

# Team (of 3) Project

## Points: 20

COSC6340.001

Dr. Minhua Huang

Out: 10/07/2022, Due 11/17/2022

### Objective:

#### Identify an English phrase on Bigram Language Model by Perplexity

You can call functions and facilities on the preprocessing procedure. You are not allowed to call functions for obtaining bigram, corpus cross entropy, and perplexity, and the test accuracy.

Given corpora  $\mathbf{D}$ , where  $\mathbf{D} = \langle \mathbf{x}[\mathbf{i}], \mathbf{y}[\mathbf{i}] \rangle$  with  $\mathbf{i} = 1 \dots n$ , s.t. each of  $\mathbf{x} = \langle \text{verb}, \text{noun}, \text{prep}, \text{prepobj} \rangle = \langle \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4 \rangle$  with a class label  $\mathbf{y} \in \{\mathbf{V}, \mathbf{N}\} = \{\mathbf{y}_1, \mathbf{y}_2\}$ . The corpora  $\mathbf{D}$  is divided into two sets which are  $\mathbf{D}_{\text{train}}$  and  $\mathbf{D}_{\text{test}}$ , specified by  $\mathbf{D}_{\text{train.csv}}$  and  $\mathbf{D}_{\text{test.csv}}$  files.

- (4.5 points, Implementation) Training procedures.
  - Compute bigram probability for  $\mathbf{j}^{\text{th}}$  attribute of  $\mathbf{i}^{\text{th}}$  feature under a class label  $\mathbf{y}$  for all  $\mathbf{i}, \mathbf{j}$  and  $\mathbf{y}$  in  $\mathbf{D}_{\text{train}}$  by MLE algorithm (and a smoothing technique), where  $\mathbf{C}$  is a counting function.

$$p(\mathbf{x}_{i,j} | \mathbf{x}_{i-1,j}, \mathbf{y}) = \frac{p(\mathbf{x}_{i-1,j}, \mathbf{x}_{i,j}, \mathbf{y})}{p(\mathbf{x}_{i-1,j}, \mathbf{y})} = \frac{C(\mathbf{x}_{i-1,j}, \mathbf{C}_{\mathbf{x}_{i,j}}, \mathbf{y})}{C(\mathbf{x}_{i-1,j}, \mathbf{y})}$$

- (4.5 points, Implementation) Testing procedures.
  - Compute the corpus cross entropy for each of the data instances in  $\mathbf{D}_{\text{test}}$ . A data instance of size  $\mathbf{m}$  is associating with a probability distribution  $\mathbf{p}$  with  $\mathbf{m}$  probabilities.

$$\mathbf{H}(\mathbf{p} | \mathbf{y}) = - \sum_{i=1}^{\mathbf{m}} p_{i,y} \log_2 p_{i,y}$$

- Compute the perplexity of the probability distribution of  $\mathbf{p}$ .

$$\mathbf{PP}(\mathbf{p} | \mathbf{y}) = 2^{\mathbf{H}(\mathbf{p} | \mathbf{y})}$$

- Assign a class label for a data instance in  $\mathbf{D}_{\text{test}}$ .

$$\mathbf{y} \leftarrow \text{argmin}_{\mathbf{y}_k} \{ \mathbf{PP}(\mathbf{p} | \mathbf{y} = \mathbf{y}_k) \}$$

- Evaluate your system by the following accuracy measurement.

$$\mathbf{ACC}_{\mathcal{D}_{\text{test}}} = \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{i=1}^{\mathbf{T}} \mathbf{L}(\hat{\mathbf{y}}^i, \mathbf{y}^i)$$

where  $\hat{\mathbf{y}}^i$  is the assigned class label by the classifier and  $\mathbf{y}^i$  is the true class label of a data instance  $\mathbf{x}^{[i]}$  in  $\mathcal{D}_{\text{test}}$  and  $\mathbf{T}$  is the number of data instances in  $\mathcal{D}_{\text{test}}$ .

$$\mathbf{L}(\hat{\mathbf{y}}^i, \mathbf{y}^i) = \begin{cases} 1 & \text{if } \hat{y}^i = y^i \\ 0 & \text{if } \hat{y}^i \neq y^i \end{cases}$$

- (7 points) Report on your design (10 pages).

Write a 10 page report. The first page should list the names of group members as well the associated tasks.

- Introduction
  - Describing your algorithms.
    - \* Preprocessing procedures
    - \* The algorithm(s) of obtaining bigrams
    - \* The algorithm(s) of obtaining corpus cross entropy and perplexity
    - \* The running time complexity of an algorithm (optional)
    - \* The missing data handling
    - \* Testing metrics
  - Experiments
    - \* Experiments and results discussions.
      - Experiment settings.
      - The results discussions and comparisons.
      - Pros and cons of the design.
  - Further improvements.
- (4 points). Presentation.
    - 13 minutes of oral presentation.
    - 2 minutes of question-answering.

### **Presentation Date**

- Nov. 21, 3 -4 groups
- Nov. 23 (reading day), 3 -4 groups
- Nov. 28, 3 -4 groups
- Nov. 30, 3 -4 groups

### **Submission:**

- Submit files
  - ReadMe.txt – describe how to operating your system.
  - Project source codes
  - Project report
  - PPT presentation slides