| Started on | Saturday, 30 September 2023, 10:42 PM |
|---|---|
| State | Finished |
| Completed on | Tuesday, 3 October 2023, 4:14 PM |
| Time taken | 2 days 17 hours |
| Marks | 5.00/5.00 |
| Grade | **1.00** out of 1.00 (**100**%) |

**Information**

# Representation of belief networks in Python

A belief (or Bayesian) network is represented by a dictionary. The keys are the names of variables. The values are dictionaries themselves. The second level dictionaries have two keys: `'Parents'` whose value is a list of the names of the variables that are the parents of the current variable, and `'CPT'` whose value is a dictionary again. The keys of the third level dictionaries are tuples of Booleans which correspond to possible assignments of values to the parents of the current node (in the order they are listed) and the values are real numbers representing the probability of the current node being <u>true</u> given the specified assignment to the parents.

## Notes

- Variable names are case sensitive.
- If a node does not have any parents, the value of `'Parents'` must be an empty list and the only key of the third level dictionary is the empty tuple.
- For simplicity, we assume that all the variables are Boolean.

## Example

The following is the representation of the *alarm network* presented in the lecture notes.

```python
network = {
    'Burglary': {
        'Parents': [],
        'CPT': {
            (): 0.001,
        }
    },

    'Earthquake': {
        'Parents': [],
        'CPT': {
            (): 0.002,
        }
    },

    'Alarm': {
        'Parents': ['Burglary','Earthquake'],
        'CPT': {
            (True,True): 0.95,
            (True,False): 0.94,
            (False,True): 0.29,
            (False,False): 0.001,
        }
    },

    'John': {
        'Parents': ['Alarm'],
        'CPT': {
            (True,): 0.9,
            (False,): 0.05,
        }
    },

    'Mary': {
        'Parents': ['Alarm'],
        'CPT': {
            (True,): 0.7,
            (False,): 0.01,
        }
    },
}
```

## Question **1**

Correct

Mark 1.00 out of 1.00

Suppose we want to predict the value of variable Y based on the values of variables X1, X2, and X3. Assuming that we want to use a Naive Bayes model for this purpose, create a belief net for the model called `network`. The probabilities must be learnt by using the dataset given below. Use Laplacian smoothing with a pseudocount of 2.

| X1 | X2 | X3 | Y |
|----|----|----|---|
| T  | T  | F  | F |
| T  | F  | F  | F |
| T  | T  | F  | F |
| T  | F  | F  | T |
| F  | F  | F  | T |
| F  | T  | F  | T |
| F  | F  | F  | T |

## Notes

- Node names are case sensitive.
- Since we are using Python syntax, you can use fraction expressions if you wish. For example you can use `3/4` or even `(2+1)/(2+1+0+1)` which will be evaluated at runtime.

**For example:**

| Test | Result |
|------|--------|
| ```
from student_answer import network
from numbers import Number

# Checking the overall type-correctness of the network
# without checking anything question-specific

assert type(network) is dict
for node_name, node_info in network.items():
    assert type(node_name) is str
    assert type(node_info) is dict
    assert set(node_info.keys()) == {'Parents', 'CPT'}
    assert type(node_info['Parents']) is list
    assert all(type(s) is str for s in node_info['Parents'])
    for assignment, prob in node_info['CPT'].items():
        assert type(assignment) is tuple
        assert isinstance(prob, Number)

print("OK")
``` | OK |

**Answer:** (penalty regime: 0, 10, ... %)

```python
 1  network = {
 2      'X1': {
 3          'Parents': ['Y'],
 4          'CPT': {
 5              (True,): 3/8,
 6              (False,): 5/7
 7          }
 8      },
 9
10      'X2': {
11          'Parents': ['Y'],
12          'CPT': {
13              (True,): 3/8,
14              (False,): 4/7
15          }
16      },
17
18      'X3': {
19          'Parents': ['Y'],
20          'CPT': {
```

