B.Tech (Computer Engineering) 8th Semester Examinations 2022
# Natural Language Processing and Information Extraction
Paper Code: CEN 807

Maximum Marks: 60                                                                              Maximum Time: 3 hr
*(Write your Roll No. on the top immediately on receipt of this question paper)*

**Note: Attempt any two parts from each question** Assume suitable data, if necessary.

| S.No. | Questions | Marks | CO |
|---|---|---|---|
| 1(a) | What is tokenization? What is the advantage of using Byte Pair Encoding scheme for tokenization? Train the BPE algorithm using the following table.<br><br><table><tr><td>low</td><td>5</td></tr><tr><td>lowest</td><td>2</td></tr><tr><td>newer</td><td>6</td></tr><tr><td>wider</td><td>3</td></tr><tr><td>new</td><td>2</td></tr></table><br>Test the BPE algorithm using the word "lower" and show its tokenization. Show all the steps for training and testing clearly. | 6 | 1 |
| 1(b) | Given the following bigram counts.<br><br><table><tr><td></td><td>i</td><td>want</td><td>to</td><td>eat</td><td>Chinese</td><td>food</td><td>lunch</td><td>spend</td></tr><tr><td>i</td><td>5</td><td>827</td><td>0</td><td>9</td><td>0</td><td>0</td><td>0</td><td>2</td></tr><tr><td>want</td><td>2</td><td>0</td><td>608</td><td>1</td><td>6</td><td>6</td><td>5</td><td>1</td></tr><tr><td>to</td><td>2</td><td>0</td><td>4</td><td>686</td><td>2</td><td>0</td><td>6</td><td>211</td></tr><tr><td>eat</td><td>0</td><td>0</td><td>2</td><td>0</td><td>16</td><td>2</td><td>42</td><td>0</td></tr><tr><td>Chinese</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td><td>82</td><td>1</td><td>0</td></tr><tr><td>food</td><td>15</td><td>0</td><td>15</td><td>0</td><td>1</td><td>4</td><td>0</td><td>0</td></tr><tr><td>lunch</td><td>2</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td></tr><tr><td>spend</td><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr></table> | 6 | 1 |

| | | | |
|---|---|---|---|
| | (1) Generate the Bigram probability matrix.<br>(2) Generate the Bigram probability matrix with add-1 smoothing. | | |
| 1(c) | Define and give examples:<br>  (1) Segmentation<br>  (2) Lemmatization<br>  (3) Stemming | 6 | 1 |
| | | | |

<br>

**2(a)** — 6 — 2

| | Doc | Words | Class |
|---|---|---|---|
| Training | 1 | Chinese, Beijing, Chinese | C |
| | 2 | Chinese, Chinese, Shanghai | C |
| | 3 | Chinese, Macao | C |
| | 4 | Tokyo, Japan, Chinese | J |
| Test | 5 | Chinese, Chinese, Tokyo, Japan | ? |

Compute the most likely class for Doc 5. Assume a multinomial naive Bayes classifier and use add-$\alpha$ La Place smoothing for the likelihoods. Suitable value of $\alpha$ may be chosen.

<br>

**2(b)** — 6 — 2

Consider the following training data:

| S.No. | Document | Class |
|---|---|---|
| 1. | Natural Language Processing | A |
| 2. | Language Model Learning | A |
| 3. | Ngram Langauge Model | A |
| 4. | Text Classification Model | A |
| 5. | Text Processing Model | A |
| 6. | Computer Vision | B |
| 7. | Image Classification Model | B |
| 8. | Object Segmentation | B |
| 9. | Image Processing | B |
| 10 | Object Recognition | B |

And Test Data:

| | | |
|---|---|---|
| 1. | Object Recognition Model | ? |
| 2. | Text Recognition Model | ? |

Predict the class for test samples using Multinomial Naïve Bayes and Bigram language model with La-Place smoothing.

<br>

| 2(c) | Define Precision, Recall and F-measure. | 6 | 2 |
|---|---|---|---|

| 3(a) | Given the following Dictionary entry for **line**.<br><br>**line²** a length of cord, rope, wire, or other material serving a particular purpose: *wring the clothes and hang them on the line | a telephone line.*<br><ul><li>one of a vessel's mooring ropes.</li><li>a telephone connection: *she had a crank on the line.*</li><li>a railroad track.</li><li>a branch or route of a railroad system: *the Philadelphia to Baltimore line.*</li></ul><br>**line³** a horizontal row of written or printed words.<br><ul><li>a part of a poem forming one such row: *each stanza has eight lines.*</li><li>(lines) the words of an actor's part in a play or film.</li><li>a particularly noteworthy written or spoken sentence: *his speech ended with a line about the failure of justice.*</li></ul><br>Which of these senses are related by homonymy, and which are related by polysemy? For any senses which are polysemous, give an argument as to how the senses are related. | 6 | 3 |
| --- | --- | --- | --- |
| 3(b) | Assume the following sentence L in which the word **line** is in focus:<br>**L** = you must wait in a long **line** at the checkout counter<br>Give a collocation feature vector (including n-gram) for in the word **line** in L, given a window size of 3 words to the left and 3 words to the right. | 6 | 3 |
| 3(c) | **C** = About three years ago, he nearly gave up because he nearly had nothing to sell; Now his shelves are full, and towels and clothes hang from a ***line*** overhead.<br><br>For the word ***line*** in the above text L, generate the bag-of-words feature vector for window size = +-2, assume **C** as the whole corpus. | 6 | 3 |

For the following term document matrix:

| 4(a) | | | | | | 6 | 4 |
| --- | --- | --- | --- | --- | --- | --- | --- |

| | Document 1 | Document 2 | Document 3 | Document 3 |
| --- | --- | --- | --- | --- |
| digital | 1 | 0 | 7 | 13 |
| computer | 114 | 80 | 62 | 89 |
| information | 36 | 58 | 1 | 4 |

| | | | | |
|---|---|---|---|---|
| | data | 20 | 15 | 2 | 3 | | |

Calculate the similarity between good and fool using (i) cosine similarity (ii) PPMI. Use add 2 smoothing if necessary.

| | | | |
|---|---|---|---|
| 4(b) | What are word vectors? Illustrate with the help of a suitable diagram. Define the cosine similarity between two word vectors. | 6 | 4 |
| 4(c) | Construct the word co-occurrence matrix, window size+-1, for the following text:<br>Document: "Roses are red. Sky is blue. | 6 | 4 |
| | | | |
| 5(a) | Define Information Extraction, Named Entity Recognition and Relation Extraction. Why is Information Retrieval task not sufficient to perform information extraction tasks? | 6 | 5 |
| 5(b) | What are the different encoding schemes for Named Entity Recognition? Illustrate with examples. | 6 | 5 |
| 5(c) | Construct a rule based e-mail extractor which can distinguish between sender and receiver e-mail addresses. | 6 | 5 |