

1. What is the best lambda smoothing parameter?

a. Development File

model 1 (0.1) : 63.73132288637212
model 2 (0.01): 43.68236600240003
model 3 (0.001): 36.08948111375689
model 4 (0.0001): 34.75148986130219
model 5 (0.00001): 34.613175504824774
model 6 (0.000001): 34.662586348571196
model 7 (0.0000001): 34.48520100538355
model 8 (0.00000001): 34.4363288307949

b. Test File

model 1 (0.1): 64.3166538296586
model 2 (0.01): 43.48238406837255
model 3 (0.001): 36.70219198874578
model 4 (0.0001): 35.13958379155097
model 5 (0.00001): 35.15101589081593
model 6 (0.000001): 34.604189112749395
model 7 (0.0000001): 35.16693870093668
model 8 (0.00000001): 35.118262408845936

- c. It seems that the best lambda is the lowest one, i.e. 0.00000001 for both dev and testing. If I hadn't used the development set to select the lambda and simply picked a lambda arbitrarily, the performance difference would have been around 30. By using the dev set for lambda selection, I was able to minimize perplexity on unseen data, as the dev set helped identify the most appropriate balance between higher-order and lower-order n-grams. The optimal lambda on the dev set was close to the optimal lambda on the test set, which suggests that the dev set was representative of the overall data distribution. The results highlight the importance of tuning hyperparameters using a dev set rather than directly on the test set, as this prevents overfitting and leads to more generalizable models. While the lambda tuning improved results, some smoothing methods like the discount model still outperformed lambda interpolation, particularly in handling unseen n-grams.

2. What is the best discount?

a. Development File (0.99, 0.9, 0.75, 0.5, 0.25, 0.1)

model 1 (0.99): 53.32040947187919
model 2 (0.9): 44.93477838459815
model 3 (0.75): 41.184766374725314
model 4 (0.5): 37.6475313961236
model 5 (0.25): 35.83315601686389
model 6 (0.1): 35.6433550049303

- b. **Test File** (0.99, 0.9, 0.75, 0.5, 0.25, 0.1)
- model 1 (0.99): 53.58513785417338
 - model 2 (0.9): 44.869933497250535
 - model 3 (0.75): 41.319168816493814
 - model 4 (0.5): 38.05546886017002
 - model 5 (0.25): 36.03787485281308
 - model 6 (0.1): 35.66690069615999

On the development set, the best discount was 0.1, which resulted in the lowest perplexity of 35.64. On the test set, the optimal discount was also 0.1, with a perplexity of 35.67. Had we not used the development set to tune the discount parameter, we might have chosen a larger discount value arbitrarily, such as 0.5 or 0.9. This would have resulted in significantly higher perplexities on both the development and test sets, with perplexities of 37.65 (for 0.5) and 44.87 (for 0.9) on the test set. This shows that careful tuning on the development set was crucial for identifying the most effective discount value. The performance improvement is consistent across both the development and test sets, with the lowest discount value (0.1) providing the best results. This suggests that giving more weight to lower-order n-grams helps capture the structure of the data better. The minimal difference in perplexity between the development and test sets (35.64 vs. 35.67) indicates that the discount model generalizes well and was not overfitted to the development set.

3. Performance

Which model is better? Provide some quantitative and/or qualitative arguments (including data or examples) of which approach is better. Make sure to clearly explain your evaluation approach and your arguments.

- The discount model performed significantly better, as evidenced by the lower average perplexities compared to the lambda model. This is because the discount model more effectively reallocates probability mass, especially to rare and unseen events, by reducing the probability of frequent n-grams and redistributing it to less frequent ones. In contrast, the lambda model tends to assign too much probability mass to unseen events without considering their context or frequency, which leads to suboptimal performance.

4. Wrap-up

Very briefly answer the following questions: how long did you spend on this assignment? What was the most fun part? least fun part? how would you improve it if I had to give it again?

- Around 15-20 hours were spent on this assignment. The most fun part was getting to work in a partnership and learn from each other. It was also a cool experience to brush up on our Java knowledge and have concepts come back to us that we'd first learned in CS 62. The least fun part was debugging the code that we thought had sound logic while still getting wrong perplexities. To improve this process, we would recommend having more mentor sessions available as well as more cohesive instructions. I think a more chronological/procedural writeup would've been more helpful in following along with what we were supposed to do.

Ethics

5. Of the six scenarios, which do you see as the most problematic? Give a couple of sentences justifying your answer.
 - I think the “Tyranny of AI Design” is the most problematic scenario. Since AIs are trained by biased data, the more we use AI, the more bias we feed into training. This potentially will create even larger disparities in an already segregated society, denying marginalized people of job opportunities, mortgage loans, health care, among others.
6. What is one scenario that is also concerning that is not listed here? It can be NLP-specific or more broadly to AI
 - Related to the Deepfake example, AI could also impersonate a person’s voice and writing style if trained enough. This could lead to scams and fraud, where, for example, scammers can use your voice to ask for ransom from someone you know, or giving written consent to something you yourself did not agree to.