```
21              (True,): 2/8,
22              (False,): 2/7
23         }
24     },
25
26 ▼    'Y': {
27          'Parents': [],
28 ▼        'CPT': {
29              (): 6/11,
30          }
31     },
32 }
```

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | `from student_answer import network`<br>`from numbers import Number`<br><br>`# Checking the overall type-correctness of the network`<br>`# without checking anything question-specific`<br><br>`assert type(network) is dict`<br>`for node_name, node_info in network.items():`<br>`    assert type(node_name) is str`<br>`    assert type(node_info) is dict`<br>`    assert set(node_info.keys()) == {'Parents', 'CPT'}`<br>`    assert type(node_info['Parents']) is list`<br>`    assert all(type(s) is str for s in node_info['Parents'])`<br>`    for assignment, prob in node_info['CPT'].items():`<br>`        assert type(assignment) is tuple`<br>`        assert isinstance(prob, Number)`<br><br>`print("OK")` | OK | OK | ✔ |

Passed all tests! ✔

(Correct)

Marks for this submission: 1.00/1.00.

---

**Information**

---

# Representation of naïve Bayes models

Naïve Bayes models can be represented with belief networks. However, since they all have a very simple topology (a directed tree of depth one where the root is the class variable and the leaves are the input features), we can use a more compact representation that is only concerned with the values of CPTs.

We assume that all the variables in a naïve Bayes network are binary. For a network with $n$ binary input features X[1] to X[n], we represent the conditional probability tables (CPTs) that are required in the network, with the following two objects:

- `prior`: a real number representing *p(Class=true)*. The probability *p(Class=false)* can be obtained by `1 - prior`.
- `likelihood`: a tuple of length $n$ where each element is a pair of real numbers such that `likelihood[i][False]` is *p(X[i]=true|C=false)* and `likelihood[i][True]` is *p(X[i]=true|C=true )*. That is, likelihood contains the *2\*n* CPTs that are required at leaf nodes.

Note: in general, indexing sequences with booleans is not ideal, however, here we are using False (for 0) and True (for 1) so that the appearance of the code is closer to the corresponding mathematical notation.

**Question 2**

Correct

Mark 1.00 out of 1.00

Write a function `posterior(prior, likelihood, observation)` that returns the posterior probability of the class variable being true, given the observation; that is, it returns *p(Class=true|observation)*. The argument `observation` is a tuple of *n* Booleans such that `observation[i]` is the observed value (True or False) for the input feature `X[i]`. The arguments `prior` and `likelihood` are as described above.

**Notes**

1. Example 9.36 in the textbook is relevant.
2. The model used in the test cases is according to this network. You can download and explore the model in the belief network applet.

**For example:**

| Test | Result |
|------|--------|
| ```from student_answer import posterior

prior = 0.05
likelihood = ((0.001, 0.3),(0.05,0.9),(0.7,0.99))

observation = (True, True, True)

class_posterior_true = posterior(prior, likelihood, observation)
print("P(C=False|observation) is approximately {:.5f}"
      .format(1 - class_posterior_true))
print("P(C=True |observation) is approximately {:.5f}"
      .format(class_posterior_true))``` | P(C=False\|observation) is approximately 0.00248<br>P(C=True \|observation) is approximately 0.99752 |
| ```from student_answer import posterior

prior = 0.05
likelihood = ((0.001, 0.3),(0.05,0.9),(0.7,0.99))

observation = (True, False, True)

class_posterior_true = posterior(prior, likelihood, observation)
print("P(C=False|observation) is approximately {:.5f}"
      .format(1 - class_posterior_true))
print("P(C=True |observation) is approximately {:.5f}"
      .format(class_posterior_true))``` | P(C=False\|observation) is approximately 0.29845<br>P(C=True \|observation) is approximately 0.70155 |
| ```from student_answer import posterior

prior = 0.05
likelihood = ((0.001, 0.3),(0.05,0.9),(0.7,0.99))

observation = (False, False, True)

class_posterior_true = posterior(prior, likelihood, observation)
print("P(C=False|observation) is approximately {:.5f}"
      .format(1 - class_posterior_true))
print("P(C=True |observation) is approximately {:.5f}"
      .format(class_posterior_true))``` | P(C=False\|observation) is approximately 0.99454<br>P(C=True \|observation) is approximately 0.00546 |
| ```from student_answer import posterior

prior = 0.05
likelihood = ((0.001, 0.3),(0.05,0.9),(0.7,0.99))

observation = (False, False, False)

class_posterior_true = posterior(prior, likelihood, observation)
print("P(C=False|observation) is approximately {:.5f}"
      .format(1 - class_posterior_true))
print("P(C=True |observation) is approximately {:.5f}"
      .format(class_posterior_true))``` | P(C=False\|observation) is approximately 0.99987<br>P(C=True \|observation) is approximately 0.00013 |

**Answer:** (penalty regime: 0, 15, ... %)

