(data_pred)[0]) == 1:
    print('Yes')
else:
    print('No')
```

Python 3: Scikit-Learn Decision Tree Classifier

```
"Yes"
```

**Naive Bayes Classifier**

```python
## libraries
from sklearn.naive_bayes import CategoricalNB

## train
clf = CategoricalNB()

clf = clf.fit(
    X = data.drop(['Approve Application'], axis = 1).values,
    y = data['Approve Application']
)

## predict
if int(clf.predict(data_pred)[0]) == 1:
    print('Yes')
else:
    print('No')
```

Python 3: Scikit-Learn Naive Bayes Classifier

```
"Yes"
```

**k-Nearest Neighbors Classifier**

$k = 1$

```python
## libraries
from sklearn.neighbors import KNeighborsClassifier

## train
clf = KNeighborsClassifier(
    n_neighbors = 1,
    metric = 'jaccard'
)

clf = clf.fit(
    X = data.drop(['Approve Application'], axis = 1).values,
    y = data['Approve Application']
)

## predict
if int(clf.predict(data_pred)[0]) == 1:
    print('Yes')
else:
    print('No')
```

Python 3: Scikit-Learn 1-NN Classifier

```
"Yes"
```

$k = 3$

```python
## libraries
from sklearn.neighbors import KNeighborsClassifier

## train
clf = KNeighborsClassifier(
    n_neighbors = 3,
    metric = 'jaccard'
)

clf = clf.fit(
    X = data.drop(['Approve Application'], axis = 1).values,
```

```python
12      y = data['Approve Application']
13 )
14
15 ## predict
16 if int(clf.predict(data_pred)[0]) == 1:
17     print('Yes')
18 else:
19     print('No')
```

Python 3: Scikit-Learn 3-NN Classifier

```
"No"
```

$k = 5$

```python
1 ## libraries
2 from sklearn.neighbors import KNeighborsClassifier
3
4 ## train
5 clf = KNeighborsClassifier(
6     n_neighbors = 5,
7     metric = 'jaccard'
8 )
9
10 clf = clf.fit(
11     X = data.drop(['Approve Application'], axis = 1).values,
12     y = data['Approve Application']
13 )
14
15 ## predict
16 if int(clf.predict(data_pred)[0]) == 1:
17     print('Yes')
18 else:
19     print('No')
```

Python 3: Scikit-Learn 5-NN Classifier

```
"Yes"
```

**Support Vector Machines Classifier**

```python
1 ## libraries
2 from sklearn.svm import SVC
3
4 ## train
5 clf = SVC()
6
7 clf = clf.fit(
8     X = data.drop(['Approve Application'], axis = 1).values,
9     y = data['Approve Application']
10 )
11
12 ## predict
13 if int(clf.predict(data_pred)[0]) == 1:
14     print('Yes')
15 else:
16     print('No')
```

Python 3: Scikit-Learn SVM Classifier

```
"Yes"
```