

3.- MLE + Smoothing – Language Models

3.1

a) Complete the program `mle.py` to estimate via MLE the parameters of a character trigram model, and write them to a file.

b) Complete the program `generate.py` to generate a random sequence of characters consistent with the loaded trigram model.

c) Run the program `smooth.py` and enter different input sentences. Discuss why some sentences have zero probability.

Modify the program `smooth.py` to perform a simple smoothing via Lidstone's or Laplace's Law. Discuss the values chosen for N and B .

3.2

a) Extend the program `mle.py` to estimate the coefficients for a linear Interpolation smoothing. Write the coefficients into the first line of the model file, followed by the trigram parameters.

Linear Interpolation: $P(z|xy) = \lambda_1 \tilde{P}(z) + \lambda_2 \tilde{P}(z|y) + \lambda_3 \tilde{P}(z|xy)$

Coefficient estimation via deleted interpolation:

$\lambda_1 = \lambda_2 = \lambda_3 = 0$

foreach trigram `xyz` with `count(xyz) > 0`

depending on the maximum of the following three values:

case `(count(xyz)-1) / (count(xy)-1)` increment λ_1 by `count(xyz)`

case `(count(yz)-1) / (count(y)-1)` increment λ_2 by `count(xyz)`

case `(count(z)-1) / (N-1)` increment λ_3 by `count(xyz)`

normalize $\lambda_1, \lambda_2, \lambda_3$

b) Extend the program `smooth.py` to load the Linear Interpolation coefficients in the first line of the file, and use them to smooth the trigram probabilities. Compare the results with the smoothing obtained in the previous exercise.