```
1  def posterior(prior, likelihood, observation):
2      """returns the posterior probability of the class variable being true.
```

```
  3              given the observation; that is, it returns p(Class=true|observation).
  4              The argument observation is a tuple of n Booleans such that observation[i]
  5              is the observed value (True or False) for the input feature X[i].
  6              The arguments prior and likelihood are as described above.
  7              """
  8              correct = prior
  9              wrong = 1 - prior
 10              for i in range(len(likelihood)):
 11                  if observation[i]:
 12                      correct *= likelihood[i][True]
 13                      wrong *= likelihood[i][False]
 14                  else:
 15                      correct *= 1 - likelihood[i][True]
 16                      wrong *= 1 - likelihood[i][False]
 17              return correct / (correct + wrong)
```

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | from student_answer import posterior<br><br>prior = 0.05<br>likelihood = ((0.001, 0.3),(0.05,0.9),<br>(0.7,0.99))<br><br>observation = (True, True, True)<br><br>class_posterior_true = posterior(prior,<br>likelihood, observation)<br>print("P(C=False\|observation) is<br>approximately {:.5f}"<br>    .format(1 - class_posterior_true))<br>print("P(C=True \|observation) is<br>approximately {:.5f}"<br>    .format(class_posterior_true)) | P(C=False\|observation) is<br>approximately 0.00248<br>P(C=True \|observation) is<br>approximately 0.99752 | P(C=False\|observation) is<br>approximately 0.00248<br>P(C=True \|observation) is<br>approximately 0.99752 | ✔ |
| ✔ | from student_answer import posterior<br><br>prior = 0.05<br>likelihood = ((0.001, 0.3),(0.05,0.9),<br>(0.7,0.99))<br><br>observation = (True, False, True)<br><br>class_posterior_true = posterior(prior,<br>likelihood, observation)<br>print("P(C=False\|observation) is<br>approximately {:.5f}"<br>    .format(1 - class_posterior_true))<br>print("P(C=True \|observation) is<br>approximately {:.5f}"<br>    .format(class_posterior_true)) | P(C=False\|observation) is<br>approximately 0.29845<br>P(C=True \|observation) is<br>approximately 0.70155 | P(C=False\|observation) is<br>approximately 0.29845<br>P(C=True \|observation) is<br>approximately 0.70155 | ✔ |

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | `from student_answer import posterior`<br><br>`prior = 0.05`<br>`likelihood = ((0.001, 0.3),(0.05,0.9),`<br>`(0.7,0.99))`<br><br>`observation = (False, False, True)`<br><br>`class_posterior_true = posterior(prior,`<br>`likelihood, observation)`<br>`print("P(C=False|observation) is`<br>`approximately {:.5f}"`<br>`        .format(1 - class_posterior_true))`<br>`print("P(C=True |observation) is`<br>`approximately {:.5f}"`<br>`        .format(class_posterior_true))` | P(C=False\|observation) is approximately 0.99454<br>P(C=True \|observation) is approximately 0.00546 | P(C=False\|observation) is approximately 0.99454<br>P(C=True \|observation) is approximately 0.00546 | ✔ |
| ✔ | `from student_answer import posterior`<br><br>`prior = 0.05`<br>`likelihood = ((0.001, 0.3),(0.05,0.9),`<br>`(0.7,0.99))`<br><br>`observation = (False, False, False)`<br><br>`class_posterior_true = posterior(prior,`<br>`likelihood, observation)`<br>`print("P(C=False|observation) is`<br>`approximately {:.5f}"`<br>`        .format(1 - class_posterior_true))`<br>`print("P(C=True |observation) is`<br>`approximately {:.5f}"`<br>`        .format(class_posterior_true))` | P(C=False\|observation) is approximately 0.99987<br>P(C=True \|observation) is approximately 0.00013 | P(C=False\|observation) is approximately 0.99987<br>P(C=True \|observation) is approximately 0.00013 | ✔ |

Passed all tests! ✔

Correct

Marks for this submission: 1.00/1.00.

**Information**

# Bayesian Spam Filter

In the next three questions, you are asked to develop the learning and classification components of a naive Bayes classifier (a spam filter).

The file spam-labelled.csv describes 200 emails, labelled as spam or non-spam by human users. Each email is specified by 12 binary attributes, indicating the presence of features such as "Lottery", "MILLION DOLLARS", significant amounts of text in CAPS, an invalid reply-to address, and so on.

The layout of the data is that each row is an example (one email), and columns correspond to attributes (features), which are binary. There are 12 input features (X1 to X12). The last (right-most) column is the class label where 1 means the example is spam (positive) and 0 means non-spam (negative).

Note that there are $2^{12}$ =4096 possible input patterns; in other words, the data set only contains a small proportion of all possible input patterns. This is a common scenario in machine learning.

The file has Unix-like line breaks. Windows users need to open the file in a proper text editor that supports different line endings. You do not need a spreadsheet to open the file.

In Python, the csv module may come in handy. You can load the content of the file as a list of tuples using the following:

