

Project Semester 1:  
Estimating a Language Model  
to Generate Wine Reviews

Joy-Anne FOSTER  
Student Number: 22203455  
Shami THIRION SEN (née SEN)  
Student Number: 22200036

Corpus Linguistics in English  
Armand STRICKER

The aim of this project is to build an algorithm that generates reviews from a language model using the probabilities of a word given the two words before it. These probabilities can be calculated by the Bayes theorem which can be explained with the following formula:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Thus, we can calculate the probability a word from the two previous words with the following logic:  $P(A|B)$  would represent the probability of knowing a word given the words before it. We then need to multiply the probabilities of A and that of B knowing A , and divide this result by  $P(B)$  in order to calculate the probability of knowing a  $w_i$  given the two words prior ( $w_1$ ,  $w_2$ ). Here we have used wine2.txt to study probability of a what precedes and follows a word.

1.a)  $P(w_i | w_{i-2}, w_{i-1}) = \text{count}(w_i, w_{i-2}, w_{i-1}) / \text{count}(w_{i-2}, w_{i-1})$

In order to calculate the probability of a word considering the two words before it, we need two words "Begin" and "Now" at the beginning of each sentence ; this will allow the language model to predict the probability of which word or words should be at the beginning of the sentence. We can test this hypothesis with the test.txt data. Some of the sentences obtained from the algorithm were the following: "I do not like chocolate pudding ." ; "I like chocolate pudding . " ; "I like chocolate ice-cream . " As can be observed, the sentences obtained are quite satisfactory for this tiny dataset. As all the sentences in the test file begin with "I", the probability that the first word in a sentence would begin with "I" is 1.0, meaning that every sentence produced in the test started with "I".

The difference between the sentences generated with the test data and wine reviews data was noticeable. Sentences generated with the wine reviews data were much longer and were grammatical, however we can see that some sentences lack meaning. An example of a sentence generated by the wine data was: *"Filled with fiji apple flavors giving way to toast , crème brûlée , green fig flavors and cream follow through on the floral-tinged finish has lots to offer gorgeous apple , sliced pear , floral scents to the kirsch and plum aromas give way to orange sherbet , yellow fruits and finds a good rosé for your money 's worth of releases to look like a turnaround in recent vintages have been standouts . "* Here, we can observe that the generator accurately makes use of punctuations, capitalisation and grammar, but we can see that as the words are generated randomly and the sentences are quite long, the semantics are lacking for the longer sentences.

```
163
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
Trebiano and Malvasia Nera. Best Va END
BEGIN NOW Boisterous, with pretty spice-accented pear and melon flavors backed with savory details. Well-structured and powerful version with
raspberry. Long. Best after 2003. Well, here? s an outstanding red. Traditional in vein, with macadamia nut and fennel aromas and easygoing '
03, Cabernet Sauvignon and END
BEGIN NOW Fruit from the well-known Hyde vineyard. Like E.T., this intense red and blackberry are shaded by notes of quince. Seductive, with
blackberries and raisins, blackberries and fresh earthiness that gives more with air, it needs cellaring. Sémillon. Drink no END
BEGIN NOW Well-meshed and mouthwatering overall. Lingering on the succulent finish. Bical, Cereal Branco, Malvasia Fina, Moscatel and Sauvign
on since t END
BEGIN NOW 4079 Daily Wine Picks found in this textbook Mosel-style '99. Complex and fascinating. Drink END
BEGIN NOW Exotic and fresh mushrooms. Drink no END
BEGIN NOW Serious dark color to the extracted blueberry, displaying modest cherry and mocha-tinged fruit, green herb and grapefruit pith on t
he tarragon-tinged finish. Seductive now, should gain finesse with cellaring. Rock-solid, with delightful aromas of fennel, green, picholine-t
o the spicy raspberry flavors shaded by tobacco leaf and crème brûlée. Drink now.30,000 cases made. Fr END
BEGIN NOW Beeswax, white currant and vanilla, raspberry ganache aromas backed by fine minerality. Modern approach focuses more on END
BEGIN NOW There-s more to Portugal than just weight in the form of sweet tapenade, maduro tobacco and fruit-filled finish has notes of damson
fruit and graham cracker, with crème brûlée. Drink or age. Sangiovese with Canaiolo and Ciliegiole. Best Valu END
BEGIN NOW Distinctively floral, lemon zest, honey, fig sauce and blackberry follow through on the END
BEGIN NOW Plenty of licorice on the sandalwood-filled finish. Boroli makes very good one, and those looking for a deliciou END
BEGIN NOW Wrapped in a light citrus highlights. Drink END
```

To conclude, it can be observed that it is possible to generate grammatical sentences from such algorithms. However, as the sentences are randomly generated by the computer, some of the longer sentences lack meaning. This is representative of a common problem of semantics in Machine Learning, or even Computational Semantics. Computers do not know the inherent meaning of words or their co-occurrence pattern for natural languages, so language generated by the machine can sometimes lack perfection.