```
with open(file_name) as in_file:
      training_examples = [tuple(row) for row in csv.reader(in_file)]
```

In the next three questions, the above file is available on the server (in the current directory). Therefore a statement like the one above, would read the file. Your function will be tested on the same file format and header but the content (the examples) may vary. Make sure your solution works with the original copy of the file, not your own format.

**Question 3**

Correct

Mark 1.00 out of 1.00

Write a function `learn_prior(file_name, pseudo_count=0)` that takes the file name of the training set and an optional pseudo-count parameter and returns a real number that is the prior probability of spam being true. The parameter `pseudo_count` is a non-negative integer and it will be the same for all the attributes and all the values.

**Notes**

- Pseudo-counts are described in the lecture notes and section 7.2.3 of the textbook.
- Although you see high values of pseudo-count in some test cases, in practice small values are mostly used.

**For example:**

| Test | Result |
|------|--------|
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv")`<br>`print("Prior probability of spam is {:.5f}.".format(prior))` | Prior probability of spam is 0.25500. |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv")`<br>`print("Prior probability of not spam is {:.5f}.".format(1 - prior))` | Prior probability of not spam is 0.74500. |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 1)`<br>`print(format(prior, ".5f"))` | 0.25743 |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 2)`<br>`print(format(prior, ".5f"))` | 0.25980 |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 10)`<br>`print(format(prior, ".5f"))` | 0.27727 |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 100)`<br>`print(format(prior, ".5f"))` | 0.37750 |
| `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 1000)`<br>`print(format(prior, ".5f"))` | 0.47773 |

**Answer:** (penalty regime: 0, 10, ... %)

```
1  import csv
2  def learn_prior(file_name, pseudo_count=0):
3      """takes the file name of the training set and an optional pseudo-count
4      parameter and returns a real number that is the prior probability of spam
5      being true. The parameter pseudo_count is a non-negative integer and
6      it will be the same for all the attributes and all the values.
7      """
8      true_case = 0
9      with open(file_name) as in_file:
10         training_examples = [tuple(row) for row in csv.reader(in_file)]
11     for num_tuple in training_examples[1:]:
12         true_case += int(num_tuple[-1])
13     return (true_case + pseudo_count) / (len(training_examples) - 1 + 2 * pseudo_
```

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv")`<br>`print("Prior probability of spam is {:.5f}.".format(prior))` | Prior probability of spam is 0.25500. | Prior probability of spam is 0.25500. | ✔ |
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv")`<br>`print("Prior probability of not spam is {:.5f}.".format(1 - prior))` | Prior probability of not spam is 0.74500. | Prior probability of not spam is 0.74500. | ✔ |
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 1)`<br>`print(format(prior, ".5f"))` | 0.25743 | 0.25743 | ✔ |
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 10)`<br>`print(format(prior, ".5f"))` | 0.27727 | 0.27727 | ✔ |
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 100)`<br>`print(format(prior, ".5f"))` | 0.37750 | 0.37750 | ✔ |
| ✔ | `from student_answer import learn_prior`<br><br>`prior = learn_prior("spam-labelled.csv", pseudo_count = 1000)`<br>`print(format(prior, ".5f"))` | 0.47773 | 0.47773 | ✔ |

Passed all tests! ✔

Correct

Marks for this submission: 1.00/1.00.

**Question 4**

Correct

Mark 1.00 out of 1.00

Write a function `learn_likelihood(file_name, pseudo_count=0)` that takes the file name of a training set (for the spam detection problem) and an optional pseudo-count parameter and returns a sequence of pairs of likelihood probabilities. As described in the representation of likelihood, the length of the returned sequence (list or tuple) must be 12. Each element in the sequence is a pair (tuple) of real numbers such that `likelihood[i][False]` is *P(X[i]=true|Spam=false)* and `likelihood[i][True]` is *P(X[i]=true|Spam=true )*.

**For example:**

| Test | Result |
|------|--------|
| `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv")`<br>`print(len(likelihood))`<br>`print([len(item) for item in likelihood])` | 12<br>[2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2] |
| `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv")`<br><br>`print("P(X1=True \| Spam=False) = {:.5f}".format(likelihood[0][False]))`<br>`print("P(X1=False\| Spam=False) = {:.5f}".format(1 - likelihood[0][False]))`<br>`print("P(X1=True \| Spam=True ) = {:.5f}".format(likelihood[0][True]))`<br>`print("P(X1=False\| Spam=True ) = {:.5f}".format(1 - likelihood[0][True]))` | P(X1=True \| Spam=False) = 0.35570<br>P(X1=False\| Spam=False) = 0.64430<br>P(X1=True \| Spam=True ) = 0.66667<br>P(X1=False\| Spam=True ) = 0.33333 |
| `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv", pseudo_count=1)`<br><br>`print("With Laplacian smoothing:")`<br>`print("P(X1=True \| Spam=False) = {:.5f}".format(likelihood[0][False]))`<br>`print("P(X1=False\| Spam=False) = {:.5f}".format(1 - likelihood[0][False]))`<br>`print("P(X1=True \| Spam=True ) = {:.5f}".format(likelihood[0][True]))`<br>`print("P(X1=False\| Spam=True ) = {:.5f}".format(1 - likelihood[0][True]))` | With Laplacian smoothing:<br>P(X1=True \| Spam=False) = 0.35762<br>P(X1=False\| Spam=False) = 0.64238<br>P(X1=True \| Spam=True ) = 0.66038<br>P(X1=False\| Spam=True ) = 0.33962 |

**Answer:** (penalty regime: 0, 15, ... %)

```
1  import csv
2  def learn_likelihood(file_name, pseudo_count=0):
3      """takes the file name of a training set (for the spam detection problem)
4      and an optional pseudo-count parameter and returns a sequence of pairs of
5      likelihood probabilities. As described in the representation of likelihood,
6      the length of the returned sequence (list or tuple) must be 12. Each elemen
7      in the sequence is a pair (tuple) of real numbers such that
8      likelihood[i][False] is P(X[i]=true|Spam=false) and likelihood[i][True]
9      is P(X[i]=true|Spam=true ).
10     """
11     with open(file_name) as in_file:
12         training_examples = [tuple(row) for row in csv.reader(in_file)]
13     likelihood = []
14
15     for i in range(len(training_examples[0]) - 1):
16         xi_true_given_class_false_count = 0 + pseudo_count
17         false_total = 0 + 2 * pseudo_count
18         xi_true_given_class_true_count = 0 + pseudo_count
19         true_total = 0 + 2 * pseudo_count
20         for num_tuple in training_examples[1: ]:
21             if int(num_tuple[-1]):  #class true
22                 true_total += 1
23                 if int(num_tuple[i]):  #xi true
24                     xi_true_given_class_true_count += 1
25             else:  #class false
26                 false_total += 1
27                 if int(num_tuple[i]):  #xi true
28                     xi_true_given_class_false_count += 1
29         likelihood += [(xi_true_given_class_false_count/false_total, xi_true_gi
30     return likelihood
```

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv")`<br>`print(len(likelihood))`<br>`print([len(item) for item in likelihood])` | 12<br>[2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2] | 12<br>[2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2] | ✔ |
| ✔ | `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv")`<br><br>`print("P(X1=True \| Spam=False) = {:.5f}".format(likelihood[0][False]))`<br>`print("P(X1=False\| Spam=False) = {:.5f}".format(1 - likelihood[0][False]))`<br>`print("P(X1=True \| Spam=True ) = {:.5f}".format(likelihood[0][True]))`<br>`print("P(X1=False\| Spam=True ) = {:.5f}".format(1 - likelihood[0][True]))` | P(X1=True \| Spam=False) = 0.35570<br>P(X1=False\| Spam=False) = 0.64430<br>P(X1=True \| Spam=True ) = 0.66667<br>P(X1=False\| Spam=True ) = 0.33333 | P(X1=True \| Spam=False) = 0.35570<br>P(X1=False\| Spam=False) = 0.64430<br>P(X1=True \| Spam=True ) = 0.66667<br>P(X1=False\| Spam=True ) = 0.33333 | ✔ |
| ✔ | `from student_answer import learn_likelihood`<br><br>`likelihood = learn_likelihood("spam-labelled.csv", pseudo_count=1)`<br><br>`print("With Laplacian smoothing:")`<br>`print("P(X1=True \| Spam=False) = {:.5f}".format(likelihood[0][False]))`<br>`print("P(X1=False\| Spam=False) = {:.5f}".format(1 - likelihood[0][False]))`<br>`print("P(X1=True \| Spam=True ) = {:.5f}".format(likelihood[0][True]))`<br>`print("P(X1=False\| Spam=True ) = {:.5f}".format(1 - likelihood[0][True]))` | With Laplacian smoothing:<br>P(X1=True \| Spam=False) = 0.35762<br>P(X1=False\| Spam=False) = 0.64238<br>P(X1=True \| Spam=True ) = 0.66038<br>P(X1=False\| Spam=True ) = 0.33962 | With Laplacian smoothing:<br>P(X1=True \| Spam=False) = 0.35762<br>P(X1=False\| Spam=False) = 0.64238<br>P(X1=True \| Spam=True ) = 0.66038<br>P(X1=False\| Spam=True ) = 0.33962 | ✔ |

Passed all tests! ✔

Correct

Marks for this submission: 1.00/1.00.

**Question 5**

Correct

Mark 1.00 out of 1.00

Write a function `nb_classify(prior, likelihood, input_vector)` that takes the learnt prior and likelihood probabilities and classifies an (unseen) input vector. The input vector will be a tuple of 12 integers (each 0 or 1) corresponding to attributes X1 to X12. The function should return a pair (tuple) where the first element is either `"Spam"` or `"Not Spam"` and the second element is the certainty. The certainty is the (posterior) probability of spam when the instance is classified as spam, or the probability of 'not-spam' otherwise. If spam and 'not spam' are equally likely (i.e. p=0.5) then choose 'not spam'.

This is a very simple function to implement as it only wraps the `posterior` function developed earlier.

Supply the following functions you developed earlier: `learn_prior` and `learn_likelihood`. Also include import statements and any other function that you may be using (e.g. `posterior`).

**For example:**

| Test | Result |
|---|---|
| ```from student_answer import learn_prior, learn_likelihood, nb_classify

prior = learn_prior("spam-labelled.csv")
likelihood = learn_likelihood("spam-labelled.csv")

input_vectors = [
    (1,1,0,0,1,1,0,0,0,0,0,0),
    (0,0,1,1,0,0,1,1,1,0,0,1),
    (1,1,1,1,1,0,1,0,0,0,1,1),
    (1,1,1,1,1,0,1,0,0,1,0,1),
    (0,1,0,0,0,0,1,0,1,0,0,0),
    ]

predictions = [nb_classify(prior, likelihood, vector)
               for vector in input_vectors]

for label, certainty in predictions:
    print("Prediction: {}, Certainty: {:.5f}"
          .format(label, certainty))``` | Prediction: Not Spam, Certainty: 0.99351<br>Prediction: Spam, Certainty: 0.57441<br>Prediction: Spam, Certainty: 0.59337<br>Prediction: Spam, Certainty: 0.83465<br>Prediction: Not Spam, Certainty: 0.99140 |
| ```from student_answer import learn_prior, learn_likelihood, nb_classify

prior = learn_prior("spam-labelled.csv", pseudo_count=1)
likelihood = learn_likelihood("spam-labelled.csv", pseudo_count=1)

input_vectors = [
    (1,1,0,0,1,1,0,0,0,0,0,0),
    (0,0,1,1,0,0,1,1,1,0,0,1),
    (1,1,1,1,1,0,1,0,0,0,1,1),
    (1,1,1,1,1,0,1,0,0,1,0,1),
    (0,1,0,0,0,0,1,0,1,0,0,0),
    ]

predictions = [nb_classify(prior, likelihood, vector)
               for vector in input_vectors]

for label, certainty in predictions:
    print("Prediction: {}, Certainty: {:.5f}"
          .format(label, certainty))``` | Prediction: Not Spam, Certainty: 0.99213<br>Prediction: Spam, Certainty: 0.57759<br>Prediction: Spam, Certainty: 0.59073<br>Prediction: Spam, Certainty: 0.83059<br>Prediction: Not Spam, Certainty: 0.98989 |

**Answer:** (penalty regime: 0, 15, ... %)

```
 1  import csv
 2  def learn_prior(file_name, pseudo_count=0):
 3      """takes the file name of the training set and an optional pseudo-count
 4      parameter and returns a real number that is the prior probability of spam
 5      being true. The parameter pseudo_count is a non-negative integer and
 6      it will be the same for all the attributes and all the values.
 7      """
 8      true_case = 0
 9      with open(file_name) as in_file:
10          training_examples = [tuple(row) for row in csv.reader(in_file)]
11      for num_tuple in training_examples[1:]:
12          true_case += int(num_tuple[-1])
```

```
13          return (true_case + pseudo_count) / (len(training_examples) - 1 + 2 * pseud
14
15  ▾ def learn_likelihood(file_name, pseudo_count=0):
16          """takes the file name of a training set (for the spam detection problem)
17          and an optional pseudo-count parameter and returns a sequence of pairs of
18          likelihood probabilities. As described in the representation of likelihood,
19          the length of the returned sequence (list or tuple) must be 12. Each elemer
20          in the sequence is a pair (tuple) of real numbers such that
21          likelihood[i][False] is P(X[i]=true|Spam=false) and likelihood[i][True]
22          is P(X[i]=true|Spam=true ).
23          """
24  ▾      with open(file_name) as in_file:
25              training_examples = [tuple(row) for row in csv.reader(in_file)]
26          likelihood = []
27
28  ▾      for i in range(len(training_examples[0]) - 1):
29              xi_true_given_class_false_count = 0 + pseudo_count
30              false_total = 0 + 2 * pseudo_count
31              xi_true_given_class_true_count = 0 + pseudo_count
32              true_total = 0 + 2 * pseudo_count
33  ▾          for num_tuple in training_examples[1: ]:
34  ▾              if int(num_tuple[-1]):   #class true
35                      true_total += 1
36  ▾                  if int(num_tuple[i]):   #xi true
37                          xi_true_given_class_true_count += 1
38  ▾              else:   #class false
39                      false_total += 1
40  ▾                  if int(num_tuple[i]):   #xi true
41                          xi_true_given_class_false_count += 1
42              likelihood += [(xi_true_given_class_false_count/false_total, xi_true_gi
43          return likelihood
44
45  ▾ def posterior(prior, likelihood, observation):
46          """returns the posterior probability of the class variable being true,
47          given the observation; that is, it returns p(Class=true|observation).
48          The argument observation is a tuple of n Booleans such that observation[i]
49          is the observed value (True or False) for the input feature X[i].
50          The arguments prior and likelihood are as described above.
51          """
52
```

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | `from student_answer import learn_prior, learn_likelihood, nb_classify`<br><br>`prior = learn_prior("spam-labelled.csv")`<br>`likelihood = learn_likelihood("spam-labelled.csv")`<br><br>`input_vectors = [`<br>`    (1,1,0,0,1,1,0,0,0,0,0,0),`<br>`    (0,0,1,1,0,0,1,1,1,0,0,1),`<br>`    (1,1,1,1,1,0,1,0,0,0,1,1),`<br>`    (1,1,1,1,1,0,1,0,0,1,0,1),`<br>`    (0,1,0,0,0,0,1,0,1,0,0,0),`<br>`    ]`<br><br>`predictions = [nb_classify(prior, likelihood, vector)`<br>`              for vector in input_vectors]`<br><br>`for label, certainty in predictions:`<br>`    print("Prediction: {}, Certainty: {:.5f}"`<br>`          .format(label, certainty))` | Prediction: Not Spam, Certainty: 0.99351<br>Prediction: Spam, Certainty: 0.57441<br>Prediction: Spam, Certainty: 0.59337<br>Prediction: Spam, Certainty: 0.83465<br>Prediction: Not Spam, Certainty: 0.99140 | Prediction: Not Spam, Certainty: 0.99351<br>Prediction: Spam, Certainty: 0.57441<br>Prediction: Spam, Certainty: 0.59337<br>Prediction: Spam, Certainty: 0.83465<br>Prediction: Not Spam, Certainty: 0.99140 | ✔ |

| | Test | Expected | Got | |
|---|---|---|---|---|
| ✔ | ```from student_answer import learn_prior, learn_likelihood, nb_classify

prior = learn_prior("spam-labelled.csv", pseudo_count=1)
likelihood = learn_likelihood("spam-labelled.csv", pseudo_count=1)

input_vectors = [
    (1,1,0,0,1,1,0,0,0,0,0,0),
    (0,0,1,1,0,0,1,1,1,0,0,1),
    (1,1,1,1,1,0,1,0,0,0,1,1),
    (1,1,1,1,1,0,1,0,0,1,0,1),
    (0,1,0,0,0,0,1,0,1,0,0,0),
    ]

predictions = [nb_classify(prior, likelihood, vector)
               for vector in input_vectors]

for label, certainty in predictions:
    print("Prediction: {}, Certainty: {:.5f}"
          .format(label, certainty))``` | Prediction: Not Spam, Certainty: 0.99213<br>Prediction: Spam, Certainty: 0.57759<br>Prediction: Spam, Certainty: 0.59073<br>Prediction: Spam, Certainty: 0.83059<br>Prediction: Not Spam, Certainty: 0.98989 | Prediction: Not Spam, Certainty: 0.99213<br>Prediction: Spam, Certainty: 0.57759<br>Prediction: Spam, Certainty: 0.59073<br>Prediction: Spam, Certainty: 0.83059<br>Prediction: Not Spam, Certainty: 0.98989 | ✔ |

Passed all tests! ✔

Correct

Marks for this submission: 1.00/1.00.